

Intelligent Automation for Scientific Workflows

Scientific facilities like SLAC face a key challenge: as experimental capabilities grow more sophisticated, the complexity of data analysis and decision-making increases dramatically. Modern AI systems, particularly Large Language Models (LLMs), offer promising capabilities for automation and assistance but face fundamental limitations when applied to demanding scientific workflows.

Here I'll try to develop the means of combination and abstraction that might be useful for composing LLMs interactions into more complex workflows. I will argue that the LLM should be treated as a basic component of a larger system that can, if built in a certain way, be much more capable than the underlying LLM.

0.1 Example: LLM-Assisted Analysis at XPP

In the last year, I developed an automated analysis pipeline for LCLS's X-ray Pump-Probe (XPP) instrument, working with my PI (Apurva Mehta) and the LCLS analytics group. The pipeline (figure 1) finds CDW signals through a contrast-enhancing transformation of the raw data and uses statistical criteria to maximize signal to noise with respect to the analysis parameters.

This automation approach reduces the user's responsibility to a single numerical input: a region of interest (ROI) on the detector. While this usage is simple, interpreting diagnostic outputs when things go wrong still requires domain expertise. To address this limitation, I experimented with a language model agent to guide users through the diagnostics. In one test, I initialized the analysis with incorrect ROI coordinates. The agent reviewed logged diagnostics and successfully identified the issue:

User: The pump-probe curve for the current run looks weird. What might be wrong?

Agent: Based on the p-value uniformity check results, the designated background ROI [80, 95, 5, 45] appears to contain signal when it should not.

This is evidenced by:

- 1 Mean p-value of 0.394 (should be 0.5)
- 2 Highly significant KS test (p=6.496e-17)

