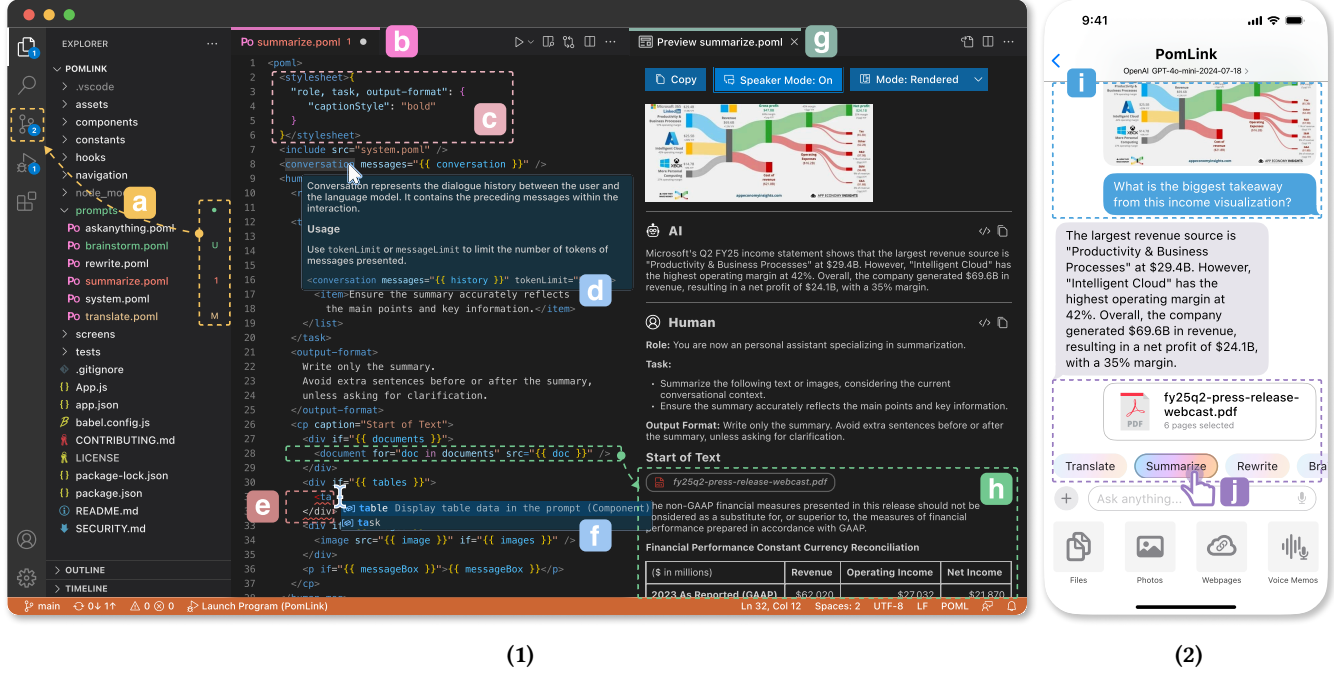


# Prompt Orchestration Markup Language

Yuge Zhang, Nan Chen, Jiahang Xu, Yuqing Yang  
semantipy@microsoft.com  
Microsoft Research  
Shanghai, China



**Figure 1: Developing the PomLink iOS prototype in VSCode – a case study of POML.** (1) The development environment in VSCode, showcasing structured prompt authoring with POML: (a) prompt version control with Git; (b) a POML file example; (c) stylesheet definition for presentation control; (d) hover documentation for inline assistance; (e) inline diagnostics for error checking; (f) context-aware auto-completion; (g) the live preview panel displaying; (h) rendered embedded data (e.g., documents and tables). (2) An application prototype developed with POML that structures multimodal inputs for LLM analysis, supporting tasks such as (i) querying visual content and (j) summarizing documents based on rich context.

## Abstract

Large Language Models (LLMs) require sophisticated prompting, yet current practices face challenges in structure, data integration, format sensitivity, and tooling. Existing methods lack comprehensive solutions for organizing complex prompts involving diverse data types (documents, tables, images) or managing presentation variations systematically. To address these gaps, we introduce POML (Prompt Orchestration Markup Language). POML employs

component-based markup for logical structure (roles, tasks, examples), specialized tags for seamless data integration, and a CSS-like styling system to decouple content from presentation, reducing formatting sensitivity. It includes templating for dynamic prompts and a comprehensive developer toolkit (IDE support, SDKs) to improve version control and collaboration. We validate POML through two case studies demonstrating its impact on complex application integration (PomLink) and accuracy performance (TableQA), as well as a user study assessing its effectiveness in real-world development scenarios.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Preprint, arXiv

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 978-1-4503-XXXX-X/2018/06  
<https://doi.org/XXXXXXX.XXXXXXX>

## CCS Concepts

• **Human-centered computing** → **Human computer interaction (HCI); Interaction design;** • **Software and its engineering** → **Markup languages; Formal language definitions;** • **Computing methodologies** → **Natural language processing.**

