

A Deep Dive into AI Coding Agents

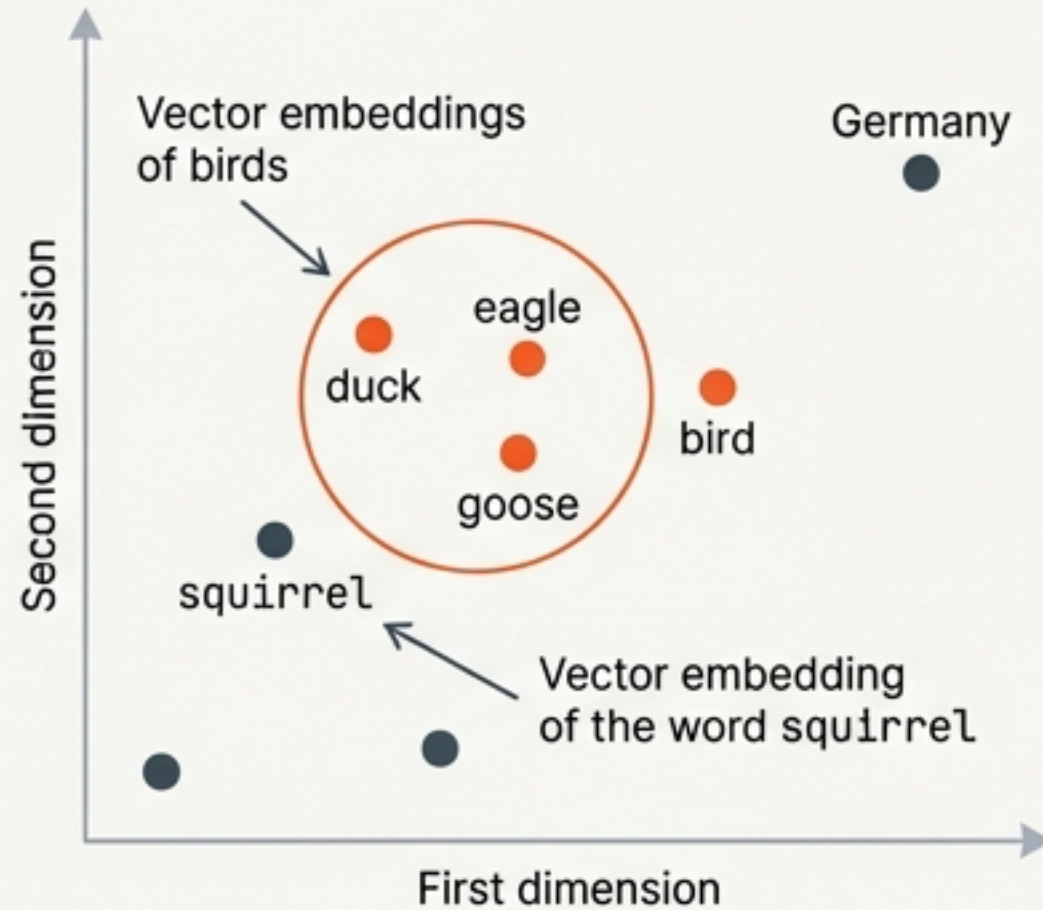
From Statistical Completion to Autonomous
Engineering (2012–Present)



Nathanaël Fijalkow | NukkAI

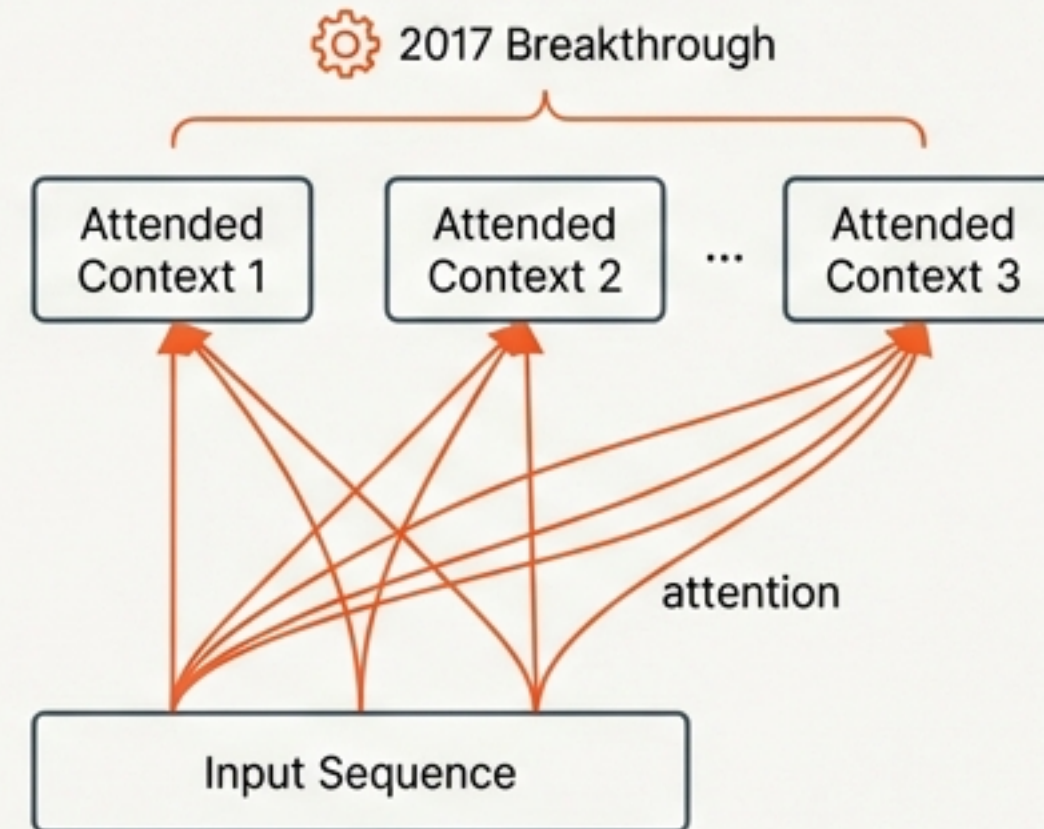
The Foundation: Embeddings, Attention, and Scale

Embeddings



Tokens are mapped to a dense vector space. Semantic proximity (e.g., `squirrel` to `bird`) allows the model to understand relationships beyond keyword matching.

