

# OLLM CLI

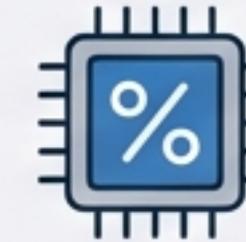
## The Anatomy of a Local-First AI Terminal

A provider-agnostic, VRAM-aware environment for intelligent development.



### Local-First

Deep integration with Ollama for privacy and offline capability.



### Resource-Aware

Mathematical context management to prevent Out-Of-Memory errors.



### Tool-Enabled

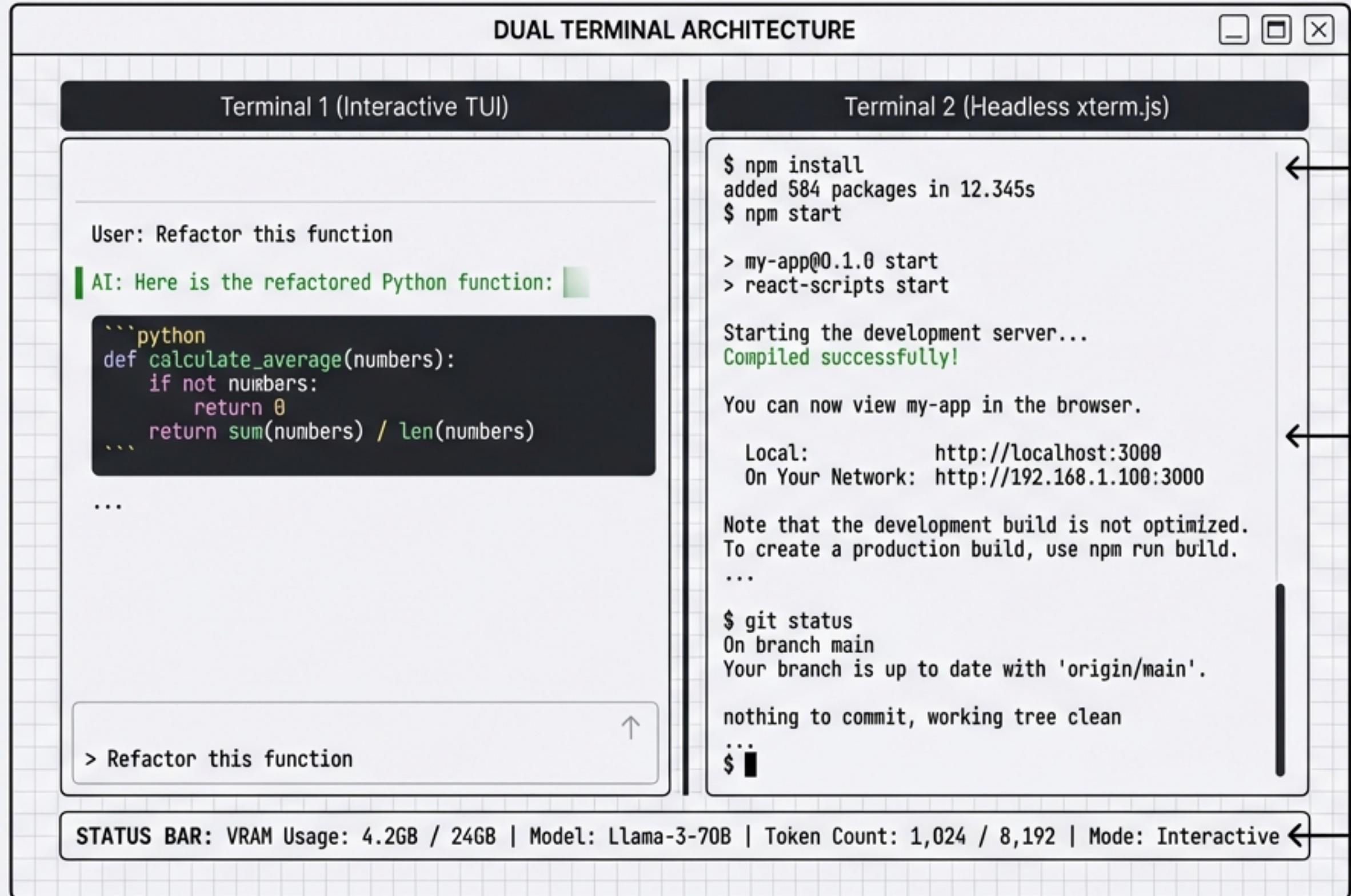
Robust file system, shell, and web access via Model Context Protocol.

```
user@ollm:~$ ollm start
System initialized.
VRAM detected: 24GB.
Model: Llama-3.1-8B loaded.
```

[STATUS: ACTIVE] [NODE: LOCAL]

A screenshot of a terminal window titled 'OLLM CLI'. It shows the command 'ollm start' being run, followed by system initialization messages: 'System initialized.', 'VRAM detected: 24GB.', and 'Model: Llama-3.1-8B loaded.'. At the bottom of the window, status indicators show '[STATUS: ACTIVE]' and '[NODE: LOCAL]'. The terminal has a dark background with light-colored text and a standard window title bar.