## Keywords

Artifact or System, Machine Learning, Programming/Development Support, Prompt Engineering, Large Language Model

## 1 Introduction

Large Language Models (LLMs) have demonstrated remarkable capabilities across diverse tasks through carefully engineered prompts. These models, exemplified by the GPT series [14, 62], LLaMA [84], and others [22, 83], have shown proficiency in parsing and acting upon instructions in various domains. Their applications span code generation [17], question answering [37], mathematical reasoning [19], and interactive agents [78, 100]. As these applications grow in complexity, so too does the sophistication of prompting techniques used to elicit desired behaviors from LLMs. Advanced prompting strategies such as chain-of-thought [97], few-shot in-context learning [14], and ReAct (reasoning and acting) [101] have emerged as critical methods for improving model performance.

Within this evolving landscape, “structured prompting” [72, 95, 104] has gained prominence as a systematic approach to organizing instructions and context data within prompts. These structured approaches offer numerous advantages: they componentize and standardize prompts for better automated editing [74], facilitate team collaboration [68], improve development efficiency [13], support fine-grained editing [24], enable reuse of proven designs [95], and potentially boost performance [44]. As LLM applications increasingly integrate multiple data types – including images [83], tables [80], and documents [45], the need for a more organized approach to prompt design becomes increasingly apparent.

Despite these advances, prompt engineering faces four significant challenges that impede effective development and deployment of LLM applications. *First*, current practices lack standardized methods for the **structured orchestration of prompts**. Projects following “structured prompting” practices often involve scattered instructions, roles, tasks, and examples [95, 104], hindering maintenance and collaboration, especially in complex projects. *Second*, **integrating diverse data sources** remains complex. Large LLM applications increasingly require references to documents [45], tables [18, 103], or images [46]. Properly integrating and presenting these varied sources within prompts poses a significant engineering burden. *Third*, LLMs exhibit significant **sensitivity to formatting and presentation** [73, 76, 89]. Research has documented “butterfly effects,” where minor textual variations can dramatically alter results. The absence of a decoupled styling layer complicates systematic testing and refinement of prompt variations [48, 75]. *Finally*, prompt development and management suffer from **inadequate tooling**. Current workflows often lack effective version control and diff visibility [24], making tracking changes and collaboration difficult. Plain-text prompts are cumbersome to manage [11, 102], and limited IDE support hinders systematic improvement [9, 79].

Several existing approaches attempt to address these challenges, but each has its specific limitations. Workflow and agent orchestration tools like LangChain [3], Microsoft Guidance [53], and PromptChainer [98] manage multi-step flows but treat prompts as plain text or minimal templates, lacking structured data presentation or advanced styling capabilities. Markup-based prompting approaches such as ChatML [77], MDXPrompt [30], PromptML

[68], and VSCode PromptTSX [90] provide component-based frameworks but lack effective data handling, a well-separated styling system, or comprehensive development support. User interfaces for prompt comparison or logging exist [9, 57], yet lack the fundamental tools to author or structure prompts before running comparisons or evaluations.

To address these multifaceted challenges comprehensively, we introduce **POML (Prompt Orchestration Markup Language)**. POML provides a novel paradigm designed to bring structure, maintainability, and versatility to advanced prompt engineering. POML is based on an HTML-like **structured markup language** featuring semantic components, such as `<role>`, `<task>`, and `<example>`, encouraging modular design and enhancing prompt reusability. To ease data integration, POML incorporates specialized **data components** (e.g., `<document>`, `<table>`, `<img>`) that seamlessly embed or reference external data sources like text files, spreadsheets, and images [18, 99]. These data components offer customizable formatting options and eliminate error-prone manual text merging. To address format sensitivity, POML introduces a **styling system** inspired by CSS [26, 94]. This decouples content from presentation, allowing developers to modify styling (e.g., verbosity, syntax format) via separate `<stylesheet>` definitions (Figure 1 (c)) and systematically experiment with format variations without altering core logic. POML also includes a built-in **templating engine** for dynamically generating complex, data-driven prompts.

Beyond the language itself, POML addresses inadequate tooling with a comprehensive **development toolkit**. This includes an Integrated Development Environment (IDE) extension (initially for Visual Studio Code, see Figure 1 (1)) providing essential development aids like real-time previews (g), inline diagnostics (e), and auto-completion (f). The toolkit also provides Software Development Kits (SDKs) for seamless integration into Python and Node.js applications [3, 54]. Overall, the toolkit streamlines the prompt engineering lifecycle, enhancing authoring efficiency, debugging, version control (Figure 1 (a)), collaboration, and the systematic evolution of prompts.

