



LOCI - The First CPU-GPU Software Execution-Aware Model

Signal Faults and Errors
No Run, Before Tests, Logs or Production

Built for Embedded Software Engineers

Meet LOCI: The GPU/CPU-Execution Aware Agent that other LLMs do not have

LOCI, an MCP and S2-based solution, provides a new layer of intelligence for developers across networking, embedded systems, AI infrastructure, inference, and performance-sensitive systems. Unlike generic AI assistants, LOCI is designed to understand how code really interacts with hardware; it does not stop on syntax- it analyzes the exe files after the project build is completed, gives feedback, and finds bugs without running the SW.

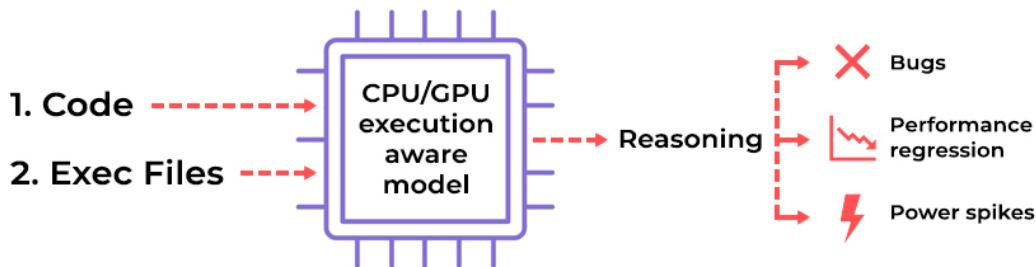
How does LOCI work?

Specialized model, trained on CPU-GPU executions in different conditions

LOCI does not run or execute the software. Instead, it reads the program's existing compiled executable as static data and analyzes its structure. Even without execution, the binary already contains the full logic of how the software can behave—its functions, branches, loops, and call relationships.

LOCI breaks this static executable into execution units, maps all possible execution paths, and reconstructs a behavioral model of the program.

To make this static structure come to life, as if it were running, LOCI applies learned proprietary performance models for CPUs and GPUs that capture real hardware execution behavior.



These models estimate how long each part of the executable would take to run, how much power it would consume, and how it would stress hardware resources. LOCI composes these estimates across full execution paths and call stacks to produce predicted timelines, flame graphs, response times, power profiles, and bottleneck analysis—all without ever running the software.

LOCI analyzes from executable files without running them:

- Throughput, Bandwidth, Latency
- AI/LLM Networking: tokens-per-second drift
- CPU execution paths and branching behavior
- GPU kernels, warp divergence, and memory coalescing
- Logical bugs such as null dereferences, race conditions, and edge-case logic
- Power anomalies
- Memory pressure, allocation patterns, and concurrency hazards

LOCI's outputs are constrained, physically meaningful

Large language models are generative systems optimized to produce plausible sequences of text or code, which can lead to hallucination when the model generates outputs that are syntactically valid but not grounded in execution reality. In contrast, the LOCI model is not generative and does not produce free-form outputs. It is trained in a supervised manner to predict a single, bounded numeric value - execution time, directly from assembly-level representations and measured ground-truth data. Because LOCI's outputs are constrained, physically meaningful, and tied to executable structure and hardware behavior, it cannot fabricate results or invent behaviors outside the defined execution domain.

LOCI in Action - Examples

LOCI automatically comments on pull requests, surfacing issues developers often miss:

- Performance drift detected (+12.4%) in `processFrame()`: nested loop increases GPU pressure.
- Power anomaly identified in `computeTransform()`: risk of IR-drop under peak load.
- Logical bug found: uninitialized variable '`bufferLen`' and potential null pointer dereference.

