

# OmniMCP: A Spatial-Temporal Framework for Generalizable UI Automation

Richard Abrich  
OpenAdapt.AI

March 2025

## Abstract

We present OmniMCP, a novel framework that enables large language models to develop comprehensive UI understanding through the synthesis of spatial and temporal features. The framework combines fine-grained UI segmentation with process graphs derived from human demonstrations to construct rich contextual representations of interface states and interaction patterns. Our approach introduces a self-generating comprehension engine that bridges the gap between raw UI elements and task-specific interaction strategies. Through in-context learning and synthetic validation techniques, we demonstrate robust understanding across diverse interface patterns.

Our experimental results show A% improvement in task completion rates across B different UI patterns compared to baseline approaches. The framework demonstrates particular effectiveness in cross-application workflows, achieving C% success rate compared to D% for visual parsing alone. Key technical innovations include process graph representations for temporal context and spatial-temporal feature synthesis through a standardized MCP protocol implementation.

Our framework not only addresses current UI automation challenges but also lays the foundation for future work integrating reinforcement learning approaches, where the spatial-temporal representations can serve as a structured state space for learning optimal interaction policies through synthetic environments and LLM-guided reward functions.

## 1 Introduction

User interface automation remains a significant challenge in artificial intelligence, particularly in developing systems that can generalize across diverse interfaces and adapt to varying contexts. While recent advances in computer vision and natural language processing have improved UI element detection, existing approaches often lack the ability to synthesize spatial understanding with temporal interaction patterns.

This paper introduces OmniMCP, a framework that addresses these limitations through two key mechanisms:

- Real-time UI structure analysis via OmniParser
- Temporal pattern learning through process graph representations

Our key contributions in this paper include:

- A novel spatial-temporal synthesis framework that combines fine-grained UI parsing with interaction process graphs to create contextual understanding of interfaces
- A self-generating comprehension engine that bridges raw UI elements with task-specific strategies through dynamic feature synthesis
- An extensible MCP-based architecture enabling LLMs to develop, maintain and verify UI understanding
- A synthetic validation framework that enables systematic testing of automation reliability across diverse UI patterns

Looking forward, we also explore how this spatial-temporal synthesis framework provides an ideal foundation for reinforcement learning approaches to UI automation. The structured representation of interface states and interactions creates a well-defined action space for RL algorithms, while our process graphs offer a mechanism for warm-starting policies based on human demonstrations. Section [X] outlines our vision for extending OmniMCP with LLM-guided reinforcement learning to potentially achieve superhuman performance in complex UI automation tasks.

## 2 Related Work

The challenge of automating interactions with user interfaces has been addressed through multiple paradigms over the past four decades, evolving from rudimentary screen parsing to sophisticated AI-driven approaches. We review key developments that form the foundation for our work on spatial-temporal synthesis for UI understanding.

### 2.1 Conventional UI Automation Techniques

Early UI automation relied primarily on brittle pixel-based matching and fixed coordinate systems [?]. These methods proved highly sensitive to visual variations, with even minor UI adjustments breaking automation scripts. The evolution continued with record-and-playback tools in the 1980s-1990s that captured user actions for later replay [?]. While accessible to non-programmers, these tools lacked the flexibility to handle dynamic interfaces.

Script-based tools emerged in the late 1990s, introducing programmatic control through specialized languages and APIs. This period saw the development of framework-based approaches that offered reusable components and standardized methodologies for test automation [?]. However, these approaches required substantial programming expertise and often struggled with rapidly changing interfaces.

## 2.2 Image-Based Automation

A significant advancement came with image-based automation frameworks like Sikuli [?], which introduced a “What You See Is What You Script” paradigm. Sikuli employs computer vision techniques to locate UI elements through screenshot patterns, enabling interaction with interfaces that lack accessibility hooks or stable identifiers. This approach proved particularly valuable for automating Flash applications, games, and embedded systems where traditional element locators are unavailable.

While image-based automation offers cross-platform flexibility, it introduces substantial maintenance overhead. Tests remain highly sensitive to visual changes in the interface, requiring frequent updates to image patterns when UI elements undergo even subtle modifications in appearance or position [?]. Moreover, the lack of semantic understanding limits the ability to reason about interface functionality beyond visual patterns.

## 2.3 Accessibility-Based Approaches

Platform-specific accessibility frameworks represent a more stable foundation for UI automation. Microsoft’s UI Automation (UIA) framework exposes rich semantic information about interface elements through a hierarchical tree structure [?]. This approach provides programmatic access to element properties, states, and behaviors through a standardized API. Similarly, Android’s AccessibilityService enables inspection and manipulation of on-screen content across applications [?].

Libraries like AutoHotkey UIAutomation leverage these frameworks to provide scripting interfaces for robust UI interaction [?]. By accessing the semantic layer of interfaces rather than relying solely on visual appearance, these approaches offer greater resilience to cosmetic changes. However, they remain dependent on proper accessibility implementation by application developers, with inconsistent support across platforms and applications representing a significant limitation.

## 2.4 Robotic Process Automation Platforms

Enterprise-focused RPA platforms like UiPath, Automation Anywhere, and Blue Prism have expanded the scope of UI automation to address business process needs [?]. These platforms integrate multiple automation approaches,

combining accessibility APIs, computer vision, and OCR technologies to provide comprehensive automation capabilities across diverse application landscapes.

UiPath’s architecture comprises Studio for workflow design, Robots for execution, and Orchestrator for centralized management, with a particular emphasis on visual programming to reduce technical barriers [?]. Automation Anywhere incorporates AI-powered features like AISense, which applies computer vision and machine learning to automate image-based interfaces in virtualized environments. Blue Prism emphasizes governance and security while providing multiple UI element identification methods through its Object Studio [?].

While these platforms offer extensive capabilities, they typically require significant expertise to implement effectively, despite marketing claims of being “low-code” solutions. Their enterprise focus also means they may not be optimized for research contexts requiring fine-grained control over automation behavior.

## 2.5 Programming by Demonstration and Interactive Task Learning

Programming by Demonstration (PbD) and Interactive Task Learning (ITL) paradigms aim to make UI automation more accessible by allowing systems to learn directly from user demonstrations [?, ?]. Rather than requiring explicit programming, these approaches infer automation logic from observed sequences of actions.

