

# EdgePrompt: Revolutionizing Education or Just Another EdTech Hype?

## A Deep Assessment of Early-stage Research of AI Guardrails in Resource-Constrained Classrooms

Irvan Putra\*, Emily Lomempow, Khrishna Malik, Riza Alaudin Syah, Yulia Nurliana, Christopher Owen, Dev Bakshi, Gustavo Demicoli, Harriet Mathew, Thang Khuat, Priyanka Paramanathan, Lalith Aditya Velayuthan Hemamalini, Milindi Kodikara, Christoforus Yoga Haryanto\*

### Abstract

We examine the early stage of EdgePrompt, a research to build an offline-capable AI framework for resource-constrained educational environments, through transdisciplinary technical, pedagogical, ethical, and philosophical lenses. Using meta-dialectical methodology, we find EdgePrompt's structured prompting and multi-stage validation demonstrably improve safety and constraint adherence on edge-deployable models, despite smaller models' JSON reliability challenges. Current state of implementation reveals critical efficiency-robustness tradeoffs (3-5x latency overhead) that potentially compromise edge deployment feasibility. Substantial implementation gaps exist in security controls, cultural responsiveness mechanisms, and hardware optimization. The architecture embodies fundamental tensions requiring contextual navigation: safety vs. capability, efficiency vs. robustness, standardization vs. contextualization, automation vs. human agency. While current implementation emphasizes Q&A validation over constructivist knowledge-building and lacks adequate cultural adaptation infrastructure, the framework's offline-capability and teacher-agency principles establish a foundation for more equitable AI deployment. EdgePrompt represents neither revolution nor mere hype but a pragmatic foundation stone requiring sustained development, community co-design, and philosophical evolution to fulfill its educational equity potential.

**Keywords:** Edge computing, artificial intelligence in education, educational equity, offline capability, teacher agency, prompt engineering, multi-stage validation, resource constraints, cultural responsiveness, pedagogical alignment, cybersecurity, AI ethics, implementation gaps, technical-educational tensions, Indonesia 3T regions

**Repository:** <https://github.com/build-club-ai-indonesia/edge-prompt>

© 2025-05-02, EdgePrompt Research Team

\* Corresponding author

## Table of Contents

[Foreword: The Promise and Peril of AI in Unequal Classrooms](#)

[Navigating This Assessment: A Guide for Diverse Readers](#)

[Bridging Technical and Educational Perspectives](#)

[The Value of Transdisciplinary Understanding](#)

[Introduction: Beyond the Hype – A Framework for Critical Assessment](#)

[The Allure of EdTech Solutions vs. Classroom Realities](#)

[Introducing EdgePrompt: Vision, Core Claims, and the 3T Context](#)

[Why a Deep, Transdisciplinary Assessment is Necessary](#)

[Methodology: Technical Analysis, Educational Ideals, First Principles, and Adversarial Thinking](#)

[Structure of the Book](#)

[Part 1: Deconstructing EdgePrompt – Vision, Architecture, and Implementation](#)

[Chapter 1: The Genesis – Necessity and Vision](#)

[Chapter 2: Architecture Blueprint – From Concept to Code](#)

[Chapter 3: The Research Engine – Validating Guardrails via Simulation](#)

[Part 2: The Crucible of Implementation – Reality Checks and Engineering Responses](#)

[Chapter 4: The Topic Consistency Crisis – Ensuring Fair Comparison](#)

[Chapter 5: The JSON Reliability Gauntlet – When LLMs Don't Follow Rules](#)

[Chapter 6: Efficiency vs. Robustness – The Edge Computing Tightrope](#)

[Part 3: Assessing EdgePrompt Through Critical Lenses](#)

[Chapter 7: Pedagogical Alignment – Tool, Tutor, or Obstacle?](#)

[Chapter 8: Human-Computer Interaction – Inside the Lived Classroom](#)

[Chapter 9: AI Ethics and Alignment – Safety, Bias, and Moral Formation](#)

[Chapter 10: Equity and Liberation – Access vs. Agency](#)

[Chapter 11: Cybersecurity in the Classroom – Beyond Standard Models](#)

[Part 4: Synthesizing the Assessment – Gaps, Tensions, and Future Paths](#)

[Chapter 12: The Gap Analysis – Vision vs. Reality Today](#)

[Chapter 13: Navigating the Tensions – Core Conflicts in Design](#)

[Chapter 14: Recommended Future Work – From Validation to Value](#)

[Chapter 15: Open Questions – The Dialectic of Further Research](#)

[Conclusion: Revolution, Hype, or Foundation Stone?](#)

[EdgePrompt as a Pragmatic Engineering Response to LLM Reality](#)

[EdgePrompt as a Socio-Technical Intervention with Moral Weight](#)

[EdgePrompt as a Work of Transition with Community Ownership and Generational Thinking](#)

[Final Thoughts: The Ongoing Work of Aligning AI with Human Flourishing in Education](#)

[References](#)

## Foreword: The Promise and Peril of AI in Unequal Classrooms

In the rapidly evolving landscape of educational technology, artificial intelligence presents both unprecedented opportunity and profound risk. This duality becomes even more pronounced when we consider the stark inequalities that characterize global education—where connectivity, resources, and technical expertise are unevenly distributed, often along lines already marked by historical disadvantage.

The 3T regions of Indonesia—"Terdepan, Terluar, Tertinggal" (Frontier, Outermost, Underdeveloped)—exemplify this reality. These regions have been formally identified through Indonesian presidential regulation as priority areas for development efforts (Widodo 2020). In these communities, the question is not just whether AI will enhance or diminish learning, but whether it will reach students at all. When reliable internet connectivity remains an unrealized promise, the cloud-based AI solutions dominating headlines and venture capital flows are effectively inaccessible. UNESCO (2020) found that during COVID-19 school closures, 43% of students worldwide had no internet access at home, with this figure rising to 82% in low-income regions, highlighting the stark digital divide affecting educational access.

EdgePrompt emerges in this context: an attempt to bring AI's educational potential to resource-constrained environments through offline-capable, teacher-controlled systems. Its ambition is significant—not merely to provide technical workarounds for connectivity barriers, but to reimagine the integration of AI into educational settings in ways that respect teacher agency, uphold safety standards, and adapt to local cultural and pedagogical realities. This approach aligns with research on teacher agency in the age of generative AI, which emphasizes the importance of maintaining teachers' power to act, make decisions, and take educational stances (Frøsig & Romero 2024).

However, we must approach such initiatives with both hope and skepticism. The history of educational technology is littered with solutions that promised revolution but delivered only incremental change, or worse, created new problems while failing to solve existing ones. Too often, technical capability has been conflated with educational value, and flashy demonstrations have obscured deeper questions about learning, agency, and equity. Holmes et al. (2022) argue for a community-wide ethical framework that acknowledges the distinction between "doing ethical things" versus "doing things ethically" in AI educational applications.

As we embark on this critical assessment of EdgePrompt, we are mindful that its development remains in early stages—more vision than reality, more potential than proof. Yet it is precisely at this formative stage that critical examination can be most valuable. By interrogating its premises, analyzing its design choices, and projecting its trajectory through various lenses, we can contribute to shaping both this specific project and the broader conversation about AI's role in education.

The stakes could not be higher. Every student deserves access to the best educational tools and approaches we can collectively create. But technological advancement alone does not guarantee more

equitable or effective learning. It is the thoughtful integration of technology into educational contexts—guided by sound pedagogy, ethical principles, and lived classroom realities—that will determine whether AI becomes a force for educational liberation or just another layer of stratification. As Swist and Gulson (2023) argue, we must view educational technologies within socio-technical assemblages that include not just the tools themselves but the social values and power structures around them.

EdgePrompt represents one early attempt to navigate this territory. Let us examine it not just for what it is today, but for what it aspires to become—and what it might teach us about creating truly accessible, equitable, and effective AI-enhanced learning environments for all students, regardless of where they live or the resources at their disposal.

## Navigating This Assessment: A Guide for Diverse Readers

### Bridging Technical and Educational Perspectives

This book intentionally operates at the intersection of multiple disciplines—computer science, education, philosophy, ethics, and more. Such a transdisciplinary approach is essential for properly evaluating AI in education, but it also creates a reading challenge. Technical specialists may find the pedagogical discussions unfamiliar, while educators may find the technical details daunting. Rather than diluting either perspective, we have chosen to maintain the necessary depth in both domains, and offer instead this guide to navigating the assessment from different starting points.

### For Educators, Policymakers, and Non-Technical Readers

If your background is primarily in education or policy rather than technology, consider the following approach to technical sections:

- **Focus on the "implications" subsections** that follow technical explanations. These translate technical findings into educational consequences.
- **Don't be intimidated by technical terminology.** Understanding conceptually what EdgePrompt attempts to do is more important than grasping every implementation detail.
- **Pay particular attention to Part 2**, where implementation challenges reveal the gap between AI theory and practice—crucial knowledge for realistic technology planning.
- **View technical constraints as design parameters** rather than merely limitations. Just as educators work within constraints of time and resources, AI systems operate within constraints of computation and data.
- **Consider Chapter 5's JSON reliability challenges** as a particularly accessible entry point to understanding technical realities, as it clearly illustrates how seemingly minor technical issues can have major educational implications.

The technical details matter not because they are interesting in themselves (though they may be), but because they shape what is possible in the classroom. A basic grasp of these constraints will make you a more effective advocate and planner for AI in education.

### For Developers, Engineers, and Technical Readers

If your expertise lies in technology rather than education, we suggest the following approach to pedagogical and philosophical sections:

- **Recognize that educational "effectiveness" cannot be reduced to metrics** alone. The philosophical discussions about meaning construction, agency, and cultural context outline what actually matters in education.
- **Pay particular attention to Chapter 7** on pedagogical alignment, which explains why technical capabilities must be evaluated against educational ideals, not just functional requirements.
- **View "tensions" as design challenges** rather than obstacles. The competing priorities in education are not bugs but features of a domain concerned with human development in all its complexity.
- **Consider Chapters 10 and 13** as particularly valuable for understanding why technically "correct" solutions might still fail educationally if they don't account for equity, liberation, and cultural context.
- **Approach the discussion of human judgment** not as a limitation to be automated away, but as a crucial component of any educational system—including those enhanced by AI.

Technical excellence in educational AI requires understanding not just what technology can do, but what education should be. The pedagogical and philosophical frameworks presented here are not abstract considerations but practical guides to building systems that actually serve educational purposes.

### The Value of Transdisciplinary Understanding

The most valuable insights in this assessment emerge precisely at the intersection of technical and educational expertise. When technical possibilities meet educational principles, we can begin to imagine AI systems that not only function reliably but contribute meaningfully to learning.

No reader is expected to become an expert in all domains covered here. Rather, we encourage a stance of respectful curiosity toward unfamiliar perspectives. The educator who appreciates the significance of technical constraints and the engineer who recognizes the complexity of educational values will both contribute more effectively to shaping AI's role in education.

As you read, consider maintaining two parallel questions: "What is technically possible?" and "What is educationally valuable?" Only by holding both questions simultaneously can we navigate the path toward AI systems that genuinely enhance learning for all students, regardless of their circumstances.

# Introduction: Beyond the Hype – A Framework for Critical Assessment

## The Allure of EdTech Solutions vs. Classroom Realities

The educational technology landscape pulses with promises of transformation. AI-powered tutors that personalize learning for every student; intelligent assistants that reduce teacher workload; adaptive systems that identify and address learning gaps in real-time. These visions captivate policymakers, funders, and technologists alike. Yet the reality in classrooms—particularly those in resource-constrained environments—often stands in stark contrast to these glossy futures. Research examining how generative AI is transforming teaching and educational practices highlights significant gaps between promises and implementation realities (Peres et al. 2023).

Teachers face immediate, practical challenges: unreliable internet, outdated hardware, minimal technical support, and the constant pressure to cover curriculum requirements while addressing diverse student needs. The gap between what EdTech solutions promise and what classroom contexts can support grows wider as AI capabilities advance, with the latest innovations typically requiring robust cloud infrastructure, high-bandwidth connections, and significant computational resources.

This is not merely a digital divide but an AI divide—one that threatens to further disadvantage already marginalized educational communities as AI becomes increasingly central to learning experiences in well-resourced settings. In such contexts, the march of technological progress risks widening existing gaps rather than closing them.

## Introducing EdgePrompt: Vision, Core Claims, and the 3T Context

EdgePrompt presents itself as a framework for enabling safe, effective AI integration in K-12 education through offline-capable systems. Its vision centers on three core principles:

1. **Offline capability:** Enabling AI-assisted learning in environments with limited or no internet connectivity, specifically targeting underserved regions like Indonesia's 3T areas.
2. **Teacher agency:** Positioning teachers as the controllers and stewards of AI tools, ensuring pedagogical decisions remain in human hands rather than delegated to algorithms.
3. **Safety guardrails:** Implementing structured prompting and multi-stage validation processes to enforce appropriate, age-suitable content generation and assessment in educational contexts.

The 3T regions (Terdepan, Terluar, Tertinggal) are formally defined in Indonesian governmental policy as frontier, outermost, and underdeveloped areas requiring special development attention (Widodo 2020). Research on teacher agency in the age of generative AI emphasizes the importance of maintaining teachers' power to act, affect matters, make decisions, and take stances when implementing AI in classrooms (Frøsig & Romero 2024). The safety guardrails approach builds on research in building

guardrails for large language models through structured prompting and validation stages to enhance safety in AI outputs (Dong et al. 2024).

Technically, EdgePrompt applies a neural-symbolic inspired approach to prompt engineering—creating structured templates and validation sequences that aim to govern LLM behavior without requiring model fine-tuning or constant cloud connectivity. This approach draws from research on augmenting neural networks with symbolic educational knowledge to create more trustworthy and interpretable AI for education (Hooshyar, Azevedo & Yang 2024). The EdgePrompt method claims to improve safety and constraint adherence on small, edge-deployable language models while maintaining acceptable alignment with cloud-quality outputs.

These claims deserve neither immediate acceptance nor dismissal, but rather careful, multidimensional evaluation. The 3T context—characterized by intermittent connectivity, limited technical infrastructure, and pressing educational needs—provides both the motivation for EdgePrompt's development and a challenging test case for its practicality. Current data shows Indonesia's internet penetration has reached 79.5%, but significant disparities remain in rural and remote areas (Adji 2024).

## **Why a Deep, Transdisciplinary Assessment is Necessary**

Educational technology exists at the intersection of multiple domains: computer science, pedagogy, ethics, psychology, sociology, and more. When we add AI to this mix, the complexity increases exponentially. No single discipline possesses the conceptual tools to fully evaluate an initiative like EdgePrompt.

A technical assessment alone might determine whether the system functions as specified but would miss critical questions about its educational value or ethical implications. A purely pedagogical evaluation might identify alignment with learning theories but overlook security vulnerabilities or implementation constraints. An ethical analysis might surface important concerns about bias or privacy but might not account for the practical realities of resource-constrained classrooms.

What we need, therefore, is a transdisciplinary approach—one that:

- Examines both the technical architecture and its educational implications
- Considers both immediate functionality and long-term sustainability
- Evaluates both stated intentions and potential unintended consequences
- Balances both opportunity and risk across diverse contexts

Such an assessment becomes even more crucial when examining early-stage work like EdgePrompt. By identifying tensions, gaps, and potential pitfalls early in the development process, we can help shape more thoughtful, effective, and equitable implementations. This approach aligns with research situating education technology within socio-technical assemblages that require examining not only the tools but the social values and power structures around them (Swist & Gulson 2023).

## Methodology: Technical Analysis, Educational Ideals, First Principles, and Adversarial Thinking

This assessment employs a multi-layered methodology that integrates diverse analytical approaches:

**Technical Analysis:** We will examine EdgePrompt's architecture, implementation choices, and research findings to understand its current capabilities, limitations, and engineering challenges. This includes analyzing the Phase 1 simulation methodology, the four-run comparison structure, and the specific technical hurdles encountered during implementation.

**Educational Lens:** We will evaluate EdgePrompt against established educational theories and principles, considering questions of pedagogical alignment, potential impact on teaching practices, and capacity to support meaningful learning experiences. This includes assessing whether its current design truly supports productive struggle, higher-order thinking, and culturally responsive pedagogy. This assessment draws on theory of culturally responsive teaching, which emphasizes incorporating students' cultural references in all aspects of learning (Gay 2018).

**First Principles Examination:** Rather than accepting conventional frameworks, we will break down EdgePrompt into its fundamental components and assumptions—examining the nature of learning, teaching, large language models, edge computing, and structured control. This approach helps uncover deeper tensions and opportunities that might be obscured by domain-specific terminology.

**Adversarial Thinking:** We will deliberately seek out potential failure modes, vulnerabilities, and unintended consequences by applying adversarial scenarios and edge cases to EdgePrompt's proposed implementation. This includes considering both technical exploits and social misuse scenarios that might emerge in actual classroom deployment.

**Meta-Dialectical Method:** Throughout our assessment, we will apply a dialectical approach that deliberately creates structured opposition to expose hidden assumptions, reveal blind spots, and navigate inherent tensions. This method involves systematically challenging positions, strengthening counter-arguments, and seeking higher-order integration of seemingly opposed perspectives.

By integrating these diverse approaches, we aim to provide a comprehensive, nuanced assessment that respects both technical complexity and educational context—one that can inform the future development of EdgePrompt while contributing to broader conversations about AI integration in resource-constrained educational environments.

## Structure of the Book

This book progresses through four major parts, each building on the previous to create a comprehensive assessment:



**Part 1: Deconstructing EdgePrompt** examines the vision, architecture, and implementation of EdgePrompt as it currently exists. We'll explore its conceptual foundations, technical design, and research methodology, establishing a clear understanding of what EdgePrompt aims to achieve and how it approaches that goal.

**Part 2: The Crucible of Implementation** focuses on the real-world challenges encountered during EdgePrompt's development. We'll analyze three critical issues—topic consistency, JSON reliability, and the efficiency-robustness tradeoff—and how the engineering responses to these challenges reveal important insights about applying AI in educational contexts.

**Part 3: Assessing EdgePrompt Through Critical Lenses** applies diverse analytical perspectives to evaluate EdgePrompt's current state and future potential. From pedagogical alignment to cybersecurity considerations, we'll examine how EdgePrompt navigates complex requirements and competing priorities across multiple domains.

**Part 4: Synthesizing the Assessment** brings together the preceding analyses to identify key gaps, essential tensions, and priority areas for future development. We'll conclude by positioning EdgePrompt within the broader landscape of educational AI and offering recommendations for its continued evolution.

Throughout this journey, we invite readers to maintain both optimism and skepticism—to see EdgePrompt not as a completed solution but as an exploratory attempt to address a crucial challenge. By critically engaging with its development at this early stage, we can contribute to shaping not just this particular initiative, but the broader field of AI-enhanced education in resource-constrained environments.

# Part 1: Deconstructing EdgePrompt – Vision, Architecture, and Implementation

## Chapter 1: The Genesis – Necessity and Vision

### The Problem: Connectivity, Equity, and Teacher Needs in 3T Regions (and beyond)

The starting point for understanding EdgePrompt is recognizing the acute challenge of educational technology access in regions where connectivity cannot be taken for granted. Indonesia's 3T regions (Frontier, Outermost, Underdeveloped) exemplify this reality, but similar conditions exist across many parts of the globe. These regions were formally identified through Presidential Regulation No. 63 of 2020, which designated specific areas requiring special development attention due to geographical isolation, limited infrastructure, and socioeconomic challenges (Widodo 2020).

In these areas, internet connectivity is often:

- **Inconsistent:** Available intermittently rather than continuously
- **Low-bandwidth:** Insufficient for real-time interaction with cloud AI services
- **Expensive:** Consuming a disproportionate share of limited educational budgets
- **Unreliable:** Subject to outages during critical educational activities

This connectivity challenge creates a paradox: the places that might benefit most from educational technology innovations are often the least able to access them. As AI increasingly drives educational tool development, this paradox threatens to widen existing educational divides. UNESCO research found that during pandemic-related school closures, 43% of students worldwide lacked internet access at home, with this figure reaching 82% in low-income regions (UNESCO 2020). While students in well-connected schools gain experience with the latest AI-enhanced learning tools, those in disconnected regions risk falling further behind—not just in technological familiarity, but in access to personalized learning experiences that could help address educational gaps.

Beyond connectivity, teachers in these regions face additional challenges:

- Limited technical expertise and support
- Significant workload pressures and resource constraints
- Cultural and linguistic considerations that may not be reflected in global AI systems
- Needs for trustworthy, predictable tools that align with local educational contexts

These challenges demand solutions that work within existing constraints rather than assuming infrastructure improvements that may be years or decades away. They call for approaches that empower rather than replace teachers, and that respect the cultural and pedagogical diversity of different educational contexts.

## The EdgePrompt Vision: Offline Capability, Teacher Agency, Safety Guardrails

In response to these challenges, EdgePrompt proposes a distinct vision for AI integration in education:

**Offline Capability:** At the core of EdgePrompt is the commitment to functioning in environments with limited or no internet connectivity. This means:

- Deploying smaller language models that can run on local, edge devices
- Enabling core educational functions without cloud dependencies
- Designing for efficient operation within hardware constraints
- Implementing intelligent synchronization for intermittent connectivity

This approach builds on educational edge computing frameworks that demonstrate how local edge servers can handle intensive tasks with minimal latency, improving throughput for remote and rural learners (Chen et al. 2022).

**Teacher Agency:** EdgePrompt positions teachers as the primary agents and decision-makers, not passive recipients of AI outputs. This principle manifests in:

- Teacher-defined content generation parameters and constraints
- Human-in-the-loop review for flagged content or edge cases
- Tools that augment rather than replace teacher expertise
- Interfaces designed for non-technical educators

Research on teacher agency in ICT use shows that structured reflection time on technology use enhances teachers' sense of agency, helping them align technology with local context and values rather than passively adopting standardized tools (Novoa-Echaurren 2024).

**Safety Guardrails:** Recognizing the particular sensitivity of K-12 educational contexts, EdgePrompt prioritizes:

- Age-appropriate content generation and filtering
- Multi-stage validation to ensure educational relevance
- Explicit constraint enforcement for generated content
- Structured prompting to maintain pedagogical alignment

This vision represents a deliberate positioning of educational AI as infrastructure rather than authority—as a tool that amplifies human educational expertise rather than supplanting it. It acknowledges that true educational transformation comes not from technology alone, but from thoughtful integration of technology into existing educational contexts, respecting both local constraints and human relationships.

## **Core Principles: Neural-Symbolic Inspiration, Prompt Engineering Focus, Educational Philosophy**

EdgePrompt draws on several foundational principles that inform its approach:

**Neural-Symbolic Inspiration:** While not formally implementing neural-symbolic computing, EdgePrompt draws inspiration from this field's integration of neural networks' learning capabilities with symbolic systems' interpretability and reasoning. This manifests in:

- Structured templates that encode symbolic constraints
- Multi-stage validation sequences that decompose complex assessment tasks
- Explicit representation of educational requirements in machine-actionable formats
- Combining statistical language models with rule-based guardrails

This hybrid approach aims to leverage the fluency and adaptability of language models while imposing the structure and safety boundaries needed for educational applications. Recent systematic reviews of neuro-symbolic AI highlight the potential of these approaches for combining the strengths of neural networks with symbolic reasoning capabilities (Colelough & Regli 2025).

**Prompt Engineering Focus:** Rather than relying on model fine-tuning (which requires significant data, expertise, and computational resources), EdgePrompt centers on prompt engineering as the primary mechanism for controlling model behavior:

- Developing explicit templates for different educational tasks
- Creating structured validation sequences with specific checking stages
- Implementing rubric formalization and transformation for assessment
- Engineering robust fallback mechanisms when models deviate from expected outputs

This approach makes EdgePrompt potentially more adaptable to different models and more accessible to implementation in resource-constrained environments. Systematic surveys of prompt engineering techniques demonstrate how different prompting strategies can guide model behavior through instructions rather than parameter updates (Sahoo et al. 2025).

**Educational Philosophy:** Beyond its technical approach, EdgePrompt embodies specific educational principles:

- Valuing productive struggle while providing appropriate scaffolding
- Supporting higher-order thinking rather than just information retrieval
- Respecting cultural and linguistic diversity in educational contexts
- Treating pedagogy as the driver of technology, not vice versa

These principles align with Vygotsky's concept of the Zone of Proximal Development, which emphasizes the gap between what learners can do independently and what they can achieve with assistance (Vygotsky 1978). They also connect to research on harnessing AI for constructivist learning,

which positions students as "active architects" of their own knowledge building rather than passive recipients (Grubaugh, Levitt & Deever 2023).

These principles suggest that EdgePrompt aims to be more than just an offline version of existing AI educational tools—it seeks to rethink how AI and education interact in ways that center human relationships, educational values, and local contexts.

### **The Initial Hypothesis (Paper vs. Refined Post-Paper)**

EdgePrompt's development began with an initial hypothesis presented in an academic paper (Syah et al.): that structured prompt templates and multi-stage validation could provide more robust safety and constraint adherence on edge-deployable language models compared to baseline prompting approaches, while maintaining acceptable alignment with cloud quality, using only prompt engineering techniques.

This hypothesis made several key assumptions:

- That smaller language models could reliably process structured output formats like JSON
- That multi-stage validation would be computationally feasible on edge devices
- That prompt engineering alone could provide sufficient guardrails without fine-tuning

As implementation progressed, this hypothesis needed refinement based on real-world findings:

**JSON Reliability Challenges:** Implementation revealed that smaller edge models struggled significantly with reliable JSON output generation—a critical issue since structured output was essential for the validation pipeline. This necessitated extensive robustness engineering through fallback mechanisms, repair utilities, and simplified templates. These challenges parallel findings from research on reliable code generation from pre-trained language models, which identified similar issues with structured output consistency (Poesia et al. 2022).

**Efficiency-Robustness Tradeoff:** The multi-stage validation approach improved safety and constraint adherence but at a significant computational cost. This highlighted the need for careful optimization to make the approach truly viable on edge devices. This tradeoff reflects broader challenges identified in research on on-device language models, which emphasizes the need for efficient architectures and compression techniques for resource-constrained deployment (Xu et al. 2024).

**Topic Consistency Issues:** Initial testing showed that baseline runs would frequently go off-topic, making fair comparisons difficult. This led to the implementation of a "shared teacher request" mechanism to ensure all test runs addressed the same task.

These refinements represent an important evolution in EdgePrompt's approach—moving from theoretical assertions to practical engineering solutions that acknowledge the messier reality of working with edge-deployed language models. This adaptation demonstrates a commitment to pragmatic

problem-solving rather than purely theoretical ideals, which will be essential for creating systems that work in challenging real-world educational environments.

## Chapter 2: Architecture Blueprint – From Concept to Code

**For Non-Technical Readers:** Think of EdgePrompt's architecture as the blueprint for a specialized school building. This section explains how the designers planned different rooms (services) that work together: a central library for storing information (database), special classrooms for creating educational content (content generation), and evaluation spaces for checking student work (validation). Unlike typical online learning tools that require constant internet access to a central headquarters, EdgePrompt's design allows each school to operate independently with its own resources. The most important design choice is putting teachers in charge of the building's operations rather than automating everything. Understanding this architecture helps you see how technical decisions directly affect who has control in the classroom and how the system works when internet connections are unreliable.

### The Conceptual Framework (Paper Fig 1, $T_c$ , $A_s$ , $R'$ , Validation Stages, Adaptation Loops)

EdgePrompt's conceptual framework, as presented in the original paper, establishes a structured approach to controlling language model behavior in educational contexts. This framework introduces several core concepts:

#### Teacher-Driven Content Generation:

- **Question Templates ( $T_c$ ):** Structured patterns for generating educational content, with placeholders for variables like topic, constraints, and learning objectives.
- **Answer Space ( $A_s$ ):** Explicit boundaries defining acceptable outputs, including parameters like word count, vocabulary level, and prohibited content.
- **Learning Objectives Mapping ( $O: T_c \rightarrow L$ ):** A formal relationship connecting templates to specific learning goals, ensuring generated content serves educational purposes.

#### Cloud/Edge Pipeline:

- **Rubric Formalization ( $R(ct, vp)$ ):** Transforming teacher criteria and validation parameters into explicit assessment rubrics.
- **Edge Transformation ( $S: R \rightarrow R'$ ):** Converting complex rubrics into forms optimized for edge deployment.
- **Grading Template ( $G(R')$ ):** Templates for evaluating student responses based on the transformed rubric.

#### Student Answer Evaluation:

- **Edge Validation ( $V(q,a)$ ):** The process of validating student answers against questions using edge-deployed models.
- **Staged Validation ( $\{v1..vn\}$ ):** Breaking validation into sequential stages addressing different aspects (relevance, safety, quality).

- **Boundary Enforcement ( $B(r)$ ):** Mechanisms to enforce constraints on generated content and evaluations.
- **Evaluation Logic ( $E(r, R')$ ):** Procedures for comparing responses against rubric criteria.
- **Scoring ( $S(e)$ ):** Converting evaluation results into numerical or categorical assessments.
- **Constraint Satisfaction ( $C(s, ct)$ ):** Verifying that scores satisfy all defined constraints.

#### Teacher Review System (Adaptation - Future Work):

- **Edge Case Detection ( $D(r, \theta)$ ):** Identifying responses requiring human review based on threshold criteria.
- **Review Triggers ( $T(d)$ ):** Mechanisms for surfacing detected edge cases to teachers.
- **Pattern Tracking ( $K(h)$ ):** Monitoring performance history to identify trends requiring intervention.
- **Rubric Adjustment ( $A: R' \rightarrow R''$ ):** Refining rubrics based on observed performance.
- **Criteria Optimization ( $O(K, \epsilon)$ ):** Using historical patterns to optimize evaluation criteria.
- **Template Refinement ( $P(Tc, h)$ ):** Evolving content templates based on performance history.