In addition, LOCI provides fixes:

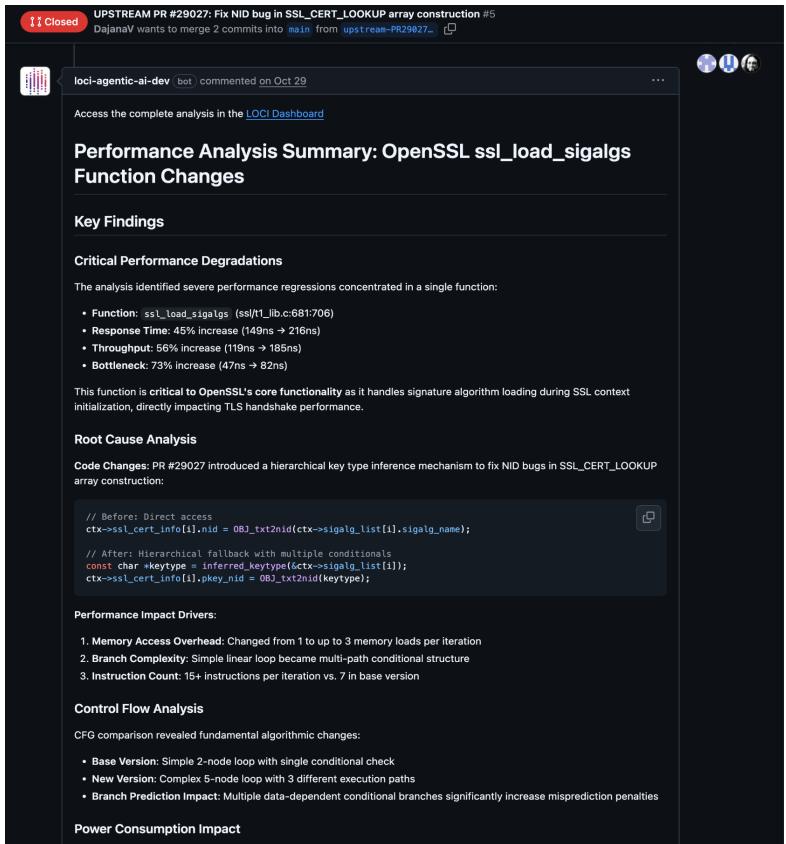
1. Move memory allocation outside the inner loop to stabilize GPU timing.
2. Align GPU and CPU execution paths to avoid unpredictable divergence.
3. Ensure all variables are initialized before branching logic.

GitHub Pull Request Example

In the screenshot below, you can see a full analysis of PR in the [GitHub OpenSSL](#) project before engaging in code review or executing any tests.

Click here to enter the project on GitHub:
<https://github.com/auroralabs-loci/openssl/pull/5>

This eliminates guesswork and provides developers with clarity before CI/CD pipelines even start.



UPSTREAM PR #29027: Fix NID bug in SSL_CERT_LOOKUP array construction #5
DajanaV wants to merge 2 commits into `main` from `upstream-PR29027...`

loci-agentic-ai-dev bot commented on Oct 29

Access the complete analysis in the [LOCI Dashboard](#)

Performance Analysis Summary: OpenSSL `ssl_load_sigalgs` Function Changes

Key Findings

Critical Performance Degradations

The analysis identified severe performance regressions concentrated in a single function:

- Function: `ssl_load_sigalgs` (`ssl/t1_lib.c:681:706`)
- Response Time: 45% increase (149ns → 216ns)
- Throughput: 56% increase (119ns → 185ns)
- Bottleneck: 73% increase (47ns → 82ns)

This function is critical to OpenSSL's core functionality as it handles signature algorithm loading during SSL context initialization, directly impacting TLS handshake performance.

Root Cause Analysis

Code Changes: PR #29027 introduced a hierarchical key type inference mechanism to fix ND bugs in `SSL_CERT_LOOKUP` array construction:

```
// Before: Direct access
ctx->ssl_cert_info[i].nid = 0B1_txt2nid(ctx->sigalg_list[i].sigalg_name);

// After: Hierarchical fallback with multiple conditionals
const char *keytype = inferenc_keytype(ctx->sigalg_list[i]);
ctx->ssl_cert_info[i].pkey_nid = 0B1_txt2nid(keytype);
```

Performance Impact Drivers:

1. Memory Access Overhead: Changed from 1 to up to 3 memory loads per iteration
2. Branch Complexity: Simple linear loop became multi-path conditional structure
3. Instruction Count: 15+ instructions per iteration vs. 7 in base version

Control Flow Analysis

CFG comparison revealed fundamental algorithmic changes:

