# Systemic Limits of AI Code Assistants in Complex Projects

**Abstract**

AI coding assistants promise faster delivery and less boilerplate. In real projects, however, the limits show up as soon as the work goes beyond trivial edits. This analysis sets out the terms, the mechanics of how assistants see code, and the consequences that follow—without anecdotes, hype, or padding.

## 1. Fundamental Concepts

To analyze the behavior of AI assistants, we must first define the operational mechanics.

- **Session:** A continuous interaction in which an assistant receives files and messages and produces answers or code.
- **Context Window:** The maximum text the model can hold at one time; while material fits, nothing is removed. Once additions would exceed that limit, earlier details cannot remain fully available in active memory.
- **Complete-Bundle:** An access pattern where the assistant is given the entire code bundle at the outset and thus starts with a unified view of the project.
- **Vibe Coding:** An access pattern where material arrives incrementally; the assistant constructs its picture of the system as the conversation touches different areas. It never truly holds the entire system in memory at once.
- **Retrieval/Search:** The auxiliary mechanisms an assistant uses to locate and load code—by filename, token, or textual pattern—into the session.

## 2. The Core Conflict: Probabilistic Generation vs. Deterministic Systems

The central limitation of AI assistants is that they generate by **probabilistic continuation rather than by proof.** They reproduce local naming and style, but correctness depends on systemic constraints that may not be represented in the active context. This leads to characteristic failures.

- **First Practical Observations:** Code that looks reasonable does not compile. Suggested changes reference configuration parameters that do not exist or call functions the project never defined. Because the presentation is uniformly confident, superficial plausibility is easy to mistake for reliability, and every suggestion requires careful review even when it appears routine.
- **Typical Misstep:** A common failure is an execution-context mix-up. A function that is correct on the server is recommended for a client-side form because the assistant recognized the label "validation" and missed where that function is permitted to run.

This demonstrates a core principle: **Availability of text is not equivalent to understanding of applicability.**

## 3. The Context Window: An Amplifier of Systemic Weakness

The context window is a hard size limit that exacerbates the core probabilistic weakness. It is critical to understand that even in the ideal scenario of a complete-bundle that fits the context window, the system remains fundamentally probabilistic; having all the information is not the same as applying it correctly. This represents the *best-case* reliability. All other variations where the full context is not available—due to truncation, summarization, or on-demand loading—can only become worse than that.

The practical effects of this limitation manifest differently across tool architectures:

- **Standard in-editor assistants**, for example, are classic examples of a highly limited, transient context. They operate based on open files and effectively "forget" a file's content moments after it is closed. Their suggestions are a series of isolated text hits rather than architecturally-aware contributions.
- More advanced, **retrieval-augmented assistants** attempt to solve this by using embeddings for whole-codebase awareness. This shifts the problem rather than solving it. Embeddings are a form of lossy

compression; they are excellent for finding semantically *similar* code but do not represent the *exact, logical* architecture. The Documentation Paradox is particularly acute here: detailed prose in an `ARCHITECTURE.md` file is converted into a vector that loses its specific, nuanced constraints. The knowledge degradation is more subtle, but it persists.

At this point, the assistant pivots toward project-wide pattern search—internal lookups akin to a powerful `grep`. It can still surface definitions and strings quickly, but the relationships between files fade. Reasoning drifts from a coherent architectural view to working off isolated text hits, which in turn encourages **repeated searching of ground that was already covered.**

This leads to the **Documentation Paradox**: Good documentation helps initially, because the assistant can quote structure and rationale. As the session grows, detailed passages are the first to be displaced from active memory or condensed into summaries. The assistant's tone remains uniformly confident while critical specifics have already slipped out of view. Attempts to recover richness by re-reading the same documents consume capacity without restoring the original architectural continuity.

### 4. Practical Consequences and Risks

This combination of probabilistic reasoning and bounded memory creates tangible risks and workflow inefficiencies.

### 4.1 Security Implications

Probabilistic generation introduces security risk. Because every snippet is presented with steady confidence, subtle flaws are easy to miss.

- Assistants can suggest unparameterized queries that invite **injection**.
- They can overlook input validation that prevents **cross-site scripting**.
- They may assemble fragile **authentication flows**.
- Training data can also encode **insecure patterns**.

As a result, AI-generated changes should be treated like contributions from an unknown external source: subject to rigorous code review and backed by deterministic scanning with tools such as CodeQL, SonarQube, or Veracode.

### 4.2 Planning Under Uncertainty

Most assistants do not disclose remaining capacity. Teams cannot tell when a session is close to the limit, and long, cross-file changes become difficult to stage. Hidden token ceilings can end a conversation abruptly, introducing operational risk and a practical form of lock-in: a central tool in the workflow may stop at an unspecified boundary, regardless of project urgency.