This conceptual framework demonstrates EdgePrompt's ambition to create a comprehensive system spanning content generation, evaluation, and adaptive refinement. However, it's important to note that the current implementation focuses primarily on the first two components (content generation and evaluation), with the adaptation components largely reserved for future development.

The framework's formalization through mathematical notation reflects its neural-symbolic inspiration—attempting to bring structured, symbolic representation to the otherwise black-box behavior of language models. This approach aims to create more transparent, controllable, and educationally aligned AI systems. This multi-stage approach parallels research on AI system evaluation frameworks that span the entire model lifecycle, implementing component-level checks and system-level validation mapped to different stakeholders and development stages (Xia et al. 2024).

#### Implemented Architecture: Backend Services, Frontend Components, Research Runner

The current implementation of EdgePrompt translates these conceptual elements into a concrete software architecture composed of three main components:

**Backend Services:** Implementing a service-oriented architecture with:

- **DatabaseService:** Centralized data access layer for SQLite operations
- **StorageService:** File management for uploaded educational materials
- **MaterialProcessor:** Content extraction and processing pipeline
- **ValidationService:** Implementation of (simplified) response validation
- **LMStudioService:** Interface to external LLM runtime

This backend-first approach ensures all LLM interactions occur server-side, preventing direct client-to-LLM communications that could bypass safety controls. The services maintain clear



separation of concerns while working together to implement the conceptual framework's core functionality.

**Frontend Components:** Organized in a feature-based hierarchy:

- **Project Components:** Managing organizational units (ProjectEditForm, ProjectPanel, etc.)
- **Teacher Components:** Supporting educator workflows (MaterialsManager, ContentGenerator, etc.)
- **Prompt Components:** Handling prompt template creation and management
- **Student Components:** Enabling learner interaction through the ResponseValidator

The frontend uses React with Context API for state management, avoiding the complexity of Redux while providing centralized access to projects, templates, and selection state. This architecture aims for simplicity and usability, particularly for non-technical educators in resource-constrained environments.

**Research Runner:** A separate Python framework implementing:

- **RunnerCore:** Orchestrating the four-run simulation strategy
- **TemplateEngine:** Processing templates with variable substitution
- **EvaluationEngine:** Executing multi-stage validation sequences
- **ConstraintEnforcer:** Implementing boundary checking on generated content
- **MetricsCollector:** Gathering performance data on latency and token usage
- **ResultLogger:** Recording experimental outputs for analysis

This research component enables systematic testing of EdgePrompt's core hypotheses through controlled experiments, comparing different approaches (baseline vs. EdgePrompt) across different execution environments (cloud vs. edge).

The division between the application (backend/frontend) and research components reflects EdgePrompt's dual nature as both a practical tool and a research project. While the research runner implements more sophisticated validation logic aligned with the paper's conceptual framework, the application currently uses a simplified approach—highlighting the gap between theoretical design and practical implementation at this stage of development.

### **Key Design Choices: Backend-First Security, SQLite, LM Studio Integration, React Context**

Several key design choices shape EdgePrompt's current implementation:

**Backend-First Security Model:** All LLM interactions occur server-side rather than directly from the client. This architectural decision:

- Prevents prompt injection or manipulation from client-side code
- Enables consistent application of safety measures and validation
- Centralizes logging and monitoring of LLM interactions

- Creates a clear boundary between user interface and AI processing

This approach represents a fundamental safety principle—ensuring that users interact with AI capabilities through well-defined interfaces rather than directly manipulating model inputs. This security-focused architecture aligns with research on cybersecurity in generative AI, which highlights the importance of controlled access patterns to prevent misuse (Gupta et al. 2023).

**SQLite for Data Persistence:** The choice of SQLite as the database technology reflects EdgePrompt's offline-first orientation:

- Self-contained database requiring no separate server process
- Minimal resource requirements suitable for edge deployment
- Simple deployment and maintenance for non-technical environments
- Reliable performance for the expected data volumes

While this choice limits certain scaling capabilities, it aligns with the goal of creating systems that can function effectively in resource-constrained environments without complex infrastructure.

**LM Studio Integration:** Rather than implementing direct model inference, EdgePrompt integrates with LM Studio via its API:

- Enables flexible model selection without code changes
- Leverages LM Studio's optimization for local model deployment
- Simplifies testing with different edge-deployable models
- Reduces development complexity by delegating inference

This integration strategy allows EdgePrompt to focus on its core value proposition (structured prompting and validation) while leveraging existing tools for the actual model execution—a pragmatic approach for early-stage development.

**React Context for State Management:** The frontend uses React Context API instead of more complex state management libraries:

- Provides sufficient centralization for the application's needs
- Reduces dependency weight and complexity
- Simplifies component interactions and data access
- Improves maintainability for a small development team

This choice exemplifies a pattern of selecting simpler, more lightweight approaches that align with both the resource constraints of the target environments and the early stage of the project's development.

Together, these design choices reflect a pragmatic approach to implementation—prioritizing simplicity, offline capability, and security while accepting certain limitations in scale and feature richness. They

create a foundation that can potentially evolve toward the more sophisticated vision presented in the conceptual framework, while delivering functional value even in its current form.

### **Alignment and Discrepancies: Where Implementation Meets (and Diverges from) the Blueprint**

Examining the relationship between EdgePrompt's conceptual blueprint and its current implementation reveals both strong alignments and significant gaps:

#### **Areas of Alignment:**

- **Backend-First Security:** The implementation successfully realizes the vision of a server-controlled AI interaction model that prevents direct user manipulation of prompts.
- **Service-Oriented Architecture:** The backend implementation follows the prescribed separation of concerns, with distinct services handling different aspects of the system.
- **Template-Based Approach:** Both the conceptual framework and implementation use structured templates to control language model outputs.
- **Research Framework Validation:** The research runner implements the multi-stage validation approach described in the conceptual framework, enabling empirical testing of the core hypotheses.

#### **Notable Discrepancies:**

- **Simplified Application Validation:** While the research framework implements sophisticated multi-stage validation, the actual web application uses a simplified, single-stage approach that lacks the nuanced checking sequence described in the blueprint.
- **Missing Adaptation Components:** The adaptive components described in the conceptual framework (edge case detection, pattern tracking, template refinement) remain largely unimplemented in the current system.
- **Limited Security Implementation:** Though the architecture specifies a comprehensive security model, the current implementation lacks key elements like user authentication, role-based access control, and data encryption.
- **Partial Offline Capability:** While designed for eventual offline use, the current implementation still functions as a standard client-server application without robust offline synchronization mechanisms.

These discrepancies highlight the gap between EdgePrompt's ambitious vision and its current reality as an early-stage system. Some of these gaps represent reasonable prioritization decisions for Phase 1 development, while others may indicate more fundamental challenges in implementing the conceptual framework.

The most significant alignment challenge appears in the validation system, where the research runner's more comprehensive approach has not yet been integrated into the application backend. This suggests

that while the multi-stage validation concept has been validated through research, translating it into a production-ready implementation remains a work in progress.

Similarly, the adaptation components represent a substantial area of future work. These components—which would enable the system to learn from usage patterns and improve over time—constitute some of the most innovative aspects of the EdgePrompt vision but also some of the most technically challenging to implement.

These alignment observations provide important context for our assessment. We should evaluate EdgePrompt not just against its current implementation but in light of the trajectory from concept to code—recognizing both the progress made and the substantial work that remains to fully realize the blueprint.

## Chapter 3: The Research Engine – Validating Guardrails via Simulation

**For Non-Technical Readers:** Imagine testing a new safety helmet design before sending it to production. This section explains how researchers systematically tested EdgePrompt's safety features through controlled experiments. They created a four-way comparison: testing both high-powered online models and smaller offline models, each with and without EdgePrompt's special safety controls. This approach is like comparing four helmet designs: premium materials with and without extra padding, and budget materials with and without the same padding. The goal was to determine whether EdgePrompt's safety methods could make the more affordable "budget" AI models (which can run without internet) perform more safely while maintaining quality. These tests help everyone understand what's actually possible with smaller AI models in disconnected classrooms, rather than relying on hopeful claims or marketing promises.

### The Phase 1 Four-Run Comparative Methodology (Rationale and Structure)

To empirically validate EdgePrompt's core claims about prompt engineering guardrails, the research team developed a structured comparative methodology implemented through a dedicated Python framework. This approach, referred to as the "Phase 1 Four-Run Comparative Methodology," enables systematic evaluation of EdgePrompt's effectiveness compared to baseline approaches.

The methodology centers on four distinct experimental configurations:

#### Run 1: Cloud Baseline

- **Executor:** CloudLLM (e.g., GPT-4o, Claude 3.7)
- **Method:** SingleTurn\_Direct (simple, unstructured prompting)
- **Purpose:** Serves as the Phase 1 Proxy Reference standard—representing high-quality outputs from state-of-the-art cloud models without EdgePrompt constraints.

#### Run 2: Cloud EdgePrompt

- **Executor:** CloudLLM (same as Run 1)
- **Method:** MultiTurn\_EdgePrompt (structured prompts, multi-stage validation)
- **Purpose:** Tests how EdgePrompt's approach performs when used with high-capability cloud models.

#### Run 3: Edge Baseline

- **Executor:** EdgeLLM (e.g., Llama 3.2 3B, Gemma 3 4B)
- **Method:** SingleTurn\_Direct (same as Run 1)
- **Purpose:** Establishes the performance of smaller, edge-deployable models with standard prompting.

#### Run 4: Edge EdgePrompt

- **Executor:** EdgeLLM (same as Run 3)
- **Method:** MultiTurn\_EdgePrompt (same as Run 2)
- **Purpose:** Demonstrates EdgePrompt's impact on smaller models—the primary focus of the research.

This structure creates a controlled experimental design with two key variables:

1. **Model Capability:** CloudLLM vs. EdgeLLM
2. **Prompting Approach:** Direct vs. EdgePrompt

By isolating these variables, the methodology enables direct comparison between Run 4 (EdgeLLM with EdgePrompt) and Run 3 (EdgeLLM baseline) to assess the specific impact of the EdgePrompt approach on edge-deployable models. Meanwhile, Run 1 (CloudLLM baseline) serves as a quality reference point for assessing how close edge outputs can get to state-of-the-art performance.

This methodological design reflects careful research planning, creating a framework for empirical validation rather than relying solely on theoretical arguments or anecdotal observations. It enables quantitative measurement of EdgePrompt's impact on key metrics like safety violations, constraint adherence, and output quality, while also capturing performance data on latency and token usage. This systematic approach parallels red-teaming methods used to identify harmful AI outputs, applying structured evaluation techniques to measure safety enhancements (Ganguli et al. 2022).

#### Operationalizing Key Concepts: CloudLLM vs. EdgeLLM Roles, Orchestrator Function

Implementing this comparative methodology required operationalizing several key concepts:

##### CloudLLM Implementation:

- Large, state-of-the-art models like GPT-4o and Claude 3.7 Sonnet
- Accessed via API calls from the research runner
- Used for two distinct purposes:
  1. Simulating teacher/student personas in generating test content
  2. Providing high-quality reference outputs and evaluations
- Configured with appropriate system prompts and generation parameters for each role

##### EdgeLLM Implementation:

- Smaller, edge-deployable models like Gemma 3 4B and Llama 3.2 3B
- Accessed through LM Studio's API for local inference
- Configured to simulate execution on edge devices
- Used to test both baseline prompting and the EdgePrompt method
- Subject to the same tasks and inputs across approaches for fair comparison

**Orchestrator Function:**

- Implemented as **RunnerCore** in the research framework
- Coordinates the entire experimental workflow:
  - Loads test cases and templates
  - Prepares shared inputs to ensure fair comparison
  - Executes each of the four runs in sequence
  - Applies constraint enforcement (e.g., word count, prohibited content)
  - Collects metrics on token usage and latency
  - Records detailed results for later analysis
- Maintains consistency across experimental conditions

**Template Processing:**

- Implemented through the **TemplateEngine** component
- Performs variable substitution in templates
- Formats prompts according to template specifications
- Applies consistent processing across all four runs
- Handles special cases like adapting to different model formats

**Multi-Stage Validation:**

- Implemented through the **EvaluationEngine** component
- Executes validation stages in the specified sequence
- Processes responses from each stage and aggregates results
- Implements robust parsing of structured outputs (with fallbacks)
- Applies the same validation logic consistently across runs

This operationalization translates the conceptual methodology into executable code, enabling systematic testing of EdgePrompt's hypotheses. The implementation includes careful attention to fair comparison (through shared inputs), robust measurement (via consistent metrics collection), and reliable output processing (with fallback mechanisms for edge cases).

It's important to note that this represents a simulation strategy rather than real-world deployment—CloudLLM simulates teacher/student personas and provides quality reference points, while EdgeLLM simulates the edge execution environment. Future phases would need to replace this simulation with actual hardware testing and human evaluation to fully validate EdgePrompt's real-world effectiveness. This approach of using more capable models to simulate different personas and evaluate outputs shares methodological elements with research on training models to follow instructions with human feedback (Ouyang et al. 2022).

## Metrics and Evaluation Criteria for Phase 1

To quantitatively assess EdgePrompt's performance, the Phase 1 research defined a clear set of metrics and evaluation criteria:

### Primary Metrics:

1. **Safety Violation Rate:** The percentage of outputs containing inappropriate content (e.g., violence, explicit material, age-inappropriate topics) as determined by both constraint checking and validation stages.
  - *Primary Comparison:* Run 4 (Edge EdgePrompt) vs. Run 3 (Edge Baseline)
  - *Expected Outcome:* Lower violation rate with EdgePrompt
2. **Constraint Adherence Rate:** The percentage of outputs meeting specified constraints such as word count limits, topic relevance, and format requirements.
  - *Primary Comparison:* Run 4 vs. Run 3
  - *Expected Outcome:* Higher adherence rate with EdgePrompt
3. **Quality Alignment:** Agreement metrics comparing outputs against the reference standard (Run 1: Cloud Baseline), using measures like:
  - Kappa scores for categorical judgments
  - Semantic similarity for content comparison
  - F1 scores for specific feature detection
  - *Primary Comparison:* (Run 4 vs. Run 1) compared to (Run 3 vs. Run 1)
  - *Expected Outcome:* Run 4 outputs closer to Run 1 quality than Run 3
4. **Efficiency Metrics:**
  - Total token usage (input + output tokens)
  - Response latency (milliseconds)
  - *Primary Comparison:* Run 4 vs. Run 3
  - *Expected Outcome:* Higher resource usage with EdgePrompt due to multi-stage processing

### Evaluation Approach:

1. **Direct Metrics:** Safety violations and constraint adherence are measured directly by the **ConstraintEnforcer** component, which checks outputs against pre-defined rules.
2. **Cross-Run Comparisons:** Quality metrics involve comparing outputs from different runs against the reference standard (Run 1), then comparing the degree of alignment between approaches.
3. **Performance Monitoring:** Efficiency metrics are collected by the **MetricsCollector** component, which records token counts from API responses and measures latency for each operation.
4. **Statistical Analysis:** Results are aggregated across multiple test cases to identify statistically significant differences between approaches.



5. **Visualization:** Key metrics are presented through comparative visualizations (bar charts, tables) to clearly communicate differences between approaches.

This evaluation framework focuses specifically on the comparative impact of the EdgePrompt approach on edge-deployable models—testing whether structured prompting and multi-stage validation can meaningfully improve safety and constraint adherence compared to baseline approaches, and at what computational cost.

The metrics deliberately prioritize quantifiable aspects of performance, recognizing that more subjective dimensions like educational effectiveness and teacher satisfaction would require human evaluation in future phases. This pragmatic scope allows for concrete initial validation while acknowledging the limitations of a simulation-based approach.

### **The Critical Role of Simulation in Isolating Variables**

The simulation-based approach of Phase 1 serves a crucial methodological purpose: isolating specific variables to test EdgePrompt's core claims while controlling for external factors. This approach offers several key advantages:

#### **Controlled Comparison Environment:**

- Ensures identical inputs across all four runs through the shared teacher request mechanism
- Eliminates variability in testing conditions that might affect results
- Provides systematic, reproducible execution of the experimental protocol
- Enables direct attribution of performance differences to the specific variables being tested

#### **Systematic Variation of Key Factors:**

- Isolates the effect of model size/capability (CloudLLM vs. EdgeLLM)
- Isolates the effect of prompting approach (Direct vs. EdgePrompt)
- Controls for other factors that might influence performance
- Creates a factorial design that reveals interaction effects between variables

#### **Efficient Initial Validation:**

- Enables testing of core hypotheses before investing in full implementation
- Provides data to guide subsequent development priorities
- Identifies critical issues early in the development process
- Creates a foundation for more targeted real-world testing

#### **Standardized Reference Points:**

- Establishes CloudLLM outputs as a consistent quality benchmark
- Creates comparable metrics across different test cases and configurations

- Provides clear baselines for measuring improvement
- Ensures fair evaluation against consistent standards

This simulation approach is particularly valuable for testing EdgePrompt's specific contribution to performance. By holding all other factors constant and varying only the model type and prompting approach, the methodology can directly measure whether EdgePrompt's structured prompting and multi-stage validation actually improve safety and constraint adherence on edge models.

However, simulation also has inherent limitations that must be acknowledged:

1. **Artificial Context:** Simulation cannot fully replicate the complexities and constraints of real classrooms and edge devices.
2. **LLM Proxy Evaluation:** Using CloudLLM for quality assessment serves as a proxy for human judgment, but cannot replace actual teacher and student evaluation.
3. **Limited Scope:** Phase 1 testing focuses narrowly on the guardrail effectiveness rather than broader educational value or usability.
4. **Idealized Conditions:** Testing occurs under more controlled conditions than would be found in actual deployment scenarios.

These limitations highlight why Phase 1 represents only an initial validation step. Future phases will need to move beyond simulation to hardware testing, human evaluation, and real-world deployment to fully validate EdgePrompt's effectiveness and value. Nevertheless, simulation provides a crucial foundation—enabling systematic testing of core hypotheses before committing to more resource-intensive evaluation approaches. This methodological approach parallels research on AI system evaluation frameworks, which emphasize the importance of component-level checks before moving to system-level validation (Xia et al. 2024).

## Part 2: The Crucible of Implementation – Reality Checks and Engineering Responses

### Chapter 4: The Topic Consistency Crisis – Ensuring Fair Comparison

#### The Problem: Baseline Runs Going Off-Topic

As the EdgePrompt team began implementing the Phase 1 research methodology, they encountered an unexpected challenge that threatened the validity of their entire experimental approach. The problem emerged when examining outputs from the four comparative runs:

**Run 1 (Cloud Baseline)** and **Run 3 (Edge Baseline)** were frequently generating content on entirely different topics than **Run 2 (Cloud EdgePrompt)** and **Run 4 (Edge EdgePrompt)**.

For example, in a test case intended to generate educational content about photosynthesis:

- The EdgePrompt runs (2 and 4) reliably produced content about the specified topic
- The baseline runs (1 and 3) sometimes produced content about completely different subjects like ocean ecosystems or weather patterns

This inconsistency created a fundamental problem: if the different approaches weren't addressing the same topic, how could their outputs be meaningfully compared? The safety violations, constraint adherence, and quality metrics would be comparing fundamentally different content, rendering the results meaningless.

Further investigation revealed the root cause: the baseline approach used more open-ended prompting without the explicit structural constraints of EdgePrompt. This gave the models more freedom to interpret the request and sometimes led them to select topics tangential to or entirely different from what was intended.

This discovery represented a critical methodological flaw. Without addressing it, the research couldn't validly claim that differences in metrics were due to the EdgePrompt approach rather than simply the topics being addressed. The very foundation of the experimental design—systematic comparison holding all variables constant except the approach—was compromised.

This crisis exemplifies the challenges of empirical AI research: apparently minor implementation details can dramatically affect results, and theoretical frameworks must often be adapted in the face of real-world behavior. It also highlights the importance of methodological rigor in validating claims about AI performance, particularly when comparing different approaches. Such challenges reflect broader findings in prompt engineering research, which highlight how different prompting strategies can significantly impact model adherence to specified tasks and constraints (Sahoo et al. 2025).

## The Fix: The "Shared Teacher Request" Mechanism

To address the topic consistency crisis, the research team implemented an elegant solution: the "shared teacher request" mechanism. This approach fundamentally changed how the experimental workflow was structured:

1. At the beginning of each test case, the orchestrator would use CloudLLM to generate a complete, detailed "teacher request" specifying:
  - The precise topic to be addressed
  - Specific learning objectives
  - Explicit constraints (word count, vocabulary level, etc.)
  - Content type expectations (e.g., explanation, question)
  - Any other relevant parameters
2. This same teacher request would then be used as input for all four experimental runs, ensuring they were all attempting the exact same task.
3. For the EdgePrompt runs (2 and 4), this request would be processed through the template system as usual, adding the structured scaffolding that defines the EdgePrompt approach.
4. For the baseline runs (1 and 3), the teacher request would be used directly as the prompt, providing clear topic guidance while still maintaining the less structured nature of the baseline approach.

This solution elegantly balanced two competing needs:

- **Consistency:** Ensuring all approaches addressed the same topic and task
- **Differentiation:** Preserving the distinction between structured (EdgePrompt) and unstructured (baseline) approaches

The "shared teacher request" mechanism created a fair comparison environment while maintaining the essential differences between approaches that the research sought to evaluate. It ensured that any measured differences in safety, constraint adherence, or quality could be more confidently attributed to the prompting approach rather than simply different topics.

This fix represents an important methodological refinement that wasn't anticipated in the original research design. It highlights the iterative nature of AI research, where initial theoretical frameworks often need practical adaptation when implemented. It also demonstrates the research team's commitment to methodological rigor—identifying and addressing a potential validity threat rather than proceeding with a flawed comparison. This approach parallels research on chain-of-thought prompting, which emphasizes the importance of clear task specification in evaluating model reasoning capabilities (Wei et al. 2022).

## **Implication: The Necessity of Controlling Context in Comparative AI Evaluation**

The topic consistency challenge and its solution reveal a broader implication for AI evaluation: the critical importance of context control in comparative studies. This insight extends beyond EdgePrompt to the evaluation of AI systems more generally:

### **Controlling Input Context is Essential:**

- Small variations in prompts can lead to dramatically different outputs
- Valid comparisons require careful standardization of inputs across approaches
- The "freedom" of less constrained approaches must be balanced with fair comparison needs
- Experimental designs must account for the high context sensitivity of language models

### **Evaluation Design Must Match Claims:**

- Claims about relative performance require controlling for all variables except the one being tested
- More structured approaches may inherently provide better topic adherence as a feature
- Evaluation protocols must separate the effects of different aspects of an approach
- Careful attention to what's actually being compared is essential for valid conclusions

### **Hidden Variables Can Undermine Comparisons:**

- Seemingly minor implementation details can have major impacts on results
- Surface-level metrics may hide important qualitative differences in outputs
- Direct side-by-side comparison is necessary to identify unexpected variances
- Assumptions about what remains constant across approaches must be verified

This experience suggests that future AI evaluation work should:

1. Implement explicit mechanisms to ensure input consistency across compared approaches
2. Include qualitative checks of outputs early in the evaluation process to identify unexpected variances
3. Design evaluation protocols that separate the effects of different aspects of an approach
4. Be transparent about methodological challenges and adaptations

For EdgePrompt specifically, this insight reinforces the value of the structured approach—part of EdgePrompt's benefit may be precisely its ability to maintain topic focus and task alignment. However, to fairly evaluate this benefit, the comparison baseline must start from the same topic and task parameters.

The shared teacher request mechanism represents a methodological innovation that could be valuable for other comparative AI evaluations. By ensuring consistent starting points while maintaining the essential differences between approaches, it enables more valid and meaningful comparisons.

This experience also highlights the importance of iterative refinement in research methodology—being willing to adapt approaches when initial implementations reveal unforeseen challenges. Such adaptability, combined with commitment to evaluation validity, is essential for developing reliable knowledge about AI system performance. These findings align with research on AI system evaluation frameworks that emphasize the importance of systematic assessment processes throughout model development lifecycles (Xia et al. 2024).

## Chapter 5: The JSON Reliability Gauntlet – When LLMs Don't Follow Rules

**For Non-Technical Readers:** This section addresses a crucial discovery that threatened the entire project: smaller AI models often struggle to consistently follow specific formatting rules when generating responses. It's similar to how a student might understand the material but fail to follow the required essay format. EdgePrompt needed the AI to produce answers in a very specific technical format (called JSON) so other parts of the system could understand and process them. However, the smaller models kept making formatting mistakes—like students forgetting paragraph breaks or citation formats. The team had to develop a series of backup plans and correction tools to handle these formatting errors. This challenge reveals an important reality about AI in education: even when AI understands the educational content, it may struggle with the structured consistency that educational systems require, necessitating significant behind-the-scenes engineering to make it work reliably.

### The "Major Issue": EdgeLLM Instability with Structured Output

As implementation of the EdgePrompt research framework progressed, the team encountered what they would later describe as the "MAJOR ISSUE" that fundamentally challenged their approach: smaller edge-deployable language models (like Gemma 3 4B and Llama 3.2 3B) proved highly unreliable at consistently generating properly formatted structured outputs, particularly JSON.

This issue struck at the heart of EdgePrompt's multi-stage validation approach, which relied on each validation stage returning structured outputs (like `{"passed": true, "score": 0.8, "feedback": "..."}` ) that could be parsed and used to drive subsequent stages. This challenge parallels findings from research on reliable code generation from pre-trained language models, which identified similar issues with structured outputs and proposed techniques like Target Similarity Tuning and Constrained Semantic Decoding to enforce syntactic and semantic constraints (Poesia et al. 2022).

The problem manifested in multiple ways:

1. **Format Inconsistency:** Models would frequently wrap JSON in triple backticks or markdown code blocks rather than returning pure JSON.
2. **Structural Errors:** Fields might be nested incorrectly, with scores appearing at the top level in some responses and nested under a "result" field in others.
3. **Missing Fields:** Key fields like "passed" or "score" would sometimes be omitted entirely.
4. **Field Name Mismatches:** Field names might vary between responses (e.g., "passed" vs. "isValid" vs. "valid").
5. **Syntactic Errors:** Invalid JSON syntax such as missing commas, unquoted keys, or unmatched brackets.

These issues occurred despite explicit instructions in prompts to return JSON in a specific format. Even when the model seemed to understand the task itself correctly, the output format remained unreliable.

This created a critical challenge: if validation outputs couldn't be reliably parsed, the entire multi-stage validation pipeline would break down. A validation stage might correctly determine that content violated safety constraints, but if that result couldn't be extracted from malformed JSON, the system would be unable to act on that determination.

The problem was particularly acute with smaller edge models, though even larger cloud models occasionally produced malformed outputs. This highlighted a fundamental challenge in relying on language models for structured output generation—while they excel at producing natural language text, they are less reliable at adhering to strict formatting requirements, especially as model size decreases.

This issue wasn't merely a minor implementation detail but a fundamental challenge to EdgePrompt's approach. If edge-deployable models couldn't reliably produce structured outputs, could EdgePrompt's multi-stage validation approach work at all on the edge?

### **Engineering Robustness: `json_utils`, Multi-Pattern Extraction, Fallbacks, Repair**

Faced with this existential challenge to their approach, the EdgePrompt team developed a comprehensive suite of robustness engineering solutions:

**1. The `json_utils` Module:** A dedicated Python module implementing multiple strategies for JSON extraction:

- Standard JSON parsing for well-formatted outputs
- Markdown code block extraction for JSON wrapped in backticks
- Regular expression matching for JSON-like structures
- Permissive parsing tolerating certain common errors
- Special case handling for known model-specific quirks

**2. Multi-Pattern Extraction:** Implementing a cascading approach that tries multiple extraction methods:

- First attempting strict parsing of the entire response
- Then looking for code blocks that might contain JSON
- Searching for JSON-like structures using regular expressions
- Applying progressively more lenient parsing rules

**3. Fallback Template System:** Creating simplified validation templates as backups:

- Reducing the complexity of JSON structures requested from models
- Using simpler field names and flatter structures
- Requesting key fields (like "passed") in predictable formats
- Explicitly instructing models about exact output formatting



**4. JSON Repair Mechanism:** Implementing advanced recovery for severely malformed outputs:

- Using stronger models to repair outputs from weaker ones
- Translating free-text responses into structured formats
- Extracting key information even from completely unstructured responses
- Maintaining graceful degradation rather than complete failure

**5. Error Recovery Flow:** Designing the system to handle parsing failures without breaking the pipeline:

- Setting sensible defaults when fields couldn't be extracted
- Limiting repair attempts to prevent infinite loops
- Preserving raw outputs for manual inspection or later recovery
- Logging detailed information about extraction attempts

This comprehensive approach transformed a potential showstopper into a manageable engineering challenge. While not eliminating the underlying issue—that smaller models struggle with structured output—it created robust mechanisms to handle the reality of model behavior rather than assuming ideal performance.

These solutions demonstrate a key principle of practical AI system development: building systems that work with actual model capabilities rather than theoretically ideal behavior. By explicitly acknowledging and addressing the limitations of current models, EdgePrompt could maintain its multi-stage validation approach even with the constraints of edge deployment. This approach aligns with research on prompt-driven safeguarding for large language models, which has investigated how safety prompts affect model behavior and proposed techniques to improve safeguarding without compromising general performance (Zheng et al. 2024).