We validate the effectiveness and utility of POML through rigorous empirical evaluation. This includes two distinct case studies and a formal user study. The first case study showcases POML’s practical application in building PomLink – an iOS agent application prototype (Figure 1 (2)). This study highlights how POML’s data components, styling, and tooling enabled the rapid development (a functional prototype in two days) of a complex application prototype requiring integration of documents, images, and tables. The second case study, focusing on Table Question Answering (TableQA), systematically explores the impact of prompt styling using POML’s styling system. It demonstrates that stylistic variations, easily managed by POML, can cause dramatic, model-specific differences in LLM performance (e.g., accuracy improvements over 9x for some models) on table-based reasoning tasks [64]. Finally, our user study involving participants with varying expertise assessed POML’s usability across five representative tasks [20, 102]. The results confirmed POML’s effectiveness in structuring prompts and handling data, with participants particularly valuing the data components and development toolkits, while also providing constructive feedback for future enhancements.

The primary contributions of this work are:

- (1) POML: A novel markup language designed specifically for prompt engineering, with specialized data components for effective integration of diverse data types, and a styling system for decoupling prompt content from presentation, enabling systematic control over formatting and mitigating LLM format sensitivity.
- (2) An integrated suite of tools including an IDE extension with features like syntax highlighting, live preview, diagnostics, interactive testing, and multi-language SDKs to improve the development workflow, improve version control, and facilitate project management and collaboration.
- (3) POML’s practical benefits and impact are empirically validated through two detailed case studies and a formal user study that assessed usability and effectiveness in realistic scenarios.

## 2 Related Works

### 2.1 Prompt Engineering and Structured Prompt

Large Language Models (LLMs) demonstrate remarkable capabilities, yet their performance is highly sensitive to the input prompt’s quality and format [72]. This sensitivity is particularly pronounced when prompts involve complex data types such as tables, charts, or images, where carefully tuned presentation is often necessary for reliable results [18, 31, 36, 99, 103]. Indeed, research highlights a “butterfly effect,” where minor variations in phrasing, formatting, markup, or content order can dramatically alter LLM outputs and accuracy [35, 48, 73, 89, 106].

Consequently, prompt engineering – the practice of designing inputs to optimize LLM outputs – has become essential [72]. Research has shown that enhancing prompts with in-context examples or iterative reasoning steps significantly improves results [14, 25, 96, 97, 101]. Within this field, structured prompting has emerged as a key practice, organizing prompts into standardized roles, tasks, or program-like formats [51, 66, 74, 95, 104]. Such approaches offer advantages including improved reusability and maintainability [13, 24, 44, 68, 74, 95]. Conversely, unstructured prompts often impede team-based workflows, making targeted modifications difficult without affecting other prompt sections [24]. The lack of clear boundaries or standardized documentation for prompt segments (e.g., roles, tasks, examples) hinders the development of shared tooling for parsing, manipulating, or reusing prompt components programmatically.

### 2.2 Prompt Development Challenges and Tooling

**Studies on Prompt Usage and Challenges** Studies on prompt usage patterns reveal significant challenges faced by users, especially non-experts. Minimalist web/chat UIs like OpenAI’s Playground and ChatGPT [59, 60] serve as short-term sandboxes, leaving a gap between experimental prototyping and large-scale and production-ready prompt designs. In industry, trial-and-error approaches are common but hinder systematic improvements, as users often make erroneous assumptions about how to communicate with LLMs [20, 102]. This continuous editing process also makes prompt history tracking difficult, complicating version control and collaboration. Text-based prompts frequently become opaque or difficult to

manage over time [11, 24]. Tagging approaches are being explored to alleviate these issues by improving differential comparisons and collaborative editing [68, 95], aligning with structured prompting practices.

**Prompt Toolkits for Development and Evaluation** Several integrated prompt toolkits have been developed to assist users in prompt engineering. Flow-based or multi-prompt systems provide pipelines for chaining multiple prompts with optional retrieval components in complex LLM applications [3, 53, 54, 98]. Other tools focus on prompt management and version control through standalone or web-based services [41, 65]. However, these systems typically manage the high-level orchestration of multiple steps or prompts, often treating individual prompts as opaque strings or simple templates, rather than providing detailed frameworks for structuring the content *within* a single complex prompt. POML can integrate with these systems as a specialized syntax for enhanced clarity and structure within individual prompt steps. Another stream of existing works is prompt evaluation interfaces, offering UIs for comparing or perturbing prompts [4, 9, 57, 67, 79]. These are largely *orthogonal* to POML’s focus on prompt authoring and structure, and could potentially use POML’s declared sub-blocks for more fine-grained testing and analysis.