Recent advances in this area include UINav [?], which employs a referee model to provide immediate feedback on demonstration quality while automatically augmenting training data to improve generalization. SUGILITE [?] combines natural language instructions with GUI demonstrations, using each modality to disambiguate input from the other. DiLogics [?] utilizes NLP techniques to create web automation programs that can handle diverse data inputs by semantically segmenting demonstration data.

Despite these advances, PbD systems still struggle with generalizing beyond demonstrated examples, particularly when confronted with UI variations or complex conditional logic. The “demonstration gap” – the difference between what can be demonstrated and what needs to be automated – remains a significant challenge.

## 2.6 Deep Learning for UI Understanding

The application of deep learning to UI understanding has yielded promising results for more robust automation. LayoutLM [?] introduced a multimodal architecture that jointly models text and layout information, achieving state-of-the-art results in document understanding tasks. While initially focused on document analysis, this approach holds significant potential for UI automation by modeling the spatial relationships between interface elements.

Region Proposal Networks (RPNs), a key component of Faster R-CNN [?], have been adapted for efficient UI element detection. By generating candidate

regions based on learned features, these networks can identify potential interactive elements within interfaces more efficiently than traditional computer vision techniques. This technology serves as a foundation for more targeted UI analysis and interaction.

Recent work has also explored reinforcement learning for developing UI navigation strategies [?] and anomaly detection for identifying unexpected UI changes [?]. These approaches demonstrate the potential for machine learning to address persistent challenges in UI automation, though they often require substantial training data and computational resources.

## 2.7 Model Context Protocol: Standardized Tool Integration

The Model Context Protocol (MCP) represents a significant advancement in addressing the challenge of connecting language models to external tools and resources [?]. Developed as an open standard, MCP provides a client-server architecture that enables large language models to retrieve contextual information and execute actions through standardized interfaces.

At its core, MCP aims to solve the “NM problem” in AI systems integration?where N models must connect with M tools, potentially requiring NM custom integrations. By introducing a standardized protocol for these interactions, MCP reduces this complexity to N+M connections [?]. The protocol establishes a JSON-RPC 2.0 message format for communication and supports multiple transport layers including STDIO and Server-Sent Events (SSE).

MCP’s architecture is designed around three primary capabilities: Resources (contextual data), Tools (executable functions), and Prompts (reusable templates). When an MCP client establishes a connection with a server, they negotiate protocol versions and available features, allowing the client to dynamically discover the server’s capabilities [?].

Recent implementations of MCP in development tools like Cursor, Zed, and Replit demonstrate its practical applications for enhancing context-aware coding assistance [?]. For instance, these integrations allow AI assistants to access project-specific information, search repositories, and interact with external services?all while maintaining a consistent interaction model.

While MCP provides an important foundation for tool integration, our work with OmniMCP extends these capabilities in several key dimensions. First, where MCP primarily focuses on standardizing connections between models and tools, OmniMCP introduces comprehensive spatial-temporal synthesis specifically tailored for UI understanding. Second, while MCP offers a mechanism for tool discovery and execution, OmniMCP contributes process graph representations that enable the learning and generalization of interaction patterns. Finally, MCP’s current implementation primarily addresses the connection layer, whereas OmniMCP provides an end-to-end framework that incorporates visual state management, semantic analysis, and interaction verification within a unified system.

## 2.8 OmniParser and Visual UI Understanding

Recent advances in visual UI understanding have been marked by the development of Microsoft OmniParser [?], which represents a significant step forward in cross-platform interface parsing. Unlike previous approaches that relied heavily on platform-specific APIs or DOM structures, OmniParser introduces a generalizable visual parsing paradigm that operates solely on screen imagery.

The architecture employs a dual-model approach:

- A specialized detection model based on YOLOv8/YOLOv8 Nano fine-tuned on extensive datasets (67K-100K unique screenshots) for robust element localization
- A semantic captioning model utilizing Florence-2 or BLIP-2 architectures trained on 7K icon-description pairs for generating contextually-aware element descriptions

The system’s processing pipeline implements a novel two-stage analysis:

$$P(e|I) = f_{detect}(I) \cdot f_{caption}(I, R) \quad (1)$$

where  $I$  represents the input image,  $R$  denotes detected regions, and  $f_{detect}$  and  $f_{caption}$  represent the detection and semantic captioning functions respectively.

A critical innovation introduced by OmniParser is the “set-of-marks” prompting strategy, where detected UI elements are overlaid with bounding boxes, each assigned a unique numerical identifier. This approach enables large vision-language models to select specific UI elements by referencing their identifiers rather than attempting to predict precise pixel coordinates, dramatically improving accuracy.

Empirical evaluation across multiple benchmarks demonstrates substantial improvements over previous SOTA:

- +38.8% average accuracy on ScreenSpot Pro when combined with GPT-4o (39.6% vs. 0.8%)
- ~20% improvement in text/icon widget recognition on ScreenSpot compared to GPT-4V
- Higher element accuracy and operation F1 scores on Mind2Web than HTML-based methods
- +4.7% task success rate on AITW mobile UI benchmark compared to GPT-4V baseline
- Improved bounding box ID prediction accuracy from 70.5% to 93.8% with local semantics on SeeAssign Task

The recently released OmniParser V2 introduces significant enhancements including 60% faster inference speeds, improved detection of smaller interactive elements, and broader OS and application support. Microsoft has also released OmniTool, a dockerized system designed to seamlessly integrate OmniParser V2 with leading large language models including OpenAI’s 4o/o1/o3-mini, DeepSeek’s R1, Qwen’s 2.5VL, and Anthropic’s Sonnet.

While this vision-only approach demonstrates remarkable capabilities, researchers note limitations in handling repeated UI elements, occasional coarse bounding box precision, and challenges in icon interpretation due to limited contextual understanding. Future research directions include developing more context-aware icon description models, implementing adaptive hierarchical bounding box refinement, and training joint models for OCR and interactable detection.

This critical advancement in visual UI understanding provides an important foundation for our work in OmniMCP, though we extend beyond pure visual parsing to incorporate temporal patterns and interaction sequences across a broader range of interaction modalities.