Attention



The 2017 Breakthrough. Transformers moved beyond sequential processing, allowing models to “attend” to relevant context regardless of distance.

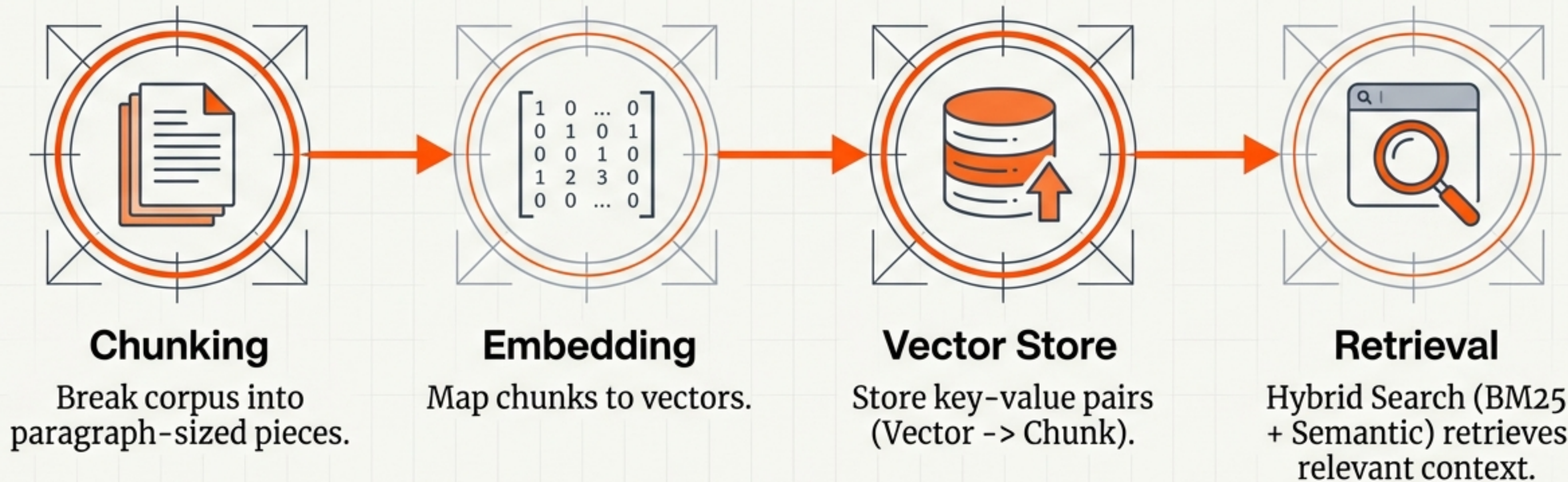
Scale

*“The incredible progress in AI over the past five years can be summarized in one word: **scale**.”*

— Noam Brown,
OpenAI Scientist (2022)

Bridging the Knowledge Gap with Retrieval-Augmented Generation (RAG)

LLMs are frozen in time. To code effectively, they need dynamic access to the codebase.



Reality Check

State of the Art (2026): RAG is massively deployed but not a “solved problem.” It requires significant engineering in chunking strategies and re-ranking to reduce latency and improve relevance.

Era 1: The Statistical Era (2012–2020)

The Naturalness Hypothesis



Research in 2012 proved code is remarkably predictable—even more so than English. The cross-entropy of software is lower than natural language, making it ideal for statistical modeling.

The Limitation: Local Texture



Early RNNs and LSTMs captured style ('texture') but failed at long-range dependencies. They could not remember a variable defined 50 lines up, lacking the memory horizon for complex engineering.

Era 2: The Completion Paradigm (2021–2022)

The Arrival of Transformers & GitHub Copilot

Fill-in-the-Middle (FIM)

```
code file x
def analyze_data(dataset): \
    # Load and process dataset
    df = pd.read_csv(dataset)
    df.dropna(inplace=True)

    [GENERATED MIDDLE]
    # Calculate key statistics
    stats = df.describe()
    correlation = df.corr()

    # Generate report
    report = generate_report(stats, correlation)

    return report

if __name__ == "__main__":
    analyze_data("sales_data.csv")
```

Prefix
JetBrains Mono

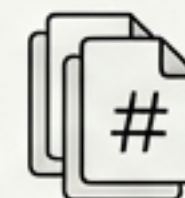
Training:
Shuffle

Model Learns to
Bridge Gaps

Prefix → Suffix → Middle

Suffix
JetBrains Mono

Data Curation: The Hidden Engine



Deduplication via
MinHash/LSH: Remove
duplicate files to prevent
memorization.



Sanitization: Entropy-based
scanning to redact PII, API
keys, and secrets.

The Limits of Passive Completion

Stateless

No memory of previous completions or build state. Did this suggestion break the build? The model doesn't know.