- Base Version: Simple 2-node loop with single conditional check
- New Version: Complex 5-node loop with 3 different execution paths
- Branch Prediction Impact: Multiple data-dependent conditional branches significantly increase misprediction penalties

Power Consumption Impact

Minimum Power consumption impact

Logic, Performance, and Power Analysis of OpenSSL PR#29027, Oct 29, 2025

AURORA LABS

Ask LOCI Anything | Via GitHub comment @LOCI-DEV or directly from your IDE

LOCI supports interactive Q&A inside the PR. Developers can ask targeted engineering questions:

- Why did execution time increase?
- Which lines introduced performance drift?
- What caused the sudden power spike?
- Does this change introduce logical or memory safety risks?

Example of LOCI response:

"Execution drift is caused by a memory allocation moved inside a high-frequency loop.

TCP_IP_Send shows a 4.7% slowdown. Variable 'packetCount' may be uninitialized.

Recommended fixes: refactor allocation, initialize variable, rebalance branching logic."

The following is LOCI's terminal App for VS Code, using its CPU/GPU execution aware modeling to find Bugs and Performance issues after project build, before running tests

llama.cpp-upstream-PR16490-branch_ggml-org-gg-graph-mamba-reuse

llama-context.cpp

```
21  llama_context::llama_context()
22
23  if (params.attention_type == LLAMA_ATTENTION_TYPE_UNSPECIFIED) {
24      cparams.causal_attn = hparams.causal_attn;
25  } else {
26      cparams.causal_attn = params.attention_type == LLAMA_ATTENTION_TYPE_CAUSAL;
27  }
28
29  cparams.flash_attn = params.flash_attn_type != LLAMA_FLASH_ATTN_TYPE_DISABLED;
30
31  // with causal attention, the batch size is limited by the context size
32  cparams.n_batch = cparams.causal_attn ? std::min(cparams.n_ctx, params.n_batch) : params.n_batch;
33
34  cparams.n_ubatch = std::min(params.n_batch, params.n_ubatch == 0 ? params.n_batch : params.n_ubatch);
35
36  cparams.op_offload = params.op_offload;
37  cparams_kv_unified = params_kv_unified;
38
39  const char * LLAMA_GRAPH_REUSE_DISABLE = getenv("LLAMA_GRAPH_REUSE_DISABLE");
40  graph_reuse_disable = LLAMA_GRAPH_REUSE_DISABLE ? (atoi(LLAMA_GRAPH_REUSE_DISABLE) != 0) : graph_reuse_disable;
41
42  if (graph_reuse_disable) {
43      LLAMA_LOG_WARN("%s: graph reuse disabled\n", __func__);
44  }
45
46  // ref: https://github.com/ggml-org/llama.cpp/pull/17046#discussion_r250388573
47  cparams.n_ctx = GGML_PAD(cparams.n_ctx, 256);
48
49  if (params_kv_unified) {
50      cparams.n_ctx_seq = cparams.n_ctx;
51  } else {
52      ...
53  }
54 }
```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

Let's get that performance back! 🔥
You: what is the impact on per second in this PR?

Ln 1, Col 1 Spaces: 4 UTF-8 LF () C++ Signed out

LOCI Terminal

LOCI flagged the most impactful degradation issue.

The top 3 critical functions are presented

llama.cpp-upstream-PR16490-branch_ggml-org-gg-graph-mamba-reuse

llama-context.cpp

```
21  llama_context::llama_context()
22
23  if (params.attention_type == LLAMA_ATTENTION_TYPE_UNSPECIFIED) {
24      cparams.causal_attn = hparams.causal_attn;
25  } else {
26      cparams.causal_attn = params.attention_type == LLAMA_ATTENTION_TYPE_CAUSAL;
27  }
28
29  cparams.flash_attn = params.flash_attn_type != LLAMA_FLASH_ATTN_TYPE_DISABLED;
30
31  // with causal attention, the batch size is limited by the context size
32  cparams.n_batch = cparams.causal_attn ? std::min(cparams.n_ctx, params.n_batch) : params.n_batch;
33
34  cparams.n_ubatch = std::min(params.n_batch, params.n_ubatch == 0 ? params.n_batch : params.n_ubatch);
35
36  cparams.op_offload = params.op_offload;
37
38  ...
39
40  #* Potential Software Bugs/Issues Causing Degradation*
41  #*# 1. Suspicious Pattern: Response Time + BUT Throughput ***
42  This is the ***biggest red flag***. The top 3 critical functions show:
43  - **Response time increased by 150-210%**
44  - **Throughput increased by 286-289%**
45
46  ***What this suggests:***
47  - Possible Bug: Functions are being called MORE frequently
48  in inefficient contexts or loops
49
50  ***Likely causes:***
51  - **Accidental N2 algorithm**: A loop that was O(n) became O(n2)
52  - **Redundant calls**: Functions being called multiple times unnecessarily
53  - **Missing caching**: Results that should be cached are being recalculated
54
55  ...
56
57  #*# 2. Critical Function #1: Tree Iterator (+218% slower)*
```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