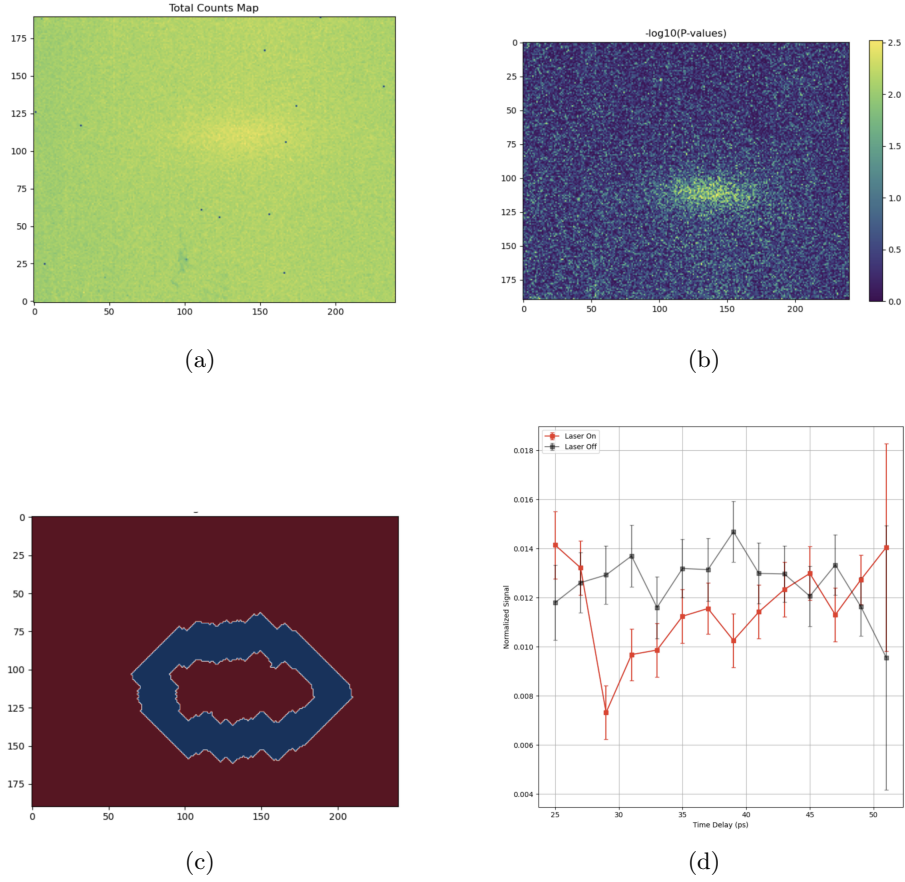


Figure 1: In low signal-to-background settings background estimation is non-trivial, as is identifying and integrating diffraction peaks (a). Our automated workflow enhances contrast using spectral distribution statistics (b), identifies a bounding mask for the signal ((c), interior red region), estimates background levels, and finally generates a pump-probe delay curve (d). Automatic optimization of the workflow parameters yields better signal to noise than the traditional (manual) analysis.

3 Large z-score deviation (9.0)

This contaminated background ROI would affect the null distribution calculation, which in turn would distort the pump-probe curve. The solution would be to select a different background ROI location that shows proper uniformity in its p-values.

While effective for modest diagnostic tasks, the system revealed fundamental limitations in more difficult situations. Attempting to use it for analyzing large collections of log files or guiding multi-step decision processes exposed a well-known underlying limitation of LLMs: difficulty maintaining analytical context across multiple interaction turns, especially when handling large datasets. This limitation is exacerbated by how LLMs perform reasoning - techniques like reflection and chain of thought show that LLM ‘thought’ is entangled with the process of token generation, and therefore each step of analysis consumes valuable context window capacity.

1 Technical Innovation: A Language-Model Architecture for Scientific Computing

The underlying challenge is that we lack principled methods for constructing natural language agent systems in scientific environments. Conventional automation at light sources depends on robust software design practices – similarly, we will need principled architectures for composing natural language agents into scientific workflows with enough scale, capability and error-tolerance.

While the last two years have seen rapid growth in LLM-based agent frameworks like LangChain and AutoGen, these systems make premature abstractions and hide too much from the user. Also rather than compensating for the limitations of LLMs, these agentic approaches either reproduce them or address them in a way that entails questionable tradeoffs. A case in point is reliance on techniques such as retrieval-augmented generation (RAG) that overcome context window size limitations at the expense of recall and accuracy.

I suggest an alternative that takes inspiration from two fundamental concepts in programming language implementation: the metacircular evaluator pattern first developed for Lisp, and staged compilation techniques used in modern compilers.

The main insight is treating natural language interaction with an LLM as a form of code execution. Suppose that we first ask the LLM to translate natural language instructions into programs written in a DSL (domain specific language). The structure of such a DSL program will represent the decomposition of a complex prompt into multiple tasks. The interpretation of the same program involves distribution of each component task to an LLM instance and the linking together of instances through a calling convention and shared memory system.

More concretely:

Algorithm 1: System Overview

Natural Language Query;
↓ [LLM Translation];
XML Task Structure (equivalent to S-expression);
↓ [Parser];
Abstract Syntax Tree;
↓ [tree traversal];
LLM execution;

In this schema, natural language queries are first translated into composite expressions made up of smaller units (atomic tasks) with the purpose to make the execution tractable while preserving semantics of the original query. A satisfying feature of the setup is that it will be self-hosting in the sense that the LLM evaluates DSL procedures generated by the LLM.

An equivalent perspective is that the framework will dynamically compile the user’s prompt into a directed acyclic graph (DAG). Each node of the DAG is dispatched to a separate LLM session, and data travels down and up the nodes of the graph in the form of environment frames and return values, respectively.