### 2.3 Web Frameworks and Prompt Markup Languages

**Inspiration from Web Frameworks and Programming** The evolution of web development offers relevant parallels and insights for prompt engineering. Both domains share concerns around clarity, dynamic content, and adaptable layouts – whether for different screen sizes or context windows [49]. Principles like the separation of structure (HTML [12, 27]), presentation (CSS [26, 94]), and behavior (JavaScript [28]) have proven crucial for managing complexity in web development, providing a blueprint for prompt frameworks. Early dynamic websites relied on templating engines to map data into presentation views [29, 34, 40, 56, 85], akin to simple prompt templating in tools like LangChain [3] or Guidance [53]. A significant advancement came with component-based frameworks [5, 70, 91], which encapsulate structure, style, and logic into reusable units, enhancing modularity and maintainability [42, 81, 87, 88]. Building upon these frameworks, UI component libraries such as Ant Design [23] and MUI [58] provide extensive sets of pre-defined, high-level components (“*widgets*”) that further accelerate development by offering ready-to-use solutions for common UI patterns, often including advanced data display capabilities. Languages like JSX [69], used extensively in React [70], exemplify this by allowing developers to define UI “*components*” using declarative, *markup-like syntax* combined with programming logic. This component-based approach improves code readability and fosters reuse. Conversely, manually constructing structured outputs via string concatenation is notoriously error-prone and hinders readability [15]. These lessons underscore the potential benefits of applying similar principles – using components, templates, and structured markup – to manage the complexity of modern prompt engineering, a core motivation for POML.

**Prompting with Markup Languages** Several markup-based or component-based prompt frameworks have emerged, though they vary in their approach and capabilities. Some approaches leverage JSX/TSX syntax within TypeScript (TS) for templating prompts [8, 10, 90]. These often prioritize adapting to context window lengths but may offer limited built-in features for complex data input handling or advanced styling mechanisms. Other frameworks adopt broader React-inspired workflows or component models for orchestrating agent interactions or generating complex outputs [32, 33]. These systems typically offer limited control over the final prompt string’s styling or the embedding of diverse data types directly via markup. MDXPrompt [30] integrates MDX syntax for prompts but requires TypeScript API usage to render dynamic content, coupling the prompt definition to the execution environment. These approaches attempt to “componentize” prompts with HTML-like syntax but are limited in “widget-level” data input encapsulation, or focus primarily on high-level agent or workflow orchestration rather than providing comprehensive prompt-level structuring and styling solutions.

Outside the JS/TS ecosystem, other markup solutions propose structured tags but typically focus more on conversational chat formats or specialized formats for LLM training [13, 77]. PromptML [68], while advocating for a domain-specific language to improve standardization and collaboration, currently lacks integrated development environment support, such as syntax highlighting, which can affect usability and adoption. SAMMO [75] operates at a higher-level abstraction, representing prompts symbolically to allow for transformations and optimization searches. While these solutions provide intention-expressive tags or enable powerful symbolic manipulations, they generally lack built-in styling layers, flexible data import features across multiple formats, or comprehensive IDE toolkits. Many existing markup or programmatic approaches are also tightly bound to specific programming languages (e.g., Python), limiting flexibility. In contrast, POML aims to unify structured markup, effective data handling, decoupled styling, and integrated templating within a standalone language specification, supported by a comprehensive toolkit and SDKs for broader compatibility and ease of use across different development environments. A table summarizing the key differences between POML and existing prompt markup languages is provided in Appendix A.

## 3 Motivations and Design Goals

### 3.1 Motivations

Insights gathered from preliminary discussions with 5 prompt engineers in our research group highlighted 4 pressing challenges in contemporary prompt engineering. These challenges directly motivate our design of POML as a comprehensive solution for prompt development.

**M1: Structured Prompt Orchestration Necessity** While the benefits of structured prompting are increasingly recognized [13, 24, 44, 68, 74, 95], a significant problem lies in the lack of a standardized framework for implementation. Current practices often result in ad-hoc approaches using unstructured plain text where instructions, roles, tasks, and examples are scattered inconsistently throughout prompts [51, 95, 104]. This absence of standardization

makes prompts difficult to maintain, optimize, and reuse, especially within collaborative development environments.

**M2: Complexity in Data Integration** Large LLM applications increasingly rely on varied data formats, requiring references to documents [45], tables [18, 103], or images [46]. Integrating these diverse data sources becomes chaotic and error-prone. Manual text shaping, multi-point editing, and frequent reformatting significantly hamper maintainability [24, 36]. Existing tools often rely on ad-hoc placeholders or separate scripts for data embedding [3, 11, 54], lacking effective, integrated solutions.

**M3: Prompt Styling and Format Sensitivity** LLMs exhibit significant sensitivity to subtle variations in input formatting, where minor presentational changes can drastically alter results [35, 73, 76, 89, 106]. This sensitivity underscores the need for structured control over prompt presentation. The absence of a mechanism to decouple presentation style from prompt content makes it difficult to systematically test, refine, or adapt prompt formats for different models or tasks [48, 75].

**M4: Prompt Development and Management Challenges** Current prompt development workflows suffer from poor diff visibility and version control [24], making it difficult to track changes and collaborate effectively. Existing prompt markup languages [13, 30, 68] often lack syntax highlighting, auto-completion, preview and debugging tools, increasing the learning curve and error rate. Many solutions [33, 90] are tightly bound to particular programming languages or frameworks, limiting their integration with existing development tools and workflows. These limitations hinder the efficiency of development and make prompt management increasingly difficult as projects scale.