**Implication: The Hidden Costs and Engineering Demands of Using Imperfect AI**

The JSON reliability challenge and the extensive engineering required to address it reveal a broader implication about working with AI systems: the substantial hidden costs and engineering demands of using imperfect AI in production systems.

**Expectation vs. Reality Gap:**

- Language models are often presented as capable of following complex instructions
- Marketing materials and demonstrations typically show ideal, cherry-picked examples
- Documentation may imply more reliability in structured output than models actually deliver
- The gap between expected and actual behavior requires significant engineering to bridge

**Engineering Complexity Beyond the Core Logic:**

- A substantial portion of the EdgePrompt codebase addresses robustness rather than core functionality

- The effort required to handle edge cases may exceed the effort to implement the primary features
- Engineering for robustness requires anticipating diverse failure modes
- Systems must be designed for graceful degradation rather than assuming perfect performance

**Quality-Resource Trade-off:**

- Smaller models (necessary for edge deployment) generally exhibit less reliable behavior
- The engineering effort required scales inversely with model capability
- Resource constraints that drive edge deployment simultaneously increase development complexity
- Models exhibiting 95% reliability still fail frequently enough to require robust handling

**Development Lifecycle Implications:**

- Initial prototypes based on ideal behavior may drastically underestimate implementation challenges
- Testing must include diverse, realistic scenarios rather than simply verifying ideal paths
- Maintenance costs increase due to the complexity of robustness mechanisms
- New model versions may introduce new failure modes requiring additional engineering

For EdgePrompt specifically, this experience points to several key lessons:

1. Implementation of AI systems in resource-constrained environments requires substantial robustness engineering beyond the core conceptual approach.
2. The reliability gap between large cloud models and smaller edge models is not just a matter of quality but can affect the fundamental feasibility of certain approaches.
3. Graceful degradation strategies are essential for real-world deployment, ensuring systems can still provide value even when components don't perform ideally.
4. The engineering effort required to make AI systems robust in the face of model limitations should be explicitly factored into development planning and resourcing.

These insights have implications beyond EdgePrompt, serving as a caution for any project deploying AI in production settings, particularly with resource constraints. The "hidden tax" of robustness engineering should be anticipated and accounted for, rather than discovered as a surprise during implementation.

The success of the EdgePrompt team in addressing these challenges demonstrates that such issues can be overcome with appropriate engineering approaches. However, it also highlights that the path from conceptual design to working implementation is rarely straightforward when working with AI systems that exhibit probabilistic rather than deterministic behavior. This aligns with findings from research on on-device language models, which highlights the various challenges in deploying LLMs on resource-constrained devices and the technical strategies needed to address them (Xu et al. 2024).

## Chapter 6: Efficiency vs. Robustness – The Edge Computing Tightrope

**For Non-Technical Readers:** Imagine having to choose between a thorough but slow grading process and a quicker but less comprehensive one when facing a stack of student essays. This section explores a similar fundamental tradeoff in EdgePrompt's design. The system's safety features require multiple rounds of checking (called multi-stage validation), which significantly increases both the time needed to generate responses and the computing power required. The research showed that adding these safety checks made the system 3-5 times slower and used 2-4 times more computational resources. In well-resourced schools with powerful computers, this might not matter much, but in the target environments with limited hardware, this difference could determine whether the system is usable at all. Understanding this tradeoff helps explain why creating AI systems for resource-constrained schools involves more complex decisions than simply using the same approaches from well-funded environments.

### Phase 1 Findings: Quantifying the Latency/Token Overhead of Multi-Stage Validation

The Phase 1 research provided crucial empirical data on one of EdgePrompt's central questions: what is the computational cost of implementing multi-stage validation on edge-deployable language models? This question is fundamental because edge deployment has inherent resource constraints, and any approach that is too computationally expensive simply won't be viable in the target environments.

The research findings quantified this cost through direct measurement of two key metrics:

#### Latency Overhead:

- **Run 4 (Edge EdgePrompt)** exhibited substantially higher response times compared to **Run 3 (Edge Baseline)**
- Multi-stage validation required multiple sequential LLM calls, each adding its own latency
- Total response times could be 3-5 times longer for the EdgePrompt approach
- Actual millisecond measurements varied based on the specific edge model and hardware profile

#### Token Usage Overhead:

- **Run 4** consumed significantly more tokens (both input and output) than **Run 3**
- Each validation stage required its own prompt context and generated its own output
- Token consumption could be 2-4 times higher for the complete EdgePrompt pipeline
- This translates directly to higher computational resource requirements

These measurements highlighted a fundamental trade-off: the multi-stage validation approach that improved safety and constraint adherence came at a substantial computational cost. This cost isn't merely a performance issue—in edge computing environments with limited computational resources, it could determine whether the approach is viable at all.

Additional analysis revealed that different validation stages contributed differently to this overhead:

- Safety checks (typically early in the validation sequence) added moderate overhead
- Semantic evaluation of content (e.g., checking relevance to the topic) was more token-intensive
- Detailed assessment against rubric criteria typically consumed the most resources
- JSON parsing and format handling added further overhead, especially with repair attempts

These findings provide crucial context for evaluating EdgePrompt's approach. While the Phase 1 research confirmed that EdgePrompt's structured prompting and multi-stage validation could improve safety and constraint adherence, it also revealed that this improvement comes at a significant efficiency cost that must be addressed for practical deployment. This efficiency-robustness tradeoff aligns with findings from research on LLMs for forecasting and anomaly detection, which notes the computational resource requirements as a significant challenge when deploying these models in resource-constrained environments (Su et al. 2024).

### **The Inherent Trade-off: Deeper Validation vs. Edge Feasibility**

The efficiency measurements highlight an inherent tension at the heart of EdgePrompt's approach: a fundamental trade-off between validation depth and edge feasibility.

#### **The Validation Depth Dimension:**

- More validation stages provide more thorough safety and quality checking
- Sequential validation allows each stage to focus on specific aspects
- Deeper validation catches more potential issues and edge cases
- Multi-stage approaches enable more nuanced assessment of different criteria

#### **The Edge Feasibility Dimension:**

- Edge devices have limited computational resources
- Response time requirements for interactive educational applications
- Battery consumption considerations for mobile/portable devices
- Thermal constraints in fanless or compact devices

This trade-off cannot be eliminated entirely—more thorough checking will always require more computation. However, it can be navigated through careful design decisions that balance protection and performance:

#### **Potential Navigation Strategies:**

- Prioritizing critical validation stages (e.g., safety) while making others optional
- Implementing more efficient prompt designs for common validation tasks
- Using non-LLM methods (e.g., keyword matching) for simpler constraints
- Applying conditional validation that adapts based on content risk assessment

- Developing early termination mechanisms that skip unnecessary stages
- Optimizing for specific edge models' token efficiency

The research findings suggest that the current implementation may be too heavily weighted toward validation depth at the expense of edge feasibility. While this makes sense for the research phase—where thoroughly validating the approach is the priority—a more balanced approach will likely be needed for practical deployment.

This tension reflects a broader challenge in edge AI: the most capable models (which provide the best quality) are also the most resource-intensive, while the models that run efficiently on edge devices often have more limited capabilities. EdgePrompt's multi-stage approach attempts to get high-quality validation from smaller models through structured prompting, but still faces this fundamental tension between capability and efficiency. This challenge is well-documented in research on educational edge computing frameworks, which demonstrate how educational applications must balance performance requirements with resource constraints (Chen et al. 2022).

The resolution of this trade-off will be crucial for EdgePrompt's future development. Finding the right balance—providing sufficient validation to ensure safety and quality while maintaining acceptable performance on target devices—represents one of the key challenges for moving from research to practical deployment.

### **Implication: The Need for Aggressive Optimization and Realistic Hardware Expectations**

The efficiency findings from Phase 1 point to a clear implication for EdgePrompt's future development: the critical need for aggressive optimization and realistic hardware expectations.

#### **Optimization Imperatives:**

##### **1. Prompt Engineering Efficiency:**

- Redesigning validation prompts to minimize token usage
- Consolidating multiple checks into single stages where possible
- Optimizing the information provided in each stage's context
- Engineering prompts specifically for token efficiency

##### **2. Validation Pipeline Streamlining:**

- Implementing early exit paths for clear pass/fail cases
- Prioritizing stages based on importance and computational cost
- Developing adaptive validation that adjusts depth based on content
- Creating more efficient information flow between stages

##### **3. Model-Specific Optimization:**

- Tuning approaches for specific edge-deployable models
- Exploiting particular models' strengths and working around weaknesses
- Benchmarking different models for validation task efficiency
- Potentially developing specialized smaller models for specific validation tasks

#### 4. **Non-LLM Approaches for Suitable Tasks:**

- Replacing LLM calls with rule-based checks for simpler constraints
- Using traditional NLP techniques for certain validation aspects
- Implementing hybrid approaches combining statistical and symbolic methods
- Reserving LLM validation for tasks requiring deeper understanding

### **Hardware Expectation Management:**

#### 1. **Baseline Performance Requirements:**

- Establishing minimum hardware specifications for acceptable performance
- Setting realistic expectations for response times on different device classes
- Defining graceful degradation paths for less capable hardware
- Communicating hardware requirements clearly to potential implementers

#### 2. **Device Class Targeting:**

- Focusing initially on more capable edge devices (e.g., NVIDIA Jetson class)
- Creating tiered functionality based on available computational resources
- Developing hardware-aware configurations that adapt to device capabilities
- Establishing a roadmap for supporting progressively less powerful devices

#### 3. **Deployment Profiling:**

- Conducting systematic testing on representative target hardware
- Measuring real-world performance under varying conditions
- Identifying specific bottlenecks in the execution pipeline
- Using profiling data to guide optimization efforts

This dual focus on optimization and hardware expectations acknowledges the reality that EdgePrompt's current approach is likely too computationally intensive for practical deployment on many edge devices. Rather than compromising on the core vision, this approach seeks to bridge the gap through a combination of technical optimization and clear expectation setting.

The need for this work highlights a broader point about AI deployment in resource-constrained environments: theoretical approaches must be rigorously tested against actual resource constraints, and development must include substantial optimization efforts specifically for target environments. This requirement adds to the already significant engineering demands of working with imperfect AI systems, further emphasizing the gap between conceptual design and practical implementation.

For EdgePrompt specifically, this points to a critical Phase 2 priority: measuring and optimizing performance on actual target hardware rather than simulated environments. Only through such testing can the team determine whether the approach can meet both the safety/quality goals and the performance requirements necessary for practical deployment in resource-constrained educational settings. This aligns with research on on-device language models, which reviews strategies for efficient deployment including architecture optimization, compression techniques, and hardware acceleration (Xu et al. 2024).

## Part 3: Assessing EdgePrompt Through Critical Lenses

### Chapter 7: Pedagogical Alignment – Tool, Tutor, or Obstacle?

#### EdgePrompt vs. Educational Ideals (Constructivism, Higher-Order Thinking)

To meaningfully assess EdgePrompt's potential educational impact, we must evaluate it against established educational ideals and theories. This examination reveals both promising alignments and significant tensions:

**Alignment with Constructivist Learning Theory:** Constructivism posits that learners actively construct knowledge through experience and reflection rather than passively receiving information. EdgePrompt shows potential alignment through:

- The vision of supporting rather than replacing teacher-guided learning
- The goal of promoting productive struggle rather than simply providing answers
- The emphasis on feedback that guides further exploration

This approach aligns with research by Grubaugh, Levitt and Deever (2023), who argue that AI tools can amplify constructivist pedagogy when designed to support learners as "active architects" of their own knowledge building rather than passive recipients of information.

However, tensions emerge in the current implementation:

- The Q&A loop focus may inadvertently reinforce transmission models of learning
- Automated validation could prioritize correctness over process or creative thinking
- The structured approach might constrain the open-ended exploration central to constructivism

These tensions reflect broader concerns about maintaining epistemic agency in educational technology, which Elgin (2013) defines as the capacity for learners to actively evaluate and construct knowledge rather than passively receive it.

**Support for Higher-Order Thinking:** Referencing Bloom's Taxonomy (which progresses from lower-order skills like remembering and understanding to higher-order skills like analyzing, evaluating, and creating), EdgePrompt's alignment varies:

- *Potential Strengths:* The framework could support analysis and evaluation through well-designed prompts and feedback
- *Current Limitations:* The implementation primarily addresses lower-order skills through direct question-answering and basic assessment
- *Future Opportunity:* Templates could be designed specifically to scaffold progression toward higher-order thinking

Research on educational question generation at different Bloom's skill levels using large language models suggests approaches for crafting questions that target varying cognitive demands, from basic recall to complex analysis and evaluation (Scaria, Dharani Chenna & Subramani 2024).

**Vygotskian Scaffolding and Zone of Proximal Development:** Vygotsky's theory emphasizes learning as a social process where novices progress with appropriate scaffolding within their Zone of Proximal Development (ZPD):

- EdgePrompt's structured templates could function as scaffolds for learning
- The multi-stage validation approach could provide graduated feedback
- However, the current system lacks mechanisms to identify individual students' ZPDs
- Without adaptation capabilities, scaffolding remains static rather than responsive

This theoretical framework, developed by Vygotsky (1978), highlights the gap between what learners can do independently and what they can achieve with assistance, emphasizing the role of social interaction and scaffolding in cognitive development.

**Cultural Responsiveness and Inclusion:** Educational ideals increasingly emphasize culturally responsive pedagogy that values diverse perspectives:

- EdgePrompt's vision includes adaptability to different cultural and linguistic contexts
- The offline capability addresses infrastructural inequities
- However, base models may contain biases that affect content generation
- True cultural responsiveness would require significant localization efforts

Culturally responsive teaching requires incorporating students' cultural references in all aspects of learning, validating and affirming diverse cultural identities, which becomes particularly important when integrating technology in multicultural classrooms (Gay 2018).

This assessment reveals that EdgePrompt's current implementation remains some distance from fully embodying educational ideals. While its vision shows alignment with constructivist principles and the potential to support higher-order thinking, the present Q&A and validation focus operates primarily at lower cognitive levels. The gap between educational ideals and current implementation highlights areas for future development—particularly in supporting more diverse learning activities, adapting to individual learners, and integrating more explicitly constructivist approaches.

### **Analyzing the Current Q&A Loop: Strengths and Limitations**

The current implementation of EdgePrompt centers on a relatively straightforward question-answer-validation loop. This approach has both strengths and limitations from a pedagogical perspective:

#### **Strengths of the Current Approach:**



1. **Structured Scaffolding:** The template-based question generation provides consistent structure that can scaffold learning activities.
2. **Immediate Feedback:** The validation system offers prompt response to student answers, potentially supporting iterative improvement.
3. **Clear Expectations:** Rubric-based assessment can make learning expectations explicit to students.
4. **Teacher Guidance:** The teacher-driven nature of content generation maintains human pedagogical direction.
5. **Offline Accessibility:** The approach makes AI-supported learning activities accessible in disconnected environments.

#### **Limitations of the Current Implementation:**

1. **Narrow Interaction Model:** The Q&A loop represents just one type of educational interaction, missing many other valuable learning activities.
2. **Potential Superficiality:** Automated validation may focus on surface features rather than deep understanding.
3. **Limited Dialogue:** The current implementation lacks true back-and-forth dialogue that characterizes rich educational interactions.
4. **Cookie-Cutter Risk:** Template-based generation could lead to formulaic, repetitive learning activities.
5. **One-Size-Fits-All:** Without adaptation capabilities, the system cannot tailor interactions to individual students.
6. **Fixed Scaffolding:** The scaffolding provided remains static rather than gradually fading as students develop mastery.

#### **Pedagogical Implications:**

From a pedagogical standpoint, the current Q&A loop resembles traditional "initiation-response-evaluation" (IRE) patterns that have been critiqued in educational research. While IRE can be effective for certain types of knowledge checking, it has limitations for developing deeper understanding, critical thinking, and student agency.

The use of templates for question generation and validation could inadvertently reinforce what Paulo Freire described as the "banking model" of education—treating students as receptacles for depositing knowledge rather than active co-creators of understanding. While this is not the intention of EdgePrompt, the technical constraints of the current implementation may push in this direction.

However, it's important to recognize that the current implementation represents just the beginning of EdgePrompt's development. The vision articulated in the conceptual framework points toward more sophisticated educational interactions, including adaptive systems that learn from patterns and more diverse interaction types beyond simple Q&A.

The gap between the current implementation and educational ideals is not a failure but rather an opportunity for guided evolution. By recognizing both the strengths and limitations of the current approach, EdgePrompt can develop toward more pedagogically rich interactions while maintaining its core commitment to offline capability, teacher agency, and safety.

### **The Risk of Shallow Learning vs. Potential for Scaffolding**

One of the most significant pedagogical tensions in EdgePrompt's current design lies between the risk of promoting shallow learning and the potential for meaningful scaffolding:

#### **The Risk of Shallow Learning:**

AI-driven educational tools, including EdgePrompt, face several inherent risks that could lead to shallow rather than deep learning:

1. **Answer-Seeking Behavior:** Students may focus on producing responses that satisfy the validation system rather than developing genuine understanding.
2. **Template Dependency:** Formulaic questions and assessments might train students to recognize patterns rather than truly engage with content.
3. **Feedback Simplification:** Automated validation may reduce complex concepts to binary correctness judgments or simplified rubrics.
4. **Process Bypassing:** The focus on final answers could inadvertently devalue the learning process, including productive struggle and iteration.
5. **Memorization Over Meaning:** Without careful design, the system could reward memorization of facts rather than conceptual understanding.

These risks are not unique to EdgePrompt but represent common challenges in educational technology. The current implementation, with its focus on Q&A and validation, may be particularly susceptible to these concerns unless deliberately designed to counter them.

#### **The Potential for Meaningful Scaffolding:**

Despite these risks, EdgePrompt also holds significant potential for providing meaningful scaffolding that supports deeper learning:

1. **Structured Cognitive Guidance:** Well-designed templates could guide students through complex thinking processes, not just prompt for answers.
2. **Progressive Challenge:** The system could implement scaffolding that gradually increases in complexity as students demonstrate mastery.
3. **Process-Oriented Feedback:** Validation could focus on identifying specific areas for improvement rather than just correctness.
4. **Metacognitive Prompting:** Questions could incorporate metacognitive elements that prompt students to reflect on their thinking.

5. **Multimodal Scaffolding:** Different template types could support various learning modalities and approaches.

For this potential to be realized, EdgePrompt would need to evolve beyond simple Q&A to incorporate more sophisticated scaffolding strategies. This might include question sequences that gradually build in complexity, feedback that highlights thinking processes rather than just outcomes, and templates specifically designed to elicit metacognition.

### **Navigating the Tension:**

The path forward involves deliberately designing EdgePrompt to maximize scaffolding potential while minimizing shallow learning risks:

1. Develop templates that emphasize process over product, including questions about reasoning and methodology
2. Implement validation approaches that assess multiple dimensions of understanding rather than binary correctness
3. Create feedback mechanisms that guide further exploration rather than simply evaluating responses
4. Build in metacognitive prompts that encourage students to reflect on their learning process
5. Ensure teacher involvement in interpreting validation results to add human judgment to automated assessment

This tension between shallow learning and meaningful scaffolding represents not just a challenge for EdgePrompt but a fundamental consideration for all AI in education. By explicitly acknowledging and addressing this tension, EdgePrompt has an opportunity to demonstrate how AI can support deeper rather than shallower learning, even within the constraints of edge deployment.

### **First Principles: Reconciling Probabilistic AI with Meaning Construction**

At a fundamental level, EdgePrompt must navigate a profound philosophical tension: how to reconcile the probabilistic, pattern-matching nature of language models with the meaning-construction processes central to genuine learning. This requires examining both the nature of learning and the nature of LLMs:

#### **The Nature of Learning as Meaning Construction:**

Learning, from first principles, involves:

- Active construction of meaning rather than passive reception of information
- Integration of new knowledge with existing cognitive structures
- Development of mental models that explain relationships and principles
- Conceptual change requiring metacognitive awareness
- Social negotiation of understanding within cultural contexts

These processes are inherently:

- Non-linear and often messy
- Deeply personal yet socially situated
- Involving both conscious and unconscious dimensions
- Requiring emotional and motivational engagement
- Embedded in identity and community

### **The Nature of Language Models as Pattern Matchers:**

Large language models, by contrast, operate through:

- Statistical prediction of token sequences based on training patterns
- No inherent understanding of meaning or truth
- No conscious experience or metacognitive awareness
- No intrinsic motivation or purpose
- No personal identity or cultural situatedness

Their outputs are:

- Probabilistic rather than reflective of understanding
- Fluent but potentially ungrounded
- Context-sensitive but without stable mental models
- Impressive but fundamentally different from human cognition

Research on the scope and limits of passive AI highlights this fundamental distinction between generating meaning through active inference (as humans do) and the pattern-matching nature of language models, with significant implications for educational applications (Pezzulo et al. 2024).

### **The Fundamental Tension:**

This creates a profound tension: how can a system based on probabilistic pattern matching effectively support the deeply human process of meaning construction? This is not merely a technical question but a philosophical one that goes to the heart of what learning is and how it can be supported.

### **Potential Reconciliation Approaches:**

EdgePrompt's architecture suggests several potential approaches to navigating this tension:

1. **Structured Symbolic Scaffolding:** Using templates and validation as a symbolic layer that gives structure to the otherwise ungrounded LLM responses.
2. **Human-AI Partnership:** Positioning the teacher as the meaning-maker who leverages AI as a tool rather than delegating meaning-making to the AI.

3. **Transparency About Limitations:** Being explicit with students about the nature of AI responses, fostering critical engagement rather than acceptance.
4. **Process Over Product:** Focusing validation on thinking processes rather than just factual correctness.
5. **Contextual Grounding:** Ensuring that AI interactions are embedded in broader learning contexts that include human relationships and real-world applications.

From first principles, EdgePrompt cannot eliminate this tension—it is inherent in using pattern-matching systems to support meaning-making processes. However, by explicitly acknowledging the tension and designing with it in mind, EdgePrompt could potentially create a productive partnership between human meaning-makers and AI pattern-matchers.

This would require moving beyond viewing EdgePrompt as simply a way to make AI work offline, toward seeing it as a fundamentally different approach to AI in education—one that acknowledges the limitations of LLMs and deliberately designs interfaces that complement human meaning-making rather than attempting to replicate it.

Such an approach would align with the vision of EdgePrompt as infrastructure rather than authority, supporting human educational processes rather than replacing them. It would require continued evolution of the system's design, with particular attention to how templates, validation, and interfaces can bridge the gap between pattern matching and meaning construction.

## Chapter 8: Human-Computer Interaction – Inside the Lived Classroom

### Evaluating the Current UI against Teacher Workload and Stress

The current EdgePrompt user interface must be evaluated against the lived reality of teaching—a profession often characterized by significant workload pressures, limited time for learning new tools, and sometimes high stress levels. While formal user studies with teachers have not yet been conducted, we can assess the current implementation against known teacher needs and constraints:

#### Interface Analysis Against Teacher Workload Factors:

1. **Setup and Configuration Demands:**
  - *Current State:* The project/template management system requires multiple setup steps before productive use.
  - *Workload Impact:* Additional setup time could create an adoption barrier, especially in time-constrained environments.
  - *Stress Factor:* Complex setup processes may increase cognitive load during already busy preparation periods.
2. **Workflow Integration:**
  - *Current State:* EdgePrompt operates as a standalone application rather than integrating with existing educational tools.

- *Workload Impact*: Teachers must manage yet another system alongside existing tools (LMS, gradebooks, etc.).
- *Stress Factor*: Context switching between different systems increases cognitive burden and potential for errors.

### 3. **Feedback Interpretation and Action:**

- *Current State*: The system provides structured feedback but requires teacher interpretation.
- *Workload Impact*: Review of AI-generated feedback adds another layer to assessment processes.
- *Stress Factor*: Determining when to override or modify AI assessments creates additional decision points.

### 4. **Technical Management Responsibilities:**

- *Current State*: Local deployment requires technical management of the application and models.
- *Workload Impact*: Technical issues could require troubleshooting time better spent on teaching.
- *Stress Factor*: Technical problems during classroom activities can create significant stress and disruption.

### 5. **Learning Curve:**

- *Current State*: The system introduces new concepts (templates, validation sequences) unfamiliar to most educators.
- *Workload Impact*: Time must be invested in learning these concepts before productive use.
- *Stress Factor*: Complex interfaces may create anxiety about correct usage, especially under time pressure.

Research on teacher agency in generative AI contexts emphasizes the importance of maintaining teachers' power to act, affect matters, make decisions, and take stances in educational settings (Frøsig & Romero 2024). Additionally, studies of teacher agency in pedagogical uses of ICT show that structured reflection time on technology use enhances teachers' sense of agency, helping them align technology with local context and values rather than passively adopting standardized tools (Novoa-Echaurren 2024).

### **Positive Aspects of the Current Design:**

1. The project organization provides a structured way to manage different educational contexts.
2. Template reuse could reduce repetitive work over time.
3. The backend-first approach shields teachers from direct LLM interaction complexity.
4. Material uploading and processing is relatively straightforward.
5. The Q&A workflow aligns with familiar educational activities.

**Areas Requiring Improvement:**

1. The multi-step process from project creation to student validation would benefit from streamlining.
2. Template creation and management appears technical and could intimidate non-technical users.
3. Material processing lacks progress indicators and transparency about what's happening.
4. The relationship between uploaded materials, generated questions, and student responses needs clearer visualization.
5. Error handling and recovery paths for when things go wrong need strengthening.

Without direct teacher feedback, this assessment remains somewhat speculative. However, it highlights a crucial consideration: EdgePrompt's success will depend not just on its technical capabilities but on how well it fits into the already demanding workflow of teachers. Reducing rather than adding to teacher workload and stress must be a central design goal as the interface evolves.

The next development phase should prioritize user research with actual teachers in target environments, focusing particularly on workflow integration, learnability, and stress reduction. This research should inform interface redesign efforts aimed at making EdgePrompt a natural extension of teaching practice rather than an additional burden.

**Simplicity vs. Control: Finding the Balance for Non-Technical Educators**

One of the most significant HCI challenges for EdgePrompt lies in balancing simplicity (making the system accessible to non-technical educators) with control (providing the flexibility and customization that effective teaching requires). This tension manifests in several key dimensions:

**Template Management:**

- *Simplicity Need:* Many teachers want ready-to-use templates without technical complexity
- *Control Need:* Teachers need to adapt templates to their specific educational contexts
- *Current State:* The template system provides control but may appear technical and complex
- *Balance Challenge:* How to provide template customization without requiring JSON editing or complex formatting

**Validation Configuration:**

- *Simplicity Need:* Teachers need straightforward ways to assess student work
- *Control Need:* Assessment criteria vary widely across subjects, grades, and pedagogical approaches
- *Current State:* Validation sequences are powerful but potentially overwhelming for novice users
- *Balance Challenge:* How to enable sophisticated assessment without exposing the full complexity of multi-stage validation

**Content Generation Parameters:**

- *Simplicity Need*: Teachers need quick content generation without too many decisions
- *Control Need*: Educational content must meet specific objectives and constraints
- *Current State*: Generation offers control but requires understanding of multiple parameters
- *Balance Challenge*: How to make parameters meaningful educationally without overwhelming users

**AI Interaction Management:**

- *Simplicity Need*: Teachers shouldn't need to understand LLM prompting details
- *Control Need*: Teachers should be able to guide AI output direction and quality
- *Current State*: Backend-first approach shields from direct prompt management but may feel like a black box
- *Balance Challenge*: How to provide intuitive control without requiring technical AI knowledge

Recent research with middle and high school teachers has uncovered their specific information needs for classroom integration of ChatGPT, revealing significant gaps in understanding how to explore AI capabilities for specific learning tasks and the need for interactive model documentation that supports teacher empowerment (Tan & Subramoniam 2024).

**Potential Approaches to Balancing Simplicity and Control:**

1. **Progressive Disclosure**: Implement interfaces that reveal complexity progressively as users become more comfortable
  - Start with simplified options covering common use cases
  - Provide "advanced" sections for users who need more control
  - Include contextual help explaining the implications of different choices
2. **Educational Metaphors**: Frame technical concepts using familiar educational terminology
  - Replace "template" with "activity type" or similar classroom language
  - Present validation in terms of rubrics and assessment criteria
  - Use educational rather than technical categorization schemes
3. **Presets with Customization**: Provide ready-made configurations that can be adjusted
  - Include curated template libraries for different subjects/grades
  - Allow modification of presets without starting from scratch
  - Enable saving of customized templates for reuse
4. **Guided Setup Processes**: Create wizard-like interfaces for complex tasks
  - Break template creation into logical, educational steps
  - Provide examples and previews at each stage
  - Include contextual suggestions based on subject/grade
5. **Visual Configuration**: Replace text-based configuration with visual alternatives
  - Implement drag-and-drop rubric builders
  - Use visual scales for setting parameters



- Provide visual previews of configuration outcomes

Finding the right balance between simplicity and control is not merely an interface design challenge but a fundamental question about EdgePrompt's relationship with teachers. The goal should be to create a system that respects teacher expertise while reducing unnecessary complexity—enhancing teacher agency rather than requiring technical skills unrelated to pedagogy.

This balance will likely evolve through iterative design informed by actual teacher usage, but the guiding principle should be clear: technical complexity should never be a barrier to pedagogical effectiveness. The interface should make the simple things easy while making the complex things possible, allowing teachers to focus on educational outcomes rather than system operation.