### 4.2a Paying for... what exactly?

This lock-in is fundamentally stronger and more severe than with traditional tools. An IDE is typically licensed by version or features, granting unlimited usage once paid for. AI assistants, however, are metered by token consumption. This introduces not just a conversational token ceiling, but a "subscription ceiling"—a hard limit on daily or monthly usage. Hitting this limit can disable a core development tool entirely, creating an unprecedented operational dependency. A team's productivity is no longer determined by its skill, but by its consumption meter, leaving it completely at the mercy of the provider's billing model.

### 4.3 Efficiency in Day-to-Day Work

The promise of speed often gives way to coordination overhead. This is made explicit in **assistants with manual context management**, where the developer is responsible for constantly curating the set of files the tool can see. The developer effectively becomes the memory manager for the AI. In all cases, developers re-supply files the assistant has aged out of memory, restate rules the assistant has partially forgotten, and verify output that looked correct but fails basic checks. **The cumulative cost—clarifying patterns, correcting assumptions, and re-running reviews—frequently approaches or exceeds the effort required to implement the change directly.**

**4.4 The Workflow Mismatch**

Modern development is iterative and stateful. Teams make cross-file edits, run tests and builds continuously, and debug against live state. Assistants operate on session snapshots with limited persistence across conversations, modest awareness of file-system events, and weak integration with test and build pipelines. This mismatch makes long-running refactors and architecture-level changes brittle, because the assistant cannot reliably carry the evolving state forward.

**5. Strategic Application and Boundaries**

Effective use requires deliberately keeping the tools within predictable lanes.

**5.1 What Today's Assistants Lack**

Dependable work on large systems requires **persistent memory across sessions, retrieval that preserves architectural relationships** rather than only matching tokens, **transparent indicators of remaining headroom**, and proper **integration with development tooling** so suggestions can be grounded in executable state. Without these, assistants remain closer to fast pattern matchers than to colleagues who hold a system steadily in mind.

**5.2 Appropriate Use Cases**

- Generating conventional scaffolding.
- Translating or explaining unfamiliar code.
- Exploring new libraries.
- Supporting greenfield projects that follow clear, standard patterns.

**5.3 High-Risk Boundaries**

- Established codebases with domain-specific conventions.
- Cross-component refactors that need sustained context.
- Architectural decisions with cross-cutting constraints.
- Performance work that demands system-wide insight.

**5.4 Recognizing When to Stop**

Certain signals indicate that continuing will cost more than it saves: references to undefined parameters or functions, code that does not compile, stylistic alignment that conflicts with domain rules, repeated re-searching of material already seen, and review time rising to match implementation time. At that point, pausing the assistant, narrowing scope, or proceeding manually is usually the rational choice.

**6. Architecting a More Reliable Future: Paths Beyond Probabilism**

The limitations of current tools are not permanent. They stem from a reliance on a single, probabilistic paradigm. A truly robust assistant would require a shift towards systems grounded in deterministic logic. This involves two distinct layers of understanding, which can then be combined into a powerful hybrid system.

**6.1 The Deterministic Foundation: A Layered Approach**

This path aims to replace statistical plausibility with verifiable proof by building upon two layers of deterministic analysis.

**6.1.1 Layer 1: Structural Correctness (AST & Language Server)** The foundation for deterministic analysis already exists in every modern IDE: the **Language Server**. It continuously parses code into an **Abstract Syntax Tree (AST)**, creating a perfect, structured "ground truth" model of the entire codebase. This is a context graph that understands:

- The exact definition of every function and variable.
- The complete call hierarchy (who calls whom).
- The full type and inheritance structure.

An assistant using this layer would have **deterministic structural knowledge**. It could never hallucinate a function that doesn't exist, because it would be checking against the definitive "map" of the AST, not guessing based on text

patterns. However, this layer only understands the *structure* of the code, not its intended *meaning* or architectural rules. It will happily let you write syntactically perfect but architecturally flawed code.

**6.1.2 Layer 2: Logical Correctness (Automated Reasoning)** This is the next layer of analysis, built **on top of** the structural model provided by the AST. Automated Reasoning (AR) applies **formal logic** to this structured representation to prove higher-level properties. The codebase is treated as a formal system with:

- **Axioms:** Facts derived from the AST (e.g., "Class `UserService` is in the 'Service' layer").
- **Rules:** User-defined architectural constraints (e.g., "A function in the 'Service' layer must not call a function in the 'Database' layer").

An AR engine can then formally **prove** whether a piece of code adheres to these rules. It moves beyond "is this code syntactically valid?" to "is this code **logically valid** according to the system's architecture?"