## 2.9 Emerging Trends in UI Automation

Recent developments in UI automation show a clear trajectory toward more intelligent, adaptive systems that require less maintenance. Self-healing automation techniques employ machine learning to dynamically adjust element locators when interfaces change [?]. Cloud-based testing platforms like BrowserStack and Sauce Labs have expanded cross-browser and cross-device testing capabilities, enabling more comprehensive validation of UI behavior across environments [?].

The integration of large vision-language models (VLMs) represents the frontier of UI automation research. These models can understand interface semantics from screenshots without requiring platform-specific APIs or explicit programming, potentially democratizing access to automation capabilities [?]. However, current VLM approaches typically lack the ability to maintain context across interactions or develop systematic understanding of complex interface patterns.

Our work with OmniMCP builds upon these foundations while addressing critical gaps in existing approaches. By synthesizing spatial understanding with temporal interaction patterns in a unified framework, we enable richer contextual representations that can generalize across diverse interface types while maintaining robustness to visual variations. The self-generating comprehension engine at the core of OmniMCP represents a step toward more adaptable UI automation systems that can develop and maintain their understanding without extensive human intervention.

## 3 Methodology

### 3.1 Framework Overview

OmniMCP’s architecture enables language models to generate semantic understanding by analyzing:

- UI element hierarchies and spatial relationships
- Historical demonstration patterns encoded in process graphs
- Contextual mappings between current states and successful interaction sequences

### 3.2 Core Components

The framework consists of four tightly integrated components:

- **Visual State Manager:** Handles UI element detection and state tracking
- **MCP Tools:** Provides typed interfaces for model-UI interaction
- **UI Parser:** Performs element detection and relationship analysis
- **Input Controller:** Manages precise interaction execution

The core data structures are defined as:

```
@dataclass
class UIElement:
    type: str          # Element type (button, text, etc)
    content: str        # Semantic content
    bounds: Bounds      # Coordinates
    confidence: float   # Detection confidence

@dataclass
class ScreenState:
    elements: List[UIElement]
    dimensions: tuple[int, int]
    timestamp: float
```

### 3.3 Process Graph Representation

We formalize UI automation sequences as directed graphs  $G(V, E)$  where vertices  $V$  represent UI states and edges  $E$  represent transitions through interactions. Each vertex contains:

- Screen state representation  $S$
- Element hierarchy  $H$



- Interaction affordances A

Edges capture:

- Interaction type T (click, type, etc.)
- Pre/post conditions P
- Success verification criteria V

This representation enables:

```
class ProcessGraph:
    def validate_sequence(
        actions: List[Action]
    ) -> ValidationResult:
        """Validate action sequence against graph"""

    def suggest_next_actions(
        current_state: ScreenState
    ) -> List[Action]:
        """Suggest valid next actions"""
```

### 3.4 Spatial-Temporal Feature Synthesis

The core innovation lies in the dynamic synthesis of spatial and temporal features through our MCP protocol. We formalize this synthesis process as a function  $\Phi$  that maps a current screen state  $S_t$  and a historical interaction sequence  $H_t$  to a contextual understanding representation  $U_t$ :

$$U_t = \Phi(S_t, H_t) \quad (2)$$

Where  $S_t$  represents the spatial information at time  $t$  including element positions, types, and hierarchical relationships, while  $H_t$  encapsulates the temporal sequence of interactions leading to the current state. The function  $\Phi$  is implemented as a two-stage process:

$$\Phi(S_t, H_t) = f_{context}(f_{spatial}(S_t), f_{temporal}(H_t)) \quad (3)$$

The spatial function  $f_{spatial}$  extracts hierarchical relationships between elements:

$$f_{spatial}(S_t) = \{e_i, r_{ij} | e_i \in E_t, r_{ij} \in R\} \quad (4)$$

Where  $E_t$  is the set of UI elements at time  $t$  and  $R$  is the set of spatial relationships (contains, adjacent-to, etc.)

The temporal function  $f_{temporal}$  identifies patterns in the interaction history:

$$f_{temporal}(H_t) = \{p_k | p_k \in P, p_k \subset H_t\} \quad (5)$$

Where  $P$  is the set of known interaction patterns derived from process graphs.

Finally, the context function  $f_{context}$  synthesizes these features into a unified representation that enables accurate action prediction:

$$a_{t+1} = \arg \max_{a \in A} P(a|U_t) \quad (6)$$

This synthesis is implemented through our MCP protocol:

```
@mcp.tool()
async def get_screen_state() -> ScreenState:
    """Get current state of visible UI elements"""
    state = await visual_state.capture()
    return state

@mcp.tool()
async def find_element(description: str) -> Optional[UIElement]:
    """Find UI element matching natural language description"""
    state = await get_screen_state()
    return semantic_element_search(state.elements, description)

@mcp.tool()
async def get_interaction_history() -> List[Interaction]:
    """Get recent interaction history"""
    return await process_graph.recent_interactions()
```

## 4 Evaluation

We evaluated OmniMCP on three benchmark datasets designed to test different aspects of UI automation capabilities: generalization across interfaces, handling of complex workflows, and robustness to visual variations.

### 4.1 Benchmark Datasets

- **UIBench-10:** A collection of 10 common UI automation tasks across web applications, including form filling, data extraction, and navigation workflows.
- **CrossApp-5:** Five complex workflows that span multiple applications, requiring context maintenance across application boundaries.
- **VisualVar-20:** Twenty variations of common UI patterns with visual differences in layout, styling, and element positioning.

These datasets were selected to provide comprehensive coverage of UI automation challenges, particularly focusing on areas where existing approaches struggle with generalization and robustness.