## **Beyond Usability: Designing for Trust, Transparency, and Relational Dynamics**

Effective educational technology must go beyond mere usability to address deeper human factors that shape its acceptance and impact. For EdgePrompt, three interconnected dimensions deserve particular attention: trust, transparency, and relational dynamics.

### **Trust Building in the Human-AI Educational Relationship:**

Trust is foundational to effective educational technology adoption, particularly for AI systems that may raise concerns about reliability, appropriateness, and agency:

#### **1. Establishing Initial Trust:**

- Clear communication about EdgePrompt's capabilities and limitations
- Explicit framing as a teacher tool rather than autonomous educator
- Credible validation of safety claims with transparent evidence
- Respecting teacher authority through visible approval mechanisms

#### **2. Maintaining Trust Through Usage:**

- Consistent behavior that aligns with educational expectations
- Predictable patterns in generation and validation
- Reliability in offline operation without unexpected failures
- Graceful handling of edge cases and errors

#### **3. Repairing Trust When Issues Arise:**

- Clear explanation of what went wrong and why
- Straightforward mechanisms to override or correct problems
- Logging of issues for systematic improvement
- Demonstration of learning from mistakes over time

Research on trust in AI-assisted decision-making has shown that cognitive forcing functions can reduce overreliance on AI by prompting users to think independently before accepting AI recommendations (Buçinca, Malaya & Gajos 2021). Additionally, work on formalizing trust in artificial intelligence defines

the prerequisites, causes, and goals of human trust in AI systems, providing frameworks for designing trustworthy AI interactions (Jacovi et al. 2021).

### **Transparency in AI-Mediated Educational Activities:**

Transparency is essential for meaningful teacher oversight and student understanding:

#### **1. Process Transparency:**

- Visibility into how content is generated and validated
- Clear indication of when AI vs. human judgment is operating
- Accessible logs of system actions and decisions
- Understandable explanations of assessment rationales

#### **2. Data Transparency:**

- Clarity about what data is stored and where
- Explicit information about data usage and sharing
- Local control over student information
- Transparency about model capabilities and limitations

#### **3. Improvement Transparency:**

- Visibility into how the system evolves over time
- Clear communication about updates and changes
- Feedback loops showing how teacher input shapes development
- Openness about ongoing challenges and limitations

### **Relational Dynamics in the Classroom:**

Technology doesn't exist in isolation but within complex educational relationships:

#### **1. Teacher-Student Relationship:**

- How does EdgePrompt affect teacher presence and authority?
- Does it strengthen or weaken teacher-student connection?
- How can it support rather than disrupt relational teaching?
- What happens to student perceptions of teacher feedback?

#### **2. Student-AI Relationship:**

- How do students conceptualize the AI component?
- What expectations do they develop about AI assistance?
- How does interacting with AI affect learning self-concept?
- What boundaries should be established for healthy interaction?

#### **3. Teacher-AI Relationship:**

- How can EdgePrompt become a trusted partner rather than a tool?
- What agency balance feels appropriate to teachers?
- How can the system respect teacher expertise while providing value?
- What communication patterns create productive collaboration?

#### 4. **Classroom Community Dynamics:**

- How does EdgePrompt affect peer collaboration and discussion?
- Can it support rather than isolate community learning?
- How might it change classroom discourse patterns?
- What new social practices might emerge around its use?

#### **Design Implications Beyond Interface:**

Addressing these dimensions requires thinking beyond traditional interface design to consider:

1. **Communication Design:** How the system explains itself and its actions
2. **Relationship Design:** How the system positions itself within educational relationships
3. **Agency Design:** How control and decision-making are distributed
4. **Identity Design:** How the system presents itself and its role
5. **Practice Design:** How the system integrates into educational activities and routines

These considerations point to a crucial insight: EdgePrompt is not merely a technical tool but a socio-technical intervention that will inevitably reshape educational relationships and practices. Designing thoughtfully for trust, transparency, and relational dynamics is essential for ensuring that this reshaping enhances rather than diminishes the human core of education.

#### **Accessibility Considerations**

For EdgePrompt to truly fulfill its mission of enhancing educational equity, accessibility must be a central consideration—ensuring that the system is usable by students and teachers with diverse abilities. Current accessibility concerns include:

#### **Interface Accessibility:**

##### 1. **Visual Accessibility:**

- *Current State:* The application uses Bootstrap for styling, which provides basic accessibility features, but lacks comprehensive screen reader optimization.
- *Needed Improvements:* Implementation of ARIA attributes, proper heading structure, focus management, and high-contrast options.
- *Impact:* Without these improvements, visually impaired teachers or students may be unable to use the system effectively.

##### 2. **Motor Accessibility:**

- *Current State:* The interface relies heavily on traditional pointer-based interaction without explicit keyboard navigation support.
- *Needed Improvements:* Complete keyboard accessibility, logical tab order, and potentially shortcuts for common operations.
- *Impact:* Teachers or students with motor impairments may struggle with basic system operations.

### 3. **Cognitive Accessibility:**

- *Current State:* The multi-step workflow and technical terminology may create barriers for users with cognitive disabilities.
- *Needed Improvements:* Simplified workflows, clear step indicators, consistent patterns, and plain language explanations.
- *Impact:* Users with cognitive disabilities may find the system confusing or overwhelming without these adaptations.

## **Content Accessibility:**

### 1. **Generated Content:**

- *Current State:* No specific mechanisms ensure that AI-generated content is accessible.
- *Needed Improvements:* Prompts that encourage accessibility in generated content (e.g., image descriptions, clear structure).
- *Impact:* Generated educational materials may inadvertently exclude students with disabilities.

### 2. **Multimedia Integration:**

- *Current State:* The system handles text well but has limited support for accessible multimedia.
- *Needed Improvements:* Support for alternative formats, captions, transcripts, and descriptions when working with multimedia content.
- *Impact:* Multimedia educational materials may not be equally accessible to all students.

### 3. **Validation Accessibility:**

- *Current State:* Validation focuses on standard text responses without accommodation for different response types.
- *Needed Improvements:* Ability to validate diverse response formats (audio, simplified text, etc.) equitably.
- *Impact:* Students who communicate in non-standard ways may be disadvantaged in assessment.

## **Contextual Accessibility:**

### 1. **Device Compatibility:**

- *Current State:* The application assumes standard computing environments without explicit support for assistive technologies.
- *Needed Improvements:* Testing with screen readers, alternative input devices, and assistive technologies commonly used in educational settings.
- *Impact:* The system may be incompatible with specialized devices used by students or teachers with disabilities.

### 2. **Environmental Adaptability:**

- *Current State:* Limited consideration of diverse usage environments in design.

- *Needed Improvements:* Adaptations for high-distraction environments, low-bandwidth settings, or special education contexts.
- *Impact:* The system may be unusable in certain educational environments where accessibility needs are greatest.

### 3. **Cognitive Load Management:**

- *Current State:* Complex workflows may create high cognitive loads.
- *Needed Improvements:* Options to simplify interfaces, break tasks into smaller steps, and provide additional scaffolding.
- *Impact:* Teachers or students with attention or executive function challenges may struggle without these accommodations.

Research on algorithmic bias in education has shown that AI systems can inadvertently disadvantage students from various groups, including those with disabilities, unless specifically designed with accessibility and inclusivity in mind (Baker & Hawn 2022).

### **Implementation Path:**

Addressing these accessibility concerns should follow a structured approach:

1. **Audit and Baseline:** Conduct a formal accessibility audit against WCAG 2.1 AA standards
2. **Prioritization:** Address critical barriers first, particularly those affecting teachers
3. **Integration:** Build accessibility into the development process rather than treating it as an add-on
4. **Testing:** Include users with disabilities in testing and feedback cycles
5. **Documentation:** Provide clear accessibility documentation for implementers and users

By prioritizing accessibility, EdgePrompt can ensure that its goal of enhancing educational equity extends to all students and teachers, regardless of ability. This aligns with both the moral imperatives of inclusive education and legal requirements in many jurisdictions, making it an essential consideration for future development.

## Chapter 9: AI Ethics and Alignment – Safety, Bias, and Moral Formation

**For Non-Technical Readers:** This section examines how well EdgePrompt's technical safety features actually work in practice and how they might shape students' values. Think of it as evaluating both the effectiveness of guardrails on a playground and how those guardrails influence children's play behavior. The research found that EdgePrompt's approach does improve safety compared to using AI without these controls, but the protection isn't perfect—there are still ways inappropriate content might slip through. More importantly, the section explores how the AI system implicitly teaches students certain values: Does it encourage critical thinking or just correct answers? Does it prioritize efficiency over exploration? Does it respect diverse cultural perspectives? These questions move beyond technical performance metrics to consider how the technology might subtly shape students' understanding of what education is and what knowledge matters, highlighting why technical design choices have profound educational and ethical implications.

### Guardrails in Practice: Effectiveness and Brittleness

The concept of guardrails—mechanisms to ensure AI systems operate within safe and appropriate boundaries—stands at the center of EdgePrompt's approach. The Phase 1 research provides initial data on the effectiveness of EdgePrompt's prompt-engineering guardrails, but also reveals potential brittleness in these mechanisms:

#### Demonstrated Effectiveness:

Phase 1 findings show promising results for EdgePrompt's guardrail approach:

1. **Safety Enforcement:** The structured prompting and multi-stage validation approach demonstrated measurable improvements in preventing inappropriate content generation compared to baseline approaches.
2. **Constraint Adherence:** EdgePrompt-guided outputs more consistently adhered to specified constraints like word count, vocabulary level, and topic relevance.
3. **Quality Maintenance:** Despite using smaller models, EdgePrompt's approach produced outputs that more closely aligned with high-quality reference standards than baseline approaches with the same models.
4. **Decomposed Checking:** Breaking validation into distinct stages allowed for more thorough assessment of different aspects of generated content.

These findings suggest that EdgePrompt's guardrail approach can effectively improve the safety and reliability of edge-deployed language models, at least within the contexts tested in Phase 1. This is significant because it indicates that prompt engineering alone—without requiring model fine-tuning or constant cloud connectivity—can meaningfully improve model behavior for educational contexts. This aligns with research on building guardrails for large language models through structured prompting and validation stages to enhance safety in AI outputs (Dong et al. 2024).

## Revealed Brittleness:

However, implementation also exposed several forms of brittleness in the current guardrail approach:

1. **JSON Parsing Vulnerability:** The difficulty in reliably extracting structured outputs from smaller models required extensive robustness engineering. This suggests that any approach relying on structured outputs from edge models may face similar challenges.
2. **Sequential Dependency:** The multi-stage validation approach creates a chain of dependencies where failure at any stage can potentially derail the entire process, requiring careful fallback handling.
3. **Context Limitations:** Edge models' more limited context windows constrain the amount of information that can be included in prompts, potentially limiting the sophistication of guardrails that can be implemented.
4. **Efficiency-Robustness Tradeoff:** More thorough checking requires more computational resources, creating tension between guardrail effectiveness and edge deployment feasibility.
5. **Unexpected Edge Cases:** Testing revealed that models sometimes produce unexpected outputs that bypass intended constraints in ways difficult to anticipate—suggesting that no guardrail system will be perfectly comprehensive.

Recent research on prompt-driven safeguarding for large language models has investigated how safety prompts affect LLM behavior in the representation space, proposing techniques to improve safeguarding without compromising general performance (Zheng et al. 2024).

## Implications for Educational Use:

These observations have important implications for EdgePrompt's application in educational settings:

1. **Graceful Degradation:** Systems must be designed to fail safely and gracefully when guardrails encounter unexpected inputs or model behaviors.
2. **Human Oversight:** Teacher review remains essential, particularly for edge cases that automated validation might miss or mishandle.
3. **Transparency About Limitations:** Both teachers and students should understand the limitations of the guardrail approach to maintain appropriate trust levels.
4. **Continuous Refinement:** Guardrails must evolve based on observed behaviors and edge cases in real educational usage.
5. **Complementary Approaches:** Prompt-based guardrails should be complemented by other safety mechanisms where possible, creating defense in depth.

The effectiveness and brittleness revealed in Phase 1 highlight both the promise and the challenges of EdgePrompt's approach. While the guardrails show meaningful improvement over baseline approaches, they remain imperfect—requiring careful design, robust implementation, and appropriate human oversight to function effectively in educational contexts.

## Data Privacy as Pedagogical Trust: Assessing the Implementation Gap

Data privacy in educational contexts is not merely a compliance issue or technical concern—it forms a crucial foundation for pedagogical trust. When students and teachers engage with technology, their willingness to participate authentically depends significantly on confidence that their data will be handled respectfully, responsibly, and transparently. For EdgePrompt, this presents both a significant challenge and a potential advantage:

### The Trust Foundation of Educational Data:

1. **Student Vulnerability:** K-12 students represent a particularly vulnerable population whose data requires special protection:
  - Legal frameworks like FERPA (US), GDPR (EU), and various national regulations
  - Ethical obligations beyond legal minimums
  - Heightened sensitivity regarding minors' data
2. **Types of Sensitive Data in Educational Contexts:**
  - Student identities and demographic information
  - Learning difficulties and accommodations
  - Performance and assessment data
  - Written responses revealing personal experiences or perspectives
  - Interaction patterns suggesting learning styles or challenges
3. **Privacy as Educational Necessity:**
  - Students need safe spaces to make mistakes without permanent records
  - Teachers need confidence that student data won't be exploited
  - Educational exploration requires freedom from excessive surveillance
  - Trust enables authentic engagement rather than strategic performance

Research on learning analytics ethics provides a socio-critical perspective on student data ethics, identifying challenges around informed consent, privacy, data ownership, and transparency, while emphasizing that students should be active agents of their data, not mere data points (Slade & Prinsloo 2013).

### EdgePrompt's Current Privacy Posture:

The existing EdgePrompt implementation shows both promising foundations and significant gaps related to data privacy:

1. **Promising Foundations:**
  - *Offline-First Design:* The core architecture reduces data transmission risks
  - *Local Storage:* Data remains primarily on local devices rather than central servers
  - *Backend-First LLM Interaction:* Creates a clear boundary for LLM-related data
  - *Design Intent:* Documentation indicates privacy as a priority
2. **Significant Implementation Gaps:**



- *Limited Authentication*: No robust user identification or authentication system
- *Minimal Access Controls*: Lack of role-based permissions for data access
- *Absence of Encryption*: No encryption for data at rest or in transit
- *Undefined Data Lifecycle*: No clear policies for data retention or deletion
- *Limited Consent Mechanisms*: No structured approach to gathering and managing consent

### 3. Architectural Considerations:

- *Database Structure*: SQLite provides a self-contained database, but lacks built-in encryption
- *Data Synchronization*: Future sync mechanisms will need careful privacy design
- *Model Interaction*: Local LLM usage reduces but doesn't eliminate privacy concerns
- *Template Management*: Templates may contain confidential educational content

Recent research on the impact of generative AI in cybersecurity and privacy has examined the implications for student data protection in educational contexts, highlighting the need for robust security and privacy measures in AI-enhanced educational tools (Gupta et al. 2023). Additionally, research on AI misuse and information privacy concerns has validated new constructs for measuring privacy concerns related to AI systems and emphasized the importance of robust safety guardrails in addressing these concerns (Menard & Bott 2024).

### The Privacy-Trust Implementation Gap:

The gap between EdgePrompt's privacy-oriented vision and its current implementation creates potential risks:

1. **Trust Erosion**: Without robust privacy implementation, the educational trust necessary for effective use could be compromised.
2. **Regulatory Exposure**: Deployment without addressing privacy gaps could create legal and compliance issues in many jurisdictions.
3. **Misaligned Incentives**: Pressure to implement features quickly might push privacy considerations to "future work" that never materializes.
4. **Security-Privacy Tension**: The focus on security as access control might overshadow privacy as appropriate data governance.
5. **Communication Gap**: Technical privacy details might not be effectively communicated to educational stakeholders.

### Implementation Priorities for Bridging the Gap:

To address these concerns, EdgePrompt's development should prioritize:

1. Implementing the comprehensive security model detailed in the documentation
2. Adding encryption for all sensitive data, particularly student responses
3. Developing clear data governance policies covering collection, retention, and deletion

4. Creating appropriate consent mechanisms for different stakeholder groups
5. Providing transparency tools showing what data is stored and how it's used
6. Designing privacy-preserving synchronization mechanisms for future offline-online bridging

By addressing these privacy implementation gaps, EdgePrompt can build the foundation of trust necessary for effective educational use. This isn't merely about compliance or security—it's about creating the conditions for authentic educational engagement in environments where students and teachers feel safe to explore, make mistakes, and grow.

### **Bias Risks: Inheritance from Base Models and Mitigation Needs**

EdgePrompt, like any system built on large language models, faces significant challenges related to bias—the tendencies of AI systems to reflect, reproduce, or sometimes amplify societal prejudices and stereotypes. These bias risks require careful assessment and mitigation, particularly in educational contexts where they can shape student understanding and self-perception:

#### **Sources of Bias in EdgePrompt:**

##### **1. Base Model Inheritance:**

- Edge-deployable models like Gemma and Llama are trained on internet-scale text that contains societal biases
- These models may inadvertently reproduce stereotypes about gender, race, ethnicity, nationality, religion, or ability
- Smaller models often undergo less extensive safety fine-tuning than their larger counterparts
- Quantization for edge deployment may sometimes impact safety layers more than core capabilities

##### **2. Template-Embedded Biases:**

- Templates created by educators may reflect their own implicit biases
- Question framing can unintentionally reinforce stereotypical perspectives
- Validation criteria might privilege certain cultural expressions or knowledge traditions
- Materials used as sources may contain historical biases or limited perspectives

##### **3. Validation Susceptibility:**

- Multi-stage validation might inconsistently enforce bias mitigation across different stages
- Different validation stages might apply inconsistent standards to different types of content
- Resource constraints for edge deployment might limit the sophistication of bias detection

##### **4. Cultural Context Gaps:**

- Models primarily trained on English-language content may perform unequally across languages
- Educational norms and examples may reflect Western educational traditions
- Culturally specific knowledge and perspectives may be underrepresented

- Indonesian 3T region contexts may be particularly underrepresented in model training

Research on algorithmic bias in education has examined how biases enter educational AI at multiple stages, connecting fairness definitions to the ML pipeline and documenting evidence of bias affecting racial, gender, and socioeconomic groups (Baker & Hawn 2022). Additionally, comprehensive surveys of bias and fairness in machine learning have reviewed techniques for identifying and mitigating various forms of bias in AI systems (Mehrabi et al. 2022).

### Current Mitigation Approaches and Limitations:

#### 1. Structured Prompting:

- *Approach*: Templates can include explicit diversity and inclusion guidance
- *Limitation*: Base model biases may still manifest despite prompt instructions

#### 2. Multi-Stage Validation:

- *Approach*: Dedicated validation stages could check for biased content
- *Limitation*: Current implementation focuses more on safety than bias detection

#### 3. Teacher Review:

- *Approach*: Human review provides a final check for biased content
- *Limitation*: Teachers may have their own implicit biases or limited time for thorough review

#### 4. Local Model Selection:

- *Approach*: Different communities could select models better aligned with their values
- *Limitation*: Limited availability of diverse, high-quality edge-deployable models

### Implications for Educational Equity:

Bias in educational AI has particularly significant implications:

1. **Representational Harm**: Students may not see themselves or their communities accurately represented in generated content
2. **Allocational Harm**: Biased assessments might unfairly evaluate students from certain backgrounds
3. **Self-Concept Impact**: Students may internalize biased perspectives about their abilities or opportunities
4. **Knowledge Distortion**: Biased content may provide inaccurate or incomplete perspectives on topics

These concerns are especially relevant for EdgePrompt's target contexts—underserved regions where educational resources are already limited and where AI-generated content might have proportionally greater influence.

### Priority Mitigation Strategies for Development:

1. **Bias-Aware Templates:** Develop template designs that explicitly counter common biases and promote inclusive representation
2. **Dedicated Bias Checking:** Implement specific validation stages focused on detecting potentially biased content
3. **Cultural Adaptation Frameworks:** Create mechanisms for adapting prompts and validation to different cultural contexts
4. **Diversity in Examples:** Ensure that examples and materials used for prompting reflect diverse perspectives
5. **Teacher Guidance:** Provide support for educators in identifying and addressing potential biases in generated content
6. **Community Input:** Engage local communities in defining appropriate bias mitigation approaches for their contexts
7. **Ongoing Monitoring:** Implement mechanisms to track and address bias issues that emerge during deployment

Addressing bias effectively will require ongoing attention throughout EdgePrompt's development—it cannot be solved through a single intervention or approach. By acknowledging these challenges explicitly and implementing layered mitigation strategies, EdgePrompt can work toward generating educational content that supports rather than undermines equity goals.

### **The Implicit Curriculum: What Values Does EdgePrompt Model?**

Beyond explicit educational content, technologies used in classrooms convey implicit messages about values, relationships, and ways of knowing. These constitute what might be called an "implicit curriculum"—the unstated lessons that students absorb through their interactions with technological systems. EdgePrompt, like any educational technology, will inevitably model certain values through its design and operation:

#### **Values Currently Modeled by EdgePrompt:**

1. **Authority and Knowledge Sources:**
  - *Current Modeling:* Knowledge comes from AI systems with teacher oversight
  - *Implicit Message:* Truth can be algorithmically determined but requires human verification
  - *Educational Implication:* May shape student understanding of what constitutes authoritative knowledge
2. **Feedback and Assessment:**
  - *Current Modeling:* Automated validation of responses against predefined criteria
  - *Implicit Message:* Learning involves meeting externally defined standards efficiently
  - *Educational Implication:* May prioritize performance over exploration or creativity
3. **Communication Patterns:**
  - *Current Modeling:* Structured question-answer-validation sequences
  - *Implicit Message:* Educational interactions follow predictable, formal patterns

- *Educational Implication*: May reinforce transmission models of education
4. **Problem-Solving Approaches:**
    - *Current Modeling*: Direct response to clearly formulated questions
    - *Implicit Message*: Problems have definite answers that can be efficiently produced
    - *Educational Implication*: May underemphasize ambiguity, complexity, and multiple perspectives
  5. **Time and Efficiency:**
    - *Current Modeling*: Rapid response generation and validation
    - *Implicit Message*: Speed and efficiency are primary virtues in learning
    - *Educational Implication*: May devalue slower, more reflective learning processes

Research on socio-technical education futures examines how educational technologies operate within socio-technical assemblages, highlighting the importance of examining social values and power structures alongside the technical tools themselves. This perspective emphasizes that technology and society co-evolve in schools, with educational technology influencing how students understand knowledge, authority, and learning itself (Swist & Gulson 2023).

### **Ethical Dimensions of the Implicit Curriculum:**

The values modeled by educational technology raise several ethical questions:

1. **Transparency**: Are students aware of the implicit messages being conveyed?
2. **Alignment**: Do the modeled values align with the stated educational philosophy?
3. **Diversity**: Does the system accommodate diverse value systems and cultural contexts?
4. **Agency**: Do students and teachers have the ability to shape or challenge the implicit curriculum?
5. **Development**: How do the modeled values influence student character and intellectual development?

Research on the ethics of AI in education has emphasized the importance of transparency in how AI makes decisions, safeguarding student autonomy and trust, and adherence to principles of fairness, accountability, and pedagogical soundness (Holmes et al. 2022).

### **Potential Value Realignment Strategies:**

To create a more intentional and positive implicit curriculum, EdgePrompt's development could:

1. **Expand Interaction Models:**
  - Move beyond simple Q&A to include exploratory, creative, and collaborative activities
  - Implement dialogue-based interactions that model more natural communication
  - Support student-initiated queries and investigations
2. **Enrich Feedback Approaches:**
  - Develop validation that emphasizes process and thinking rather than just correctness

- Implement feedback that asks questions rather than just providing judgments
  - Create space for multiple valid approaches to problems
3. **Model Intellectual Humility:**
    - Explicitly acknowledge limitations and uncertainty in AI-generated content
    - Include language that models appropriate tentativeness about complex topics
    - Provide multiple perspectives on issues where reasonable disagreement exists
  4. **Emphasize Human Relationships:**
    - Position AI as facilitating rather than replacing human educational relationships
    - Create workflows that enhance teacher-student and peer interactions
    - Avoid anthropomorphizing the AI in ways that suggest it as a relationship substitute
  5. **Support Metacognition:**
    - Include prompts that encourage reflection on learning processes
    - Develop templates focusing on the "how" and "why" rather than just the "what"
    - Create space for students to evaluate both their own thinking and the AI's responses

By attending to these implicit curriculum dimensions, EdgePrompt development can move beyond the technical challenges of implementation to consider the deeper ethical questions of what kind of learning environment the system helps to create. This attention to values alignment represents an essential aspect of responsible AI development for educational contexts—particularly when working with young, impressionable K-12 students in diverse cultural settings.

### **First Principles: Aligning AI Behavior with Human Educational Values**

From first principles, we must consider a fundamental question: how can a system driven by pattern matching (LLMs) be aligned with human educational values that transcend patterns? This requires examining both the nature of human educational values and the mechanisms through which AI behavior can be guided:

#### **Core Educational Values:**

Educational systems worldwide, while diverse in implementation, often share certain foundational values:

1. **Human Flourishing:** Education aims to enhance human capabilities and potential
2. **Critical Thinking:** Developing the ability to evaluate claims and evidence independently
3. **Autonomy:** Fostering the capacity for self-directed learning and decision-making
4. **Respect:** Treating each learner with dignity and acknowledging their uniqueness
5. **Truth-Seeking:** Commitment to accuracy, honesty, and intellectual integrity
6. **Cultural Heritage:** Preserving and transmitting cultural knowledge and perspectives
7. **Justice:** Creating fair opportunities and addressing systemic barriers to learning

These values often exist in tension with one another and are interpreted differently across cultural contexts, but they provide essential orientation for educational practice.

Research on artificial intelligence, values, and alignment provides philosophical frameworks for understanding how AI systems can be aligned with human values, offering approaches for thinking about how technologies like EdgePrompt can be designed to uphold educational values (Gabriel 2020).

### **Alignment Mechanisms in EdgePrompt:**

EdgePrompt employs several mechanisms that could potentially align AI behavior with these educational values:

#### **1. Explicit Value Encoding:**

- *Mechanism:* Templates and validation stages can explicitly encode educational values
- *Example:* Prompts instructing models to encourage critical thinking or respect diverse perspectives
- *Limitation:* Models may simulate value alignment without truly embodying it

#### **2. Structural Constraints:**

- *Mechanism:* System architecture can enforce boundaries aligned with educational values
- *Example:* Multi-stage validation creating checks and balances against potential misalignment
- *Limitation:* Constraints may address symptoms rather than underlying misalignment

#### **3. Human Oversight:**

- *Mechanism:* Teacher review and intervention can correct misaligned outputs
- *Example:* Review of flagged responses that may violate educational values
- *Limitation:* Relies on human capacity and judgment, which may be inconsistent or limited by time

#### **4. Alignment Through Use:**

- *Mechanism:* System design can encourage value-aligned usage patterns
- *Example:* Templates structured to promote exploration rather than just information delivery
- *Limitation:* Users may adapt tools in ways that undermine intended alignment

Research on AI moral enhancement examines the socio-technical systems of moral engagement, proposing modular approaches to AI ethics that complement multi-stage validation architectures like EdgePrompt's (Volkman & Gabriels 2023).

### **Alignment Challenges from First Principles:**

Several fundamental challenges complicate the alignment of AI with educational values:

1. **Value Pluralism:** Educational values vary across cultures, communities, and individual educators
2. **Ambiguity:** Values often require contextual interpretation resistant to algorithmic implementation
3. **Evolution:** Educational values shift over time, requiring adaptive alignment approaches
4. **Measurement Gap:** Many important educational values resist clear metrics or validation

5. **Emergence:** Complex interactions between components may produce unintended misalignment

### **Toward Principled Alignment Engineering:**

To better align EdgePrompt with human educational values, several first-principles approaches warrant consideration:

1. **Value Articulation:** Create explicit processes for educational stakeholders to articulate their values in forms that can guide system development
2. **Multi-Level Alignment:** Implement alignment mechanisms at multiple levels—from hardware constraints to user interfaces—creating defense in depth against misalignment
3. **Contextual Adaptation:** Design systems to adapt alignment approaches based on cultural, institutional, and individual educational contexts
4. **Reflective Practices:** Build in mechanisms for ongoing reflection on alignment effectiveness, including feedback loops from actual educational usage
5. **Epistemic Humility:** Acknowledge the inherent limitations of aligning pattern-matching systems with human values, maintaining appropriate caution about capabilities

From first principles, the alignment challenge is not primarily technical but philosophical: it involves reconciling fundamentally different modes of operation (statistical pattern completion versus value-driven human judgment) in service of complex educational goals. EdgePrompt's approach—using human-authored templates, staged validation, and teacher oversight—represents one potential path toward this reconciliation, but it requires continued refinement and critical assessment to ensure that the technology truly serves human educational values rather than subtly reshaping those values to fit technological constraints.