The architecture consists of three main components:

1.1 Execution Model

Two mutually recursive procedures, `eval` and `apply`, work with each other to evaluate DSL expressions:

```
1 ; Evaluates a task in the given environment; returns direct result or ↵
   decomposed tasks
2 (define (eval-task task env)
3   (cond
4     ; For atomic tasks, try direct application; if it fails, decompose
5     ((atomic? task)
6      ; amb tries the first option; if it fails, backtracks to the second
7      (amb (apply-proc task '() env) ; Direct application
8            (eval-task (decompose task) env))) ; Or decomposition
9
10    ; For compound tasks: evaluate arguments, apply procedure
11    (else
12     (let ((proc (task-proc task))
13           (args (map (lambda (arg) (eval-task arg env))
14                      (task-args task))))
15       (apply-proc proc args env))))
16
17 ; Applies a procedure to evaluated arguments in given environment
18 (define (apply-proc proc args env)
19   (cond
20     ; For primitives, try direct execution; if it fails, decompose and retry
21     ((primitive? proc)
22      ; amb tries the first option; if it fails, backtracks to the second
23      (amb (execute-llm proc args env) ; Direct execution
24            (eval-task (decompose proc args) env))) ; Or decomposition
25
26     ; For compound procedures: create new environment, evaluate body
27     (else
28      (let ((new-env (extend-environment proc args env)))
29        (eval-task (procedure-body proc) new-env))))
```

Listing 1: Scheme sketch of the evaluation procedure

When task execution ends with an error (e.g. context window overrun, output validation failure), the executor can retry by generating and evaluating an alternate procedure – for example, a decomposition of the task into multiple subtasks.

The creation of the execution data context `'env'` is mediated by the memory subsystem.

1.2 Associative memory system

The memory system explicitly separates storage and working contexts through a hierarchical design:

- Long-term memory for data and procedures
- Working memory for active computations

- Context frames that capture execution environments (including working memory)

Working memory is instantiated from long-term storage using an associative retrieval mechanism that is itself an (atomic) LLM procedure whose purpose is to match an atomic task to a contextually relevant subset of the data in long-term memory.

1.3 Task Expression Framework

The expression system supports nested procedures and basic functional patterns:

1	AtomicTask	-> Direct LLM execution
2	NestedTask	-> Compositional evaluation
3	MapExpression	-> Parallel processing
4	ReduceExpression	-> Result combination

Listing 2: Task Expression Types

These expressions, which can be extended, provide formal semantics for the DSL.

2 Implementation Plan

The implementation strategy builds on recent work and collaborations at SLAC. Working with LCLS beamline scientist Lingjia Liu and Frederic Poitevin from the LCLS analytics group, I developed the previously mentioned analysis approach for charge density wave dynamics in pump-probe experiments. The project aligns with broader LCLS initiatives, led by Jana Thayer and others, to develop real-time analysis capabilities for experimental beamlines.

Building on these experiences and collaborations, the implementation will proceed in two phases:

First, we will develop the core architectural components: the memory system for managing analysis contexts, the execution model for task decomposition, and initial task libraries. These libraries will include specialized agents for software architecture, code generation, analysis refinement, and experimental log interpretation. We'll work with the LCLS analytics group to ensure the framework complements their efforts, and with beamline scientists to get feedback from the end-user point of view.

Second, we will focus on development across scientific workflows based on reusable task patterns that combine automated processing with human-in-the-loop guidance. The framework will support diverse needs, from real-time experiment optimization to offline analysis and documentation. The goals will include accelerated analysis turnaround and reduced downtime during beamtimes.

3 PL concepts and concrete examples

3.1 Compilation

Staged compilation traditionally refers to breaking down compilation into distinct phases, where each stage transforms the program into a new representation closer to the target execution form:

Source code \rightarrow Parse tree \rightarrow AST \rightarrow Intermediate code \rightarrow Machine code

Our system uses the LLM to parse natural language into structured expressions:

*TODO this is incomplete. There should be a cycle connecting 'Task executable' back to 'AST', to represent dynamic / incremental reparsing. See also: the other TODOs, *dynamic reparsing*;*

Source code (English) \rightarrow Parse tree (XML) \rightarrow AST (Python) \rightarrow Task data + executable (XML)