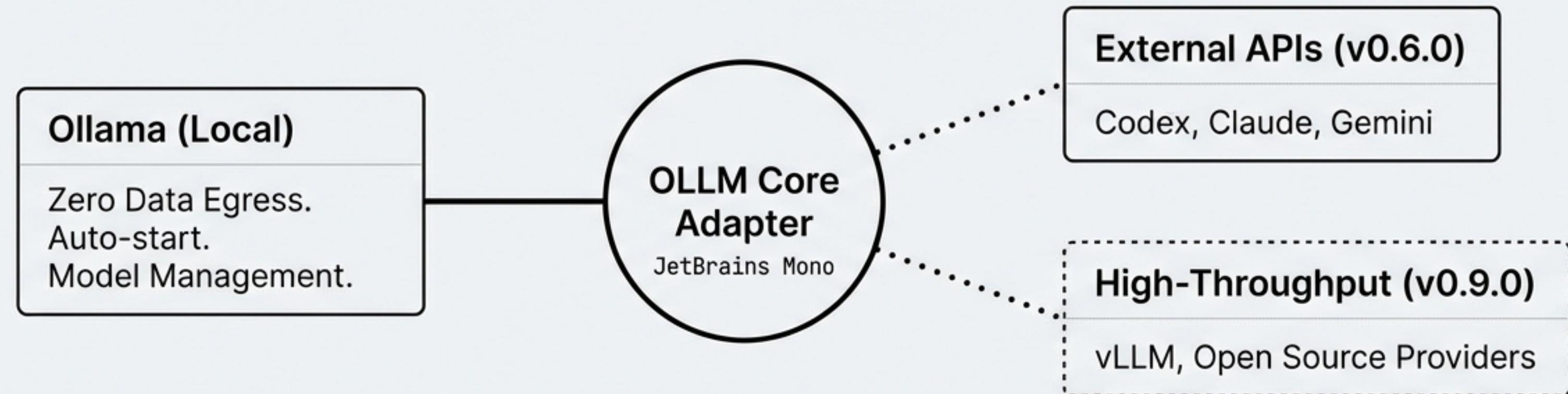
# The Interface: React + Ink Rendering



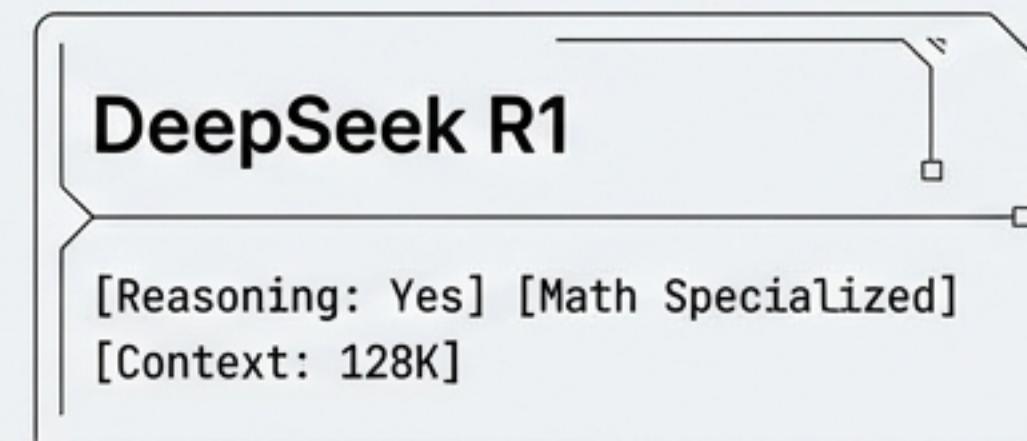
## NAVIGATION CONTROLS

- J / K : Vim-style Scrolling
- Ctrl+1 : Chat Tab
- Ctrl+2 : Tools Tab
- Ctrl+4 : Files Tab
- Ctrl+8 : MCP Tab