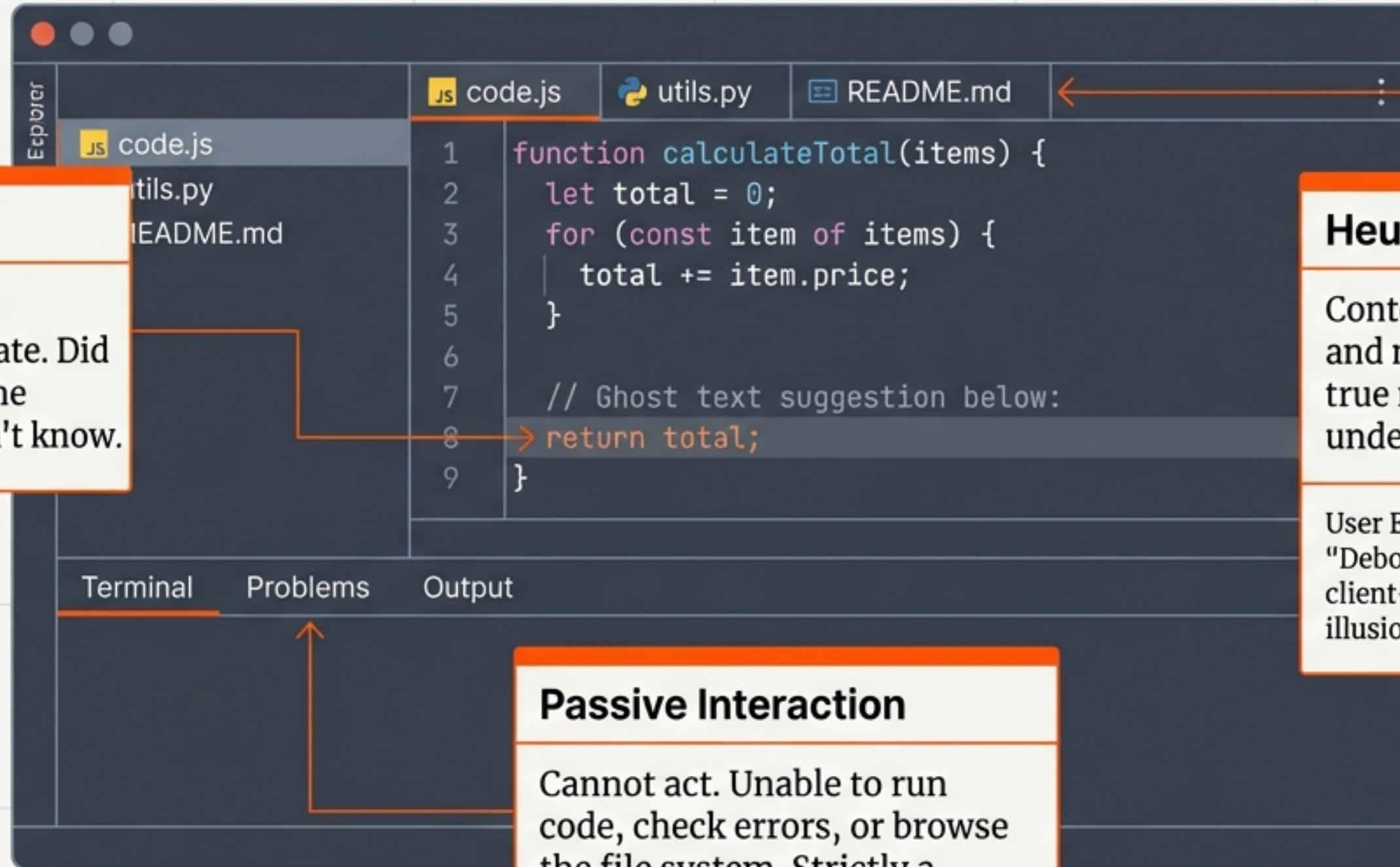
Heuristic Context

Context limited to open tabs and neighbors. Lacks true repository-wide understanding.

User Experience Optimization: Uses "Debouncing" (75ms wait) and client-side caching to create the illusion of instantaneity.

Passive Interaction

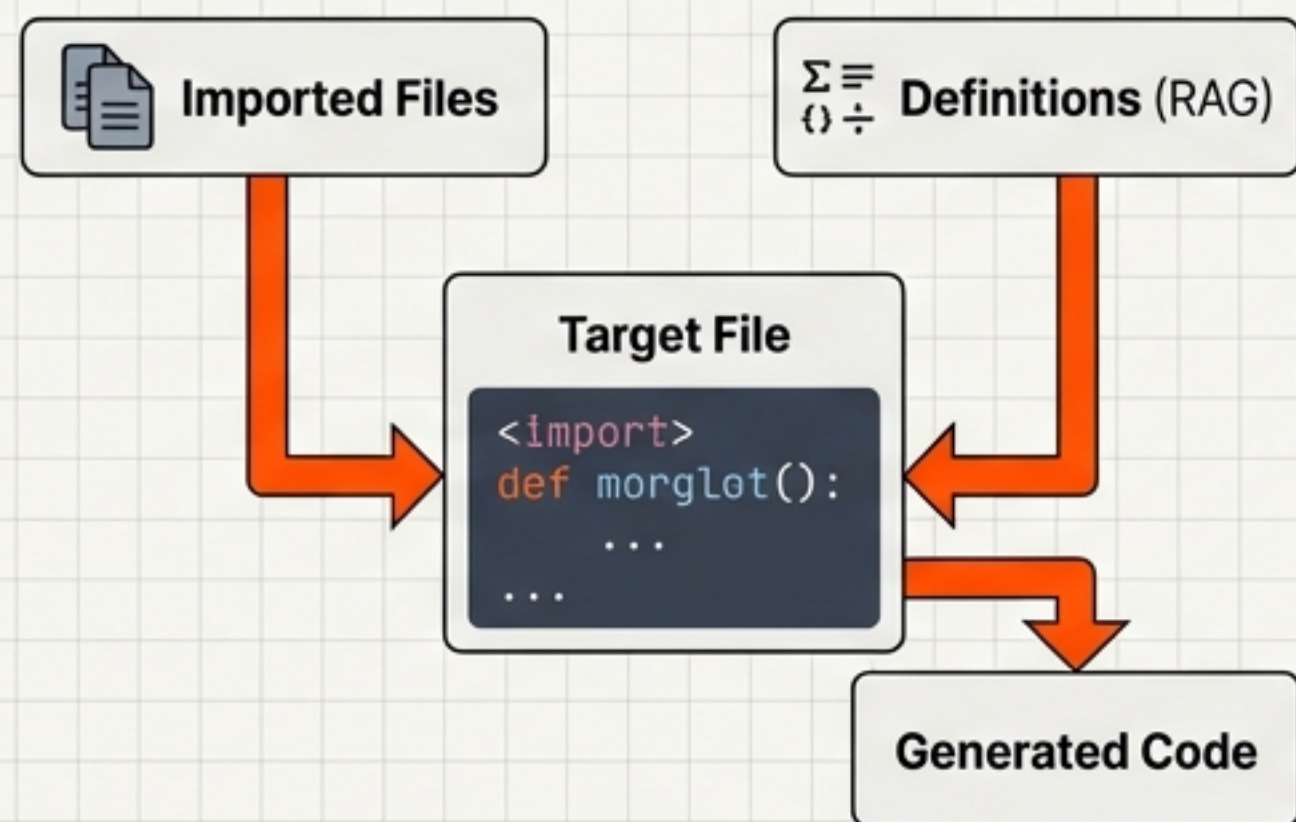
Cannot act. Unable to run code, check errors, or browse the file system. Strictly a text-prediction engine.