In Python the first three steps (Source code (English) \rightarrow Parse tree (XML) \rightarrow AST (Python)) are orchestrated by the class Compiler. The remaining portion is generated by the interaction between Compile and the evaluation loop (see Evaluator, next section):

```
1
2 @dataclass
3 class Operator:
4     """Represents an operation to be performed"""
5     type: str # "atomic" or "compound"
6     task: str # The actual task description
7     params: Optional[Dict] = None # Additional parameters if needed
8
9
10 @dataclass
11 class ASTNode:
12     """Represents a node in the Abstract Syntax Tree"""
13     operator: Any
14     args: List['ASTNode']
15
16 @dataclass
17 class Operator:
18     """Represents an operation to be performed"""
19     type: str # Maps to TaskType values
20     task: str # Task description
21     params: Optional[Dict[str, Any]] = None
22
23 class TaskSystem(Protocol):
24     """Interface for task management and prompt generation"""
25
26     def get_decomposition_prompt(self, task: str, resource: ResourceType) <->
27         -> str:
28         """Generate prompt for decomposing a task that exceeded resources
```

```

28
29     Args:
30         task: The task description that failed
31         resource: Which resource was exhausted
32
33     Returns:
34         Prompt for LLM to generate decomposed task structure
35         """
36         ...
37
38     def get_alternative_prompt(self, task: str, failure_reason: str) -> str:
39         """Generate prompt for alternative approach after task failure
40
41     Args:
42         task: The task description that failed
43         failure_reason: Why the task failed
44
45     Returns:
46         Prompt for LLM to generate alternative approach
47         """
48         ...
49
50     class TaskType(Enum):
51         """Types of tasks in decomposition"""
52         ATOMIC = "atomic" # Direct LLM execution
53         MAP = "map" # Parallel subtasks
54         REDUCE = "reduce" # Combine results
55         SEQUENCE = "sequence" # Sequential steps
56
57     @dataclass
58     class TaskStructure:
59         """Parsed representation of a task from XML"""
60         type: TaskType
61         description: str
62         subtasks: list['TaskStructure'] = None
63         parameters: Optional[Dict[str, str]] = None
64
65
66     class Compiler:
67         def __init__(self, task_system: TaskSystem):
68             self.task_system = task_system
69             self.parser = XMLTaskParser()
70
71         def compile(self, query: str) -> ASTNode:
72             """Initial compilation of natural language to AST
73
74     Args:
75         query: Natural language query from user
76
77     Returns:
78         ASTNode: Root node of compiled AST
79         """

```



```

80     # Generate structured task description
81     xml = self.llm_translate(query)
82
83     # Parse into intermediate task structure
84     task_structure = self.parser.parse_task(xml)
85
86     # Convert to AST
87     return self.parser.task_to_ast(task_structure)
88
89     def reparse(self,
90                 failed_task: str,
91                 error: ExecutionError,
92                 env: Environment) -> ASTNode:
93         """Generate new AST nodes when execution fails
94
95         Creates new task decomposition based on error type and context.
96
97         Args:
98             failed_task: Description of task that failed
99             error: Details about the failure
100            env: Current execution environment
101
102            Returns:
103                ASTNode: New AST subtree to try executing
104
105            Raises:
106                ExecutionError: If reparsing itself fails
107
108            """
109            # TODO the reparsing prompt shouldn't be exposed to the compiler.
110            # Instead it should construct a specialized Operator for reparsing
111            # And then execute it via Evaluator.execute_llm() (or whatever we're
112            # calling the method that has that responsibility
113
114            # Get specialized prompt for reparsing
115            reparse_prompt = self._get_reparse_prompt(failed_task, error)
116
117            # Generate new task structure via LLM
118            xml = self.llm_translate(reparse_prompt)
119
120            # Parse into intermediate structure
121            task_structure = self.parser.parse_task(xml)
122
123            # Convert to AST
124            node = self.parser.task_to_ast(task_structure)
125
126            # Validate before returning
127            self._validate_reparse(node, error)
128
129            return node
130
131     def _get_reparse_prompt(self, task: str, error: ExecutionError) -> str:
132         """Get appropriate reparse prompt based on error type"""