### 6.2 The Hybrid System: Synthesizing All Layers

This is the most promising and practical path forward. It doesn't seek to replace the LLM, but to constrain it by combining its generative strengths with the deterministic validation of both layers.

The workflow would be a **"Generate-then-Validate"** loop:

1. **The LLM as the "Creative Proposer":** The developer interacts with the LLM in natural language. The LLM generates a plausible code suggestion.
2. **The Deterministic Validator:** Before the suggestion is shown, it undergoes a two-stage internal check:
   - **Layer 1 Check:** Does the generated code parse and fit into the existing **AST**? Do all types match and do all functions exist? If not, reject.
   - **Layer 2 Check:** If the structure is valid, does this change violate any of the formal architectural rules defined in the **AR engine**? If so, reject.
3. **The Result:** Only suggestions that pass both the structural and logical checks are presented to the developer.

### 6.3 Why These Approaches Are Not Yet Standard

If this hybrid system is so superior, why do current tools still rely on unreliable text retrieval?

1. **Latency:** The "Validate" step, especially the AR part, is computationally expensive. Today's tools are optimized for the "magic" of sub-second text completion. A reliable but slower suggestion is often perceived as a worse user experience.
2. **The "Messy State" Problem:** This is the biggest hurdle. An AST requires syntactically correct code. Developers, however, spend most of their time writing code that is temporarily broken. A robust validator must be incredibly sophisticated to handle these partial, invalid states, which is a massive engineering challenge.
3. **Formalization Overhead:** Implementing Automated Reasoning (Layer 2) requires developers to formally define their architectural rules in a machine-readable language. Most projects do not have this level of formal specification.

The current approach is a pragmatic compromise: it sacrifices reliability for speed and the ability to function in a messy, informal environment.

### 7. Conclusion

Whether an assistant begins with a **Complete-Bundle** snapshot or constructs its view through **Vibe Coding**, two fundamentals govern outcomes: **bounded context** and **probabilistic generation**. The first constrains how much detail can be carried forward; the second makes superficially correct code easy to produce and easy to trust. Complete-Bundle offers the best start and still cannot guarantee correct application of project-specific rules. Vibe Coding scales to large codebases but tends toward project-wide pattern search once the window is pressured, trading architectural coherence for fast text matches.

Used deliberately—with static analysis, rigorous review, and realistic scope—assistants remain useful for exploration, explanation, and routine code. Expecting them to handle domain logic, sustained multi-file changes, or security-critical work reliably is where the vibes fade. A future grounded in the hybrid architecture of probabilistic generation checked by deterministic validation may finally overcome these limits, but today's tools require careful, critical application.

**Further samples**

1. Hard Coupling Without Thinking

- Error: Hard-coded TableRegistry into the function
- Junior Behavior: Just wrote import without considering testability
- Correct: Dependency injection from the start

2. Overengineering Error Handling

- Error: Complex error arrays instead of simple throw
- Junior Behavior: "More code = better" mentality
- Correct: KISS - just throw on problems

3. Inconsistency and Opinion Changes

- Error: Yesterday: "TableRegistry is perfect!" → Today: "TableRegistry is bad!"
- Junior Behavior: No clear architectural vision
- Correct: Stick to decisions

4. Misusing TypeScript Types

- Error: typeof AliasedTableRegistry instead of interface
- Junior Behavior: Using concrete types for flexibility
- Correct: Use interfaces for abstraction

5. Tests as Afterthought

- Error: Code first, then "oh shit, how do I test this?"
- Junior Behavior: Not thinking about testability from the start
- Correct: Test-driven design

6. Sloppy Implementation

- Error: Forgot test parameters, incomplete refactors
- Junior Behavior: Not following through on changes
- Correct: Systematically update all places

Argument: "But it works"

---

## 8. The Compiler Feedback Gap: A Case Study in Systematic Inefficiency

**Abstract:** A straightforward refactoring task (replacing a component prop across 8 files) took 10x longer with AI assistance than manual implementation. This reveals three critical gaps in current AI-assisted workflows.

### 8.1 The Three Gaps

### 8.1.1 The Implementation Verification Gap

AI tools report edit operations as "successful" without post-execution validation. In practice:

- Changes claimed as complete were not written to files
- Changes written contained incorrect logic despite correct display
- User must manually verify every claimed change

**Core Issue:** Edit operation returns "success" → File may remain unchanged → Discovery only when subsequent steps fail.

**Impact:** Step-by-step verification by user cannot compensate for unreliable implementation layer. Trust erosion creates massive overhead.