### 3.2 Design Goals

Based on these motivations, we established four core design goals for POML, each intended to address one of the identified challenges. These goals focus on standardizing syntax, supporting effective data handling, separating styling from content, and providing developer-focused tooling.

**DG1: Reusable and Maintainable Prompt Markup** A primary goal is to provide a standardized, JSX-like markup language that inherently supports structured prompting practices (M1). The language *should* enhance prompt clarity and ensure that prompt logic is easily refactorable and shareable across projects or teams. It *should* support a hierarchical nesting structure to reflect logical relationships and improve flexibility. Furthermore, the system *should* offer commonly used elements like roles, tasks, or few-shot examples, promoting modularity and consistency in line with structured prompting approaches [95, 104].

**DG2: Comprehensive Data Handling** To address the complexities of data integration (M2), the system *should* offer effective and comprehensive data handling capabilities, including dynamic prompt generation. Specialized components *should* be provided to seamlessly embed commonly used static data sources into prompts, including documents, tables, and images [18, 45, 99]. These components *should* preserve important data details while offering easy-to-use interfaces for adjusting presentation formats (e.g., syntaxes for



tables) [36], reducing the burden of manual text formatting. Furthermore, the language itself *should* support integrated templating logic (variables, loops, conditionals) as most data are only available at the runtime. Support for fallback text for non-textual data *should* also be included to accommodate LLMs with varying multi-modal capabilities and enhance prompt robustness across models.

**DG3: Decoupled Presentation Specifications** To manage LLM sensitivity to formatting (M3), a core design goal *is* to provide flexible control over prompt presentation while simultaneously decoupling these presentation specifications (styling) from the core prompt content. This separation *should* allow styling rules (e.g., bullet style, syntax format, verbosity) to be defined centrally with minimized modifications to prompt contents. The system *should* support format configurations at multiple levels (e.g., both overall syntax and local syntax) to adapt prompts easily to different LLM sensitivities or task requirements [35, 48]. This decoupling will ideally enable rapid testing and optimization while maintaining content integrity.

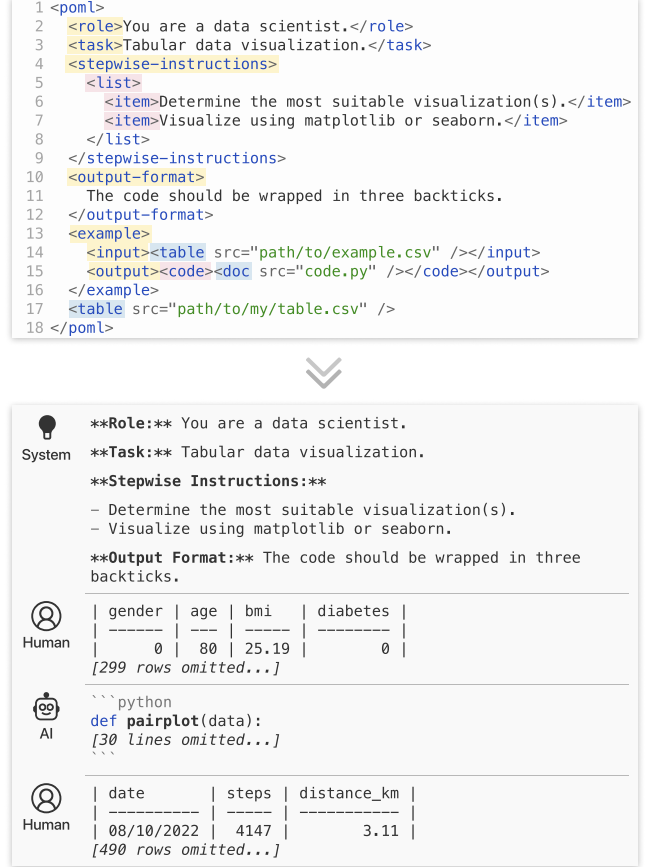
**DG4: Enhanced Tooling for Development** To address development and management challenges (M4), the framework *must* be accompanied by enhanced tooling. It *should* utilize standalone prompt files with a clear structure, facilitating integration with standard version control systems like Git and improving readability and collaborative editing. Dedicated IDE support *should* be provided with essential features like syntax highlighting, context-aware auto-completion, hover documentation, inline diagnostics, and a live preview. Finally, the system *should* maintain interoperability with popular LLM frameworks (e.g., LangChain [3] and Guidance [53]) through integration SDKs, allowing seamless integration into existing development workflows.