```

```

132     # TODO the reparsing prompt shouldn't be exposed to the compiler.
133     # Instead it should construct a specialized Operator for reparsing.
134     # If we really need to track the error type, then it should be an
135     # attribute of that specialized Operator
136     if error.type == ErrorType.RESOURCE_EXHAUSTED:
137         return self.task_system.get_decomposition_prompt(
138             task,
139             error.details["resource"]
140         )
141     elif error.type == ErrorType.TASK_FAILED:
142         return self.task_system.get_alternative_prompt(
143             task,
144             error.details["reason"]
145         )
146     else:
147         raise ExecutionError.invalid_reparse(
148             task=task,
149             validation_error="Unsupported error type for reparsing"
150         )
151
152     def _validate_reparse(self, node: ASTNode, error: ExecutionError) -> ↔
153     None:
154         """Validate reparsed AST is acceptable
155
156         TODO: Implement validation logic:
157         - Check structure is valid
158         - Verify resources are bounded
159         - Ensure termination properties
160         """
161         pass
162
163     def llm_translate(self, prompt: str) -> Element:
164         """Translate natural language to XML via LLM
165
166         TODO: Implement LLM interaction:
167         - Format prompt with XML expectations
168         - Handle LLM response parsing
169         - Validate XML structure
170         """
171         pass

```

Example task structure in XML:

```

1 <task>
2   <description>analyze peak patterns across detectors</description>
3   <inputs>
4     <input name="detector1_data">
5       <task>
6         <description>load and preprocess detector 1 data</description>
7         <expected_output>
8           Preprocessed detector 1 data in standard format:

```

```

9         - Intensity values
10        - Peak positions
11        - Background levels
12    </expected_output>
13    </task>
14 </input>
15 <input name="detector2_data">
16     <task>
17         <description>load and preprocess detector 2 data</description>
18         <expected_output>
19             Preprocessed detector 2 data in standard format:
20             - Intensity values
21             - Peak positions
22             - Background levels
23         </expected_output>
24     </task>
25 </input>
26 </inputs>
27 <expected_output>
28     Comparative peak analysis:
29     - Peak correlations between detectors
30     - Intensity pattern matching
31     - Anomaly detection
32 </expected_output>
33 </task>
34 <!--
35 # TODO the above xml is just an example, but we need to clearly define a ↵
36 # mapping between the
37 # structure of xml generated by the llm (when it does <reparsing>) and ↵
38 # AST subtrees. (We could
39 # potentially simplify this by constraining the llm <reparsing> process ↵
40 # so that it can only generate
41 # one ASTNode at a time. In the case of composite nodes, this would mean ↵
42 # generating xml for an outer
43 # task (e.g. a reduction function) and inner task(s) (e.g. the individual ↵
44 # tasks whose outputs the
45 # reduction is operating on)
46 -->

```

In summary:

1. The XML stage lets the LLM express task composition and input / output conventions in a structured way
2. XML provides the interface between LLM and the evaluator
3. Every AST node follows a uniform structure representing operations and their arguments
4. Composite task behavior emerges from AST semantics and atomic task behavior

5. Atomic task behavior emerge from natural-language operator definitions, which are *not* exposed at this level
6. Environments handle local variable bindings and global working memory for LLM operations (see next section)
7. New atomic task patterns can be added without changing the evaluator or execution system

3.2 Environment

An environment represents the complete context needed to evaluate expressions. In traditional programming languages, this mostly means lexical scope—the set of variables accessible from inside a given stack frame.