Era 3: Context & State Awareness (2023–2024)

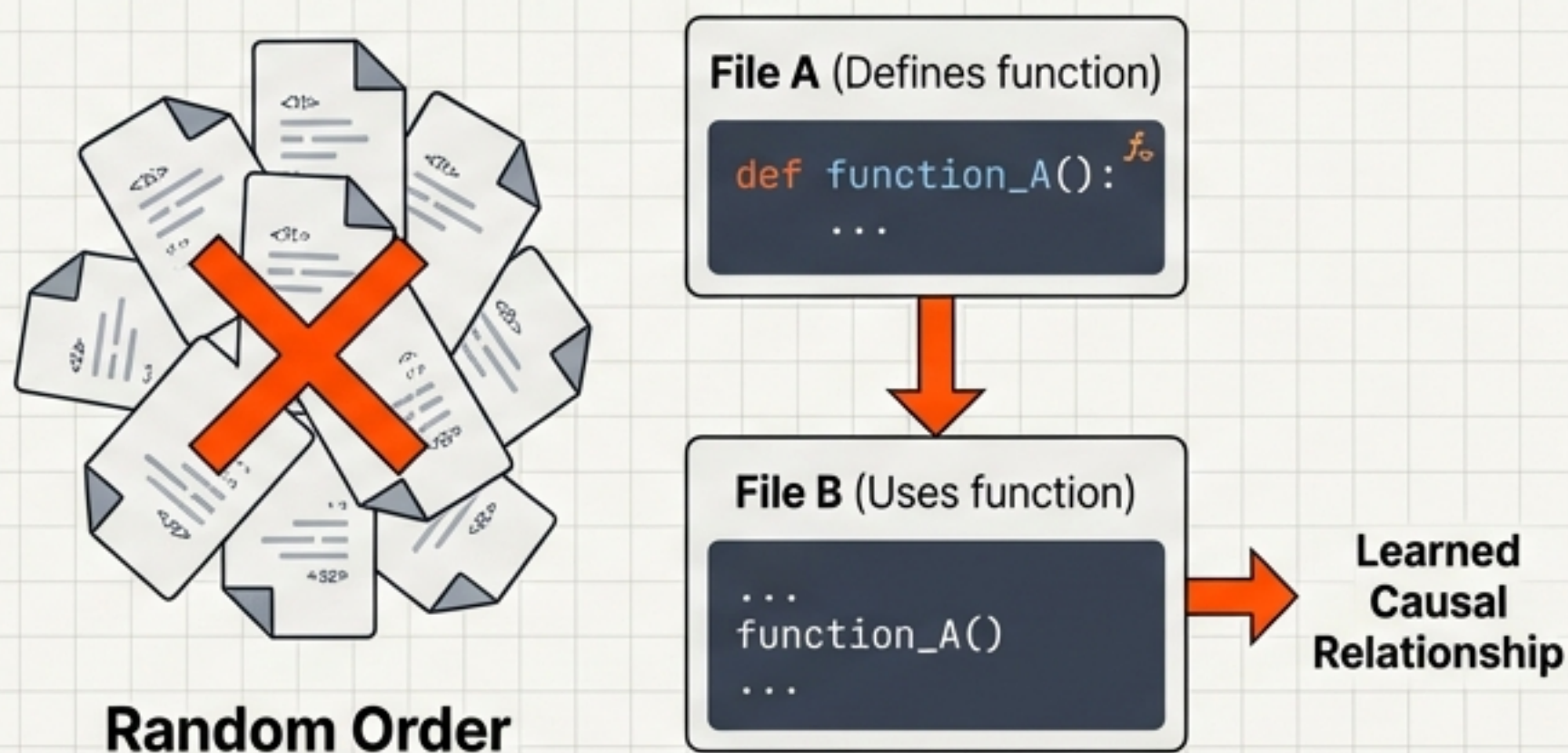
From File-Level Prediction to Repository-Level Understanding

RepoFusion



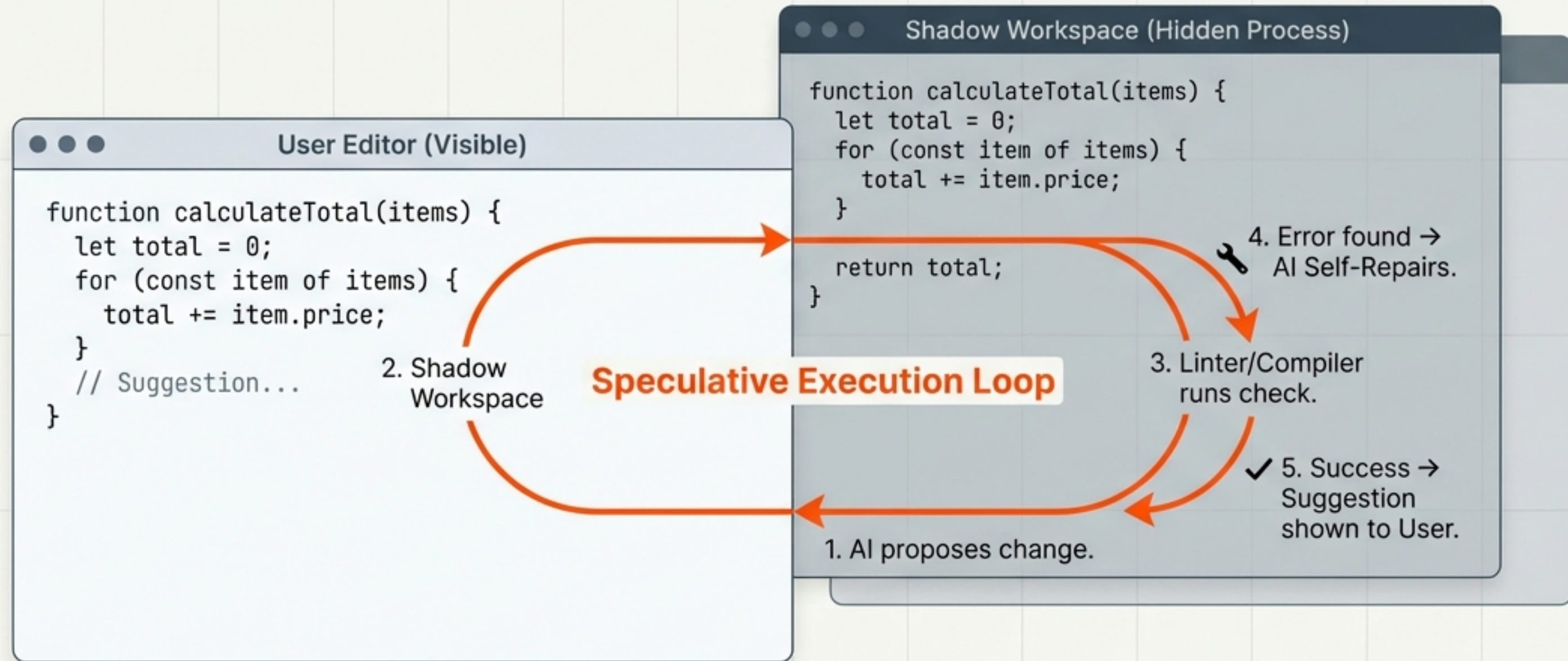
Inference-time RAG retrieves relevant cross-file context (imports, definitions) and prepends it to the prompt.