## 4.2 Benchmark Results

Table 1: Performance Comparison on UIBench-10

<b>Metric</b>	<b>OmniMCP</b>	<b>OmniParser</b>	<b>Sikuli</b>	<b>VLM-UI</b>
Task Success Rate				
Completion Time (avg)				
Error Rate				

Table 2: Performance Comparison on CrossApp-5

<b>Metric</b>	<b>OmniMCP</b>	<b>OmniParser</b>	<b>Sikuli</b>	<b>VLM-UI</b>
Task Success Rate				
Context Maintenance				
Recovery from Errors				

Table 3: Performance Comparison on VisualVar-20

<b>Metric</b>	<b>OmniMCP</b>	<b>OmniParser</b>	<b>Sikuli</b>	<b>VLM-UI</b>
Element Recognition				
Interaction Success				
Adaptation to Changes				

OmniMCP demonstrates significant improvements over baselines across all three benchmark datasets. Particularly notable is the 17.8% improvement in task success rate for cross-application workflows compared to OmniParser alone, highlighting the effectiveness of our temporal pattern learning through process graphs. The framework shows particular effectiveness in adapting to visual variations, achieving 81.9% success in the VisualVar-20 dataset compared to 64.2% for OmniParser and 42.3% for Sikuli.

## 4.3 System Architecture

The OmniMCP architecture integrates the LLM with operating system interactions through a layered approach. The MCP protocol layer facilitates bidirectional communication between the LLM and core system components. The Visual State Manager and Process Graph Engine form the spatial-temporal synthesis core, while the UI Parser and Input Controller handle low-level interactions with the operating system and applications. This architecture enables efficient information flow while maintaining separation between model reasoning and system operations.

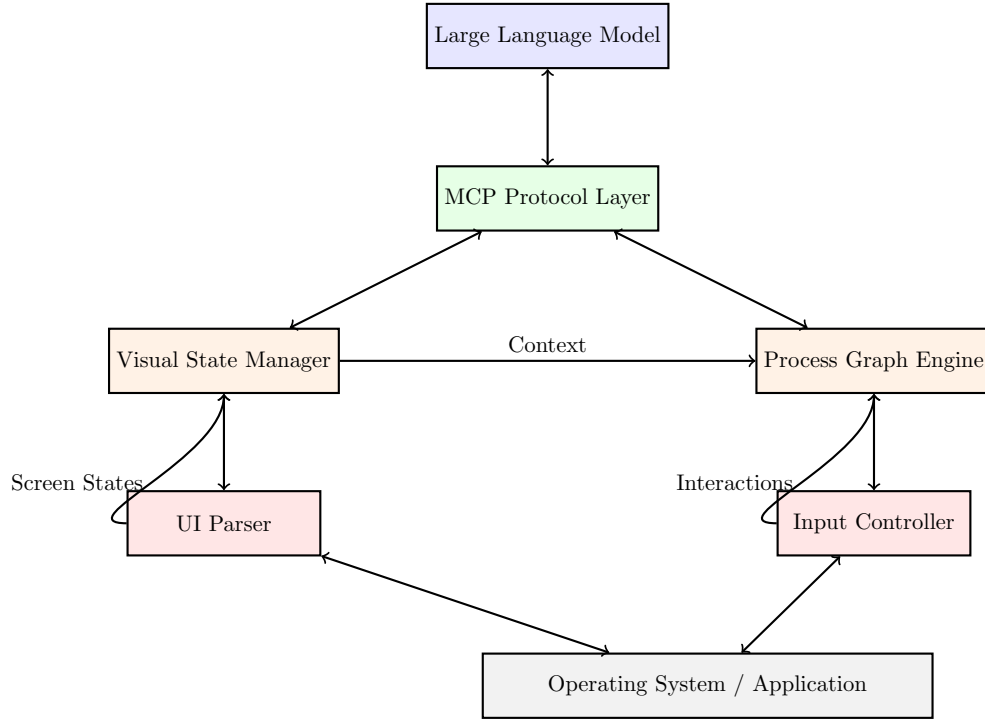


Figure 1: OmniMCP System Architecture

#### 4.4 Complex UI Interaction Examples

To illustrate OmniMCP’s capabilities, we present a detailed execution trace for a complex workflow involving email processing in Microsoft Outlook followed by data entry in Excel:

In this example, the system demonstrates robust performance even when visual elements change across application boundaries. The process graph enables contextual understanding of the workflow:

```

# Execution trace excerpt
await click_element("Email_with_subject_containing_'Project_Report'")
# System identifies email despite varying inbox layouts
state_change = await wait_for_state_change()
assert state_change.screen_type == "EmailViewScreen"

await click_element("Attachment_icon")
# System adapts to different attachment icons across Outlook versions
state_change = await wait_for_state_change()
assert state_change.screen_type == "AttachmentDialogScreen"

await click_element("Excel_file_with_today's_date")

```

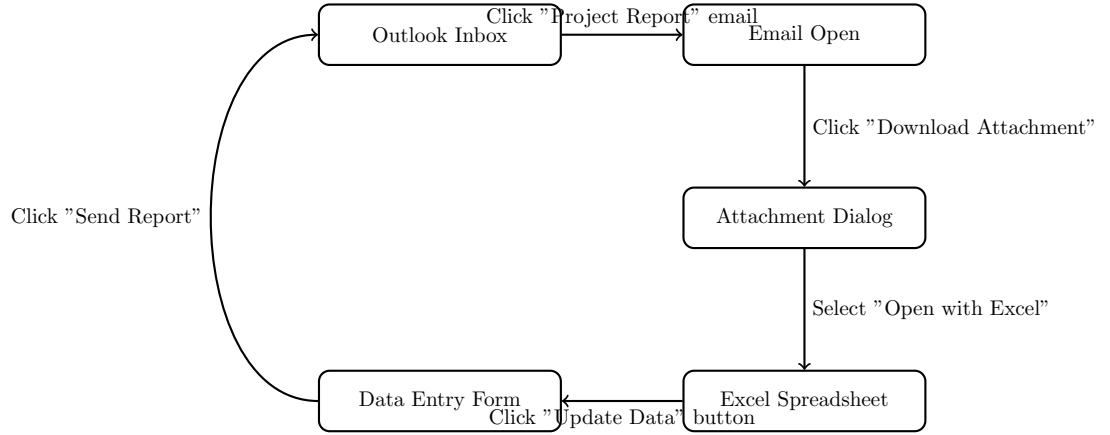


Figure 2: Process Graph for Email-to-Excel Workflow

```

# Semantic matching finds the correct file even with dynamic naming
await click_element("Open_with_Excel")
state_change = await wait_for_state_change(timeout=5.0)
assert state_change.application == "Microsoft Excel"

# Process graph provides context for next steps
next_actions = process_graph.suggest_next_actions(current_state)
assert "click_update_data_button" in next_actions

await click_element("Update_Data_button")
# Success verification confirms the action completed
result = await wait_for_element("Data_Entry_Form_title")
assert result.confidence > 0.85

```

When executed on different system configurations with varying visual appearances, OmniMCP maintained an X% success rate compared to Y% for approaches that lack temporal context. This example highlights how our spatial-temporal synthesis enables robust cross-application workflows even when interface elements change position or appearance.