## Chapter 10: Equity and Liberation – Access vs. Agency

### The Offline Premise: A Genuine Step Towards Equity?

EdgePrompt's core premise—bringing AI-enhanced education to offline or connectivity-challenged environments—positions it as an equity intervention. But does this offline capability truly advance educational equity? This question requires critical examination:

#### The Equity Challenge:

Educational equity faces a multi-layered digital divide:

1. *Basic connectivity*: Many communities lack reliable internet access
2. *Hardware access*: Resource disparities limit device availability
3. *Digital literacy*: Technical knowledge varies widely among educators
4. *Content relevance*: Available digital resources often lack cultural relevance
5. *Language barriers*: Educational technology predominantly favors major languages

These divides risk creating a two-tiered educational future—one where well-resourced schools benefit from AI advances while under-resourced schools fall further behind. EdgePrompt's offline capability directly addresses the first challenge, potentially enabling more equitable access to AI-enhanced education across connectivity barriers. This approach aligns with research on AI for Sustainable Development Goals, which examines how artificial intelligence can be applied to address educational equity and other development challenges (Singh et al. 2024).

#### Potential Equity Benefits:

EdgePrompt's offline approach offers several potential equity advantages:

1. **Geographic Inclusion**: Enables participation by rural, remote, and infrastructure-limited regions
2. **Continuity of Access**: Provides consistent educational tools regardless of connectivity fluctuations
3. **Cost Reduction**: Eliminates ongoing connectivity costs that burden resource-constrained schools
4. **Local Control**: Enables communities to maintain ownership and governance of educational AI
5. **Adaptation Potential**: Creates foundation for local customization addressing cultural and linguistic needs

These benefits suggest that EdgePrompt's offline capability could indeed represent a meaningful step toward more equitable AI distribution in education—addressing a critical barrier that currently excludes many communities from emerging educational technologies.

**Critical Limitations:**

However, offline capability alone faces several limitations as an equity intervention:

1. **Hardware Requirements:** Edge-deployed LLMs still require moderately powerful devices, potentially beyond the reach of the most resource-constrained schools
2. **Technical Support Needs:** Local deployment requires technical capacity that may be limited in target communities
3. **Update Challenges:** Offline systems risk becoming outdated without connectivity for updates
4. **Content Constraints:** Edge models may have more limited capabilities than cloud alternatives
5. **Implementation Gaps:** Current implementation lacks many features needed for equitable deployment

These limitations suggest that while offline capability is necessary for equity, it is not sufficient—additional considerations around hardware accessibility, technical support, content relevance, and implementation completeness must be addressed.

**Beyond Technical Access:**

True educational equity requires looking beyond technical access to deeper questions:

1. **Pedagogical Alignment:** Does the system support pedagogical approaches relevant to diverse communities?
2. **Cultural Representation:** Does the content reflect diverse cultural perspectives and knowledge systems?
3. **Language Support:** Can the system function effectively across multiple languages, including local ones?
4. **Teacher Empowerment:** Does the system enhance rather than undermine teacher agency and expertise?
5. **Community Ownership:** Do communities have meaningful control over how the system is used?

EdgePrompt's vision acknowledges many of these dimensions, but its current implementation primarily addresses the technical access barrier rather than these deeper equity considerations.

**Assessment:**

EdgePrompt's offline capability represents a genuine and significant step toward more equitable AI distribution in education, directly addressing a critical barrier that currently excludes many communities. However, it remains an initial step rather than a complete solution—true equity will require addressing hardware access, technical support, content relevance, cultural adaptation, and community ownership alongside offline capability.

The potential equity impact is substantial but depends heavily on how the other dimensions are addressed in future development. EdgePrompt opens a door that has been closed to many communities, but walking through that door will require attention to the full spectrum of equity considerations beyond mere technical access.

### **Potential Pitfalls: Amplifying Existing Gaps (Hardware, Digital Literacy)**

While EdgePrompt aims to bridge educational divides through offline capability, there's a risk that it could inadvertently amplify other existing gaps. Several potential pitfalls warrant careful consideration:

#### **Hardware Disparities:**

1. **Minimum Requirements Gap:** Despite offline design, EdgePrompt still requires moderately powerful hardware to run edge-deployed LLMs effectively:
  - *Potential Impact:* Schools unable to afford suitable devices would be excluded despite the offline capability
  - *Risk Amplification:* Could create a secondary divide between schools that can afford edge-capable devices and those that cannot
  - *Current Status:* Phase 1 efficiency findings suggest meaningful performance demands, but actual hardware requirements remain untested
2. **Device Maintenance Challenges:** Edge deployments require ongoing maintenance and occasional hardware replacement:
  - *Potential Impact:* Schools without technical support or replacement budgets face sustainability challenges
  - *Risk Amplification:* Schools in better-resourced areas might maintain systems longer, creating widening capability gaps over time
  - *Current Status:* Deployment guide lacks guidance on sustainable hardware maintenance for resource-constrained environments
3. **Energy Infrastructure Dependencies:** Edge devices require reliable power sources to function:
  - *Potential Impact:* Schools in areas with unreliable electricity may face additional barriers
  - *Risk Amplification:* Power stability often correlates with other resource advantages, potentially reinforcing existing divides
  - *Current Status:* Limited attention to power management and low-energy operation modes

Research on on-device language models has highlighted the various challenges in deploying LLMs on resource-constrained devices and proposed strategies including efficient architectures, compression techniques, and hardware acceleration (Xu et al. 2024).

#### **Digital Literacy Divides:**

1. **Teacher Technical Capabilities:** EdgePrompt assumes basic technical literacy for system management:

- *Potential Impact:* Teachers lacking technical confidence may struggle with deployment and troubleshooting
  - *Risk Amplification:* Schools with technically experienced staff gain additional advantage
  - *Current Status:* Current UI requires moderate technical understanding; limited scaffolding for low-tech-literacy users
2. **Deployment Complexity:** Setting up and maintaining the system involves multiple technical steps:
    - *Potential Impact:* More technically proficient schools implement more quickly and effectively
    - *Risk Amplification:* Early adopters gain experience advantages that compound over time
    - *Current Status:* Deployment documentation assumes technical background knowledge
  3. **Troubleshooting Barriers:** When systems fail, technical knowledge determines recovery speed:
    - *Potential Impact:* Schools with limited technical capacity experience more downtime
    - *Risk Amplification:* Reliability differences can create self-reinforcing adoption patterns
    - *Current Status:* Limited robust error handling or non-technical recovery documentation

### Content and Usage Sophistication:

1. **Template Creation Skills:** Creating effective templates requires both technical and pedagogical knowledge:
  - *Potential Impact:* Schools with more experienced teachers can create more sophisticated learning activities
  - *Risk Amplification:* Template quality differences could create substantial learning experience gaps
  - *Current Status:* Limited template examples or guidance for educators new to AI-enhanced teaching
2. **AI Literacy Disparities:** Understanding AI capabilities and limitations affects effective usage:
  - *Potential Impact:* Teachers with greater AI literacy leverage the system more effectively
  - *Risk Amplification:* Existing technology exposure correlates with other advantage markers
  - *Current Status:* Limited onboarding for teachers without prior AI experience
3. **Integration Sophistication:** Effective educational technology integration requires pedagogical expertise:
  - *Potential Impact:* Schools with stronger instructional leadership integrate more meaningfully
  - *Risk Amplification:* Surface-level adoption in some schools versus transformative use in others
  - *Current Status:* Limited integration guidance addressing diverse pedagogical contexts

Research on educational 5G edge computing has demonstrated how edge computing architectures can support educational applications in resource-constrained environments, highlighting the potential for

local edge servers to improve throughput and latency for remote learners when properly implemented (Chen et al. 2022).

### **Mitigation Strategies:**

To address these potential pitfalls, EdgePrompt development could:

1. **Tiered Hardware Support:** Define multiple deployment profiles optimized for different hardware capabilities, from minimal to ideal
2. **Simplified Deployment Paths:** Create step-by-step deployment guides with minimal technical assumptions, including visual guides
3. **Teacher Support Materials:** Develop non-technical onboarding resources specifically designed for educators with limited technology experience
4. **Template Libraries:** Provide extensive pre-built templates across subjects to reduce creation barriers
5. **Community Support Structures:** Design implementation models that leverage community knowledge-sharing to address technical limitations
6. **Progressive Interface Complexity:** Implement interfaces that adapt to user technical proficiency, starting with basics
7. **Robustness Engineering:** Design for high fault tolerance and self-recovery to minimize technical support needs

By acknowledging and explicitly addressing these potential pitfalls, EdgePrompt can work to ensure that its offline capability genuinely advances equity rather than inadvertently amplifying existing gaps. This requires viewing equity not as a single-dimension challenge (connectivity) but as a multi-faceted consideration that shapes every aspect of system design, implementation, and support.

### **Culturally Responsive Adaptation: Moving Beyond Vision to Implementation**

EdgePrompt's vision emphasizes adaptability to different cultural and linguistic contexts—a crucial consideration for equitable deployment across diverse educational settings. However, translating this vision into implementation requires moving beyond general principles to specific architectural and design choices.

### **Current State Assessment**

- **Vision-Implementation Gap:** While cultural responsiveness appears as a priority in vision documents, the current implementation provides only minimal adaptation mechanisms—primarily a language selection flag and template customization potential.
- **Target Context Identification:** The project specifically identifies Indonesia's 3T regions as the initial focus but lacks region-specific implementation features beyond basic offline capability.

- **Architectural Foundation:** The template-based approach creates a theoretical foundation for adaptation, but current templates reflect primarily Western educational assumptions without explicit cultural variation mechanisms.
- **Documentation versus Code:** A significant discrepancy exists between the cultural adaptation commitments in documentation and their limited implementation in the actual codebase.
- **User Interface Considerations:** The interface currently lacks comprehensive localization or culturally adaptive elements, offering minimal support for non-Western educational contexts.

The implementation prioritizes basic functionality over deep contextual adaptation, creating a substantial gap between EdgePrompt's cultural responsiveness vision and its current reality.

### Language Support Implementation

- **Beyond translation to localization:** Moving from language translation to comprehensive localization requires incorporating region-specific educational terminology, assessment approaches, and feedback conventions. The current implementation enables basic language selection but doesn't yet support deeper localization.
- **Dialect and regional variation support:** Many regions use local dialects or variants of national languages that impact student comprehension. Implementation must support these variations through adaptable prompting structures that can incorporate regional expressions and terminology.
- **Character set and text direction:** Current implementation lacks robust support for diverse writing systems, bidirectional text, or special character requirements of many languages used in target regions.
- **Multilingual validation capabilities:** The validation system requires adaptation to evaluate responses in multiple languages while maintaining effectiveness—a capability not yet implemented in the current system.
- **Interface language consistency:** The implementation needs mechanisms ensuring consistent language use across all interface elements and generated content, avoiding mixed-language experiences that create cognitive barriers.

The current implementation offers language selection but lacks the comprehensive linguistic adaptation infrastructure needed for genuine cultural responsiveness in multilingual educational environments.

### Content Adaptation Implementation

- **Culturally relevant examples in generation:** Implementation challenges include ensuring generated content references culturally familiar contexts, examples, and scenarios that resonate with local students—avoiding imported contexts that create unnecessary cognitive barriers.
- **Locally appropriate metaphors and references:** Educational explanations often rely on metaphors and references that vary significantly across cultures. Current implementation lacks mechanisms for adapting these cultural touchpoints to local contexts.

- **Context-sensitive educational scenarios:** Effective educational activities often rely on scenarios familiar to students. The implementation needs capabilities for generating region-appropriate contexts for problems and examples.
- **Multiple knowledge traditions representation:** Different cultures have distinct knowledge traditions and epistemological frameworks. Implementation requires mechanisms for respecting and incorporating diverse approaches to knowledge beyond Western academic traditions.
- **Bias detection and mitigation:** EdgePrompt's reliance on LLMs trained predominantly on Western content necessitates implementation of bias detection and mitigation mechanisms specifically for multicultural educational contexts.

The current template system theoretically allows for cultural customization but lacks explicit adaptation mechanisms or culturally diverse template collections. Research on culturally responsive teaching emphasizes incorporating students' cultural references in all aspects of learning, validating and affirming diverse cultural identities (Gay 2018). Additionally, work on culturally adaptive thinking in education for AI highlights the importance of localizing AI curricula and tools to align with local cultural norms and languages to avoid cultural bias and resistance (Samuel et al. 2023).

### **Pedagogical Alignment Implementation**

- **Different instructional approaches across cultures:** Educational philosophies and approaches vary significantly across cultural contexts. Implementation must support diverse pedagogical approaches rather than embedding assumptions from specific educational traditions.
- **Varied assessment practices and feedback styles:** Assessment and feedback conventions differ culturally in directness, emphasis, and structure. Current implementation reflects primarily Western assessment assumptions without adaptation mechanisms.
- **Cultural values around learning relationships:** Educational interactions reflect cultural values regarding authority, relationship, and knowledge transmission. Implementation must adapt to these diverse relationship expectations.
- **Age-appropriate interactions by culture:** Expectations for age-appropriate interactions vary across cultures, requiring implementation of culturally calibrated scaffolding and interaction patterns.
- **Knowledge organization frameworks:** Different cultures organize knowledge domains in distinct ways. Implementation must accommodate these varied organizational frameworks rather than imposing standardized knowledge structures.

The template architecture theoretically allows pedagogical flexibility but current templates embed specific pedagogical assumptions without explicit cultural adaptation mechanisms.

## Technical Contextualization Implementation

- **Device ecosystem awareness across regions:** Implementation must account for prevailing device types in target regions—whether Android-dominant, Windows-based, or mixed ecosystems—ensuring the framework functions effectively on locally available hardware.
- **Connectivity pattern adaptations:** Recognition that "offline" exists on a spectrum from completely disconnected environments to intermittent connectivity with varying characteristics. Implementation must adapt synchronization strategies to match local connectivity patterns.
- **Storage and memory optimization:** Adaptations accounting for regional variations in available resources—implementing compression techniques, efficient caching strategies, and tiered storage models that optimize performance on entry-level devices prevalent in target areas.
- **Power infrastructure considerations:** Implementation must address inconsistent electricity access through robust crash recovery, efficient battery usage patterns, and graceful degradation during low-power states—tested under conditions mirroring actual usage environments.
- **Data management aligned with cultural practices:** Implementation needs data retention, sharing, and privacy features that reflect local cultural perspectives on information ownership and appropriate data stewardship models, which may differ from Western assumptions.

The current implementation provides basic offline capability but lacks the sophisticated technical contextualization needed for optimal functioning across diverse infrastructure environments.

## Process Adaptation Implementation

- **Workflow adaptation to regional teaching practices:** Different regions employ diverse teaching workflows, from planning to assessment. Implementation must allow flexible workflows that align with established practices rather than imposing standardized processes.
- **Teacher decision points reflecting local autonomy models:** The implementation of where and how teachers can intervene in the AI process should reflect local educational governance models and teacher authority structures.
- **Integration with existing educational rhythms:** From scheduling to assessment cycles, educational processes follow region-specific patterns. Implementation must support integration with these established rhythms rather than disrupting them.
- **Community involvement mechanisms:** In many regions, education involves broader community participation. Implementation should include options for appropriate community engagement aligned with local practices, from parental review to community elder input.
- **Contextual interpretation of assessment results:** Implementation must avoid universal scoring or feedback standards, instead providing interpretation frameworks that align with local educational expectations and values.

The process adaptation gap remains particularly wide, with minimal implementation of region-specific process adaptations beyond basic workflow structuring.



## Implementation Path Forward

- **Co-design methodology development:** Creating structured approaches for engaging target communities in adaptation design rather than imposing external solutions.
- **Cultural reference architecture:** Developing a comprehensive framework for cultural adaptation across all system components.
- **Adaptation infrastructure prioritization:** Building core adaptation mechanisms before expanding feature set to ensure cultural responsiveness is fundamental rather than supplemental.
- **Community ownership of adaptation:** Creating mechanisms for local communities to control and evolve cultural adaptations rather than depending on outside developers.
- **Continuous cultural assessment:** Implementing ongoing evaluation of cultural appropriateness and effectiveness rather than one-time adaptations.

Bridging the substantial gap between EdgePrompt's cultural responsiveness vision and current implementation requires systematic engagement with educators and technical experts from target regions, with cultural adaptation treated as core functionality rather than optional enhancement.

## From Equity (Access) to Liberation (Epistemic Agency): Does EdgePrompt Empower or Domesticate?

The ultimate test of EdgePrompt's contribution to educational equity transcends mere technological access. True equity demands advancing epistemic agency—the capacity for learners to engage as knowledge creators rather than mere consumers, particularly for historically marginalized communities. This section examines whether EdgePrompt's approach empowers or potentially domesticates learners and educators.

### Epistemic Agency vs. Technological Dependency:

EdgePrompt's dual nature creates tension between empowerment and dependency:

- **Teacher agency preservation:** EdgePrompt's philosophy of maintaining teacher control over prompts, validation criteria, and the educational process theoretically preserves teacher epistemic agency. However, the current implementation's complexity may subtly shift actual control to technical experts who understand the system's internal workings.
- **Student agency implications:** The current generation-validation approach primarily positions students as responders to AI-generated questions rather than co-creators of knowledge. This pedagogical model may reinforce traditional epistemic hierarchies rather than fostering student agency.
- **Framework adaptability vs. prescribed patterns:** While EdgePrompt theoretically enables adaptation to diverse educational approaches, the structured prompt-validation architecture potentially normalizes specific interaction patterns that originate from particular educational philosophies, potentially marginalizing alternative knowledge construction approaches.

- **Knowledge validation power dynamics:** The multi-stage validation system, while technically impressive, centralizes evaluative authority in pre-defined criteria and AI assessment—potentially reinforcing dominant epistemologies rather than creating space for diverse ways of knowing.

Research on epistemic agency in education emphasizes the importance of viewing learners as active agents capable of evaluating and constructing knowledge, shifting educational focus from passive reception to active engagement (Elgin 2013). Additionally, socio-critical perspectives on learning analytics ethics highlight that students should be active agents of their data, not mere data points, requiring policies that protect student identity and rights (Slade & Prinsloo 2013).

The critical question remains: does EdgePrompt's technical approach to AI guardrails inadvertently constrain epistemic exploration in its pursuit of safety and structure? The answer likely depends on implementation choices that either emphasize control or prioritize flexible knowledge co-construction.

### **Beyond Binary Access Models:**

Advancing from access to liberation requires moving beyond binary thinking about technology availability:

- **Spectrum of meaningful access:** Implementation must recognize that meaningful access exists on a spectrum beyond the binary "online/offline" distinction. Even with offline capability, EdgePrompt must address additional access barriers including language, digital literacy, and knowledge prerequisites.
- **Capability vs. agency balance:** The current focus on what EdgePrompt can do may inadvertently overshadow examination of how it shifts power over knowledge construction. Future implementation should explicitly address how features either enhance or potentially limit user agency.
- **Access to what, for whom, and why:** Implementation must continuously examine the values embedded in what content is made accessible, which users are prioritized, and what educational purposes are centered. These choices can either challenge or reinforce existing inequities.
- **From consumption to creation:** True liberation requires implementation to evolve beyond students merely consuming AI-generated content toward students actively shaping the knowledge construction process, potentially through student-authored prompts, co-designed validation criteria, or customizable interaction models.

EdgePrompt's current implementation primarily addresses physical access barriers through offline capability but has yet to fully engage with these deeper dimensions of accessibility that connect to epistemic liberation.

### Pedagogical Colonization Risks:

Any educational technology, regardless of good intentions, carries risks of pedagogical colonization that must be actively addressed:

- **Hidden curriculum analysis:** EdgePrompt's structured approach to education encodes particular values about what constitutes valid knowledge and appropriate learning processes. These implicit assumptions must be examined to prevent unintentional imposition of dominant educational paradigms.
- **Indigenous and alternative knowledge systems:** Implementation must create space for knowledge systems that may operate through different structures, validation approaches, or epistemological foundations than those assumed in the current prompt-validation architecture.
- **Potential for educational extraction:** Even offline frameworks potentially extract educational data and practices from communities while providing limited reciprocal value. Implementation must incorporate robust data governance that ensures value flows back to the communities providing educational contexts.
- **Teacher deskilling prevention:** Implementation must guard against gradual erosion of teacher expertise through over-reliance on AI-generated content and validation, which could diminish the role of human judgment in educational processes.

The current EdgePrompt implementation focuses primarily on technical guardrails with limited explicit attention to these pedagogical colonization risks. Future development requires active engagement with these concerns through participatory design with diverse educational stakeholders.

### Liberatory Technology Implementation:

For EdgePrompt to advance true liberation, implementation must embody specific characteristics:

- **Community ownership mechanisms:** Moving beyond "designed for" communities toward development "by and with" communities through concrete governance mechanisms, open-source accessibility, and adaptation rights that ensure technology serves community-defined needs.
- **Constructionist knowledge building:** Implementation should evolve toward supporting students and teachers as creators rather than consumers, including capabilities for collaborative prompt design, custom validation criteria development, or community-specific template creation.
- **Critical technological consciousness:** Future implementations should incorporate transparency mechanisms that make AI limitations, potential biases, and decision processes visible to users, fostering critical engagement rather than passive acceptance.
- **Contextual flexibility prioritized over standardization:** Implementation should embrace adaptability even at the cost of standardization, allowing diverse educational philosophies to shape how the technology functions rather than imposing a universal approach.

These liberatory characteristics remain aspirational for EdgePrompt. While the offline capability and teacher control philosophy align with liberatory principles, the current implementation's focus on technical validation and standardized processes has not yet fully embodied these characteristics.

The path from EdgePrompt's current state to a truly liberatory educational technology requires not just technical refinement but philosophical evolution. The project must continuously examine whether it empowers communities to define and create their educational futures or subtly reinforces technological dependency under the guise of democratization.

## Chapter 11: Cybersecurity in the Classroom – Beyond Standard Models

**For Non-Technical Readers:** This section explains why protecting a classroom AI system requires different security approaches than those used for business computers. Traditional cybersecurity assumes controlled environments with individual users at separate devices, but classrooms feature shared computers, visible screens, curious students, and limited technical support. For example, while a business might worry about external hackers, a classroom system needs to prevent students from manipulating the AI to complete assignments for them or generate inappropriate content. The current version of EdgePrompt has detailed security plans on paper but minimal actual security protections in the software—like a house with blueprints for sophisticated locks but currently no doors installed. This gap represents a critical issue that must be addressed before real classroom use, as students will naturally explore and test boundaries of any technology they use, requiring special protection approaches designed specifically for educational contexts.

The deployment of AI systems in educational settings introduces unique cybersecurity challenges that transcend traditional enterprise security models. K-12 classrooms—particularly in resource-constrained environments like Indonesia's 3T regions—represent complex security environments with distinctive threat models, unusual attack vectors, and security requirements that standard frameworks often fail to address. This chapter examines EdgePrompt's security posture against these classroom-specific security realities.

### Why Classrooms Defy Traditional Security Assumptions

Traditional cybersecurity models rest on assumptions that fundamentally misalign with K-12 educational environments, creating a security paradigm mismatch that EdgePrompt must navigate:

#### Physical Environment Realities:

- **Shared physical access:** Unlike corporate environments where physical access is controlled, classroom devices are often shared among multiple students and sometimes across multiple classes, creating inherent physical security risks.
- **Visible screens and inputs:** Classroom layouts typically feature visible screens that can be observed by multiple students simultaneously, compromising confidential interactions and potentially exposing sensitive information.
- **Intermittent supervision:** Teachers cannot maintain constant supervision of all devices simultaneously, creating windows of opportunity for misuse that rarely exist in enterprise environments.
- **Device mobility:** In many resource-constrained settings, devices move between classrooms or even travel home with students, crossing physical security boundaries in ways corporate equipment rarely does.

EdgePrompt's security model acknowledges these physical realities in its specification documents but lacks concrete implementation of mitigations in the current codebase. The security model document details a theoretical approach but the application code reveals minimal physical security controls.

### User Behavior Patterns:

K-12 students exhibit fundamentally different security behaviors than enterprise users:

- **Exploratory mindset:** Students, particularly adolescents, naturally test boundaries and experiment with systems in ways adult professional users typically don't, creating unique attack vectors.
- **Peer dynamics:** Peer pressure and social status may incentivize security circumvention for entertainment or status rather than material gain, shifting attack motivations.
- **Limited security awareness:** Students have developing rather than mature security consciousness, with security often treated as an obstacle rather than a protection.
- **Deliberate constraint testing:** Unlike professional environments where compliance is expected, educational settings may inadvertently reward creative system circumvention as a demonstration of technical skill.

EdgePrompt's documentation acknowledges these behavioral challenges conceptually but provides limited concrete implementation guidance. The security engineering document outlines authentication models but doesn't specifically address student-specific behavior patterns or metrics for measuring security effectiveness in classroom contexts yet.

### Resource Constraints Impact:

Resource-constrained educational settings introduce additional security complexities:

- **Limited update capabilities:** Irregular connectivity may prevent timely security updates, creating persistent vulnerability windows.
- **Minimal logging infrastructure:** Resource limitations often preclude comprehensive security logging and monitoring that enterprise environments take for granted.
- **Restricted security expertise:** Schools rarely have dedicated security personnel or the expertise to implement complex security controls effectively.
- **Device longevity:** Educational devices typically remain in service far longer than enterprise equipment, increasing exposure to vulnerabilities in aging software and hardware.

EdgePrompt acknowledges these constraints philosophically, particularly in its offline-first design, but the current implementation lacks robust mechanisms for security operation under these constraints. The architecture enables offline functionality but doesn't yet address how security will be maintained under prolonged disconnection or resource limitation.

### Analysis: EdgePrompt's Security Foundation:

EdgePrompt's security model theoretically recognizes educational environments' unique characteristics. The security model document outlines an architecture designed for offline operation, minimal resource requirements, and appropriate user authentication. However, this theoretical model has minimal implementation in the current codebase. The project's security posture reflects a vision-implementation gap typical of early-stage development—sophisticated in conception but nascent in execution.

Research on the impact of generative AI in cybersecurity and privacy has examined security implications for educational contexts, highlighting the importance of protecting student data and implementing appropriate safeguards (Gupta et al. 2023). Additionally, work on regulating large generative AI models provides context for understanding emerging regulatory frameworks that might affect AI safety requirements in educational applications (Hacker, Engel & Mauer 2023).

### **Adversarial Scenarios: Cheating, Misuse, Sabotage, Social Harm**

K-12 environments present distinctive adversarial scenarios that traditional security models rarely address. These scenarios reflect the unique motivations, capabilities, and opportunities present in educational settings.

#### **Academic Integrity Threats:**

Students may engage with AI systems in ways that undermine academic integrity, creating challenges specific to educational security:

- **Generative AI for assignment completion:** Students might attempt to use EdgePrompt to complete assignments by generating inappropriate levels of assistance, undermining learning objectives.
- **Validation circumvention:** Students could potentially craft inputs designed to manipulate the validation system into providing complete answers under the guise of feedback.
- **Answer sharing networks:** Academic collusion could leverage EdgePrompt outputs shared across student networks, potentially via screenshots or copied text.
- **Credential sharing:** Students might share access credentials to bypass individual usage monitoring or restrictions, particularly in environments with limited device access.

The current EdgePrompt implementation lacks specific protections against these academic integrity threats. The validation model focuses primarily on content safety rather than academic integrity protection, and the authentication system remains largely unimplemented.

#### **Social-Emotional Manipulation:**

Classroom dynamics create unique opportunities for technology misuse that causes social or emotional harm:

- **Impersonation attacks:** Students could potentially craft inputs that cause the system to generate content appearing to come from other students or faculty.

- **AI-generated bullying content:** A significant risk involves students using the generative capabilities to create convincing but harmful content targeted at peers, such as embarrassing fictional scenarios.
- **Rumor amplification:** Students might prompt the system to elaborate on or appear to validate harmful rumors, leveraging AI's perceived authority.
- **Authority undermining:** Strategic system manipulation could be used to challenge teacher authority by generating content that appears to contradict educational materials.

EdgePrompt's multi-stage validation architecture theoretically addresses content safety concerns but lacks specific mechanisms for detecting and preventing social-emotional manipulation. The current validation stages focus on educational appropriateness rather than social harm prevention.

### Technical Sabotage Vectors:

Educational environments present unique technical attack patterns:

- **Deliberate resource exhaustion:** Students might intentionally craft inputs designed to consume excessive resources, potentially as disruption or competition.
- **Data integrity attacks:** Motivated students might attempt to corrupt local databases or submission records, particularly regarding their own performance data.
- **Template manipulation:** In systems allowing template sharing or modification, students might create misleading or inappropriate templates.
- **Synchronization exploitation:** Systems with intermittent connectivity create potential attack vectors during synchronization events.

EdgePrompt's architecture acknowledges some of these technical risks conceptually, but the current implementation provides minimal protection. The SQLite database lacks robust integrity protections, and the template system doesn't yet incorporate strong verification mechanisms against manipulation.

### Privacy Compromises:

Privacy threats in educational settings have distinct characteristics:

- **Peer privacy violations:** Students may attempt to access other students' work, potentially aided by physical proximity and shared devices.
- **Educational profiling risks:** AI systems that track student performance could create sensitive profiles that require protection beyond standard data security.
- **Context-inappropriate data collection:** Systems must balance assessment needs against age-appropriate privacy protections, particularly for younger students.
- **Cross-boundary data exposure:** Educational data inappropriately shared across classroom, administrative, and home boundaries creates unique privacy risks.



The current EdgePrompt implementation lacks the robust user separation and role-based access control described in the security specification. While the architecture conceptually supports these protections, the implementation reflects early-stage development priorities focused on core functionality over comprehensive security controls.

Research on AI misuse and information privacy concerns has validated new constructs for measuring privacy concerns related to AI systems and emphasized the importance of robust safety guardrails in addressing these concerns (Menard & Bott 2024).