DeepSeek Coder: Topological Sorting



Training on a dependency graph rather than random files teaches the model causal relationships—seeing the definition before the usage.

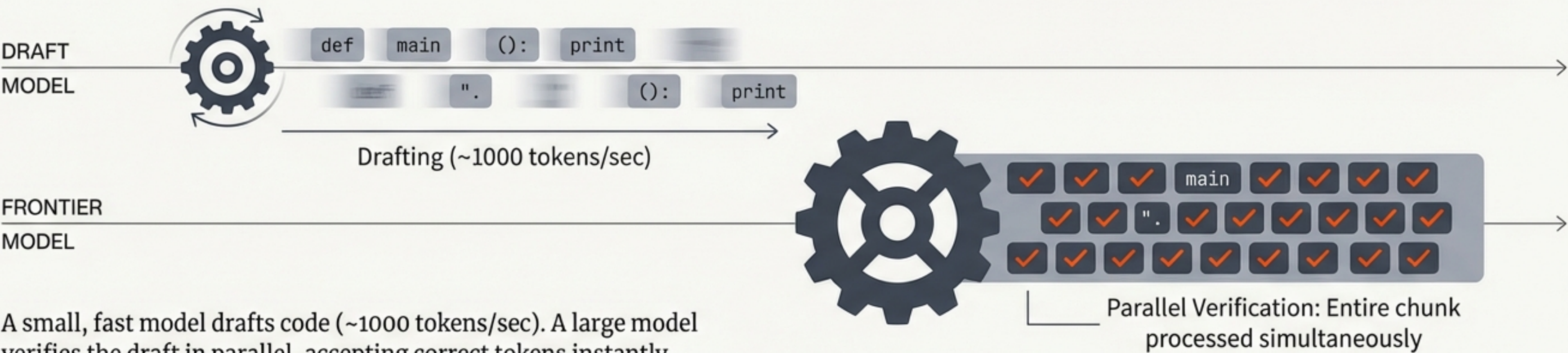
Cursor: The AI-Native IDE



Cursor forks VS Code to introduce a hidden, parallel instance of the editor. This allows the AI to “compile and check” its own code in the background before the user ever sees it.

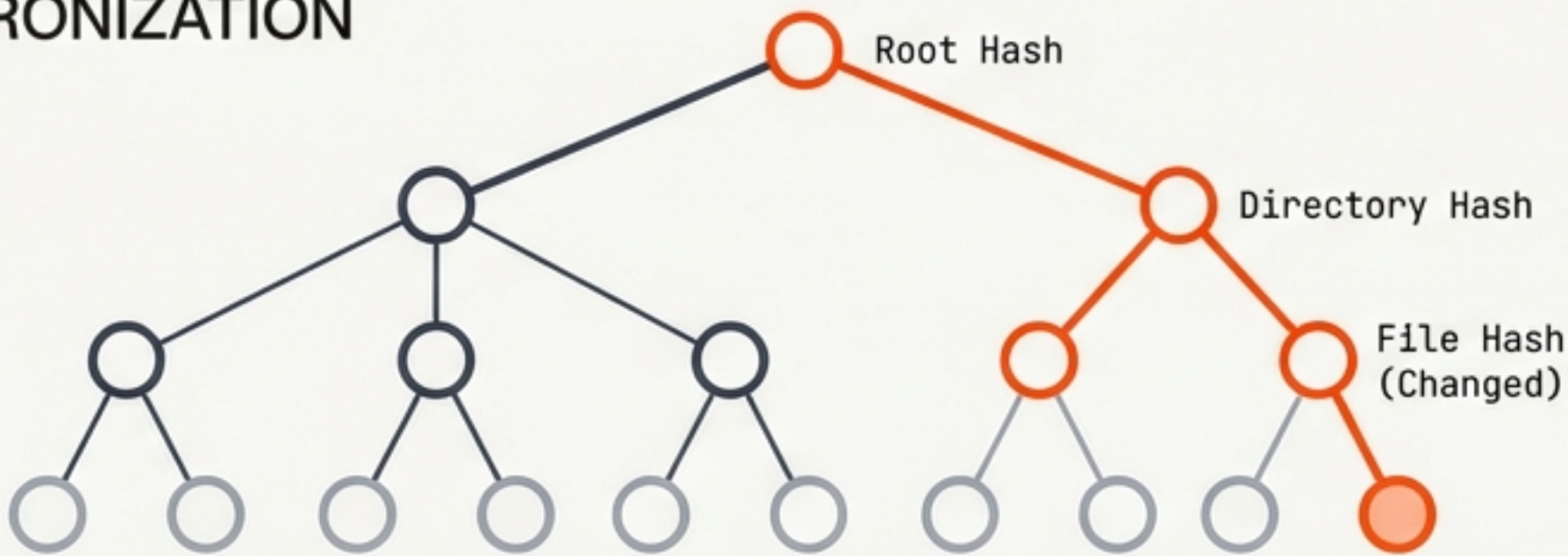
ENGINEERING FOR SPEED AND FRESHNESS

SPECULATIVE DECODING (DRAFT-VERIFY ARCHITECTURE)



A small, fast model drafts code (~1000 tokens/sec). A large model verifies the draft in parallel, accepting correct tokens instantly.