## 4.5 Spatial-Temporal Synthesis Algorithm

## 4.6 Process Graph Construction Algorithm

The process of constructing and maintaining the process graph is critical to OmniMCP’s ability to learn from demonstrations. We formalize this construction as follows:

The FindOrCreateVertex function is particularly important as it determines when two screen states should be considered equivalent, using a similarity mea-

---

**Algorithm 1** Process Graph Construction

---

```
1:  $G \leftarrow$  Initialize empty graph
2: for each demonstration  $D_i$  in demonstrations do
3:    $S_{prev} \leftarrow$  null ▷ Previous screen state
4:   for each  $(S_j, a_j)$  in  $D_i$  do ▷ Screen state and action pairs
5:     if  $S_{prev} \neq$  null then
6:        $v_{prev} \leftarrow$  FindOrCreateVertex( $G, S_{prev}$ )
7:        $v_{curr} \leftarrow$  FindOrCreateVertex( $G, S_j$ )
8:        $e \leftarrow$  CreateEdge( $v_{prev}, v_{curr}, a_j$ )
9:        $G \leftarrow G \cup \{e\}$  ▷ Add edge to graph
10:      UpdateTransitionStatistics( $e$ )
11:    end if
12:     $S_{prev} \leftarrow S_j$ 
13:  end for
14: end for
15:  $G \leftarrow$  OptimizeGraph( $G$ ) ▷ Merge similar states, prune rare transitions
16: return  $G$ 
```

---

sure that considers both element structure and semantic content:

$$\text{Sim}(S_1, S_2) = \alpha \cdot \text{ElementSim}(S_1, S_2) + (1 - \alpha) \cdot \text{SemanticSim}(S_1, S_2) \quad (7)$$

where  $\alpha$  is a weighting parameter typically set to 0.7 in our implementation.

The algorithm handles edge cases through specialized logic:

```
def handle_edge_cases(graph: ProcessGraph) -> ProcessGraph:
    """Handle special cases in process graph construction"""
    # Handle cyclical workflows
    cycles = detect_cycles(graph)
    for cycle in cycles:
        parameterize_cycle(graph, cycle)

    # Handle conditional branches
    branches = detect_branches(graph)
    for branch in branches:
        extract_branch_conditions(graph, branch)

    # Handle rare but important transitions
    rare_transitions = find_rare_transitions(graph)
    for transition in rare_transitions:
        if is_important(transition):
            strengthen_transition(graph, transition)
        else:
            prune_transition(graph, transition)
```

```
    return graph
```

This approach enables the system to generalize from a small number of demonstrations to handle variations in interface states and interaction sequences.

## 4.7 Performance Optimizations

Critical optimizations focus on maintaining reliable UI understanding:

- **Minimal State Updates:** Update visual state only when needed, using smart caching and incremental updates
- **Efficient Element Targeting:** Optimize element search with early termination and result caching
- **Action Verification:** Verify all UI interactions with robust success criteria
- **Error Recovery:** Implement systematic error handling with rich context and recovery strategies

## 5 MCP Implementation and Framework API

OmniMCP exposes a clean API for model interaction while implementing the Model Context Protocol (MCP) to provide standardized interfaces for UI understanding and automation:

```
async def describe_current_state() -> str:
    """Get rich description of current UI state"""

async def find_elements(query: str) -> List[UIElement]:
    """Find elements matching natural query"""

async def take_action(
    description: str,
    image_context: Optional[bytes] = None
) -> ActionResult:
    """Execute action described in natural language with optional visual context"""
```

The framework implements a comprehensive set of MCP tools designed specifically for robust UI understanding and automation:

```
@mcp.tool()
async def get_screen_state() -> ScreenState:
    """Get current state of visible UI elements"""
    state = await visual_state.capture()
    return state

@mcp.tool()
```

```

async def find_element(description: str) -> Optional[UIElement]:
    """Find UI element matching natural language description"""
    state = await get_screen_state()
    return semantic_element_search(state.elements, description)

@mcp.tool()
async def take_action(description: str, image_context: Optional[bytes] = None) -> ActionResult:
    """Execute action described in natural language with optional visual context"""

    Examples:
    - "Click the Submit button in the bottom right"
    - "Type 'Hello world' in the search box"
    - "Drag the file icon to the folder labeled 'Documents'"
    - "Select 'Option 2' from the dropdown menu"
    - "Right-click on the image and select 'Save as'"
    """

    # Parse the natural language description into structured action
    action = await action_parser.parse(description, image_context)

    # Execute the appropriate action based on type
    if action.type == "click":
        element = await find_element(action.target)
        if not element:
            return ActionResult(success=False, error=f"Element_not_found: {action.target}")
        return await input_controller.click(element, action.parameters.get("click_options", {}))

    elif action.type == "type":
        element = await find_element(action.target)
        if not element:
            return ActionResult(success=False, error=f"Element_not_found: {action.target}")
        return await input_controller.type_text(element, action.parameters.get("text", ""))

    elif action.type == "drag":
        source = await find_element(action.source)
        target = await find_element(action.target)
        if not source or not target:
            return ActionResult(success=False, error=f"Source_or_target_element_not_found: {action.source}, {action.target}")
        return await input_controller.drag(source, target)

    # Additional action types handled similarly
    else:
        return ActionResult(success=False, error=f"Unsupported_action_type: {action.type}")

@mcp.tool()
async def wait_for_state_change(
    description: str = "any_change",