## 4 POML: Prompt Orchestration Markup Language

POML, standing for Prompt Orchestration Markup Language, introduces a novel structured paradigm designed for authoring, managing, and testing prompts for LLMs. This section provides an overview of POML’s foundational elements, including its core syntax based on HTML principles, its features for integrating diverse data types, its styling system for formatting prompt content, and a templating engine for dynamic prompt generation. These features address common challenges in prompt engineering. POML’s hierarchical structure enhances clarity and reusability (DG1). Its data components streamline the integration of varied information sources (DG2). Its styling system manages LLM formatting sensitivities (DG3). Its overall design combined with tooling (introduced in § 5) supports an enhanced development experience (DG4).

### 4.1 Structured Prompting Markup

The fundamental syntax of POML adopts an HTML-like structure, closely mirroring HTML conventions. This design choice stems from the observation that prompt engineering shares similarities with web design; both involve crafting interfaces — textual for LLMs, visual for humans — that require clear hierarchy, consistent presentation, and adaptability. Using an HTML-inspired approach



**Figure 2: An example illustrating POML’s structured markup and rendering (§ 4.1). Top: POML source code demonstrating core components: intention components (e.g., `<role>`, `<task>`, `<example>`), data components (e.g., `<table>`, `<doc>`), and basic structural components (e.g., `<list>`, `<item>`). Bottom: The corresponding rendered output, demonstrating how the structured markup translates into a clear prompt for an LLM. Note the rendering of intention components (e.g., `<role>` becomes the title “**\*\*Role:\*\***”) and data components (e.g., the first `<table>` is rendered as Markdown from its CSV source).**

also leverages developers’ familiarity with web technologies, reducing the learning curve [49]. Nested tags (called **components** in our case) form the overall structure, enabling the hierarchical composition of complex prompts from smaller, modular parts. Within these tags, attributes provide metadata or configure the behavior of individual components, similar to how attributes like `src` or `href` function in HTML. POML supports standard HTML comment syntax, `<!-- ... -->`, allowing developers to embed annotations within the prompt source code without affecting rendering. POML also maintains compatibility with plain text, allowing for gradual adoption. The root `<poml>` tag is optional, akin to the optional `<html>` tag in web development, which minimizes overhead for simple prompts where only specific features like data embedding might initially be needed.

Functionally, POML components fall into several key categories, conceptually similar to how complex user interface widgets are built upon fundamental HTML elements.

- **Basic Structural Components** provide fundamental text formatting and structural grouping capabilities, analogous to common HTML tags. Elements like `<b>` (bold), `<i>` (italic), `<p>` (paragraph), and `<div>` (division) allow for control over text presentation and logical grouping within larger content blocks, enhancing readability and structure. Alongside these static elements, POML includes components and syntax for dynamic content generation through its integrated templating engine, such as constructs for loops (`for`), conditionals (`if`), variable definition (`<let>`), and substitution (`{{ . . . }}`). These templating features are detailed further in § 4.4.
- **Intention Components** define the core logic and overall structure of the prompt. These components implement the “structured prompting” paradigm [30, 51, 68, 95, 104], establishing the interaction’s purpose and guiding the LLM’s response. Examples include the `<role>` component, used to define the persona the LLM should adopt; the `<task>` component, which specifies the objective the LLM needs to achieve; and the `<example>` component, crucial for providing few-shot demonstrations to guide the model’s behavior. Unlike data components, intention components organize the logical blocks of the prompt rather than presenting complex data directly. When rendered, components like `<role>` typically appear as formatted titles (e.g., `**Role:**`), with the presentation adjustable via styling rules (§ 4.3). These intention components orchestrate the overall flow and content of the prompt, as illustrated in Figure 2.
- **Data Components** are designed to handle the integration of diverse external data formats into the prompt context. Elements such as `<document>`, `<img>`, and `<table>` provide methods to embed content from files or data structures. These components, detailed in § 4.2, are essential for grounding LLM responses in specific information or enabling tasks that operate on external data.

**Rationale** This hierarchical and descriptive markup enforces a logical organization, offering significant advantages over unstructured plain text. For instance, as illustrated in Figure 2, the explicit separation of conceptual parts — using intention components like `<role>` and `<task>`, structuring instructions with `<stepwise-instructions>`, defining outputs with `<output-format>`, and providing few-shot demonstrations [14] via `<example>` (which itself nests `<input>` containing data components like `<table>` and `<output>` referencing an external `<doc>`) — significantly improves prompt clarity and understandability, reducing ambiguity for both LLMs and human readers. These components, such as the tables shown integrated directly or within examples, can then be easily reused across different prompts or modified in isolation, fostering systematic prompt engineering (DG1). Compared to simpler formats like PromptML [68], which primarily offer basic task and example definitions, POML provides richer semantic components crucial for handling diverse content types (DG2). Meanwhile, the modularity allows individual components to be treated as self-contained units. Modifying components independently minimizes the risk of

unintended side effects elsewhere in the prompt, facilitating faster and safer iterations, supporting more agile development workflows. The combination of a standardized, text-based format and clearly named elements simplifies prompt management within standard version control systems like Git (Figure 1 (a)). This enhances collaboration, as changes become easier to track, merge, and discuss (DG4).