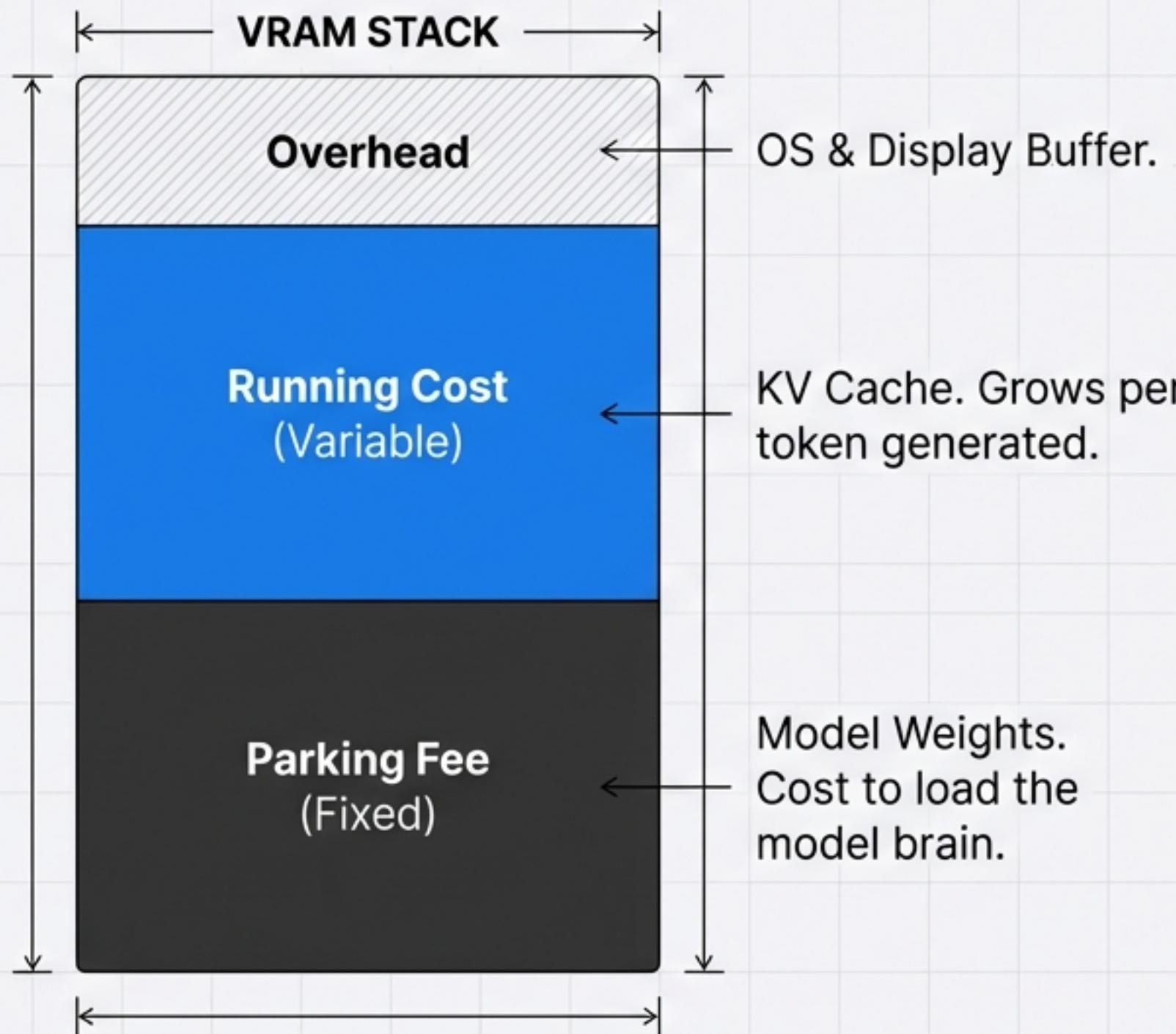
# The Brain: Provider-Agnostic Intelligence



## Model Database (35+ Profiles)



# Resource Logic: Math, Not Guesswork



## The Formula

$$\text{Total VRAM} = \text{Model Weights} + (\text{Context Length} \times \text{KV Cache per Token})$$

## Auto-Sizing Mechanism

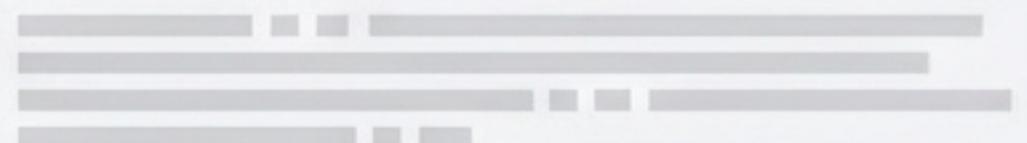
- **Target:** ~85% VRAM Utilization
- **Action:** Automatic context compression when limit reached.
- **Example:** 8GB Card → ~8k Context | 24GB Card → ~32k+ Context.