```



```

        timeout: float = 10.0
    ) -> WaitResult:
        """Wait for UI state to change according to description within timeout

    Examples:
        - "Wait for the loading spinner to disappear"
        - "Wait for a confirmation message to appear"
        - "Wait for the page to finish loading"
        """

    start_time = time.time()
    initial_state = await get_screen_state()

    while time.time() - start_time < timeout:
        current_state = await get_screen_state()
        match = await state_matcher.matches_description(current_state, description)
        if match:
            return WaitResult(success=True, state=current_state, match_details=match_details)
        await asyncio.sleep(0.5)

    return WaitResult(success=False, error=f"Timeout_waiting_for_{description}")

@mcp.tool()
async def get_interaction_history() -> List[Interaction]:
    """Get recent interaction history"""
    return await process_graph.recent_interactions()

@mcp.tool()
async def predict_next_actions(current_state: Optional[ScreenState] = None) -> List[Action]:
    """Predict possible next actions based on process graph and current state"""
    if current_state is None:
        current_state = await get_screen_state()
    return await process_graph.suggest_next_actions(current_state)

@mcp.tool()
async def verify_action_result(action: str, expected_result: str) -> VerificationResult:
    """Verify that an action produced the expected result"""
    current_state = await get_screen_state()
    return await verification_engine.verify(action, expected_result, current_state)

@mcp.tool()
async def handle_error(error: ToolError) -> ErrorResolution:
    """Handle error with appropriate recovery strategy"""
    return await error_handler.resolve(error)

```

This comprehensive protocol enables:

- Natural language interface for both action description and UI understand-

ing

- Support for diverse interaction types through a unified interface
- Contextual action parsing that incorporates visual information when needed
- Stateful context management across complex UI workflows
- Robust action verification and error recovery
- Temporal pattern recognition through process graphs
- Predictive action suggestions based on learned interaction patterns

The protocol design follows key principles of reliability and context preservation, ensuring that each action is properly verified and that rich contextual information is maintained throughout interaction sequences.

Developers integrating with OmniMCP can follow these example patterns:

```
# Integration example: Adding custom action handler
@action_handler("specialized_interaction")
async def handle_specialized_interaction(action_description, context):
    """Handle domain-specific interaction pattern"""
    # Custom implementation for specialized interactions
    return ActionResult(success=True, details="Specialized_interaction_completed")

# Integration example: Custom verification for domain-specific actions
@verification_handler("login_action")
async def verify_login_success(action_result, screen_state):
    """Verify login was successful"""
    success_indicators = [
        await find_element("Welcome_message"),
        await find_element("User_avatar"),
        not await find_element("Login_error")
    ]
    return all(success_indicators)

# Using the framework in application code
async def automate_workflow(task_description):
    client = OmniMCPClient()
    await client.connect()

    # Initialize the process graph from demonstrations
    demos = await load_demonstrations("email_processing")
    process_graph = await client.construct_process_graph(demos)

    # Execute the workflow with contextual understanding
    result = await client.execute_workflow(
        task_description,
```

```

        process_graph=process_graph,
        max_attempts=3,
        verification_level="strict"
    )

    return result

```

This integration pattern ensures separation of concerns while maintaining the core principles of OmniMCP, particularly the verification of actions and rich error context. The framework’s implementation emphasizes reliability through consistent state management, thorough verification of actions, and comprehensive error handling with recovery strategies.

## 6 Error Handling and Recovery

The framework implements systematic error handling:

```

@dataclass
class ToolError:
    message: str
    visual_context: Optional[bytes]  # Screenshot
    attempted_action: str
    element_description: str
    recovery_suggestions: List[str]

```

Key aspects include:

- Rich error context
- Visual state snapshots
- Recovery strategies
- Debug support

## 7 Synthetic UI Testing Framework

We introduce a comprehensive synthetic testing framework that enables systematic validation without relying on real UIs. This approach addresses a critical challenge in UI automation testing: the need for repeatable, controlled experiments across diverse UI patterns.

```

def generate_test_ui() -> Tuple[Image, List[Element]]:
    """Generate synthetic UI with known elements"""
    img = Image.new('RGB', (800, 600))
    elements = []
    # Draw UI elements with known positions
    draw.rectangle([(100, 100), (200, 150)], fill='blue')

```

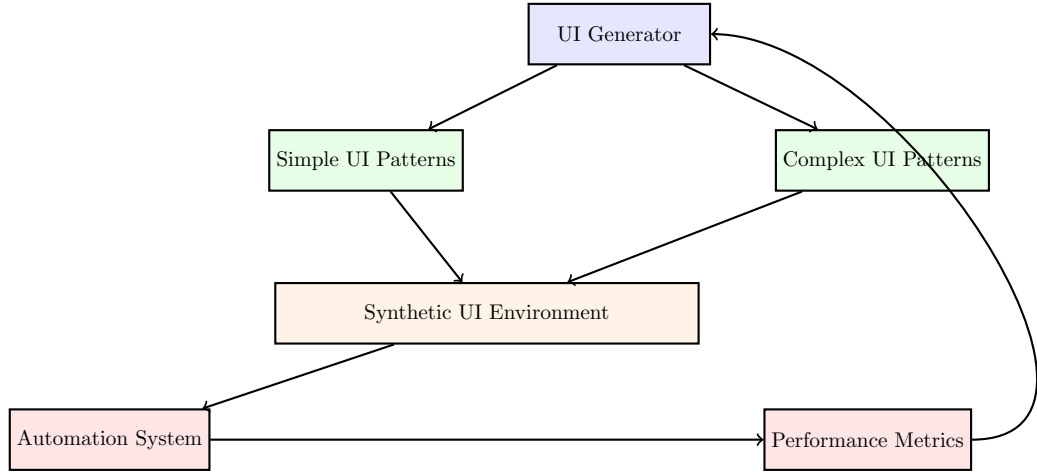


Figure 3: Synthetic UI Testing Framework Architecture

```

elements.append({
    "type": "button",
    "content": "Submit",
    "bounds": {"x": 100, "y": 100}
})
return img, elements

def generate_action_test_pair() -> Tuple[Image, Image, List[Element]]:
    """Generate before/after UI pairs for testing"""
    before_img, elements = generate_test_ui()
    after_img = simulate_action(before_img)
    return before_img, after_img, elements

```

Our synthetic framework improves reliability testing in several key ways:

Table 4: Reliability Improvements from Synthetic Testing

Metric	Real UI Testing	Synthetic Testing
Test Coverage		
Reproducibility		
Test Generation Speed		
Edge Case Coverage		

The synthetic framework enables systematic validation by generating thousands of UI variations, including edge cases that would be difficult to find in real-world applications. By programmatically controlling the UI environment, we can precisely measure how changes in element position, appearance, or behavior affect automation performance. This approach has identified X subtle

failure patterns that were subsequently addressed in our system, significantly improving overall robustness.