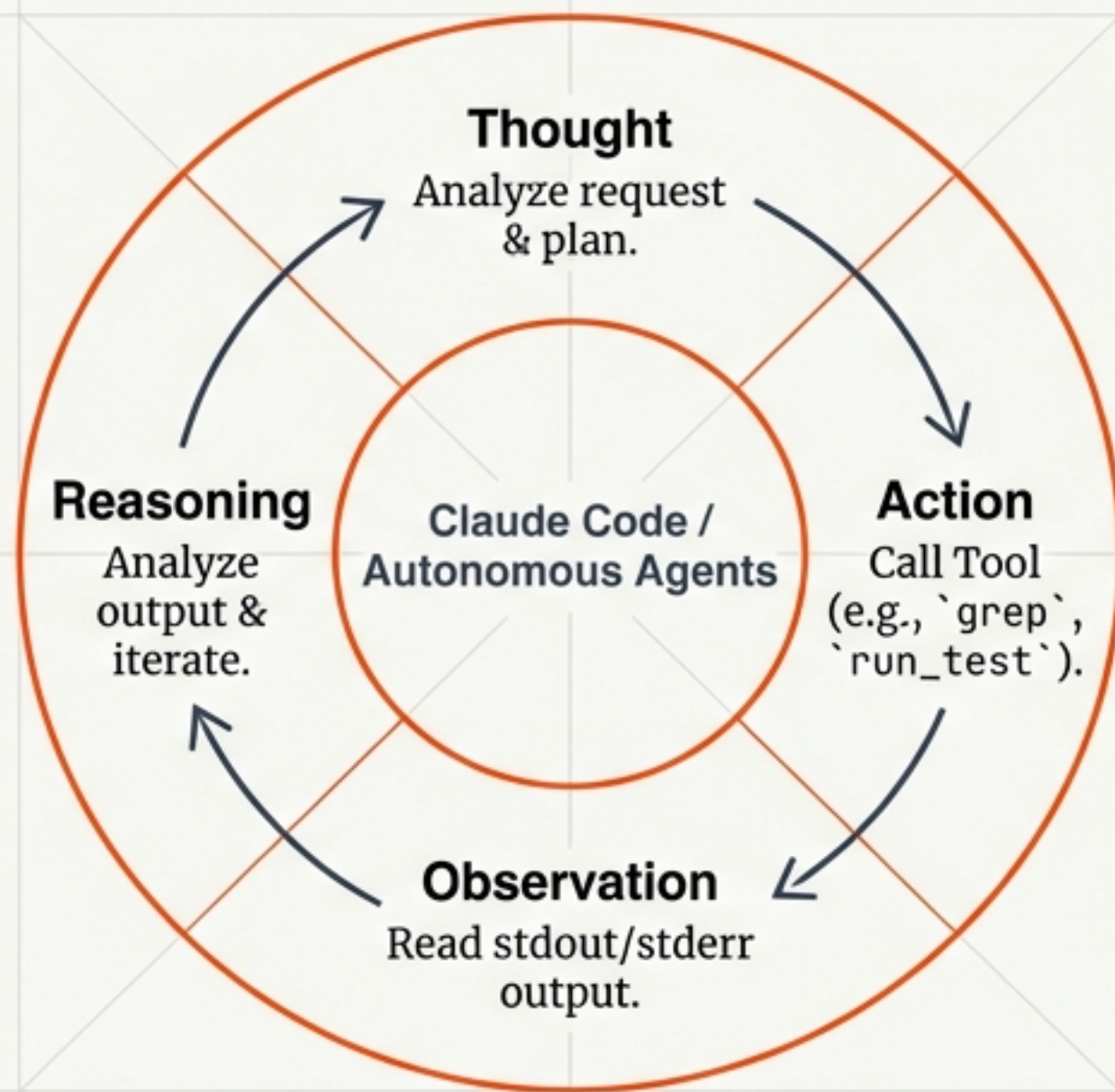
MERKLE TREE SYNCHRONIZATION



To keep vector stores fresh, files are hashed into a Merkle Tree. When typing occurs, only the branches with changed hashes are re-indexed, enabling efficient differential sync.

Era 4: The Agentic Era (2024–Present)