Our environments support traditional variable scoping while also managing the short-term memory component of the LLM execution context. The basic aspects are:

```

1 class Environment:
2     def __init__(self):
3         """Environment holds bindings and LLM working context"""
4         self.bindings = {}    # Current variable bindings
5         self.context = {}    # LLM working memory
6
7     def extend(self, names, values):
8         """Create new environment with additional bindings"""
9         new_env = Environment()
10        new_env.bindings = dict(zip(names, values))
11        # The full implementation will update the context using associative
12        # matching via the long term - short term memory system instead of
13        # just cloning it
14        new_env.context = self.context.copy()
15        return new_env
16
17    def lookup(self, name):
18        """Look up value in current bindings"""
19        return self.bindings.get(name)

```

This is sufficient for us to:

1. Track context through nested evaluations
2. Pass relevant state between task executions
3. Allow the evaluator to pass around execution contexts and create new ones using the associative memory procedure

3.3 Metacircular Evaluator

A metacircular evaluator is an interpreter implemented using similar fundamental operations to the ones it aims to interpret. In our context, we implement a domain-specific language (DSL) evaluator using LLM operations as one basic component of the evaluation machinery, and this evaluator in turn coordinates and executes higher-level LLM tasks.

The architecture has two key aspects. First, the environment must be a first-class data structure that can be explicitly introspected by both the evaluator and the LLM. The environment captures not just variable bindings (as in a conventional programming language implementation) but the complete context needed for task interpretation and decomposition. (In contrast, in a traditional language implementation, the execution environment is entangled with the parsing and code generation process in a rather opaque way.)

The second aspect is the self-hosting property mentioned above. The LLM provides the evaluator with primitive capabilities for generating structured output from natural language (as XML task descriptions), and for executing atomic tasks. The evaluator combines these primitive operations to implement higher-level functionality: managing execution environments, parsing and dispatching structured task descriptions, and collecting execution outputs.

Here’s the core evaluator pseudo-implemented in Python:

```
1 from dataclasses import dataclass
2 from typing import List, Any, Optional
3 from enum import Enum
4
5 # Shared type definitions
6 class OperatorType(Enum):
7     ATOMIC = "atomic" # Direct LLM execution
8     MAP = "map" # Process multiple inputs
9     REDUCE = "reduce" # Combine results
10    SEQUENCE = "sequence" # Execute in order
11
12 # TODO only ATOMIC operators should have (or populate) the task attribute,
13 # because the other operator types aren't directly llm-executable.
14 # Note an asymmetry between llm parsing and llm execution: when parsing,
15 # the llm is locally aware of the AST structure, but when executing the
16 # llm only sees one atomic node at a time. <reparsing>
17 @dataclass
18 class Operator:
19     type: OperatorType
20     task: str # Task description/prompt
21     params: Dict = None # Optional parameters
22
23 class Evaluator:
24     def eval(self, node: Node, env: 'Environment') -> Any:
25         """Evaluate a node in the given environment"""
26         # TODO <dynamic reparsing> if LLM execution fails (typically bc ↔
27         # of resource exhaustion or verification failure)
```

```

27     # we need to recover by reparsing the AST node into one or more ↵
    rewritten subtrees
28     # (i.e. either a new atomic expression or a decomposed, compound ↵
    expression). See the
29     # Scheme description to get the general idea. also see <errors>
30     # </dynamic reparsing>
31     if self.is_atomic(node.operator):
32         return self.execute_llm(node.operator, env)
33
34     # For compound expressions, evaluate args then apply
35     evaluated_args = [self.eval(arg, env) for arg in node.args]
36     return self.apply(node.operator, evaluated_args, env)
37
38     def is_atomic(self, node: ASTNode) -> bool:
39         """Check if node represents direct LLM execution"""
40         return len(node.args) == 0 and node.operator.type == "atomic"
41
42     def execute_llm(self, operator: Any, env: 'Environment') -> Any:
43         """Execute atomic task with LLM"""
44         return self.llm_execute(operator.task, env)
45
46     def apply(self, operator: Any, args: List[Any], env: 'Environment') ↵
    -> Any:
47         """Apply compound operator to evaluated arguments"""
48         new_env = env.extend(operator.params, args)
49         return self.eval(operator.body, new_env)

```

Note that the dynamic reparsing is spread out as an interaction between Evaluator and Compiler, but it is conceptually equivalent to this simple expression from the Scheme version:

```

1  ...
2      ; For atomic tasks, try direct application; if it fails, decompose
3      ((atomic? task)
4       ; amb tries the first option; if it fails, backtracks to the second
5       (amb (apply-proc task '()) env)                ; Direct application
6            (eval-task (decompose task) env)))         ; Or decomposition
7  ...