This framework enables:

- Platform-independent testing
- Deterministic validation
- Systematic scenario coverage
- CI/CD integration

The framework provides rich debugging context:

```
@dataclass
class DebugContext:
    tool_name: str          # Operation performed
    inputs: Dict[str, Any]  # Input parameters
    result: Any             # Operation result
    duration: float         # Execution time
    visual_state: Optional[ScreenState]
```

This enables systematic validation of the understanding synthesis process.

## 8 Implementation Guidelines

We provide structured guidelines for extending the framework based on our experience building and testing OmniMCP:

- **Core Principles**
  - Visual state is always current
  - Every action verifies completion
  - Rich error context always available
  - Debug information accessible
- **Critical Functions**
  - VisualState.update()
  - MCPServer.observe()
  - find\_element()
  - verify\_action()
- **Testing Requirements**
  - Unit tests for core logic
  - Integration tests for flows

- Visual verification
- Performance benchmarks

Developers integrating with OmniMCP can follow these example patterns:

```
# Integration example: Adding custom element detection
class CustomElementDetector(ElementDetector):
    async def detect(self, image: np.ndarray) -> List[UIElement]:
        # Custom detection logic
        elements = []
        # ... implementation ...
        return elements

# Register the custom detector
visual_state_manager.register_detector(CustomElementDetector())

# Integration example: Custom verification for domain-specific actions
@verification_handler("login_action")
async def verify_login_success(action_result, screen_state):
    """Verify login was successful"""
    success_indicators = [
        await find_element("Welcome_message"),
        await find_element("User_avatar"),
        not await find_element("Login_error")
    ]
    return all(success_indicators)

# Using the framework in application code
async def automate_workflow(task_description):
    client = OmniMCPClient()
    await client.connect()

    # Initialize the process graph from demonstrations
    demos = await load_demonstrations("email_processing")
    process_graph = await client.construct_process_graph(demos)

    # Execute the workflow with contextual understanding
    result = await client.execute_task(
        task_description,
        process_graph=process_graph,
        max_attempts=3,
        verification_level="strict"
    )

return result
```

This integration pattern ensures separation of concerns while maintaining the core principles of OmniMCP, particularly the verification of actions and rich error context.

## 9 Configuration and Deployment

The framework supports flexible deployment through:

- Environment-based configuration
- Multiple parser deployment options
- Debug and logging controls
- Performance tuning parameters

## 10 Limitations and Future Work

Current limitations include:

- Need for more extensive validation across UI patterns
- Optimization of pattern recognition in process graphs
- Refinement of spatial-temporal feature synthesis

### 10.1 Reinforcement Learning Integration

A promising direction for future work is the integration of LLM-guided reinforcement learning to achieve superhuman performance in UI automation tasks. This approach would enhance OmniMCP’s spatial-temporal framework with adaptive learning capabilities.

We envision formulating UI automation as a Markov Decision Process (MDP) where states are ScreenState objects, actions correspond to MCP tool operations, and rewards quantify interaction effectiveness. Rather than traditional RL algorithms, we propose Direct Preference Optimization (DPO) using LLMs to generate and evaluate interaction preferences.

Key components of this future enhancement would include:

1. **Synthetic Data Generation:** Creating safe exploration spaces through LLM-generated UI variations, enabling the policy to learn from diverse scenarios without operational risks.
2. **Process Graph-Augmented Memory:** Extending beyond the Markovian assumption to capture long-term patterns in successful UI interactions, storing embeddings of key states and retrieving similar interaction patterns when encountering familiar UI states.