# Adaptive System Prompts

## Tier 1 (Minimal)

2K-4K Context | ~200 Token Budget

### System Prompt



Scales with Hardware

## Tier 5 (Ultra)

128K Context | ~1500 Token Budget

### System Prompt



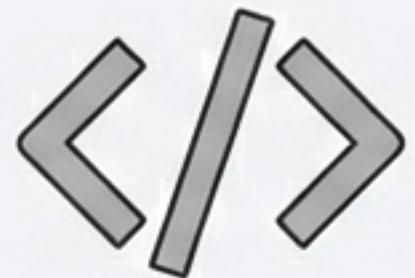
**Focus:** Essential guardrails only.  
High efficiency for low-end hardware.

**Focus:** Expert-level sophistication,  
mentoring, advanced reasoning patterns.

Maximizes quality on high-end rigs. Maximizes workspace on low-end laptops.

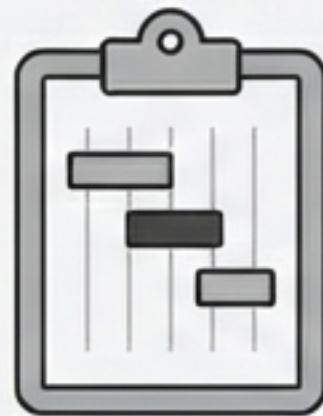
# Operational Modes

Context-aware personas for specific tasks.



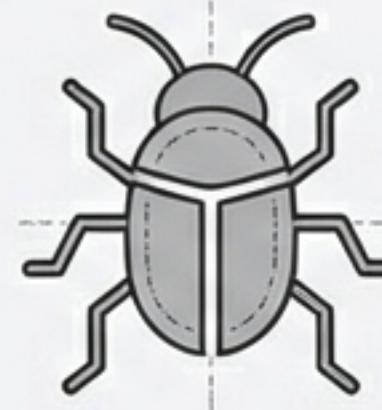
## Dev Mode

- **Focus:** SOLID principles, Code Quality, Testing, Testing
- **Output:** Production-ready code.



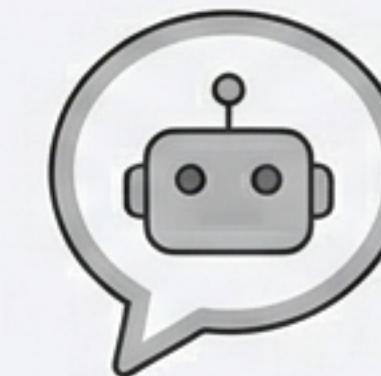
## Plan Mode

- **Focus:** Risk Assessment, Dependencies, Estimation
- **Output:** Actionable plans.



## Debug Mode

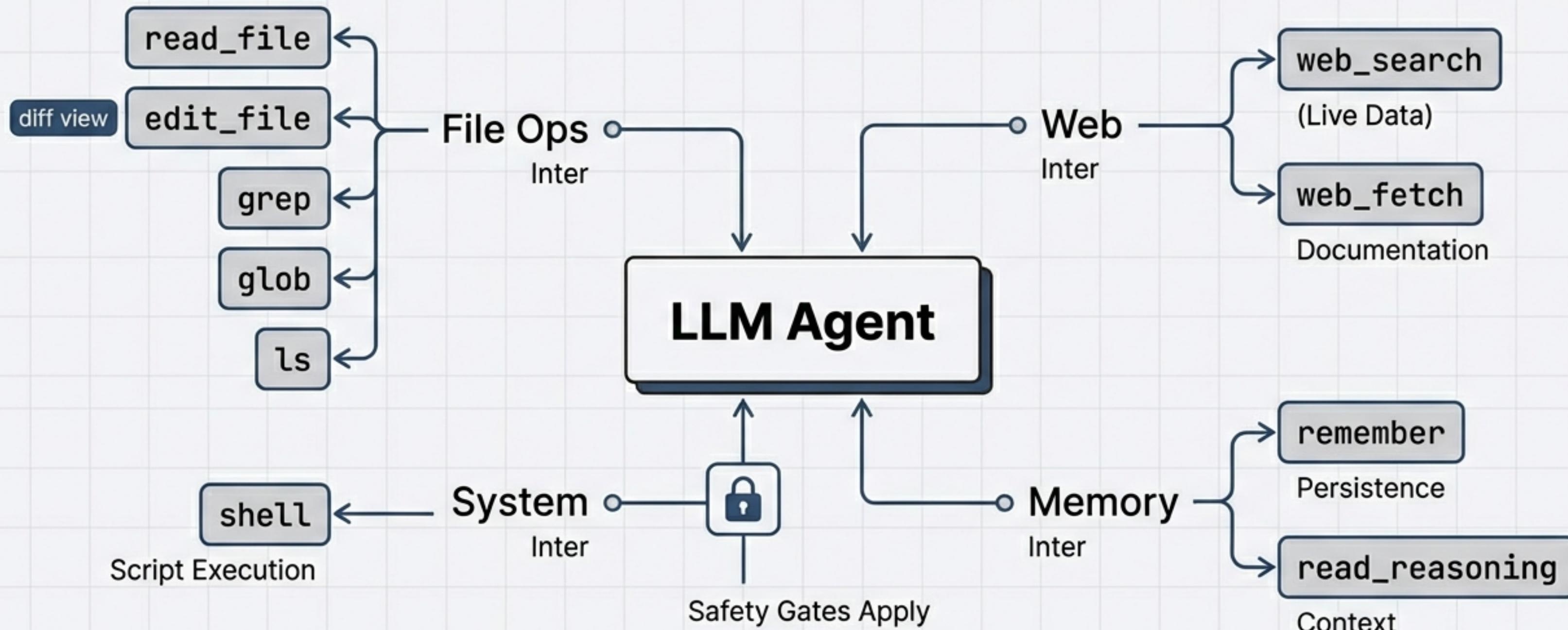
- **Focus:** Root Cause Analysis, Reproduction, Hypothesis Testing
- **Output:** Systematic fixes.