```

Listing 3: nondeterministic evaluation using the amb (ambiguous) operator

3.4 End-to-End Example

Let's examine how the system handles a user request that benefits from natural language understanding:

'Review our XRD analysis and check if we chose a good background region.'

The LLM first translates this into nested tasks:

```

1  <task>

```

```

2   <description>analyze background region quality</description>
3   <inputs>
4     <input name="region_stats">
5       <task>
6         <description>extract statistics from experiment logs</description>
7         <parameters>
8           Extract and analyze:
9           - Background region coordinates
10          - Statistical test p-values
11          - Distribution uniformity metrics
12        </parameters>
13        <expected_output>
14          Structured statistics including:
15          - ROI coordinates
16          - P-value series
17          - Distribution metrics
18        </expected_output>
19      </task>
20    </input>
21  </inputs>
22  <expected_output>
23    Quality assessment including:
24    - Statistical validity evaluation
25    - Potential signal contamination check
26    - Recommendations for improvement if needed
27  </expected_output>
28 </task>

```

The compiler builds an AST:

```

1  node = ASTNode(
2    operator=Operator(
3      type="atomic",
4      task="analyze_region_quality"
5    ),
6    args=[
7      ASTNode(
8        operator=Operator(
9          type="atomic",
10         task="extract_stats",
11         params={"instruction": "Find background..."}
12       ),
13       args=[]
14     )
15   ]
16 )

```

The evaluator processes this with environment handling:

```

1  # Initialize environment with both logs and documentation

```

```

2 env = Environment(context={
3   "log_contents": ""
4   2024-04-06 10:15:32 INFO: Starting analysis of run 123
5   2024-04-06 10:15:33 DEBUG: Background ROI set to [80,95,5,45]
6   2024-04-06 10:15:34 DEBUG: Method: local linear
7   2024-04-06 10:15:35 WARNING: High fit residuals
8   2024-04-06 10:15:36 DEBUG: P-values: [0.394, 0.412, 0.378]
9   ...""
10
11  "analysis_docs": ""
12  Background Region Quality Assessment Guide:
13  - P-values should follow uniform distribution (mean approximately 0.5)
14  - KS test should show p > 0.05
15  - Region should be at least 20 pixels from any peak
16  - Common failure modes:
17    * Signal contamination causes p-value clustering
18    * Edge effects near beam stop distort background
19  ...""
20 })
21
22 # Evaluate full expression
23 result, final_env = evaluator.eval(node, env)
24
25 # The evaluation proceeds:
26 # 1. Inner extract_stats task receives filtered environment:
27 #   env.context = {
28 #     "log_contents": "..." # Only the log data needed for extraction
29 #   }
30 #   -> Returns structured data like:
31 #     {"roi": [80,95,5,45],
32 #      "pvalues": [0.394, 0.412, 0.378]}
33
34 # 2. Outer analyze_region_quality task receives full context:
35 #   env.context = {
36 #     "analysis_docs": "...", # Documentation for analysis guidance
37 #     "extraction_results": {"roi": [80,95,5,45], ...}
38 #   }
39 #   -> Uses docs to guide analysis:
40 #     "Background region shows signs of signal contamination.
41 #     P-value clustering (mean=0.394) matches known failure mode
42 #     described in analysis guide."

```

In this example, each task gets minimal context needed for its operation and the outer tasks can access both the log-parsing results (as direct input) and reference documentation (as context / short-term memory).