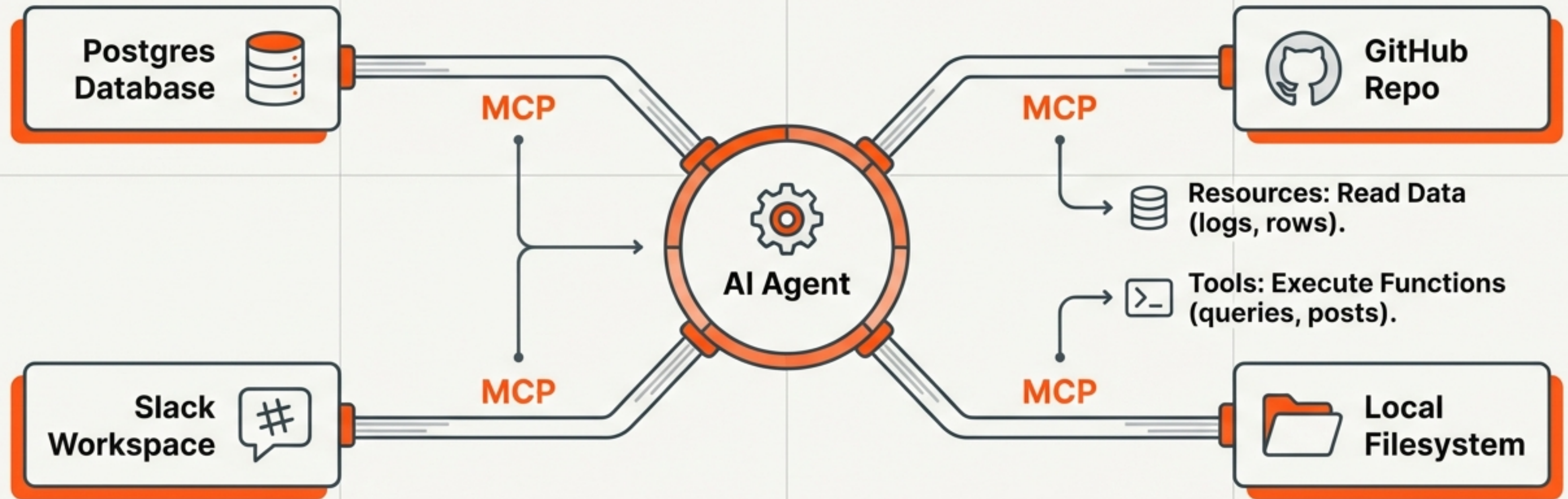
From “Helping” to “Acting”—The Autonomous Loop



Agents do not just predict text; they execute autonomous loops of reasoning, using terminal tools to explore, fix, and verify code.

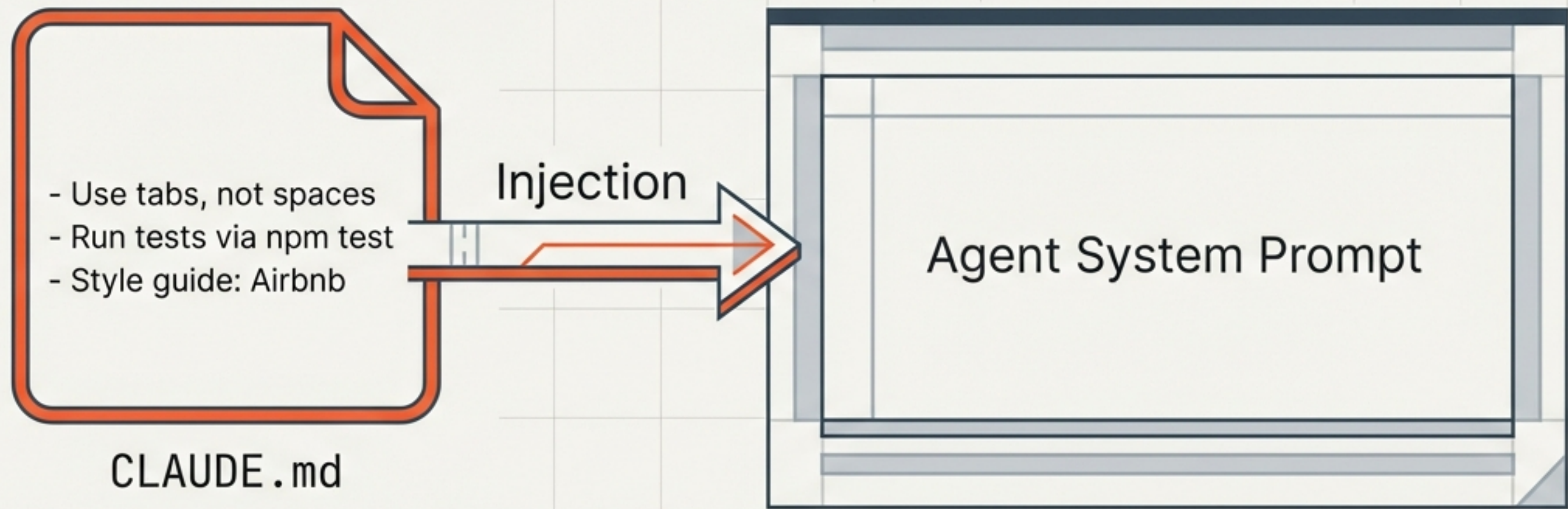
Standardizing Tools: Model Context Protocol (MCP)

Introduced by Anthropic (2024), MCP is an open standard that transforms the coding agent from a text editor into a systems integrator, allowing instant connectivity to any data source or tool.



Impact: This transforms the coding agent from a file-editor into a systems integrator.