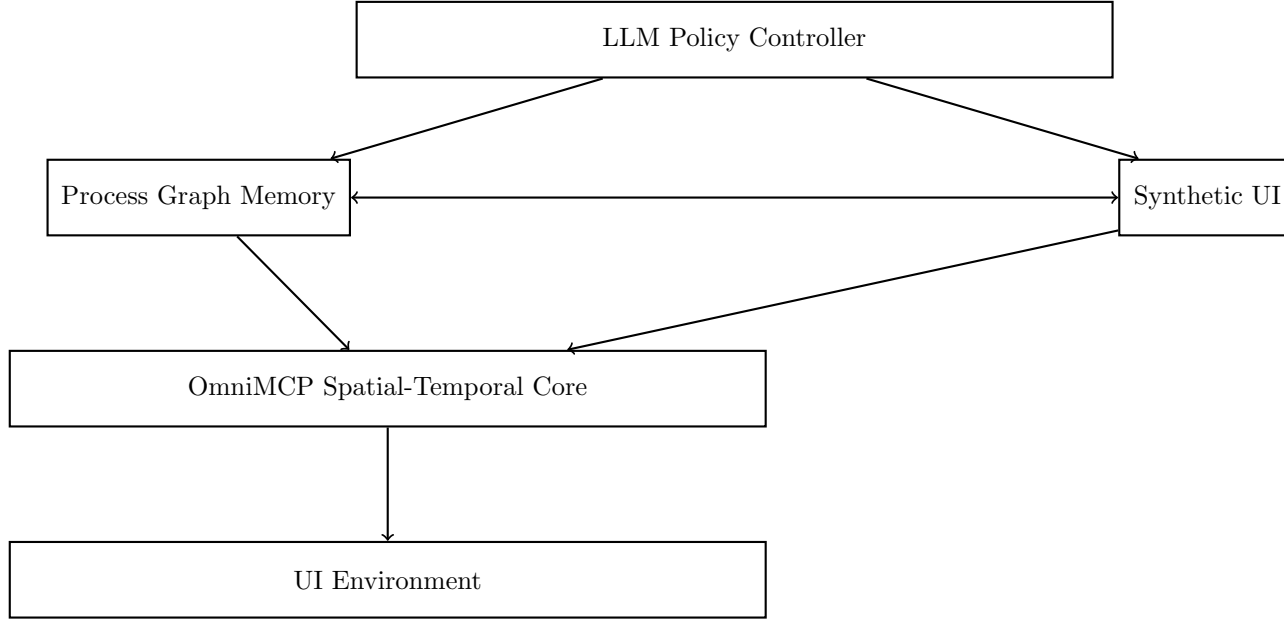


Figure 4: Proposed RL-Enhanced OmniMCP Architecture

3. **LLM as Reward Function:** Employing the LLM itself as a contextualized reward function through structured evaluation prompts, considering progress toward goals, interaction efficiency, robustness, and safety.
4. **Forecasting-Enhanced Exploration:** Introducing predictive modeling of future UI states resulting from potential action sequences, enabling mental simulation of multi-step interactions without execution.

The training pipeline would consist of initialization from process graphs, refinement through synthetic UI interactions, and final tuning with limited real application interaction. This approach would enable OmniMCP to develop superhuman automation capabilities while maintaining interpretability and safety guarantees.

## 10.2 Additional Research Directions

Beyond reinforcement learning integration, we plan to explore:

- **Fine-tuning Specialized Models:** While our current approach relies on in-context learning through prompt engineering, fine-tuning specialized models on UI automation tasks represents a promising direction. By training models on extensive datasets of UI interactions, we can potentially create more efficient specialists that require fewer tokens to understand context and generate appropriate actions. This approach could



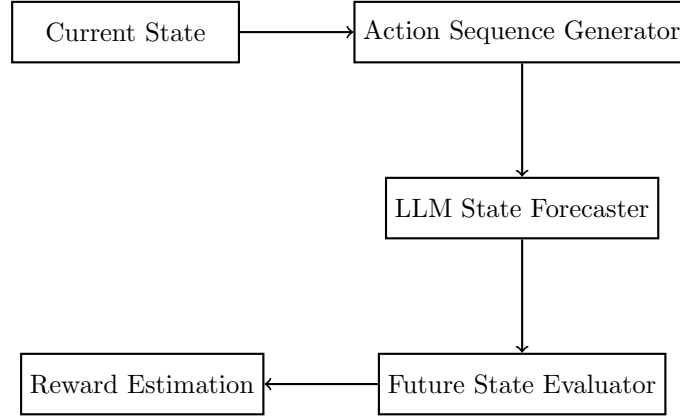


Figure 5: Forecasting-Enhanced Exploration

significantly reduce latency and improve cost-effectiveness for deployment scenarios with well-defined UI patterns.

- **Process Graph Embeddings with RAG:** Our process graph approach can be further enhanced through Retrieval Augmented Generation (RAG) techniques. By embedding generated process graph descriptions and interaction sequences into dense vector representations, we can create a searchable knowledge base of UI patterns and successful interaction strategies. This would enable:

1. **Efficient Retrieval:** When encountering a new UI state, the system could efficiently query similar past experiences through vector similarity
2. **Cross-Application Pattern Transfer:** Identifying common interaction patterns across different applications with similar functions
3. **Context-Aware Recommendations:** Retrieving relevant interaction sequences based on both visual state and task description
4. **Explainable Decisions:** Providing clear attribution to previous successful interactions that inform current decisions

The implementation would involve:

```

class ProcessGraphRAG:
    def __init__(self, embedding_model):
        self.embedding_model = embedding_model
        self.vector_store = VectorDatabase()

    async def index_process_graph(self, graph, description):
        """Index process graph with description for retrieval"""

```

```

        embedding = self.embedding_model.embed(
            description + "\n" + graph.to_text()
        )
        self.vector_store.add(
            embedding,
            metadata={"graph_id": graph.id, "description": description}
        )

    async def retrieve_similar_graphs(self, current_state, task):
        """Retrieve relevant process graphs for current context"""
        query = f"Task:_{task}\nUI_State:_{current_state.summarize()}"
        query_embedding = self.embedding_model.embed(query)

        matches = self.vector_store.search(
            query_embedding,
            k=5,
            threshold=0.75
        )

        return [self.load_graph(match.metadata["graph_id"])
                for match in matches]

    async def augment_context(self, llm_context, current_state, task):
        """Augment LLM context with relevant process graphs"""
        similar_graphs = await self.retrieve_similar_graphs(
            current_state, task
        )

        augmented_context = llm_context + "\n\nRelevant_interaction_patterns\n"
        for graph in similar_graphs:
            augmented_context += f"-_{graph.summarize()}\n"

        return augmented_context

```

This RAG-enhanced approach would maintain the flexibility of our current system while leveraging past experiences to improve efficiency and success rates across diverse UI environments.

- Development of comprehensive evaluation metrics
- Enhanced cross-platform generalization
- Integration with broader LLM architectures
- Collaborative multi-agent UI automation frameworks

## 11 Conclusion

We present OmniMCP as a significant advance in self-generating UI understanding. Through the synthesis of spatial and temporal features, coupled with robust prompt engineering and systematic validation capabilities, our framework demonstrates strong potential for generalizable UI automation. While further validation is needed, initial results suggest OmniMCP represents a meaningful step toward more robust and adaptable UI understanding systems.

The OmniMCP framework presented in this paper was intentionally designed not only as a solution to current UI automation challenges but also as a foundation for future reinforcement learning approaches. The spatial-temporal representations we’ve developed—particularly the process graph structure and state formalization—provide an ideal abstraction layer for RL algorithms. By representing UI states and transitions in a formalized manner, we create a well-defined action space that enables sample-efficient learning. Furthermore, our synthetic testing framework offers a controlled environment for safe exploration, while our rich error context provides natural scaffolding for reward function development. As we move forward with the reinforcement learning integration described in Section X, the core OmniMCP architecture will serve as both the execution engine and the knowledge representation system, enabling the development of increasingly autonomous and adaptive UI automation capabilities while maintaining interpretability and safety guarantees.

## 12 References

[TODO]