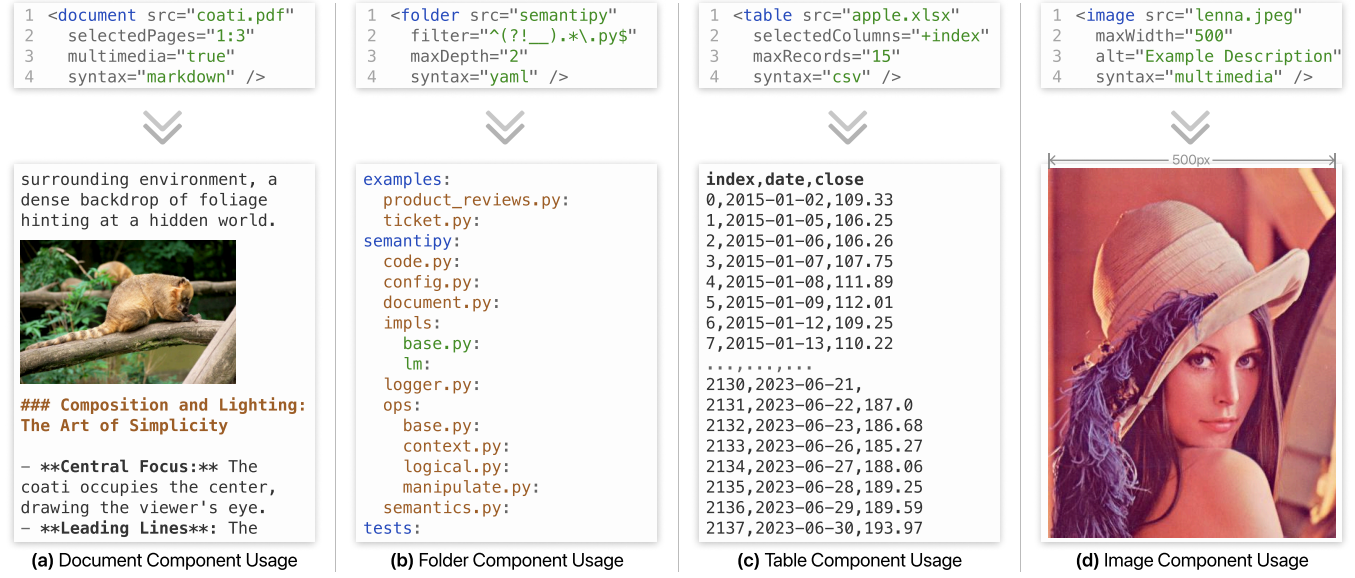
## 4.2 Data Components

POML provides specialized data components to systematically integrate diverse data modalities directly within the prompt structure, as exemplified in Figure 3. This capability directly addresses the critical need to combine LLM reasoning with varied external data sources, a key requirement for many advanced applications (DG2). POML’s modular data components reduce the inherent complexity associated with constructing prompts that incorporate multiple data inputs. These built-in components distinguish POML by offering systematic handling for common data types within a unified framework, analogous to component libraries in modern UI frameworks [23, 58]. Currently, POML supports **7 distinct data components**: document, table, folder, image, conversational messages, audio clip, web page. We detail **4 representative examples** here; the remaining components operate similarly and are demonstrated in the case studies (§ 7).

**Document** The document component (Figure 3 (a)) provides a mechanism for referencing and embedding content from external text-based files, supporting formats like `.txt`, `.docx`, and `.pdf`. This is particularly crucial for tasks involving large amounts of background information or context, such as in retrieval-augmented generation (RAG) systems [16, 45]. It offers granular control through attributes, allowing developers to specify character encoding, whether to preserve original formatting, whether to include or discard embedded images within documents like PDFs, and how to trim content — for instance, by selecting specific page ranges (e.g., `selectedPages="1:3"`) — thereby helping manage the LLM’s limited context window.

**Folder** For scenarios involving interaction with file systems or code repositories, such as file system analysis, code reviews, and project navigation [39, 47, 86], the `<folder>` component (Figure 3 (b)) provides a structured way to represent directory hierarchies. It allows customization of the representation through attributes controlling the maximum depth of the displayed tree, filtering options to include or exclude files based on extensions or name patterns (e.g., `filter="^(?!__).*\.py"`), and optional automatic summarization for very large directories to conserve context space. It can also display metadata like file sizes or modification dates and supports multiple output formats, including classic tree views using box-drawing characters, YAML, or JSON representations.