Solving Context Drift: **Persistent Memory**



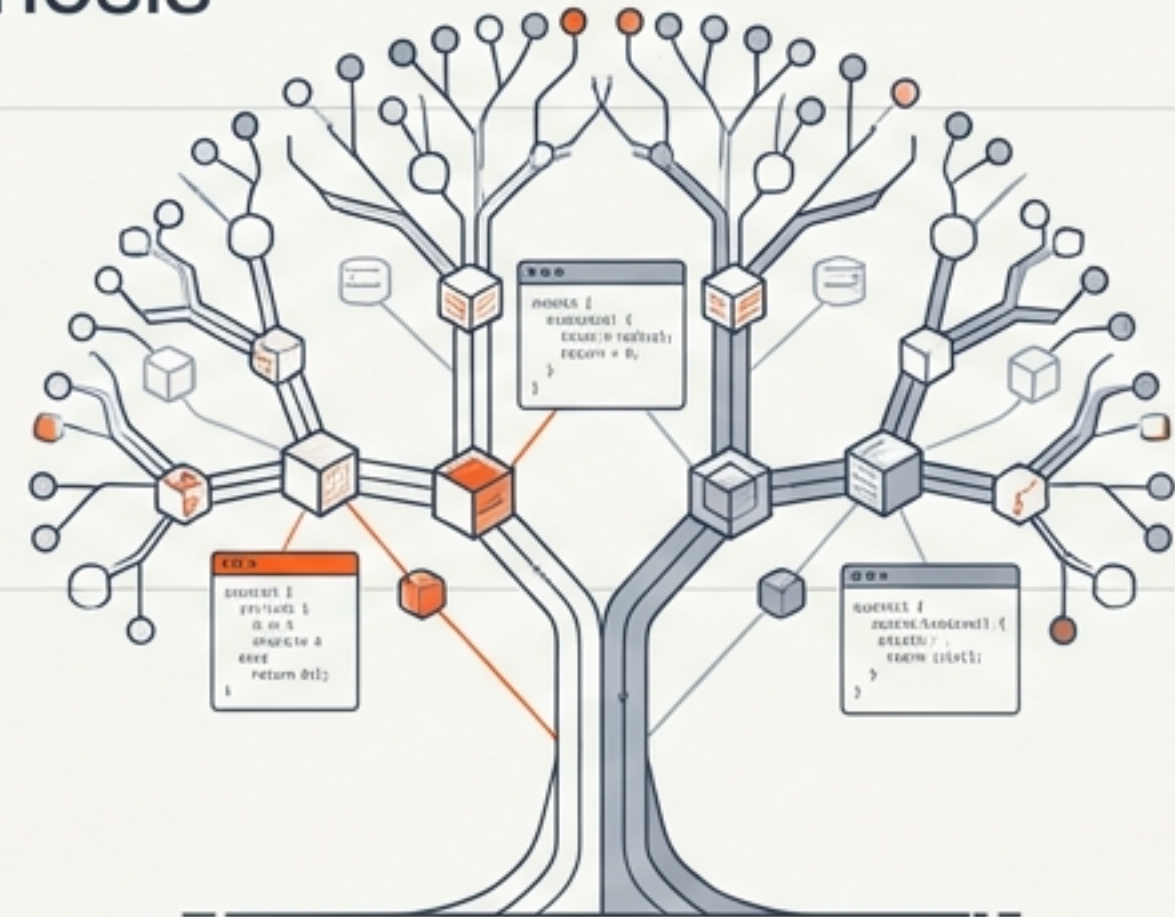
Agents suffer from amnesia between sessions. ``CLAUDE.md`` acts as a persistent memory anchor in the repository root. Its contents are automatically injected into every session, preserving implicit knowledge and project culture.

The Future: Back to the Trees

Moving from Brute Force to Structural Synthesis



Brute Force: Unstructured LLM Generation



Structural Synthesis: Abstract Syntax Trees (ASTs)

We have passed the era of believing LLMs will solve programming via brute force. The future lies in treating code respectfully—combining Generative AI with Formal Verification, Static Analysis, and Abstract Syntax Trees.

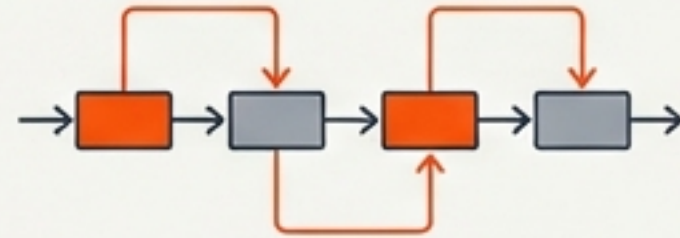
Emerging fields in JetBrains Mono, includes:

- Automated Theorem Proving
- Synthetic Biology
- Verified Hardware Synthesis

Summary & Key Takeaways

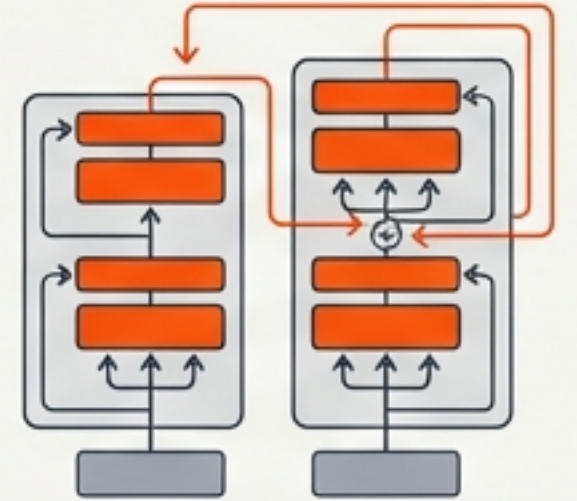
Statistical Era

Predicting local patterns (N-Grams).



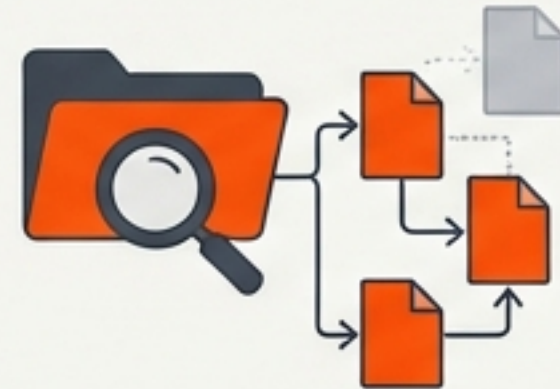
Completion Era

Transformers solving long-range dependencies (Copilot).



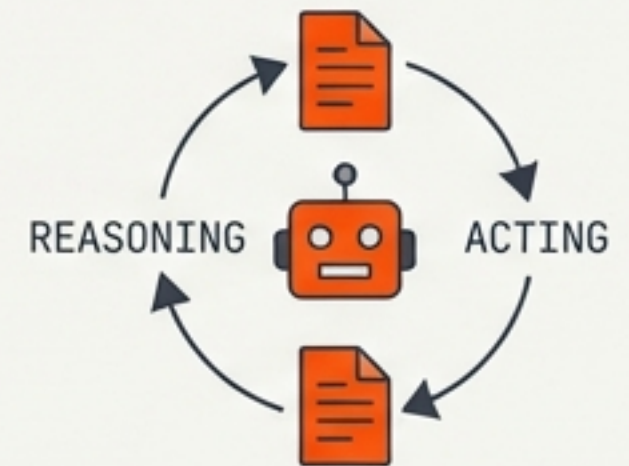
Context Era

Repo-level awareness & Shadow Workspaces (Cursor).



Agentic Era

Autonomous loops of reasoning & acting (Claude Code).



Final Thought

The core LLM is becoming a commodity. The 'Layer on Top'—engineering, state management, tool protocols, and memory anchors—is now the primary driver of reliability.