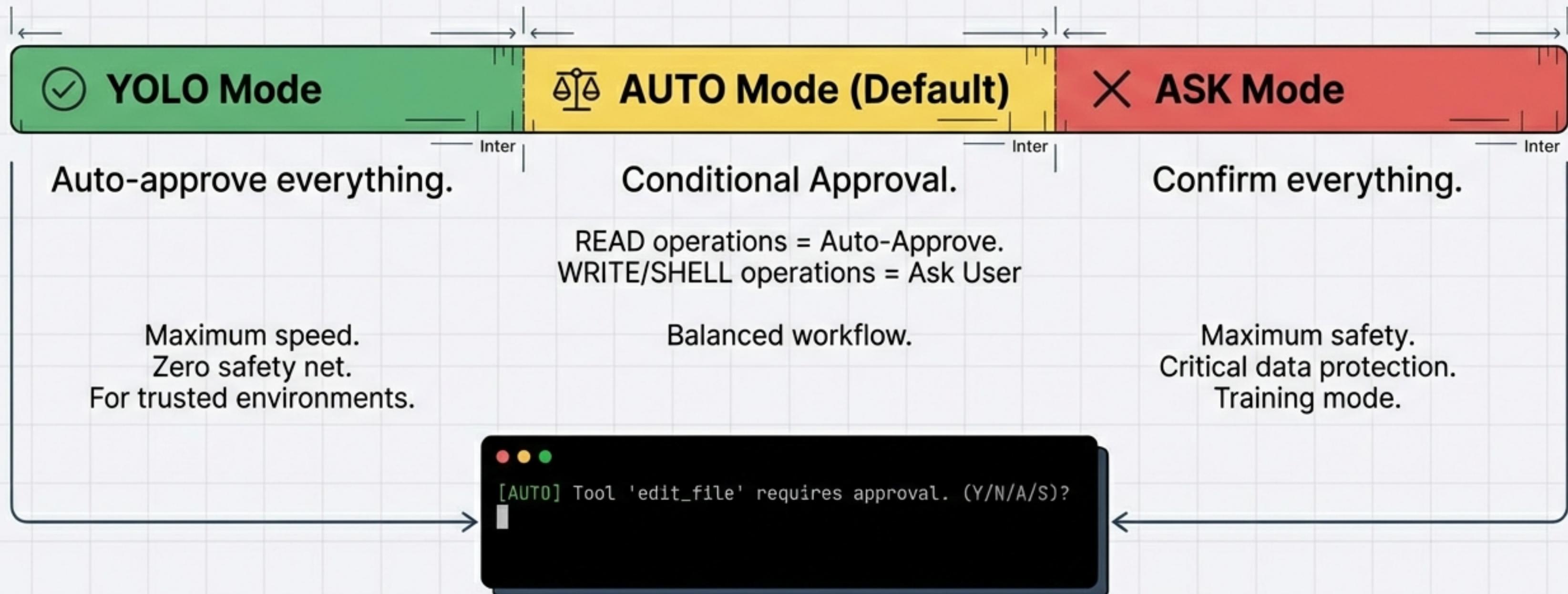
## Assistant Mode

- **Focus:** General Help, Communication, User Preferences
- **Output:** Clear explanations.