### 8.1.2 The Context Fragmentation Anti-Pattern

Selective file reads (offset/limit parameters) intended to reduce token usage paradoxically increase:

- String-matching error rate: +400%
- Total token consumption: +200% (due to re-reads)
- Edit attempts per file: 2-4x

**Mechanics:**

```
Read lines 50-80 only
Attempt edit of string spanning lines 45-52
→ FAILS (missing context)
Re-read lines 40-85
Retry edit
→ Success on 2nd or 3rd attempt
```

**Principle:** Context completeness reduces errors more than it costs in tokens. Full file reads are more efficient than selective reads for files <2000 lines.

### 8.1.3 The Compiler Integration Gap

**Manual Developer Workflow:**

1. Make breaking change
2. Run: `npm run check`
3. Compiler lists ALL affected files with errors
4. Fix systematically in batch
5. Verify → Done (15 minutes)

**AI-Assisted Workflow (Observed):**

1. Make breaking change
2. [No compiler check]
3. Discover File A needs update (manually)
4. Edit File A
5. [No compiler check]
6. Discover File B needs update (user reports)
7. Edit File B
8. Repeat until all 8 files discovered one-by-one
9. Done (90 minutes)

**Efficiency Loss:** 6x from incremental discovery vs. systematic resolution.

**Core Principle:** The compiler is not an adversary - it's the most reliable guide for systematic changes. Development without compiler feedback is development without instruments.

### 8.2 Secondary Manifestations

### 8.2.1 Wrong Logical Order in Composition

Incorrect precedence in composing query parameters (SQL ORDER BY):

- Default sort added BEFORE user sort (wrong)
- User sort never takes effect (overridden by default)
- Affects 3 API functions

**Root Cause:** Misunderstanding of precedence semantics. No compiler/integration test to catch logic error.

### 8.2.2 Incomplete Parameter Handling

Function signature claims to accept parameter, implementation ignores it:

```
// Signature
async loadData(where?: Condition, orderBy?: Sort): Promise<T[]>

// Implementation
const request = { where: { /* hardcoded only */ } }
// where parameter silently ignored
```

**Pattern:** Copy-paste from reference without understanding - signature copied but not implementation logic.

### 8.2.3 Pattern Inconsistency

Failed to recognize and follow established pattern:

- Codebase had `ApiLoadFunc<T>` type definition
- Needed equivalent `SortFunc<T>` type
- Instead: inline types in 8 files → inconsistency → required second pass

**Root Cause:** Not following established patterns + no compiler feedback to identify inconsistency early.

### 8.3 The Incremental Discovery Anti-Pattern
**Observed Pattern:**

```
Change component X
→ Discover file A uses X → fix
→ Discover file B uses X → fix
→ Discover file C uses X → fix
[continue discovering...]
```

**Correct Pattern:**

```
grep -r "X" → ALL 8 affected files identified upfront
→ Plan changes for all 8 files
→ Batch implement
→ Verify with compiler
```

**Impact:** Linear time complexity $O(n)$ instead of constant $O(1)$ planning. Each discovery requires context switch and separate edit attempt.

### 8.4 Metrics from Real Session

**Task Complexity:** Low (standard refactoring) **Expected Duration:** 20-30 minutes (manual) **Actual Duration:** 90-120 minutes (AI-assisted) **Efficiency:** 11% (89% wasted time)

**Error Distribution by Root Cause:**

- 40% - No compiler integration
- 30% - Selective read inefficiency
- 20% - Implementation verification failures
- 10% - Incremental discovery

**File Operation Metrics:**

- Files requiring multiple edits: 5
- Total edit attempts: 20+
- String-matching failures: 6+
- Average attempts per file: 2.5x

**Time Distribution:**

- Effective work: ~10 minutes (11%)
- Verification failures: ~30 minutes (25%)
- Context fragmentation: ~30 minutes (25%)
- Incremental discovery: ~40 minutes (39%)

**8.5 Required Workflow Changes**

**Before ANY Implementation:**

1. `grep -r "pattern"` → identify ALL affected files
2. List all required changes
3. Get user confirmation
4. THEN implement

**After Every Breaking Change:**

1. Run: `npm run check`
2. Parse TypeScript errors
3. Group by error type
4. Batch fix all instances
5. Re-run compiler
6. Repeat until green

**File Operations:**

1. Read full file (not offset/limit) unless >2000 lines
2. Make edit
3. Read back edited section to verify
4. Report actual vs intended change

**Post-Edit Verification Protocol:** Every edit must be verified that change was actually applied to file. Tool success report is insufficient.

**8.6 Meta-Lesson: Novice vs. Expert Workflow**

**What Failed (Novice Pattern):**

- Make change
- See what breaks
- Fix that
- Repeat

**What Works (Expert Pattern):**

- Analyze impact (grep, compiler)

- Plan all changes
- Execute systematically
- Verify completion

**Implication:** AI tools need guardrails to enforce expert workflows, not freedom to replicate novice patterns. Without these constraints, automation becomes slower than manual work.