### **Analysis: EdgePrompt's Adversarial Readiness:**

EdgePrompt's adversarial scenario preparedness reveals a significant gap between security vision and implementation reality. The security specifications outline a sophisticated model with identity management, role-based access control, and comprehensive data protection, but the current implementation reflects early development priorities focused on validating core functionality rather than security hardening. This prioritization is reasonable for a research-stage project but represents a critical area for focused development before classroom deployment.

### **Assessing the Current (Minimal) Security Posture vs. the Specification**

EdgePrompt presents a compelling security vision in its specification documents, but implementation reality reflects early development priorities. This section examines the gap between specified security controls and current implementation.

#### **Authentication & Authorization:**

##### **Security Specification:**

- Detailed user identity management with credential creation, validation, and rotation
- Role-based access control with fine-grained permissions
- Multi-factor authentication options
- Token-based session management

##### **Current Implementation:**

- No user authentication or authorization system implemented
- No role separation between teacher and student interfaces
- No credential management or session controls
- No implementation of the specified identity layer

The authentication and authorization gap represents the most significant security implementation shortfall. While the specification details a sophisticated identity model, the application currently operates without user authentication, creating a fundamental security vulnerability.

## **Data Protection:**

### **Specification:**

- Envelope encryption for sensitive data
- Key management framework with master key encryption
- Data compartmentalization by user
- Secure storage backends for offline operation

### **Current Implementation:**

- Unencrypted SQLite database for all storage
- No implementation of the specified encryption framework
- No data compartmentalization or user isolation
- Basic file storage without additional protection

The data protection implementation gap creates significant risks for sensitive educational data, particularly student responses and assessment results. The absence of the specified encryption framework leaves all data in plaintext, accessible to anyone with device access.

## **Input Validation & Content Filtering:**

### **Specification:**

- Multi-stage validation of all generated content
- Content filtering for age-appropriate outputs
- Constraint enforcement on inputs and outputs
- Robust input sanitization

### **Current Implementation:**

- Basic input validation for API routes
- Simplified validation logic implemented
- Limited constraint enforcement in the research framework
- Partial implementation of the multi-stage validation in research code but not in the application

The input validation gap represents a more nuanced implementation status. While the research framework implements aspects of the specified validation architecture, the web application includes only simplified validation, creating inconsistent protection levels.

**Offline Security:****Specification:**

- Secure offline authentication flow
- Credential caching with local validation
- Synchronization security for intermittent connectivity
- Resource-constrained token validation

**Current Implementation:**

- Basic offline operation capability
- No implementation of secure offline authentication
- No credential caching or secure synchronization
- No implementation of resource-optimized security primitives

The offline security gap reflects the project's current focus on basic functionality over security hardening. While the application supports basic offline operation, it lacks the sophisticated security mechanisms described in the specification for maintaining security during disconnected operation.

**Analysis: Security Implementation Gap:**

EdgePrompt's current security implementation reflects early-stage development prioritization rather than fundamental architectural flaws. The detailed security specifications provide a robust foundation for future implementation, but the current codebase represents a minimal viable product focused on validating core functionality rather than comprehensive security implementation.

This gap between specification and implementation creates significant security risks for any near-term classroom deployment. The absence of authentication, data protection, and robust validation in the application would expose sensitive educational data and create potential for misuse. These gaps must be addressed before considering actual classroom deployment beyond controlled research contexts.

**Towards Classroom-Specific Defenses: Zero Trust Edge, Robust Filtering, Teacher Oversight**

Addressing EdgePrompt's classroom-specific security challenges requires targeted solutions that go beyond generic security patterns. This section outlines a classroom-specific security framework that aligns with EdgePrompt's vision while addressing its current implementation gaps.

**Zero Trust Edge Architecture:**

Traditional security models often assume protected perimeters and trusted internal users—assumptions that fail in classroom environments. A classroom-specific adaptation of Zero Trust principles would include:

- **Continuous revalidation:** Security architecture should assume devices regularly change users and continuously revalidate identity through appropriate mechanisms (potentially including non-credential factors in younger grades).
- **Least privilege by default:** Every function should operate with minimal necessary permissions, with particular attention to content generation and validation capabilities.
- **Trust boundaries around individuals, not devices:** Security models must assume shared devices with multiple potential adversaries, establishing trust boundaries at the user level rather than device level.
- **Offline-capable verification:** Authentication and authorization must function robustly without connectivity, potentially through cached credentials with appropriate expiration and local validation.
- **Tamper-evident operations:** All significant operations should generate tamper-evident audit records suitable for teacher review, documenting who performed what actions when.

Implementation of these Zero Trust adaptations would significantly enhance EdgePrompt's security posture, particularly in addressing the shared access challenges inherent in classroom environments.

### **Content Filtering Beyond Keywords:**

Educational environments require more sophisticated content filtering than typical enterprise systems:

- **Age-appropriate filtering calibration:** Content filtering requires precise calibration to student developmental stages, balancing protection against educational exposure needs.
- **Context-sensitive evaluation:** Filtering must consider educational context—content appropriate in a history lesson may be inappropriate in another context.
- **Multi-dimensional assessment:** Rather than binary allow/block decisions, educational filtering should employ graduated responses appropriate to content severity.
- **Local cultural alignment:** Filtering standards must adapt to community standards and cultural contexts rather than imposing universal definitions of appropriateness.
- **Teacher review mechanisms:** Filtering should include efficient teacher review processes for edge cases, potentially with pre-prepared alternative content.

EdgePrompt's multi-stage validation architecture theoretically supports these capabilities, but current implementation reflects simplified approaches. Future development should enhance filtering sophistication beyond simple keyword matching or constraints.

### **Teacher-Centered Security Controls:**

Effective classroom security places teachers at the center of the security model:

- **Teacher dashboard with security overview:** Implementation should include teacher interfaces for monitoring system usage, reviewing flagged content, and assessing security status.

- **Graduated intervention options:** Security controls should offer teachers a range of intervention options beyond binary allow/block decisions, supporting teachable moments.
- **Classroom-wide policy management:** Teachers need efficient tools to manage security policies across multiple students simultaneously, adapting to lesson requirements.
- **Transparent but efficient oversight:** Security interfaces must balance complete information with teacher cognitive load, presenting actionable security insights rather than overwhelming details.
- **Security event analysis tools:** Teachers need support interpreting security events through educational rather than technical lenses, focusing on learning interventions rather than purely technical responses.

The current EdgePrompt implementation lacks these teacher-centered security controls, representing a significant gap between educational security needs and implementation reality.

### **Student-Appropriate Security Messaging:**

Security communications must adapt to student developmental stages:

- **Age-calibrated security explanations:** Security messages should explain restrictions in developmentally appropriate language rather than technical terms.
- **Learning-oriented feedback:** Security blocks should be framed as learning opportunities rather than punitive measures, where appropriate.
- **Positive security framing:** Security measures should be presented as enabling safe exploration rather than merely imposing limitations.
- **Progressive security awareness:** Interfaces should increase security transparency as students develop, gradually building security understanding.

The current implementation lacks student-facing security messaging entirely, creating both immediate risks and missed educational opportunities around digital citizenship.

### **Classroom-Specific Threat Monitoring:**

Traditional threat detection requires adaptation for educational contexts:

- **Behavioral baseline adaptation:** Threat detection must accommodate natural classroom behavior patterns that might trigger false positives in enterprise systems.
- **Educational context awareness:** Security monitoring should incorporate lesson context to distinguish legitimate experimentation from actual threats.
- **Peer-interaction patterns:** Monitoring should identify concerning patterns across student interactions, not just individual actions.
- **Resource constraint adaptation:** Threat detection must operate efficiently on resource-constrained devices without requiring constant connectivity.

EdgePrompt's security specifications acknowledge monitoring needs conceptually but lack specific implementations adapted to classroom realities.

### **Implementation Recommendations:**

Addressing EdgePrompt's classroom security needs requires targeted enhancements to the current implementation:

1. **Prioritize authentication implementation:** The most critical security gap is basic user authentication and role separation, which should be implemented before any broader deployment.
2. **Implement data compartmentalization:** Even before full encryption, implementing basic data isolation between users would significantly enhance security.
3. **Enhance validation robustness:** Expanding the application's validation logic to match the research framework's multi-stage approach would improve content safety.
4. **Develop teacher security dashboard:** Creating basic security monitoring tools would leverage teacher oversight as a security control.
5. **Implement audit logging:** Adding tamper-evident operation logging would enhance accountability even in shared-device environments.

These targeted enhancements would address the most critical security gaps while aligning with EdgePrompt's resource-constrained deployment vision.

## Part 4: Synthesizing the Assessment – Gaps, Tensions, and Future Paths

Having examined EdgePrompt through multiple specialized lenses—from pedagogical alignment to cybersecurity—we now synthesize these insights into a comprehensive assessment that identifies core gaps, navigates inherent tensions, and outlines future development paths. This synthesis aims not to simply enumerate shortcomings but to constructively map the journey from EdgePrompt's current promising foundation to its ambitious vision of educational AI that serves all learners equitably.

### Chapter 12: The Gap Analysis – Vision vs. Reality Today

EdgePrompt embodies an ambitious vision: bringing AI-enhanced education to resource-constrained environments through offline-capable systems that maintain teacher control while ensuring content safety. The current implementation represents meaningful progress toward this vision, but significant gaps remain between aspiration and reality. This chapter catalogs these gaps across functional, technical, and philosophical dimensions.

#### Functional Gaps: From Vision to Working System

The functional gaps between EdgePrompt's conceptual vision and current implementation reflect both development prioritization and technical challenges:

- **Simplified vs. Multi-Stage Validation:** Perhaps the most significant functional gap exists between the comprehensive multi-stage validation pipeline described in the documentation and implemented in the research framework versus the simplified validation in the current application. While the research code demonstrates the feasibility of the EdgePrompt approach, the application implements only a basic approximation of this sophisticated architecture.
- **Missing Teacher Review System:** The envisioned teacher review system ( $D(r, \theta)$ ,  $T(d)$ ,  $K(h)$ ) for detecting edge cases, triggering reviews, and tracking patterns remains entirely unimplemented. This system, essential for maintaining human oversight, exists only as a conceptual framework.
- **Absent Adaptation Mechanisms:** The adaptation loops ( $A$ ,  $O$ ,  $P$ ) that would enable automatic rubric adjustment, criteria optimization, and template refinement based on performance history are absent from the current implementation, limiting the system's ability to improve over time.
- **Limited Offline Synchronization:** While the system supports basic offline operation, the sophisticated synchronization mechanisms needed for effective operation in intermittently connected environments remain unimplemented. Current functionality focuses on local operation without addressing the complexities of data synchronization during connectivity windows.
- **Missing Security Implementation:** Despite detailed security specifications, the current implementation lacks user authentication, authorization, data protection, and other basic

security controls. This represents a fundamental functional gap that must be addressed before any real-world deployment.

These functional gaps reflect reasonable development prioritization for early-stage research software—validating core concepts before implementing sophisticated features. However, they represent critical limitations for any near-term deployment in actual educational settings.

A comprehensive survey of large language models for education highlights the various applications, challenges, and future research directions in this domain, providing context for understanding EdgePrompt's development gaps within the broader educational AI landscape (Wang S. et al. 2024).

### **AI Capability Gaps: Model Reality vs. Conceptual Needs**

EdgePrompt's design reveals tensions between what current AI models can reliably deliver and what the conceptual framework requires:

- **Edge LLM Reliability Limitations:** The most significant AI capability gap involves the unreliability of edge-deployable models in producing structured outputs like JSON—a limitation that required extensive engineering workarounds. This fundamental gap between expected and actual model behavior has significant implications for validation robustness.
- **Context Window Constraints:** Current edge-deployable models have limited context windows that restrict the complexity of prompts, templates, and validation instructions that can be effectively used. This constraint forces compromises between template sophistication and processing efficiency.
- **Reasoning Depth Limitations:** Edge models exhibit limitations in complex reasoning tasks, constraining validation sophistication and potentially requiring more explicit step-by-step decomposition than initially envisioned.
- **Multi-Stage Efficiency Costs:** The token and latency costs of multi-stage validation using current edge models significantly exceed initial expectations, creating performance challenges for real-time classroom use.
- **Inherent Pedagogical Understanding:** Current models lack deep understanding of pedagogical concepts, requiring explicit encoding of educational principles that might ideally emerge more naturally from the models' knowledge base.

These AI capability gaps reflect the fundamental challenge of deploying sophisticated AI capabilities in resource-constrained environments using current technology. While engineering workarounds can address some limitations, others may require fundamental advances in model architecture or training approaches.

Research on implementing Constitutional AI with smaller models like LLaMA 3-8B has found that increasing harmlessness often comes at the cost of helpfulness, with smaller models showing signs of collapse—highlighting the technical challenges of implementing safety features in edge-deployable models (Zhang 2025).



## Deployment Gaps: From Simulation to Real-World Implementation

EdgePrompt's transition from simulation to real-world implementation faces several critical gaps:

- **Hardware Performance Validation:** While the research framework observes performance metrics during simulation, the system lacks validation on actual target hardware like Jetson Nano or similar edge devices. This creates uncertainty about real-world performance under genuine resource constraints.
- **User Interface Maturity:** The current user interface supports basic functionality but lacks the refinement, error handling, and user experience optimizations needed for non-technical teachers in challenging environments.
- **Deployment Packaging:** The system currently lacks packaging for easy deployment, update management, and maintenance in resource-constrained settings. The development environment configuration differs significantly from deployment requirements.
- **Teacher Training Resources:** Despite the focus on teacher agency, the system lacks comprehensive training materials that would enable teachers to effectively leverage its capabilities without technical support.
- **Integration Mechanisms:** Current implementation doesn't address integration with existing educational systems, content repositories, or classroom workflows, potentially creating adoption barriers.

These deployment gaps highlight the distance between a functioning research prototype and classroom-ready educational technology. Bridging this gap requires significant investment in implementation refinement, user experience optimization, and deployment infrastructure.

Research on on-device language models has reviewed strategies for deploying LLMs on resource-constrained devices, including efficient architectures, compression techniques, and hardware acceleration—approaches that could help address EdgePrompt's deployment challenges (Xu et al. 2024).

## Pedagogical Gaps: Design Intentions vs. Learning Realities

Several gaps exist between EdgePrompt's pedagogical vision and its current implementation:

- **Teacher Agency vs. System Complexity:** While designed to enhance teacher agency, the system's current complexity may actually reduce agency for less technical teachers who cannot effectively manipulate templates or understand validation processes.
- **Constructivist Intent vs. Assessment Focus:** Despite valuing constructivist learning principles, the current implementation focuses heavily on question generation and answer validation rather than knowledge construction or exploratory learning.
- **Student-Centered Vision vs. Tool-Centric Reality:** The vision emphasizes student-centered learning, but the implementation centers around tool capabilities rather than directly supporting diverse student learning processes.

- **Pedagogical Flexibility vs. Template Rigidity:** The template-based approach, while providing safety guardrails, potentially limits pedagogical flexibility and creative teaching approaches.
- **Metacognitive Development vs. Content Delivery:** The current focus on content generation and assessment provides limited support for metacognitive development or reflection, despite their importance for deep learning.

These pedagogical gaps require attention not just to technical implementation but to fundamental design questions about how EdgePrompt can genuinely support varied learning processes rather than simply automating traditional educational activities.

Research on harnessing AI for constructivist learning argues that modern AI tools can amplify constructivist pedagogy when designed to support learners as "active architects" of their own knowledge building (Grubaugh, Levitt & Deever 2023). Additionally, Vygotsky's work on the Zone of Proximal Development emphasizes the importance of appropriate scaffolding that supports learning in the gap between independent and assisted achievement (Vygotsky 1978).

### **Ethical Gaps: Stated Values vs. Implemented Safeguards**

EdgePrompt's ethical vision requires implementation of concrete safeguards that remain partially realized:

- **Privacy Protection Claims vs. Implementation:** Despite emphasizing student privacy, the current implementation lacks basic privacy protections like data compartmentalization, encryption, or consent mechanisms.
- **Equity Vision vs. Accessibility Reality:** While designed to enhance educational equity, the current implementation addresses only connectivity barriers without integrated solutions for language, ability, or digital literacy barriers.
- **Cultural Responsiveness Goal vs. Western Defaults:** The vision emphasizes cultural responsiveness, but implementation relies on western educational assumptions embedded in templates and validation approaches.
- **Safety Focus vs. Limited Content Filtering:** Despite centering safety, current content filtering capabilities remain basic compared to the sophisticated multi-stage approach described in the documentation.
- **Transparency Intention vs. Opaque Processes:** The system aims for transparent AI use but provides limited visibility into validation decisions or reasoning processes in the current implementation.

These ethical gaps reflect the challenge of translating ethical principles into concrete technical implementations. Addressing them requires both enhanced technical controls and deeper integration of ethical considerations into the development process.

Research on the ethics of AI in education emphasizes the importance of transparency in AI decision-making, safeguarding student autonomy and trust, and adherence to principles of fairness,

accountability, and pedagogical soundness (Holmes et al. 2022). Additionally, philosophical work on artificial intelligence, values, and alignment provides frameworks for understanding how AI systems can be designed to uphold human values (Gabriel 2020).

### **Infrastructure Gaps: Vision Requirements vs. Implementation Realities**

EdgePrompt's infrastructure vision faces several implementation challenges:

- **Offline Capability vs. Synchronization Complexity:** While conceptually supporting offline operation, the system hasn't yet addressed the complex synchronization challenges that arise in intermittently connected environments.
- **Resource Efficiency Goals vs. Multi-Stage Costs:** The efficiency cost of multi-stage validation potentially conflicts with the goal of operation on minimal hardware, requiring careful optimization.
- **Secure Operation vs. Basic Implementation:** The absence of robust security mechanisms creates a fundamental infrastructure gap between secure operation requirements and current capabilities.
- **Extensibility Vision vs. Limited Modularity:** Despite emphasizing extensibility, the current implementation lacks clearly documented extension points or modularity for community adaptation.
- **Hardware Target Clarity vs. Development Environment:** The system lacks clear documentation of minimum hardware requirements or optimization for specific target platforms.

These infrastructure gaps highlight the challenges of developing sophisticated systems for resource-constrained environments. While the conceptual architecture acknowledges these challenges, the implementation requires additional refinement to fully address them.

Research on educational edge computing frameworks demonstrates how local edge servers can handle intensive tasks with minimal latency, improving throughput for remote and rural learners—approaches that could inform EdgePrompt's infrastructure optimization (Chen et al. 2022).

### **Nature of the Gaps: Technical Limits, Engineering Challenges, Socio-Technical Complexity**

Understanding the nature of these gaps informs prioritization and resolution approaches:

- **Technical Limitation Gaps:** Some gaps reflect fundamental limitations of current technology, particularly regarding edge LLM capabilities and performance. These gaps may require advances in model architecture or training before full resolution is possible.
- **Engineering Implementation Gaps:** Many gaps represent engineering tasks that are conceptually understood but not yet implemented due to development prioritization. These gaps can be addressed through focused development efforts without requiring fundamental research advances.

- **Socio-Technical Design Gaps:** The most complex gaps involve the intersection of technical capabilities with social, educational, and ethical requirements. These socio-technical gaps require not just implementation but iterative co-design with educational stakeholders.
- **Resource Constraint Gaps:** Several gaps reflect the inherent challenge of implementing sophisticated AI capabilities within significant resource constraints. These gaps require creative optimization and potential compromise rather than straightforward implementation.

This nuanced understanding of gap types highlights that resolution requires diverse approaches—from technical research to stakeholder engagement to resource optimization—rather than simply more development time.

Research on AI system evaluation frameworks has proposed comprehensive approaches that span the entire model lifecycle, implementing component-level checks and system-level validation mapped to different stakeholders and development stages (Xia et al. 2024).

### **Prioritization Framework: Critical Path vs. Enhancement Opportunities**

Not all gaps require immediate resolution for EdgePrompt to deliver value. A prioritization framework might categorize gaps as follows:

- **Critical Path Gaps:** Gaps that fundamentally prevent safe, effective use in educational settings. These include basic security implementation, validation robustness, and hardware performance validation. Without addressing these gaps, classroom deployment would be premature.
- **Core Value Gaps:** Gaps that significantly limit the system's ability to deliver on its core value proposition, such as multi-stage validation implementation in the application, teacher review mechanisms, and basic offline synchronization.
- **Enhancement Opportunity Gaps:** Gaps that represent opportunities for expanding value but don't prevent effective use, such as adaptation mechanisms, advanced cultural responsiveness, or sophisticated integration capabilities.
- **Visionary Future Gaps:** Gaps between current capabilities and the long-term vision that may require substantial technological advances or pedagogical innovation beyond immediate implementation, such as advanced metacognitive support or fully personalized learning.

This prioritization framework enables strategic planning for moving EdgePrompt from promising research prototype to classroom-ready educational technology in a measured, responsible manner.

## Chapter 13: Navigating the Tensions – Core Conflicts in Design

Beyond specific implementation gaps, EdgePrompt embodies fundamental tensions—inherent trade-offs with no perfect resolution. These tensions don't represent flaws but rather core design conflicts requiring continual, contextual navigation rather than one-time resolution. This chapter examines these tensions and explores potential navigation strategies.

### Safety vs. Capability

Perhaps the most fundamental tension in EdgePrompt's design is between ensuring safety through constraints and enabling powerful educational capabilities:

- **The Tension:** Stronger safety guardrails typically limit the system's expressive capabilities and pedagogical flexibility. Conversely, increasing capability often requires relaxing constraints that ensure safety and alignment.
- **Manifestation in EdgePrompt:** This tension appears throughout the architecture, from the structured template design that ensures safety at the cost of flexibility to the multi-stage validation that improves safety but increases processing time and complexity. The JSON reliability challenges exemplify this tension—structured outputs enhance safety but create reliability issues with current models.
- **Trade-off Landscape:** This tension isn't binary but exists on a complex landscape where the optimal balance varies by:
  - Student age and developmental needs
  - Subject matter sensitivity
  - Educational setting formality
  - Teacher technical expertise
  - Cultural and community context
- **Navigation Strategies:** Potential approaches for navigating this tension include:
  - Contextual constraint adjustment based on educational parameters
  - Graduated safety controls with teacher configuration capabilities
  - Transparent constraint explanation for educational context
  - Intelligent fallback mechanisms when safety requirements conflict with educational needs

This tension cannot be "solved" but must be continually navigated through deliberate design choices that acknowledge the inherent trade-offs between safety and capability in educational AI.

Research on building guardrails for large language models through structured prompting and validation stages provides approaches for enhancing safety in AI outputs while maintaining usability (Dong et al. 2024). Additionally, studies on implementing Constitutional AI with smaller models highlight the tradeoffs between harmlessness and helpfulness, noting that increasing safety constraints can sometimes lead to reduced model utility (Zhang 2025).

## Efficiency vs. Robustness

EdgePrompt's deployment in resource-constrained environments creates inherent tension between computational efficiency and system robustness:

- **The Tension:** More robust systems typically require additional computation for validation, error checking, and graceful handling of edge cases. However, resource-constrained environments demand minimal computational overhead.
- **Manifestation in EdgePrompt:** This tension appears most clearly in the multi-stage validation approach, which creates significant token and latency overhead compared to simpler approaches. The robust JSON parsing mechanisms similarly add computational complexity to address LLM output unreliability.
- **Trade-off Dimensions:** This tension varies across:
  - Hardware capability targets
  - Response time expectations
  - Validation criticality for specific content
  - Connectivity patterns affecting synchronization robustness
- **Navigation Strategies:** Potential approaches for navigating this tension include:
  - Tiered validation based on content risk assessment
  - Asynchronous processing for non-time-critical validations
  - Computational resource allocation proportional to validation importance
  - Template optimization for specific hardware targets
  - Progressive enhancement based on available resources

This tension requires continuous balancing between ensuring system reliability and maintaining acceptable performance in resource-constrained environments, with no single optimal balance point across all contexts.

Research on educational edge computing frameworks demonstrates approaches for balancing performance requirements with resource constraints in educational applications, showing how edge setups can improve throughput and latency for remote learners (Chen et al. 2022). Additionally, comprehensive reviews of on-device language models highlight strategies for efficient deployment on resource-constrained devices, including architecture optimization and compression techniques (Xu et al. 2024).

## Standardization vs. Contextualization

EdgePrompt faces tension between standardizing approaches for consistency and contextualizing for diverse educational settings:

- **The Tension:** Standardized templates, validation approaches, and interactions enhance predictability, maintainability, and systematic improvement. However, effective education typically requires contextualization to specific cultural, linguistic, and pedagogical environments.

- **Manifestation in EdgePrompt:** This tension emerges in template design (standardized vs. customizable), validation criteria (universal vs. context-specific), and interaction models (consistent vs. culturally adapted).
- **Contextual Variations:** This tension varies across:
  - Cultural contexts and linguistic requirements
  - Educational system structures
  - Subject-specific pedagogical approaches
  - Teacher preferences and teaching styles
  - Community educational values
- **Navigation Strategies:** Potential approaches for navigating this tension include:
  - Core standardization with contextual extension points
  - Layered templates with cultural adaptation capabilities
  - Configurable validation standards reflecting local educational norms
  - Participatory design processes for context-specific implementations
  - Flexible architecture supporting diverse interaction patterns

This tension requires ongoing negotiation between standardization for scalability and contextualization for effectiveness, ideally finding balance through modular design that enables customization within consistent architectural patterns.

Research on culturally adaptive thinking in education for AI highlights the importance of localizing AI curricula and tools to align with local cultural norms and languages to avoid cultural bias and resistance (Samuel et al. 2023). Additionally, work on culturally responsive teaching emphasizes incorporating students' cultural references in all aspects of learning, validating and affirming diverse cultural identities (Gay 2018).

### **Automation vs. Human Agency (Teacher & Student)**

EdgePrompt embodies the fundamental tension between automating educational processes for efficiency and preserving human agency for pedagogical effectiveness:

- **The Tension:** Automation can increase efficiency, consistency, and access to educational resources but potentially reduces human agency, judgment, and educational relationship quality. Preserving human agency maintains educational authenticity but limits scalability and resource efficiency.
- **Manifestation in EdgePrompt:** This tension appears in decisions about automation levels for content generation, validation, and feedback, as well as in the balance between system-defined and teacher-defined parameters.
- **Agency Dimensions:** This tension varies across:
  - Teacher vs. student agency considerations
  - Process types (assessment vs. exploration)
  - Educational goals (procedural knowledge vs. creative expression)

- Resource availability affecting teacher time constraints
- Stakeholder technical capabilities
- **Navigation Strategies:** Potential approaches for navigating this tension include:
  - Explicit agency design showing where humans retain control
  - Task allocation based on comparative advantages (human vs. AI)
  - Transparent intervention options at multiple process points
  - Configurable automation levels based on context
  - Co-creative rather than replacement-oriented automation

This tension requires ongoing reflection about which aspects of education benefit from automation and which require human judgment, with the understanding that this balance may shift across contexts and over time.

Research on teacher agency in the age of generative AI emphasizes the importance of maintaining teachers' power to act, affect matters, make decisions, and take stances in educational settings (Frøsig & Romero 2024). Additionally, work on epistemic agency in education highlights the value of viewing learners as active agents capable of evaluating and constructing knowledge rather than passive recipients (Elgin 2013).

### Openness vs. Sustainability

EdgePrompt's vision as "public infrastructure" creates tension between openness for community ownership and sustainability for long-term operation:

- **The Tension:** Open, community-owned systems enhance equity, adaptation, and educational sovereignty. However, sustainable operation requires resources, maintenance, and potentially commercial or institutional support structures.
- **Manifestation in EdgePrompt:** This tension appears in decisions about licensing, development models, deployment approaches, and long-term support strategies. The project's community-oriented vision conflicts with the resource requirements of sophisticated AI infrastructure.
- **Contextual Factors:** This tension varies across:
  - Deployment contexts (resource levels, institutional support)
  - Technical expertise availability for maintenance
  - Funding models and priorities
  - Community capacity for ownership
  - Educational system structures
- **Navigation Strategies:** Potential approaches for navigating this tension include:
  - Tiered sustainability models based on context
  - Community-institutional partnerships for resource sharing
  - Modularity enabling component-level sustainability approaches
  - Transparent value exchange mechanisms



- Progressive ownership transfer strategies
- Documentation and training for community maintenance capacity

This tension requires creative approaches that balance openness for equitable access with sustainability mechanisms for continued operation and improvement, potentially through hybrid models that combine community ownership with institutional support.

### Technical Optimization vs. Human Judgment

An additional tension emerges between optimizing for technical metrics and prioritizing human educational judgment:

- **The Tension:** Technical optimization improves measurable performance metrics like token efficiency, validation accuracy, or generation quality. However, educational effectiveness ultimately depends on human judgment about learning value that may not align with technical optimization targets.
- **Manifestation in EdgePrompt:** This tension appears in decisions about what to optimize—technical metrics like token efficiency or validation accuracy versus less measurable factors like educational relevance or cultural appropriateness.
- **Judgment Dimensions:** This tension varies across:
  - Stakeholder perspectives (technical vs. educational)
  - Optimization targets (efficiency vs. effectiveness)
  - Measurement approaches (quantitative vs. qualitative)
  - Value frameworks (technical vs. pedagogical)
- **Navigation Strategies:** Potential approaches for navigating this tension include:
  - Mixed-metric evaluation frameworks
  - Participatory evaluation design with educators
  - Complementary technical and educational assessment
  - Contextual weighting of technical vs. human judgment
  - Transparent rationale for optimization priorities

This tension requires maintaining the primacy of educational judgment even while pursuing technical optimization, ensuring that computational efficiency serves pedagogical effectiveness rather than becoming an end in itself.