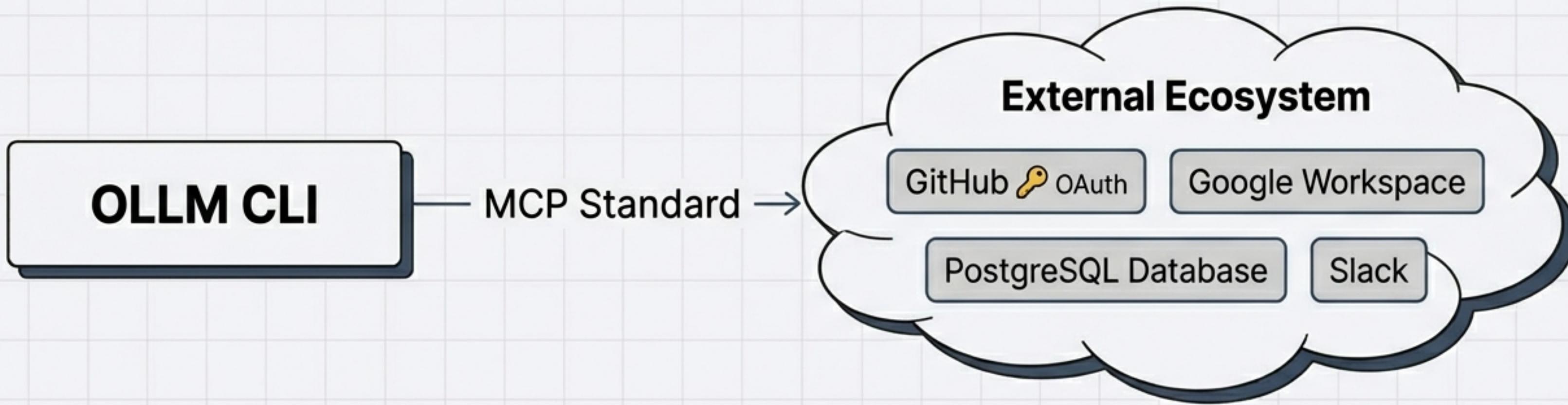
# The Tool System: Agency & Action



# Safety Spectrum: The Approval Protocol



# Extensibility: Model Context Protocol (MCP)



## Standardized Protocol

Connects AI to external tools, resources, and prompt templates.

## Transports

- Supports `Stdio` (Local Processes)
- Supports `SSE` (Remote Streams)

## Capabilities

- Marketplace browsing (/mcp)
- Built-in OAuth token management
- Automatic Health Monitoring