ans-loci/llama.cpp-upstream-PR16490-branch_ggml-org-gg-graph-mamba-reuse .py --pull-url https://github.com/auroralabs/llama.cpp-upstream-PR16490-branch_ggml-org-gg-graph-mamba-reuse@454ab90

Ln 1, Col 1 Spaces: 4 UTF-8 LF () C++ Signed out

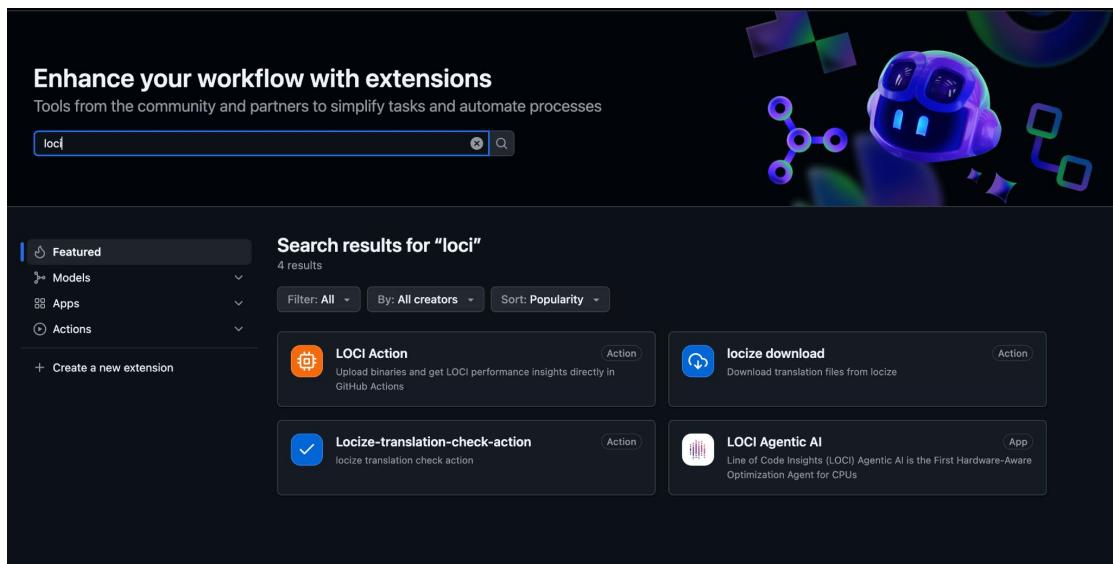
Real ROI in 2-3 Weeks

Engineering teams using LOCI consistently report measurable improvements within 2-3 weeks:

- 20%-40% reduction in test cycle time due to earlier detection
- Fewer regressions reaching integration and system-level tests
- Faster code reviews with factual, hardware-aware insights
- Stronger predictability in release cadence
- Reduced engineering cost through minimized rework

LOCI reduces the costliest factor in modern engineering: the late discovery of regressions.

Developers can see LOCI live today on open-source workflows such as OpenSSL, where LOCI interacts directly with contributors.



LOCI | Software Engineering Intelligence for Embedded Software Engineers.

Catch regressions before testing. Build with confidence

Next Steps:

1. See LOCI in action in [GitHub](#)
2. [Book a discovery meeting with a demo >>](#)
3. Contact us at info@auroralabs.com or visit www.auroralabs.com