### Productive Navigation: From Resolution to Balance

These tensions cannot be permanently resolved but require continuous, contextual navigation. Several principles might guide this navigation:

- **Explicit Acknowledgment:** Recognizing and explicitly discussing tensions rather than obscuring them enables more thoughtful decision-making.

- **Contextual Calibration:** Different balances may be appropriate for different educational contexts, age groups, or subject areas.
- **Stakeholder Participation:** Involving diverse stakeholders, particularly educators and students, in navigating tensions ensures decisions reflect educational priorities.
- **Transparency of Trade-offs:** Making trade-off rationales transparent helps users understand why certain balances were chosen.
- **Adaptation Capability:** Building systems that allow recalibration of these balances based on experience and changing needs.
- **Values-Based Decision Frameworks:** Developing explicit frameworks connecting design decisions to educational values.

Research on socio-technical education futures emphasizes examining educational technology within broader socio-technical assemblages that include social values and power structures, acknowledging that technology and society co-evolve in educational settings (Swist & Gulson 2023). Additionally, work on new cybernetics provides systems thinking approaches for understanding educational environments as complex systems requiring adaptive and responsive design (Bell 2022).

EdgePrompt's long-term success depends not on finding perfect resolution to these tensions but on building architecture that supports their thoughtful, contextually appropriate navigation across diverse educational settings. This requires moving beyond technical optimization toward values-based design that centers educational priorities even while acknowledging technical constraints.

## Chapter 14: Recommended Future Work – From Validation to Value

Having analyzed EdgePrompt's current state, gaps, and inherent tensions, this chapter outlines a strategic roadmap for moving the project from its promising research foundation toward real-world educational value. This roadmap focuses not on aspirational features but on concrete next steps that address critical gaps while navigating the core tensions identified in previous chapters.

### **Prioritized Roadmap: Security/Privacy, Efficiency/Hardware Testing, Full Validation Implementation, Human Evaluation, Pedagogical Enhancement, Offline Architecture**

A strategic roadmap for EdgePrompt development would prioritize work in the following areas:

#### **1. Security and Privacy Foundation (Critical)**

Current security implementation gaps represent fundamental barriers to responsible deployment. Priority security enhancements should include:

- **User Authentication Implementation:** Integrate the identity management architecture described in the security specification, with appropriate adaptations for educational settings.
- **Role-Based Access Control:** Implement teacher/student role separation with appropriate permissions and boundaries.
- **Data Compartmentalization:** Ensure separation between user data even before implementing full encryption.
- **Basic Audit Logging:** Implement tamper-evident logging of significant system actions for accountability.
- **Input/Output Validation:** Enhance input validation and content filtering to prevent injection and inappropriate content generation.

These security enhancements represent minimum requirements for responsible classroom testing beyond controlled research environments. While full implementation of the sophisticated security specification may be deferred, these basic protections should be prioritized.

#### **2. Efficiency Optimization and Hardware Testing (Critical)**

The gap between simulation observations and real-world performance requires targeted optimization and testing:

- **Validation Sequence Optimization:** Review and refine validation stages to reduce token usage and latency while maintaining effectiveness.
- **Real Hardware Benchmarking:** Test performance on actual target hardware platforms (Jetson, entry-level laptops) to validate feasibility.
- **Tiered Validation Implementation:** Develop context-sensitive validation levels that adapt to content risk and available resources.

- **Template Efficiency Refinement:** Optimize templates for token efficiency while maintaining educational effectiveness.
- **Resource Monitoring Integration:** Implement real-time resource monitoring to support dynamic optimization strategies.

These optimizations are critical for determining whether EdgePrompt's multi-stage validation approach is genuinely feasible on target hardware or requires fundamental reconsideration.

### 3. Full Validation Implementation in Application (High)

The gap between the research framework's sophisticated validation and the application's simplified approach requires addressed:

- **Multi-Stage Validation Integration:** Implement the full multi-stage validation sequence from the research framework in the application backend.
- **Robust JSON Handling:** Integrate the robust json\_utils and fallback mechanisms developed in research into the application codebase.
- **Validation Visualization:** Develop UI components showing validation stage outcomes for transparency.
- **Teacher Review Interface:** Implement basic review capabilities for flagged content or borderline validations.
- **Validation Configuration:** Create teacher-accessible controls for adjusting validation parameters based on educational context.

These enhancements would bring the application's capabilities in line with the validated research approach, enabling more accurate assessment of the full EdgePrompt methodology in classroom settings.

### 4. Human Evaluation Framework (High)

Moving beyond simulation requires structured human evaluation by actual educators:

- **Teacher Evaluation Protocol:** Develop a structured protocol for teacher evaluation of EdgePrompt's educational effectiveness.
- **Student Experience Assessment:** Create student-appropriate feedback mechanisms for understanding the learning experience.
- **Comparative Testing Framework:** Implement approaches for comparing EdgePrompt with alternative educational approaches.
- **Longitudinal Assessment Design:** Plan for measuring impact over extended usage periods, not just initial impressions.
- **Mixed-Methods Evaluation:** Combine quantitative metrics with qualitative insights for comprehensive assessment.

Human evaluation represents a critical step beyond simulation for understanding EdgePrompt's real-world educational value and guiding further development.

## 5. Pedagogical Enhancement (Medium)

Addressing pedagogical gaps requires targeted enhancements to better align with educational principles:

- **Constructivist Interaction Models:** Develop templates and workflows supporting more exploratory, constructivist learning approaches.
- **Student Agency Enhancement:** Create capabilities for student-initiated interactions rather than only teacher-defined activities.
- **Metacognitive Support:** Implement reflection prompts and metacognitive scaffolding beyond direct content delivery.
- **Culturally Responsive Templates:** Develop template adaptations for specific cultural and educational contexts beyond translation.
- **Pedagogical Flexibility:** Enhance system adaptability to diverse teaching approaches and philosophies.

These enhancements would help bridge the gap between EdgePrompt's educational vision and its current implementation, moving beyond content delivery toward deeper learning support.

## 6. Robust Offline Architecture (Medium)

Fully realizing the offline vision requires addressing synchronization and deployment challenges:

- **Offline Data Synchronization:** Implement robust synchronization mechanisms for intermittent connectivity environments.
- **Conflict Resolution Strategies:** Develop approaches for handling data conflicts during synchronization.
- **Deployment Packaging:** Create packaging for easy deployment, updates, and maintenance in resource-constrained settings.
- **Resource-Adaptive Operation:** Implement capabilities for adapting to varying resource availability across environments.
- **Self-Diagnostics:** Develop system self-assessment capabilities for identifying issues in offline environments.

These enhancements would strengthen EdgePrompt's core value proposition of offline operation in resource-constrained environments.

## 7. Teacher Ownership Enhancement (Medium)

Deepening teacher control and ownership requires specific capability development:

- **Template Customization Tools:** Create user-friendly interfaces for template adaptation without technical expertise.
- **Collaborative Template Libraries:** Develop mechanisms for sharing and reusing templates across educator communities.
- **Transparent AI Operation:** Enhance visibility into AI decision processes for teacher understanding.
- **Progressive Technical Control:** Implement layered interfaces providing increasing technical control as teacher expertise develops.
- **Local Governance Tools:** Create mechanisms supporting community oversight of AI deployment in educational settings.

These enhancements would strengthen EdgePrompt's commitment to teacher agency and community ownership.

## 8. Systematic Technical Hardening (Ongoing)

Beyond specific feature enhancements, systemic technical improvements include:

- **Comprehensive Testing:** Develop systematic testing across diverse scenarios and edge cases.
- **Performance Optimization:** Continuously refine core operations for efficiency on resource-constrained hardware.
- **Error Handling Robustness:** Enhance error detection, reporting, and recovery throughout the system.
- **Documentation Enhancement:** Develop comprehensive documentation for developers, administrators, and users.
- **Accessibility Compliance:** Ensure interfaces meet accessibility standards for diverse users.

These ongoing technical improvements support overall system quality and sustainability.

## The Role of Transdisciplinary Input (Education Management, Cybernetics, Philosophy)

EdgePrompt's success requires not just technical development but transdisciplinary input across several domains:

### Education Management Perspectives

Educational deployment requires management insights beyond technical implementation:

- **Integration with Educational Systems:** Understanding how EdgePrompt fits within existing educational structures, policies, and workflows.
- **Professional Development Models:** Designing effective teacher training and support for technology adoption.
- **Resource Allocation Frameworks:** Developing sustainable approaches to hardware, maintenance, and support resource allocation.

- **Policy Alignment Strategies:** Ensuring alignment with educational policies at institutional, regional, and national levels.
- **Impact Assessment Frameworks:** Creating comprehensive approaches for evaluating educational technology effectiveness.

Education management perspectives help bridge the gap between technical capabilities and educational implementation realities.

### Cybernetic Systems Thinking

EdgePrompt's architecture can benefit from cybernetic insights about system behavior:

- **Feedback Loop Design:** Structuring effective feedback mechanisms between system components and with human users.
- **Adaptation Mechanisms:** Designing robust approaches for system adaptation to changing conditions and requirements.
- **Emergent Behavior Analysis:** Understanding and managing emergent behaviors in complex human-AI educational systems.
- **System Boundary Management:** Defining appropriate boundaries between technical systems and human educational processes.
- **Homeostatic Balance:** Maintaining system stability while enabling appropriate adaptation and growth.

Cybernetic perspectives help ensure that EdgePrompt functions as an effective component within broader educational systems rather than an isolated technological artifact. Research on new cybernetics provides systems thinking approaches for understanding educational environments as complex systems requiring adaptive and responsive design (Bell 2022).

### Philosophical and Ethical Frameworks

EdgePrompt's development requires grounding in philosophical and ethical considerations:

- **Value Alignment Methodologies:** Approaches for ensuring system alignment with human educational values.
- **Epistemic Agency Frameworks:** Understanding how AI systems affect human knowledge construction and agency.
- **Justice and Equity Considerations:** Frameworks for evaluating and enhancing technology's contribution to educational justice.
- **Cultural-Philosophical Perspectives:** Understanding how diverse cultural philosophies should shape educational technology.
- **Human-Technology Relationship Models:** Frameworks for conceptualizing productive human-AI educational partnerships.

Philosophical perspectives help ensure that technical development remains grounded in deeper human values and educational purposes.

### **Addressing the Long-Term Vision: Towards the Learning Companion**

Beyond addressing immediate gaps, future work should advance EdgePrompt toward its vision as a true "learning companion":

#### **Memory and Longitudinal Relationship**

Developing meaningful educational relationships requires memory across interactions:

- **Student Progress Tracking:** Implementing mechanisms for tracking individual student progress over time.
- **Interaction History Maintenance:** Building appropriate memory of past interactions to inform future responses.
- **Pattern Recognition Across Sessions:** Developing capabilities for identifying learning patterns across multiple sessions.
- **Forgetting Mechanisms:** Implementing appropriate information decay to maintain relevance and privacy.
- **Context-Sensitive Recall:** Developing nuanced approaches for when and how to reference past interactions.

These capabilities would transform EdgePrompt from a stateless tool to a system capable of meaningful educational relationships.

#### **Adaptive Personalization**

Moving beyond standardized interactions requires adaptive personalization:

- **Learning Style Adaptation:** Developing capabilities for adapting to individual learning styles and preferences.
- **Difficulty Calibration:** Implementing mechanisms for adjusting challenge levels based on demonstrated capabilities.
- **Interest-Based Engagement:** Creating approaches for connecting content to individual student interests.
- **Metacognitive Scaffolding:** Developing personalized support for individual metacognitive development.
- **Intervention Targeting:** Implementing capabilities for identifying and addressing specific learning challenges.

These capabilities would enhance EdgePrompt's educational effectiveness through more tailored support.



## Teacher Augmentation Rather than Replacement

Fully realizing the teacher-centric vision requires:

- **Teacher Insight Tools:** Developing mechanisms for surfacing patterns and insights for teacher consideration.
- **Workflow Integration:** Creating seamless connections to existing teacher workflows rather than parallel processes.
- **Cognitive Load Reduction:** Identifying and automating appropriate administrative tasks while preserving pedagogical judgment.
- **Co-Teaching Capabilities:** Developing models for effective human-AI educational collaboration.
- **Teacher Learning Support:** Creating mechanisms that enhance teacher professional development alongside student learning.

These capabilities would strengthen EdgePrompt's role as teacher augmentation rather than replacement.

## Community Ownership and Adaptation

Advancing the "public infrastructure" vision requires:

- **Local Extension Mechanisms:** Developing clear, accessible extension points for community customization.
- **Knowledge Sovereignty:** Implementing approaches ensuring community control over educational content and data.
- **Participatory Governance:** Creating frameworks for community participation in system governance.
- **Value Exchange Transparency:** Ensuring clear understanding of costs, benefits, and data usage.
- **Progressive Transfer:** Developing pathways for increasing community ownership and control over time.

These capabilities would help realize EdgePrompt's vision as community-owned educational infrastructure rather than merely a technological product.

## The Path Forward: From Good Intentions to Transformative Impact

EdgePrompt represents a promising approach to bringing AI-enhanced education to resource-constrained environments. Moving from current research to transformative impact requires:

- **Technical Discipline:** Addressing critical implementation gaps, particularly in security and validation robustness.
- **Educational Integration:** Ensuring technical development remains grounded in educational purpose and effectiveness.

- **Participatory Development:** Engaging educators and students as partners rather than subjects in ongoing development.
- **Ethical Vigilance:** Continuously examining how implementation choices affect equity, agency, and educational values.
- **Sustained Investment:** Acknowledging the significant resources required to bridge the gap between vision and reality.

This path involves not just technical development but ongoing negotiation of the inherent tensions identified in EdgePrompt's design. Success will be measured not merely by functional completeness but by genuine contribution to educational equity and effectiveness in the communities it aims to serve.

A comprehensive survey of large language models for education provides context for understanding EdgePrompt's place within the broader educational AI landscape, highlighting applications, challenges, and future research directions in this domain (Wang S. et al. 2024). Additionally, research on AI system evaluation frameworks offers approaches for comprehensive assessment spanning the entire model lifecycle (Xia et al. 2024).

## Chapter 15: Open Questions – The Dialectic of Further Research

### Beyond Answers to Deeper Questions

As we conclude our assessment of EdgePrompt, we find ourselves not with neat conclusions but with richer questions. The transdisciplinary examination has revealed not just specific implementation gaps but fundamental tensions that require ongoing navigation rather than one-time resolution. This chapter frames critical research directions not as a simple list but as a dialectical journey—deliberately exploring opposing forces whose productive tension may generate deeper insights than either perspective alone could provide.

Our approach here is intentionally Socratic, positioning questions as tools for thinking rather than merely gaps to fill. We organize these questions around core tensions identified throughout our assessment, recognizing that the most valuable research will not merely solve technical problems but navigate complex trade-offs between competing values, approaches, and priorities.

Each section presents a fundamental tension, explores its manifestations across EdgePrompt's domains, and offers structured research questions whose investigation might help navigate—though never fully resolve—these tensions. This dialectical approach reflects the reality that EdgePrompt exists not in a space of simple solutions but at the intersection of multiple valid perspectives that must be held in productive balance.

### The Security-Accessibility Dialectic: Protection Without Exclusion

**The Tension:** EdgePrompt faces a fundamental tension between robust security that protects educational stakeholders and accessibility that ensures equitable participation. More sophisticated security mechanisms often create barriers to participation, particularly for those with limited technical resources or expertise. Yet inadequate security creates vulnerabilities that may harm the very communities the system aims to serve.

**The Questions:** How might we conceptualize security not as a barrier to access but as an enabler of appropriate participation? What would security look like if designed from the perspective of the least-resourced participants rather than the most technically sophisticated? Can security and accessibility be reconceptualized not as opposing forces but as mutually reinforcing aspects of a trustworthy system?

**Research Directions:**

Research Question	Investigation Approach	Expected Outcome
How can authentication remain secure and usable in intermittent connectivity environments?	Simulation of authentication flows under varying connectivity patterns; analysis of offline authentication mechanisms across domains	Authentication framework balancing security, usability, and offline functionality
What classroom-specific threat model best captures unique K-12 adversarial scenarios?	AI-driven analysis of documented educational technology incidents; scenario simulation of student-specific exploitation vectors	Education-specific threat model with probability-impact matrices calibrated to classroom contexts
How can content safety be maintained on edge-deployed LLMs across diverse cultural contexts?	Adversarial prompting tests; comparative analysis of filtering approaches; cultural variation modeling	Multi-layered safety framework with performance optimization for edge deployment
What synchronization protocol best secures educational data during brief connectivity windows?	Simulation of connectivity scenarios with varying constraints; bandwidth-security optimization modeling	Resource-efficient synchronization protocol with conflict resolution and security guarantees
How can zero trust principles be adapted to educational environments with shared devices?	Modeling trust boundaries for educational workflows; usability-security trade-off analysis	Education-specific zero trust framework addressing classroom device sharing realities

**Opposing Perspectives:** Consider how different stakeholders conceptualize this tension: Security professionals might prioritize comprehensive protection against all threats, regardless of usability

impact. Community advocates might prioritize removing barriers to participation, even at some security cost. Teachers might seek a middle ground where security remains invisible until genuinely needed. What new approaches emerge when we refuse to accept either perspective as complete?

Research on the impact of generative AI in cybersecurity and privacy has examined security implications for educational contexts, highlighting the importance of protecting student data and implementing appropriate safeguards (Gupta et al. 2023). Additionally, work on AI misuse and information privacy concerns has emphasized the importance of robust safety guardrails in addressing these concerns (Menard & Bott 2024).

### **The Technical-Pedagogical Dialectic: Capability vs. Educational Value**

**The Tension:** A persistent tension exists between technical capabilities (what AI systems can do) and pedagogical value (what enhances learning). Technical optimization often prioritizes measurable metrics like efficiency, accuracy, and reliability. Educational effectiveness, however, sometimes benefits from productive struggle, ambiguity, and personalized approaches resistant to standardization. How should EdgePrompt navigate between technical excellence and pedagogical wisdom?

**The Questions:** What happens when we view technical performance as subordinate to educational purpose rather than as an end in itself? Conversely, what educational possibilities emerge when we fully embrace new technical capabilities rather than fitting them into existing pedagogical models? How might the productive friction between technical and educational perspectives generate approaches that transcend both?

### **Research Directions:**

<b>Research Question</b>	<b>Investigation Approach</b>	<b>Expected Outcome</b>
How can AI-driven activities better support constructivist learning rather than knowledge transmission?	Analysis of interaction patterns that foster knowledge construction; prototype development of constructivist-aligned templates	Design patterns for AI educational activities that support authentic knowledge construction
What approaches enable edge-deployed AI to support higher-order thinking despite resource constraints?	Cognitive task analysis of higher-order thinking activities; template optimization for metacognitive support	Framework for designing validation templates that assess and encourage higher-order thinking

Research Question	Investigation Approach	Expected Outcome
How can EdgePrompt better support formative assessment rather than simply summative evaluation?	Analysis of formative feedback patterns; simulation of feedback mechanisms under resource constraints	Assessment design patterns that emphasize learning process over correctness
What balance between structure and exploration best supports productive struggle?	Modeling of scaffolding approaches across learning contexts; prototype testing of variable support mechanisms	Adaptive scaffolding framework that maintains productive challenge while preventing frustration
How can teacher epistemological beliefs be better reflected in AI-assisted educational activities?	Analysis of template customization patterns; modeling of pedagogical value expression in prompts	Design guidelines for aligning AI-assisted activities with diverse teaching philosophies

**Opposing Perspectives:** Technical and educational stakeholders often speak different languages about the same system. Engineers might focus on optimizing token efficiency and response accuracy, viewing education as an application domain. Educators might focus on learning outcomes and student engagement, viewing technology as merely a tool. What happens when we refuse to privilege either perspective and instead seek approaches that transform both our understanding of education and our implementation of technology?

Research on harnessing AI for constructivist learning argues that modern AI tools can amplify constructivist pedagogy when designed to support learners as "active architects" of their own knowledge building (Grubaugh, Levitt & Deever 2023). Additionally, work on epistemic agency in education emphasizes the importance of viewing learners as active agents capable of evaluating and constructing knowledge rather than passive recipients (Elgin 2013).

### **The Global-Local Dialectic: Universal Technology vs. Contextual Knowledge**

**The Tension:** EdgePrompt embodies a tension between global technological approaches (standardized AI models, architecture patterns, interface designs) and local knowledge systems (cultural practices, pedagogical traditions, linguistic patterns). Technology typically seeks universality and standardization

for efficiency and scale. Education, however, derives much of its power from cultural relevance, contextual appropriateness, and local knowledge traditions.

**The Questions:** Is it possible to create technological systems that genuinely respect and incorporate diverse knowledge systems rather than merely accommodating them? What would it mean to design educational AI not from assumed universal principles but from the specific knowledge traditions of target communities? Can the tension between global and local approaches generate new hybrid knowledge systems that honor both technological capabilities and cultural wisdom?

**Research Directions:**

Research Question	Investigation Approach	Expected Outcome
How can EdgePrompt address hardware disparities beyond connectivity barriers?	Hardware availability analysis across target regions; progressive enhancement modeling	Technical architecture supporting functionality tiers across diverse hardware capabilities
What mechanisms best ensure knowledge sovereignty and community ownership?	Governance model analysis across cultural contexts; community involvement simulation	Community ownership framework with appropriate governance structures and control mechanisms
How can cultural knowledge systems be meaningfully represented in AI-driven educational activities?	Knowledge organization analysis across cultures; epistemological framework modeling	Design patterns for representing diverse knowledge systems in educational AI
What adaptation mechanisms allow effective operation across linguistic variations beyond standard languages?	Linguistic variation analysis in target regions; prompt adaptation simulation for dialects	Linguistic adaptation framework supporting regional dialects and language variations

Research Question	Investigation Approach	Expected Outcome
How can community involvement in design and implementation be structured for authentic participation?	Participatory design pattern analysis; community engagement model simulation	Co-design methodology adapted to various community contexts and technical accessibility

**Opposing Perspectives:** Consider the perspective of global technology developers seeking universal solutions versus local communities with specific needs and traditions. Technologists might argue that fundamental patterns transcend cultural differences and enable efficiency through standardization. Cultural advocates might counter that meaningful education must emerge from specific cultural contexts and knowledge traditions. What approaches emerge when we refuse to see either universal or contextual approaches as inherently superior?

Research on culturally adaptive thinking in education for AI highlights the importance of localizing AI curricula and tools to align with local cultural norms and languages to avoid cultural bias and resistance (Samuel et al. 2023). Additionally, work on culturally responsive teaching emphasizes incorporating students' cultural references in all aspects of learning, validating and affirming diverse cultural identities (Gay 2018).

### **The Efficiency-Robustness Dialectic: Resource Constraints vs. System Reliability**

**The Tension:** EdgePrompt's deployment in resource-constrained environments creates inherent tension between computational efficiency and system robustness. More thorough safety mechanisms, validation processes, and adaptation capabilities require additional computational resources. Yet these same environments offer the least resource availability, creating a paradoxical situation where those most in need of reliable systems have the least capacity to support resource-intensive robustness.

**The Questions:** Is this tension truly inevitable, or does it reflect particular technical approaches that could be reconceptualized? What if robustness were reimagined not as resource-intensive checking but through fundamentally different architectural approaches? Can we transcend the apparent trade-off between efficiency and robustness by reconsidering what each means in educational contexts?



**Research Directions:**

Research Question	Investigation Approach	Expected Outcome
What optimization approaches enable multi-stage validation within edge hardware constraints?	Performance profiling across validation stages; resource allocation modeling	Optimization framework balancing validation robustness with performance requirements
How can structured output generation be made more reliable on small edge-deployable models?	Analysis of output formatting patterns; error recovery mechanism testing	Robust framework for extracting structured data from edge model outputs
What synchronization approaches best handle extended offline periods with eventual connectivity?	Data synchronization pattern analysis; conflict resolution modeling	Offline-first synchronization architecture robust to extended disconnection
How can model size and capabilities be optimally balanced for specific educational tasks?	Task requirement analysis; model capability benchmarking across sizes	Decision framework for model selection based on educational requirements
What technical acceleration approaches most effectively improve edge performance without compromising safety?	Acceleration technique benchmarking; safety impact analysis	Edge optimization framework maintaining safety guardrails while improving performance

**Opposing Perspectives:** Consider how different priorities shape approaches to this tension. Resource optimization proponents might argue for the minimal viable system that functions within constraints, even if less robust. Safety advocates might argue that compromising on validation creates

unacceptable risks, regardless of resource costs. What approaches emerge when we refuse to accept either efficiency or robustness as the dominant priority?

Research on educational edge computing frameworks demonstrates approaches for balancing performance requirements with resource constraints in educational applications, showing how edge setups can improve throughput and latency for remote learners (Chen et al. 2022). Additionally, comprehensive reviews of on-device language models highlight strategies for efficient deployment on resource-constrained devices, including architecture optimization and compression techniques (Xu et al. 2024).

### **The Human-AI Dialectic: Automation vs. Agency**

**The Tension:** Perhaps the most profound tension in EdgePrompt lies between automation (using AI to handle educational tasks for efficiency and consistency) and human agency (preserving teacher and student control over the educational process). This tension manifests across multiple dimensions: teacher vs. system control, standardized vs. personalized approaches, algorithmic vs. human judgment, and efficiency vs. relationship.

**The Questions:** What educational activities genuinely benefit from automation, and which require irreducible human judgment? How might we reconceptualize the relationship between human and artificial intelligence not as a zero-sum competition but as a complementary partnership with appropriate domains for each? What happens when we move beyond both techno-optimistic and techno-pessimistic positions to explore nuanced approaches to human-AI collaboration?

### **Research Directions:**

<b>Research Question</b>	<b>Investigation Approach</b>	<b>Expected Outcome</b>
What interface design best reduces cognitive load for non-technical educators in resource-constrained settings?	Cognitive walkthrough analysis; simulation of teacher workflows under time constraints	Interface design patterns optimized for teachers with limited technical experience and time
How can trust-building mechanisms be integrated into the system to foster appropriate teacher confidence?	Analysis of trust development patterns in educational technology; transparency mechanism testing	Trust-building framework with calibrated transparency and appropriate system explanations

Research Question	Investigation Approach	Expected Outcome
What progressive disclosure approach best balances simplicity with control for diverse teacher technical capabilities?	Usability testing simulations; complexity progression modeling across user expertise	Interface progression framework allowing growing technical control as teacher expertise develops
How can workflow integration minimize disruption to established teaching practices?	Analysis of teaching workflow patterns across cultural contexts; integration point modeling	Integration guidelines for embedding AI capabilities within existing educational workflows
What accessibility approaches best address the diverse needs of teachers and students in target regions?	Accessibility requirement analysis across deployment contexts; adaptive interface simulation	Comprehensive accessibility framework addressing regional variation in ability needs

**Opposing Perspectives:** Automation advocates might argue that AI can handle routine tasks more efficiently, freeing humans for higher-value activities. Agency defenders might counter that seemingly routine educational tasks often contain subtle judgment opportunities essential to the learning relationship. What approaches emerge when we refuse to see automation and human agency as necessarily opposed, seeking instead their appropriate balance in specific educational contexts?

Research on teacher agency in the age of generative AI emphasizes the importance of maintaining teachers' power to act, affect matters, make decisions, and take stances in educational settings (Frøsig & Romero 2024). Additionally, work on trust in AI-assisted decision-making has shown that cognitive forcing functions can reduce overreliance on AI by prompting users to think independently before accepting AI recommendations (Buçinca, Malaya & Gajos 2021).

### **The Theory-Practice Dialectic: Research vs. Implementation**

**The Tension:** Throughout our assessment, we've observed tension between theoretical frameworks and practical implementation reality. EdgePrompt's conceptual architecture represents a sophisticated vision for educational AI, yet implementation challenges—from JSON reliability to teacher interface

complexity—reveal the gap between theory and practice. This gap isn't merely an implementation failure but reflects genuine difficulty in translating abstract principles into functioning systems.

**The Questions:** How might we reconceptualize the relationship between research and implementation beyond the traditional "research then implement" model? What if implementation challenges were viewed not as obstacles to theoretical purity but as essential knowledge that should reshape theory itself? Can we develop research methodologies that more effectively bridge theoretical insight and practical reality?

**Research Directions:**

Research Question	Investigation Approach	Expected Outcome
How can research-implementation feedback loops be structured for more rapid iteration in resource-constrained environments?	Analysis of development processes across resource contexts; modeling of feedback mechanisms	Framework for integrating research and implementation in accelerated, context-sensitive cycles
What co-design methodologies best integrate theoretical expertise and practical wisdom in educational AI development?	Comparative analysis of co-design approaches; stakeholder engagement pattern modeling	Co-design methodology balancing theoretical insight with practitioner knowledge
How can implementation friction be used productively to refine theoretical models rather than being treated as mere obstacles?	Case study analysis of implementation challenges; friction-to-insight pattern modeling	Framework for using implementation challenges as theory-building opportunities
What forms of practitioner knowledge remain systematically undervalued in educational AI research, and	Analysis of knowledge types across stakeholder groups; expertise integration modeling	Methodology for identifying and integrating diverse knowledge forms in educational AI

Research Question	Investigation Approach	Expected Outcome
how might they be better integrated?		
How can prototype testing methodologies be adapted for communities with limited prior exposure to AI technologies?	Comparative analysis of testing approaches across technological familiarity; adaptation pattern modeling	Testing methodology appropriate for communities with limited prior AI exposure

**Opposing Perspectives:** Theoretical researchers might emphasize conceptual clarity, architectural elegance, and systematic approaches. Implementation practitioners might prioritize functional systems, user acceptance, and contextual adaptation. What approaches emerge when we refuse to privilege either theoretical understanding or practical functionality, instead seeking their productive integration?