# Memory Architecture & Persistence

## Checkpoint Aging Strategy

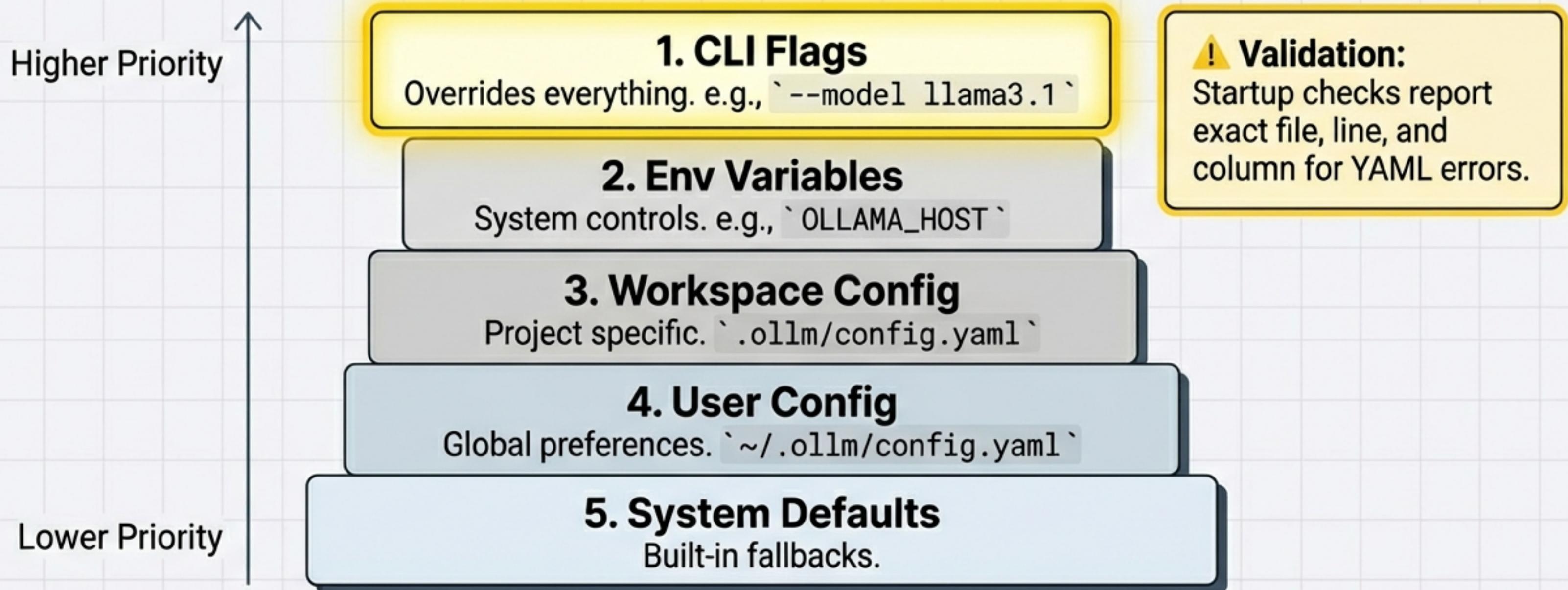


## Persistence Layer

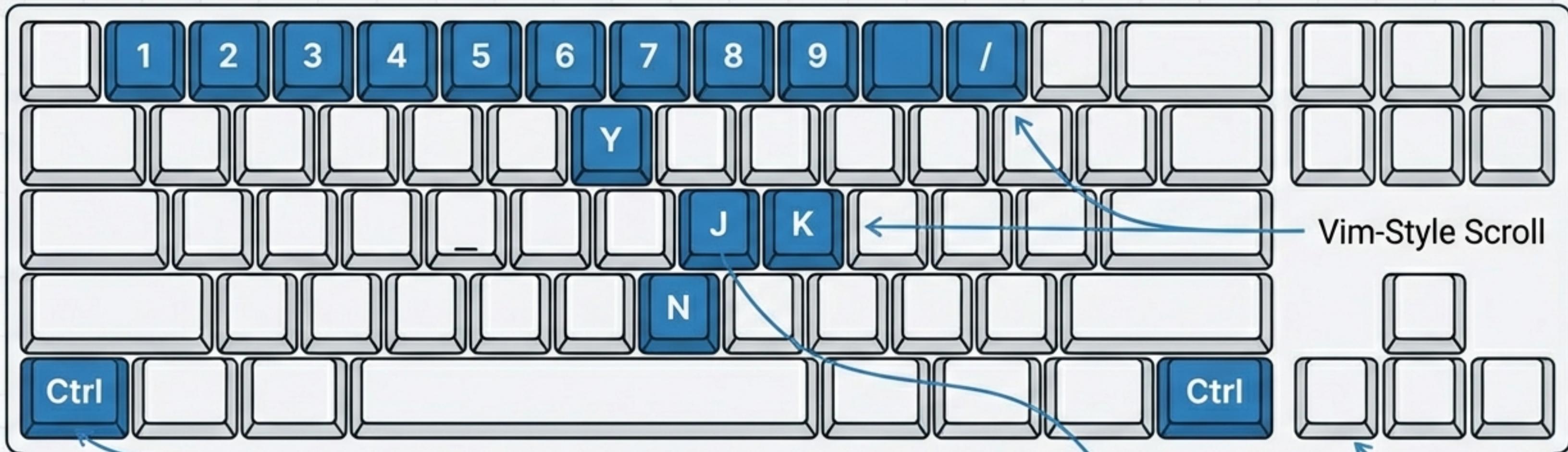
- **Session Management:** Save/Resume (`/session save`). Compression history.
- **Cross-Session Memory:** stored in `~/.o11m/memory.json`. Remembers user preferences and facts.

# The Precedence Stack

Fully tunable behavior via layered configuration.