### Integrating the Dialectics: Meta-Questions for EdgePrompt's Future

These dialectical tensions cannot be resolved through a single synthesis but require ongoing, contextual navigation. Several meta-questions emerge for EdgePrompt's future development:

1. **From Binary to Spectrum Thinking:** How might we move beyond binary oppositions (online/offline, secure/accessible, technical/educational) toward spectrum thinking that acknowledges the contextual nature of appropriate balance points?
2. **From Fixed Solutions to Adaptation Frameworks:** Instead of seeking universal answers, how might EdgePrompt develop frameworks that enable contextual adaptation across its tension dimensions?
3. **From External Design to Community Ownership:** How can EdgePrompt transition from a system designed for communities to one owned, governed, and evolved by communities themselves?
4. **From Technical Metrics to Educational Values:** What would development prioritization look like if driven primarily by educational values rather than technical capabilities or implementation ease?
5. **From Problem-Solving to Capacity-Building:** How might EdgePrompt's approach shift from solving educational problems to building community capacity for technological sovereignty?

These meta-questions suggest that EdgePrompt's most important contribution may ultimately lie not in its technical architecture but in its exploration of how AI can be thoughtfully integrated into diverse educational contexts. By explicitly acknowledging and navigating these tensions rather than seeking to

eliminate them, EdgePrompt can contribute to a deeper understanding of technology's appropriate role in human learning—a question that transcends any single implementation while being concretely manifested in each design decision.

The research directions outlined in this chapter offer not merely topics for investigation but opportunities for navigating these fundamental tensions. By approaching these questions with both technical rigor and philosophical depth, EdgePrompt can continue evolving from promising concept toward meaningful educational contribution.

Research on the ethics of AI in education emphasizes the importance of transparency in AI decision-making, safeguarding student autonomy and trust, and adherence to principles of fairness, accountability, and pedagogical soundness (Holmes et al. 2022). Additionally, work on socio-technical education futures examines educational technology within broader socio-technical assemblages, highlighting the importance of involving diverse stakeholders in "technical democracy" (Swist & Gulson 2023).

## Conclusion: Revolution, Hype, or Foundation Stone?

As we return to the question posed in our title—"EdgePrompt: Revolutionizing Education or Just Another EdTech Hype?"—we find ourselves confronting a more nuanced reality than either revolution or hype alone would suggest. Our comprehensive assessment reveals EdgePrompt as neither a fully realized revolution nor merely empty promise, but rather a foundation stone—an initial building block with the potential to support more equitable, teacher-centered AI integration in education if developed with appropriate care, resources, and community engagement.

### EdgePrompt as a Pragmatic Engineering Response to LLM Reality

First and foremost, EdgePrompt represents a pragmatic engineering response to the inherent tensions between AI's potential educational benefits and its significant challenges in resource-constrained environments. Rather than merely theorizing about ideal systems, the project has confronted messy implementation realities—from the JSON reliability challenges of edge LLMs to the efficiency costs of multi-stage validation—and developed concrete technical approaches to address them.

This pragmatic orientation distinguishes EdgePrompt from both revolutionaries who promise transformation without addressing implementation barriers and skeptics who dismiss AI's potential value without exploring practical solutions. The project demonstrates that bringing AI capabilities to resource-constrained educational environments is neither trivially simple nor fundamentally impossible but rather a complex engineering challenge requiring sustained, focused effort.

The research framework's successful validation of the core EdgePrompt approach—structured prompts and multi-stage validation improving safety and adherence compared to baseline approaches—provides evidence that prompt engineering can create meaningful guardrails without model fine-tuning. This validation, while limited to simulation rather than real-world deployment, represents genuine progress toward the goal of safe, offline-capable educational AI.

However, the significant gap between the research framework and the current application implementation highlights that moving from validated approach to deployed solution requires substantial additional engineering effort. The project has demonstrated concept feasibility but not yet developed a classroom-ready implementation, particularly regarding security, full validation implementation, and hardware optimization.

### EdgePrompt as a Socio-Technical Intervention with Moral Weight

Beyond its technical dimensions, EdgePrompt represents a socio-technical intervention with significant moral implications. The project's focus on educational equity, teacher agency, and offline capability reflects a particular vision of how AI should integrate with education—one that prioritizes democratization, local control, and accessibility over centralization, automation, and dependency.

This moral orientation shapes not just the project's goals but its architectural decisions, from the backend-first security model that prevents direct prompt manipulation to the template system that preserves teacher control over content generation. These architectural choices aren't merely technical but embed specific values about appropriate human-AI relationships in educational contexts.

The project's approach to cultural responsiveness—while still primarily aspirational in implementation—similarly reflects a moral commitment to adaptation rather than standardization. This orientation toward pluralism rather than universalism represents an important alternative to educational AI approaches that implicitly impose dominant cultural perspectives through standardized interactions.

However, the gap between these moral commitments and current implementation highlights the challenge of translating values into technical reality. Achieving true educational equity requires more than offline capability alone but comprehensive engagement with language barriers, digital literacy variations, and diverse cultural contexts. Similarly, realizing teacher agency demands more than theoretical control—it requires interfaces accessible to non-technical educators and workflows that genuinely enhance rather than complicate teaching practice.

EdgePrompt thus stands as a morally serious attempt to navigate AI's educational possibilities in a principled way, but one whose implementation has not yet fully realized its moral vision. The project demonstrates that ethical integration of AI in education is possible but not automatic—it requires deliberate design choices and sustained development effort guided by clear values.

## **EdgePrompt as a Work of Transition with Community Ownership and Generational Thinking**

Perhaps most importantly, EdgePrompt represents a work of transition—a project situated at the intersection of AI's rapidly evolving capabilities, education's enduring human purposes, and the practical realities of resource-constrained environments. This transitional nature demands particular attention to context, community ownership, and generational thinking.

The project's focus on specific contexts—particularly Indonesia's 3T regions—represents an important shift from universal technological solutions toward contextually grounded interventions. This approach acknowledges that educational technology's effectiveness depends not just on technical capabilities but on alignment with specific educational realities, cultural contexts, and infrastructure conditions. Future development must deepen this contextual grounding through genuine co-design with educators and communities in target regions.

The vision of EdgePrompt as "public infrastructure" rather than commercial product similarly points toward a transitional model of educational technology development—one that prioritizes community ownership, adaptation rights, and knowledge sovereignty over proprietary control. Realizing this vision requires development of concrete governance mechanisms, sustainability models, and knowledge transfer approaches that enable genuine community ownership beyond rhetorical commitment.



Finally, EdgePrompt implicitly adopts a generational thinking timeframe appropriate to its educational mission. Unlike technologies designed for immediate commercial impact, educational infrastructure must consider developmental timescales—how it shapes not just immediate learning activities but long-term educational trajectories and capabilities. This perspective demands patience, sustainability, and commitment to foundational development rather than merely visible features.

As a work of transition, EdgePrompt should be judged not just by its current capabilities but by the architectural foundation it provides for future development. The project's explicit identification of gaps, tensions, and development priorities demonstrates self-awareness about its transitional status and provides a roadmap for evolution toward its fuller vision.

## **Final Thoughts: The Ongoing Work of Aligning AI with Human Flourishing in Education**

EdgePrompt exists within a broader context of ongoing efforts to align artificial intelligence with human flourishing in education. This alignment is not a one-time achievement but an ongoing process requiring continuous adaptation to evolving technological capabilities, educational needs, and ethical understandings.

The project's focus on teacher control, offline capability, and safety guardrails represents one particular vision of this alignment—one that prioritizes human oversight, equitable access, and ethical boundaries. This vision stands in productive tension with alternative approaches emphasizing automation, cloud integration, or capability maximization. The educational AI landscape benefits from this diversity of approaches, with different models appropriate for different contexts and purposes.

What distinguishes EdgePrompt is not just its particular alignment choices but its transparency about them. Through detailed documentation of architecture, validation approaches, and design principles, the project makes its alignment mechanisms explicit rather than implicit. This transparency enables critical engagement, potential adaptation, and informed debate about appropriate AI integration approaches for various educational contexts.

The ultimate test of EdgePrompt will not be whether it revolutionizes education broadly—no single technology can or should claim such impact—but whether it contributes meaningfully to equitable, effective, and ethical AI integration in the specific contexts it targets. This contribution depends not just on technical development but on ongoing engagement with educators, students, and communities to ensure the technology genuinely serves their educational needs rather than imposing external priorities.

In this light, EdgePrompt is best understood neither as revolution nor hype but as a serious contribution to the ongoing work of aligning artificial intelligence with human flourishing in education—work that requires technical innovation, ethical reflection, and community engagement across generations. The project has laid a foundation; building upon it responsibly requires continued commitment to both technological excellence and educational purpose.

## References

### AI Guardrails and Prompt Engineering

**Dong Y, Mu R, Jin G, Qi Y, Hu J, Zhao X, Meng J, Ruan W and Huang X** (2024) 'Building Guardrails for Large Language Models', [arXiv:2402.01822](https://arxiv.org/abs/2402.01822)

**Ganguli D, Lovitt L, Kernion J, Askell A, Bai Y, Kadavath S, Mann B, Perez E, Schiefer N, Ndousse K et al.** (2022) 'Red Teaming Language Models to Reduce Harms: Methods, Scaling Behaviors, and Lessons Learned', [arXiv:2209.07858](https://arxiv.org/abs/2209.07858)

**Gil De Zúñiga H, Goyanes M and Durotoye T** (2024) 'A Scholarly Definition of Artificial Intelligence (AI): Advancing AI as a Conceptual Framework in Communication Research', *Political Communication*, 41(2):317-334, [doi:10.1080/10584609.2023.2290497](https://doi.org/10.1080/10584609.2023.2290497).

**Gupta A, Akiri C, Aryal K, Parker E and Praharaj L** (2023) 'From ChatGPT to ThreatGPT: Impact of Generative AI in Cybersecurity and Privacy', *IEEE Access*, 11:80218-80245, [doi:10.1109/ACCESS.2023.3300381](https://doi.org/10.1109/ACCESS.2023.3300381).

**Frøsig TB and Romero M** (2024) 'Teacher agency in the age of generative AI: towards a framework of hybrid intelligence for learning design', [arXiv:2407.06655](https://arxiv.org/abs/2407.06655).

**Hacker P, Engel A and Mauer M** (2023) 'Regulating ChatGPT and other Large Generative AI Models', *2023 ACM Conference on Fairness, Accountability, and Transparency*, ACM, pp. 1112-1123, [doi:10.1145/3593013.3594067](https://doi.org/10.1145/3593013.3594067).

**Ouyang L, Wu J, Jiang X, Almeida D, Wainwright CL, Mishkin P, Zhang C, Agarwal S, Slama K, Ray A et al.** (2022) 'Training language models to follow instructions with human feedback', [arXiv:2203.02155](https://arxiv.org/abs/2203.02155)

**Sahoo P, Singh S, Kala CSR, Mukherjee A, Shrimal M, Hota C and Raychoudhury V** (2025) 'A Systematic Survey of Prompt Engineering in Large Language Models: Techniques and Applications', [arXiv:2402.07927](https://arxiv.org/abs/2402.07927)

**Wei J, Wang X, Schuurmans D, Bosma M, Ichter B, Xia F, Chi E, Le Q and Zhou D** (2022) 'Chain-of-Thought Prompting Elicits Reasoning in Large Language Models', [arXiv:2201.11903](https://arxiv.org/abs/2201.11903)

**White J, Fu Q, Hays S, Sandborn M, Olea C, Gilbert H, Elnashar A, Spencer-Smith J and Schmidt DC** (2023) 'A Prompt Pattern Catalog to Enhance Prompt Engineering with ChatGPT', [arXiv:2302.11382](https://arxiv.org/abs/2302.11382)

**Xia B, Lu Q, Zhu L and Xing Z** (2024) 'An AI System Evaluation Framework for Advancing AI Safety: Terminology, Taxonomy, Lifecycle Mapping', [arXiv:2404.05388](https://arxiv.org/abs/2404.05388)

**Yao S, Yu D, Zhao J, Shafran I, Griffiths TL, Cao Y and Narasimhan K** (2023) 'Tree of Thoughts: Deliberate Problem Solving with Large Language Models', [arXiv:2305.10601](https://arxiv.org/abs/2305.10601)

**Zheng C, Yin F, Zhou H, Meng F, Zhou J, Chang KW, Huang M and Peng N** (2024) 'On Prompt-Driven Safeguarding for Large Language Models', [arXiv:2401.18018](https://arxiv.org/abs/2401.18018).

**Zhang X** (2025) 'Constitution or Collapse? Exploring Constitutional AI with Llama 3-8B', [arXiv:2504.04918](https://arxiv.org/abs/2504.04918).

## **Educational AI Applications and Question Generation**

**Bastani H, Bastani O, Sungu A, Ge H, Kabakcı Ö and Mariman R** (2024) 'Generative AI Can Harm Learning', *The Wharton School Research Paper*, [doi:10.2139/ssrn.4895486](https://doi.org/10.2139/ssrn.4895486).

**Bhowmick AK, Jagmohan A, Vempaty A, Dey P, Hall L, Hartman J, Kokku R and Maheshwari H** (2023) 'Automating Question Generation From Educational Text', in Bramer M and Stahl F (eds) *Artificial Intelligence XL*, Springer Nature Switzerland, pp. 437–450, [doi:10.1007/978-3-031-47994-6\\_38](https://doi.org/10.1007/978-3-031-47994-6_38)

**Bubeck S, Chandrasekaran V, Eldan R, Gehrke J, Horvitz E, Kamar E, Lee P, Lee YT, Li Y, Lundberg S et al.** (2023) 'Sparks of Artificial General Intelligence: Early experiments with GPT-4', [arXiv:2303.12712](https://arxiv.org/abs/2303.12712).

**Furze L, Perkins M, Roe J and MacVaugh J** (2024) 'The AI Assessment Scale (AIAS) in action: A pilot implementation of GenAI-supported assessment', *Australasian Journal of Educational Technology*, 40(2), [doi:10.14742/ajet.9434](https://doi.org/10.14742/ajet.9434)

**Fui-Hoon Nah F, Zheng R, Cai J, Siau K and Chen L** (2023) 'Generative AI and ChatGPT: Applications, challenges, and AI-human collaboration', *Journal of Information Technology Case and Application Research*, 25(3):277-304, [doi:10.1080/15228053.2023.2233814](https://doi.org/10.1080/15228053.2023.2233814).

**Grubaugh S, Levitt G and Deever D** (2023) 'Harnessing AI to Power Constructivist Learning: An Evolution in Educational Methodologies', *Journal of Effective Teaching Methods*, 1(3), [doi:10.59652/jetm.v1i3.43](https://doi.org/10.59652/jetm.v1i3.43)

**Hang CN, Wei Tan C and Yu PD** (2024) 'MCQGen: A Large Language Model-Driven MCQ Generator for Personalized Learning', *IEEE Access*, 12:102261–102273, [doi:10.1109/ACCESS.2024.3420709](https://doi.org/10.1109/ACCESS.2024.3420709)

**Magana Vsevolodovna RI and Monti M** (2025) 'Enhancing Large Language Models through Neuro-Symbolic Integration and Ontological Reasoning', [arXiv:2504.07640](https://arxiv.org/abs/2504.07640).

**Niknazar M, Haley PV, Ramanan L, Truong ST, Shrinivasan Y, Bhowmick AK, Dey P, Jagmohan A, Maheshwari H, Ponoth S et al.** (2024) 'Building a Domain-specific Guardrail Model in Production', [arXiv:2408.01452](https://arxiv.org/abs/2408.01452)

**Novoa-Echaurren Á** (2024) 'Teacher Agency in the Pedagogical Uses of ICT: A Holistic Perspective Emanating from Reflexive Practice', *Education Sciences (MDPI)*, 14(3):254, [doi:10.3390/educsci14030254](https://doi.org/10.3390/educsci14030254)

**Peres R, Schreier M, Schweidel D and Sorescu A** (2023) 'On ChatGPT and beyond: How generative artificial intelligence may affect research, teaching, and practice', *International Journal of Research in Marketing*, 40(2):269-275, [doi:10.1016/j.ijresmar.2023.03.001](https://doi.org/10.1016/j.ijresmar.2023.03.001).

**Riza LS, Firdaus Y, Sukanto RA, Wahyudin and Abu Samah KAF** (2023) 'Automatic generation of short-answer questions in reading comprehension using NLP and KNN', *Multimedia Tools and Applications*, 82(27):41913–41940, [doi:10.1007/s11042-023-15191-6](https://doi.org/10.1007/s11042-023-15191-6)

**Scaria N, Dharani Chenna S and Subramani D** (2024) 'Automated Educational Question Generation at Different Bloom's Skill Levels Using Large Language Models: Strategies and Evaluation', in Olney AM, Chounta IA, Liu Z, Santos OC and Bittencourt IG (eds) *Artificial Intelligence in Education*, Springer Nature Switzerland, pp. 165–179, [doi:10.1007/978-3-031-64299-9\\_12](https://doi.org/10.1007/978-3-031-64299-9_12)

**Tan M and Subramoniam H** (2024) 'More than Model Documentation: Uncovering Teachers' Bespoke Information Needs for Informed Classroom Integration of ChatGPT', *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, Article No. 269, pp. 1-19, [doi:10.1145/3613904.3642592](https://doi.org/10.1145/3613904.3642592).

**Wang S, Xu T, Li H, Zhang C, Liang J, Tang J, Yu PS and Wen Q** (2024) 'Large Language Models for Education: A Survey and Outlook', [arXiv:2403.18105](https://arxiv.org/abs/2403.18105).

**Wang X, Fan S, Houghton J and Wang L** (2022) 'Towards Process-Oriented, Modular, and Versatile Question Generation that Meets Educational Needs', *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Association for Computational Linguistics, Seattle, USA, pp. 291–302, [doi:10.18653/v1/2022.naacl-main.22](https://doi.org/10.18653/v1/2022.naacl-main.22)

## Resource-Constrained Environments and Indonesian Context

**Adji LK** (2024) 'Indonesia's internet penetration hits 79.5 percent, trend continues', *Antara News*, <https://en.antaranews.com/news/304593/indonesias-internet-penetration-hits-795-percent-trend-continues>

**Arifin AA and Lennerfors TT** (2022) 'Ethical aspects of voice assistants: a critical discourse analysis of Indonesian media texts', *Journal of Information, Communication and Ethics in Society*, 20(1):18-36, [doi:10.1108/JICES-12-2020-0118](https://doi.org/10.1108/JICES-12-2020-0118).

**Chen Q, Zhao W, Wang Z, Luo T, Jian X and Zheng L** (2022) 'Educational 5G Edge Computing: Framework and Experimental Study', *Electronics (MDPI)*, 11(17):2727, [doi:10.3390/electronics11172727](https://doi.org/10.3390/electronics11172727)

**Kementerian Desa** (2025) 'Official Website of the Ministry of Villages', <https://www.kemendesa.go.id>

**Singh A, Kanaujia A, Singh VK and Vinuesa R** (2024) 'Artificial intelligence for Sustainable Development Goals: Bibliometric patterns and concept evolution trajectories', *Sustainable Development*, 32(1):724–754, [doi:10.1002/sd.2706](https://doi.org/10.1002/sd.2706)

**UNESCO** (2020) 'Startling digital divides in distance learning emerge', UNESCO Press Release, 21 April 2020, <https://en.unesco.org/news/startling-digital-divides-distance-learning-emerge>

**Widodo J** (2020) 'Peraturan Presiden (PERPRES) Nomor 63 Tahun 2020 Penetapan Daerah Tertinggal Tahun 2020-2024', *Database Peraturan | JDIH BPK*, <https://peraturan.bpk.go.id/Details/136563/perpres-no-63-tahun-2020>

## AI Alignment, Educational Philosophy, Ethics, and Trust

**Baker RS and Hawn A** (2022) 'Algorithmic Bias in Education', *International Journal of Artificial Intelligence in Education*, 32(4), [doi:10.1007/s40593-021-00285-9](https://doi.org/10.1007/s40593-021-00285-9)

**Bell G** (2022) 'The New Cybernetics: Systems Thinking for the 21st Century', ANU School of Cybernetics, [URL: https://cybernetics.anu.edu.au/news/2022/03/28/the-new-cybernetics-systems-thinking-for-21st-century/](https://cybernetics.anu.edu.au/news/2022/03/28/the-new-cybernetics-systems-thinking-for-21st-century/)

**Buçinca Z, Malaya MB and Gajos KZ** (2021) 'To Trust or to Think: Cognitive Forcing Functions Can Reduce Overreliance on AI in AI-assisted Decision-making', *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW1):1–21, [doi:10.1145/3449287](https://doi.org/10.1145/3449287)

**Dwivedi YK, Kshetri N, Hughes L, Slade E, Jeyaraj A, Kar AK, Baabdullah AM, Koohang A, Raghavan V, Ahuja M et al.** (2023) 'Opinion Paper: "So what if ChatGPT wrote it?" Multidisciplinary perspectives on opportunities, challenges and implications of generative conversational AI for research, practice and policy', *International Journal of Information Management*, 71:102642, [doi:10.1016/j.ijinfomgt.2023.102642](https://doi.org/10.1016/j.ijinfomgt.2023.102642).

**Elgin CZ** (2013) 'Epistemic Agency', *Theory and Research in Education*, 11(2):135-152, [doi:10.1177/1477878513485171](https://doi.org/10.1177/1477878513485171)

**Gabriel I** (2020) 'Artificial Intelligence, Values, and Alignment', *Minds and Machines*, 30(3):411–437, [doi:10.1007/s11023-020-09539-2](https://doi.org/10.1007/s11023-020-09539-2)

**Gay G** (2018) 'Culturally Responsive Teaching: Theory, Research, and Practice', Teachers College Press, [URL: https://www.tcpres.com/culturally-responsive-teaching-9780807758762](https://www.tcpres.com/culturally-responsive-teaching-9780807758762)

**Hendrycks D, Mazeika M and Woodside T** (2023) 'An Overview of Catastrophic AI Risks', [arXiv:2306.12001](https://arxiv.org/abs/2306.12001)

**Holmes W, Porayska-Pomsta K, Holstein K, Sutherland E, Baker T, Buckingham Shum S, Santos OC, Rodrigo MMT, Cukurova M, Bittencourt II and Koedinger KR** (2022) 'Ethics of AI in Education: Towards a Community-Wide Framework', *International Journal of Artificial Intelligence in Education*, 32(4), [doi:10.1007/s40593-021-00295-8](https://doi.org/10.1007/s40593-021-00295-8)

**Jacovi A, Marasović A, Miller T and Goldberg Y** (2021) 'Formalizing Trust in Artificial Intelligence: Prerequisites, Causes and Goals of Human Trust in AI', *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, ACM, Virtual Event Canada, pp. 624–635, [doi:10.1145/3442188.3445923](https://doi.org/10.1145/3442188.3445923)

**Kumar D and Suthar N** (2024) 'Ethical and legal challenges of AI in marketing: an exploration of solutions', *Journal of Information, Communication and Ethics in Society*, 22(1):124-144, [doi:10.1108/JICES-05-2023-0068](https://doi.org/10.1108/JICES-05-2023-0068).

**Mehrabi N, Morstatter F, Saxena N, Lerman K and Galstyan A** (2022) 'A Survey on Bias and Fairness in Machine Learning', *ACM Computing Surveys*, 54(6):1–35, [doi:10.1145/3457607](https://doi.org/10.1145/3457607)

**Menard P and Bott GJ** (2024) 'Artificial intelligence misuse and concern for information privacy: New construct validation and future directions', *Information Systems Journal*, [doi:10.1111/isj.12544](https://doi.org/10.1111/isj.12544).

**Samuel Y, Mvogo A, Akinyemi T, Segun F, Smith E, Ajilore I and Okolo CT** (2023) 'Cultivation of Human Centered Artificial Intelligence: Culturally Adaptive Thinking in Education (CATE) for AI', *Frontiers in Artificial Intelligence*, 6:1198180, [doi:10.3389/frai.2023.1198180](https://doi.org/10.3389/frai.2023.1198180)

**Slade S and Prinsloo P** (2013) 'Learning Analytics: Ethical Issues and Dilemmas', *American Behavioral Scientist*, 57(10), [doi:10.1177/0002764213479366](https://doi.org/10.1177/0002764213479366)

**Swist T and Gulson KN** (2023) 'Instituting Socio-Technical Education Futures: Encounters with/through Technical Democracy, Data Justice, and Imaginaries', *Learning, Media and Technology*, 48(2), [doi:10.1080/17439884.2023.2205225](https://doi.org/10.1080/17439884.2023.2205225)

**Volkman R and Gabriels K** (2023) 'AI Moral Enhancement: Upgrading the Socio-Technical System of Moral Engagement', *Science and Engineering Ethics*, 29(2):11, [doi:10.1007/s11948-023-00428-2](https://doi.org/10.1007/s11948-023-00428-2)

**Vygotsky LS** (1978) 'Mind in Society: The Development of Higher Psychological Processes', Harvard University Press, [URL: https://www.hup.harvard.edu/catalog.php?isbn=9780674576292](https://www.hup.harvard.edu/catalog.php?isbn=9780674576292)

## Neural-Symbolic Approaches and Cognitive Models

**Binz M and Schulz E** (2023) 'Turning large language models into cognitive models', [arXiv:2306.03917](https://arxiv.org/abs/2306.03917)

**Colelough BC and Regli W** (2025) 'Neuro-Symbolic AI in 2024: A Systematic Review', [arXiv:2501.05435](https://arxiv.org/abs/2501.05435)

**Hooshyar D, Azevedo R and Yang Y** (2024) 'Augmenting Deep Neural Networks with Symbolic Educational Knowledge: Towards Trustworthy and Interpretable AI for Education', *Machine Learning and Knowledge Extraction (MDPI)*, 6(1), [doi:10.3390/make6010028](https://doi.org/10.3390/make6010028)

**Kaswan KS, Dhatteval JS, Malik K and Baliyan A** (2023) 'Generative AI: A Review on Models and Applications', *2023 International Conference on Communication, Security and Artificial Intelligence (ICCSAI)*, IEEE, pp. 699-704, [doi:10.1109/ICCSAI59793.2023.10421601](https://doi.org/10.1109/ICCSAI59793.2023.10421601).

**Lewis P, Perez E, Piktus A, Petroni F, Karpukhin V, Goyal N, Küttler H, Lewis M, Yih WT, Rocktäschel T, Riedel S and Kiela D** (2020) 'Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks', *Advances in Neural Information Processing Systems*, 33:9459-9474.

**Mazzaglia P, Verbelen T, Çatal O and Dhoedt B** (2022) 'The Free Energy Principle for Perception and Action: A Deep Learning Perspective', *Entropy*, 24(2):301, [doi:10.3390/e24020301](https://doi.org/10.3390/e24020301)

**Pezzulo G, Parr T, Cisek P, Clark A and Friston K** (2024) 'Generating meaning: active inference and the scope and limits of passive AI', *Trends in Cognitive Sciences*, 28(2):97–112, [doi:10.1016/j.tics.2023.10.002](https://doi.org/10.1016/j.tics.2023.10.002)

**Rozenblit L and Keil F** (2002) 'The misunderstood limits of folk science: an illusion of explanatory depth', *Cognitive Science*, 26(5):521–562, [doi:10.1207/s15516709cog2605\\_1](https://doi.org/10.1207/s15516709cog2605_1)

**Su J, Jiang C, Jin X, Qiao Y, Xiao T, Ma H, Wei R, Jing Z, Xu J and Lin J** (2024) 'Large Language Models for Forecasting and Anomaly Detection: A Systematic Literature Review', [arXiv:2402.10350](https://arxiv.org/abs/2402.10350).



## Implementation Approaches and Retrieval-Augmentation

**Lewis P, Perez E, Piktus A, Petroni F, Karpukhin V, Goyal N, Küttler H, Lewis M, Yih WT, Rocktäschel T, Riedel S and Kiela D** (2020) 'Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks', *Advances in Neural Information Processing Systems*, 33:9459-9474, <https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html>

**Poesia G, Polozov O, Le V, Tiwari A, Soares G, Meek C and Gulwani S** (2022) 'Synchromesh: Reliable code generation from pre-trained language models', [arXiv:2201.11227](https://arxiv.org/abs/2201.11227)

**Wang H, Huang W, Deng Y, Wang R, Wang Z, Wang Y, Mi F, Pan JZ and Wong KF** (2024) 'UniMS-RAG: A Unified Multi-source Retrieval-Augmented Generation for Personalized Dialogue Systems', [arXiv:2401.13256](https://arxiv.org/abs/2401.13256)

**Xu J, Li Z, Chen W, Wang Q, Gao X, Cai Q and Ling Z** (2024) 'On-Device Language Models: A Comprehensive Review', [arXiv:2409.00088](https://arxiv.org/abs/2409.00088).

**Yazaki M, Maki S, Furuya T, Inoue K, Nagai K, Nagashima Y, Maruyama J, Toki Y, Kitagawa K, Iwata S et al.** (2024) 'Emergency Patient Triage Improvement through a Retrieval-Augmented Generation Enhanced Large-Scale Language Model', *Prehospital Emergency Care*, 1–13, [doi:10.1080/10903127.2024.2374400](https://doi.org/10.1080/10903127.2024.2374400)