# Developer-Centric Workflow



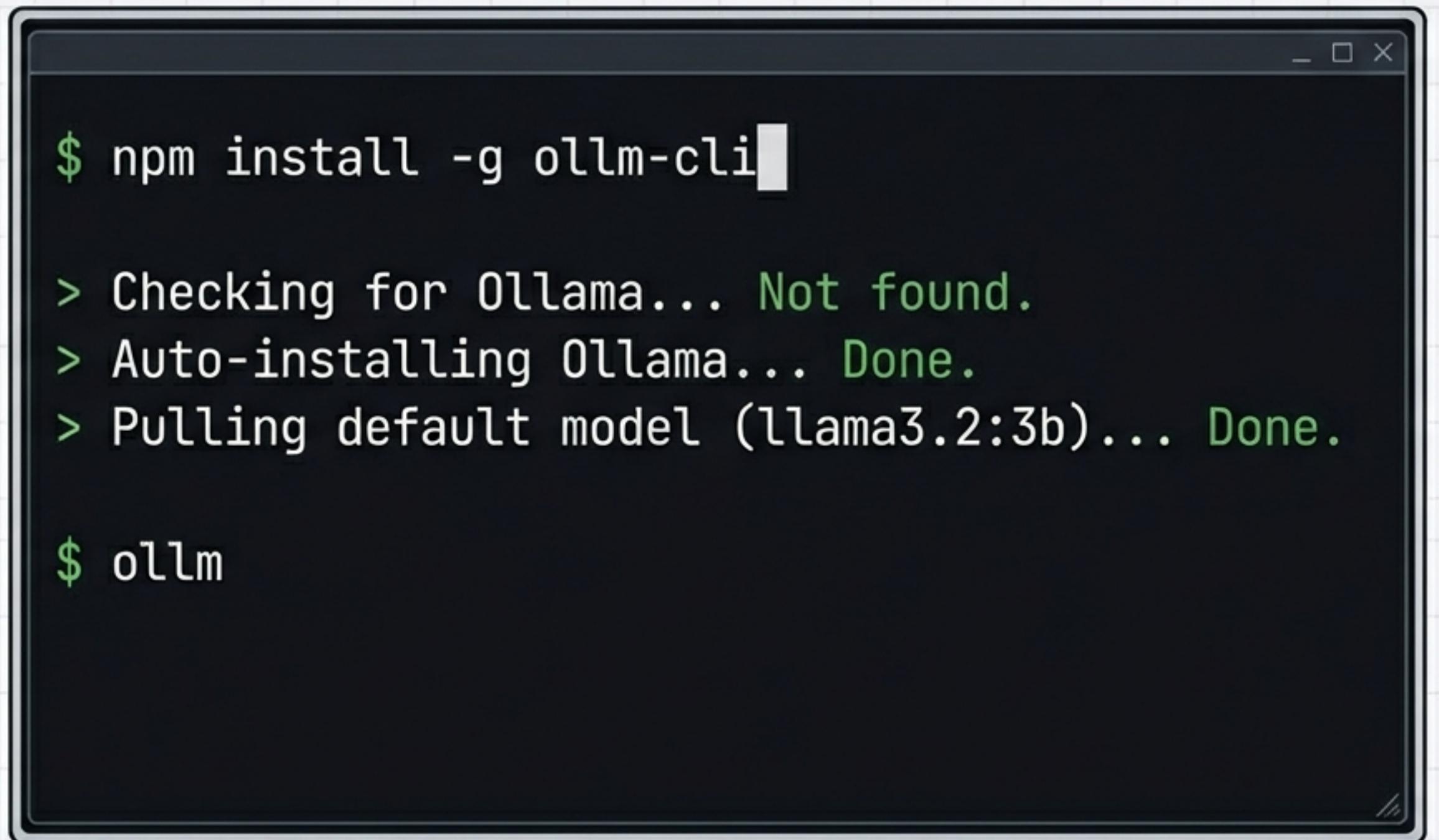
## Power Features

- **Quick Open Files** (**Ctrl+P**) - in JetBrains Mono.
- **Command Palette** (`/model`, `/context`).
- **Review Mode** for rapid code diff approvals.

# Development Roadmap (Alpha Phase)



# Quick Start: Zero to AI in Minutes



A terminal window showing the installation of the Ollama CLI and its initial run. The terminal has a dark background with white text.

```
$ npm install -g ollm-cli
> Checking for Ollama... Not found.
> Auto-installing Ollama... Done.
> Pulling default model (llama3.2:3b)... Done.

$ ollm
```

## Key Commands

- **/help:** Show all commands.
- **/model list:** Switch models.
- **/context stats:** View VRAM/Context usage.

# Capabilities at a Glance

CAPABILITY	STATUS	DESCRIPTION
Local-First		Offline capable, Privacy-focused, Zero data egress.
VRAM-Aware		Mathematical context management prevents OOM crashes.
Tool-Enabled		Native Files, Shell, and Web access with safety gates.
Provider-Agnostic		Ready for vLLM, OpenAI-compatible backends.
Developer-Centric		TUI, Vim keys, Configurable hierarchy.

## Build the future of local AI tooling.

[github.com/o11m/o11m-cli](https://github.com/o11m/o11m-cli)