**Table** Tables are vital for complex question-answering over structured data [64], generating data visualizations [18], or performing data manipulation tasks [103]. The `<table>` component (Figure 3 (c)) supports a variety of input sources, including CSV, TSV, Excel spreadsheets, and JSON arrays of objects. Crucially, it allows specifying the output representation (e.g., Markdown, HTML, XML, plain CSV) and controlling content details such as the inclusion of



**Figure 3: Examples of POML data components demonstrating integration of diverse data types (§ 4.2). (a) `<document>` rendering selected PDF pages with multimedia; (b) `<folder>` displaying a filtered directory structure as YAML; (c) `<table>` extracting and formatting spreadsheet data as CSV; (d) `<img>` inserting a referenced image with resizing.**

headers or indices, and selective presentation of rows or columns for large tables, offering flexibility in presentation. Employing a consistent and well-defined table format via this component can significantly enhance an LLM’s ability to accurately parse and reason over the tabular data, applicable to scenarios like [36, 80, 103].

**Image** Acknowledging the rapid advancement of multimodal LLMs capable of processing visual information [62, 83], POML includes the `<img>` component (Figure 3 (d)) for direct image insertion into prompts. Inspired by web accessibility principles [92–94], this component supports standard attributes like `alt` text, allowing prompts to be easily adapted to LLMs with varying capabilities, including those without vision, thereby enhancing robustness. Layout can be influenced using a position attribute (e.g., `before`, `after`, `here`) to control placement relative to surrounding text. Attributes like `maxWidth` and `maxHeight` allow developers to suggest rendering constraints, potentially influencing the number of tokens consumed when sending image data to the model. Compared to interactive uploads in chat interfaces [21, 60] or less structured parameter passing in raw API calls [62, 83], inserting images via a dedicated component within the markup offers superior programmatic control and integration.

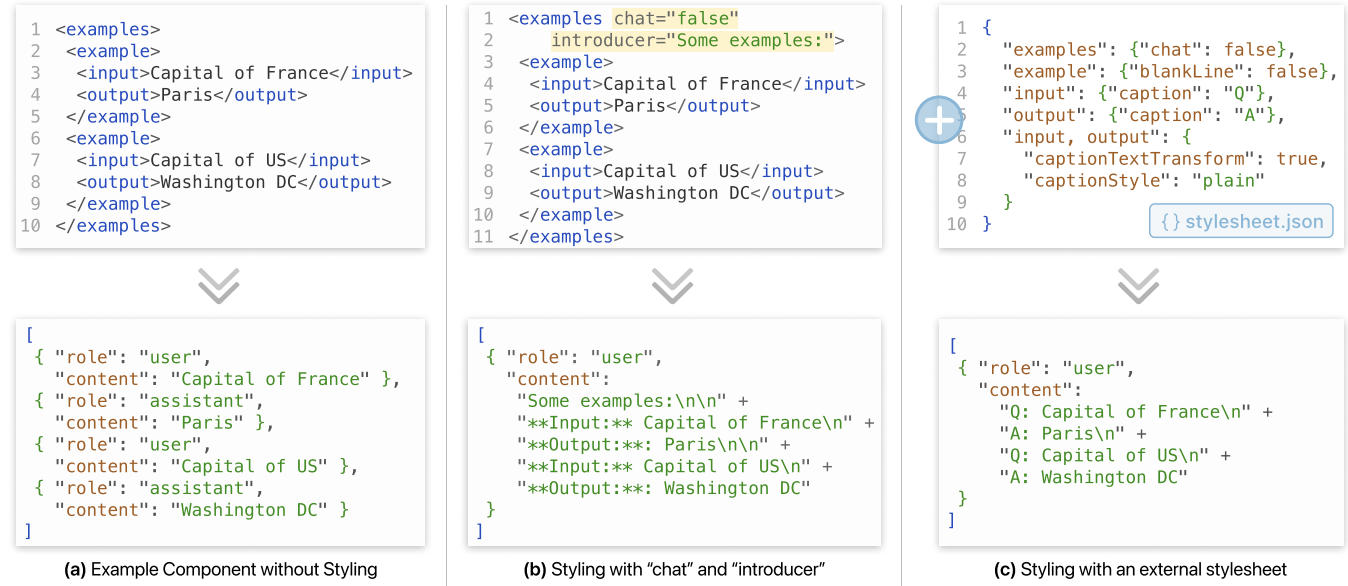
**Extensions** POML’s architecture is designed to be extensible. As LLMs evolve to handle new modalities, POML can incorporate additional data components for types like video segments and node-edge graphs, maintaining a consistent framework for diverse data integration.

### 4.3 Styling System

A critical challenge in prompt engineering stems from the documented sensitivity of LLMs to subtle variations in input formatting

[35, 73, 76, 89, 106]. Minor changes in presentation can significantly impact model performance, necessitating precise control over how prompt content is rendered (DG3). POML addresses this through a dedicated styling system designed to manage presentation effectively. The styling options provided, such as control over syntax formats, layout adjustments (e.g., chat vs. block formatting), list styles, and caption presentation, are informed by findings in prompt engineering research exploring these sensitivities [48, 76]. Inspired by the separation of concerns in web development (HTML for structure, CSS for style [26, 94]), POML deliberately decouples presentation rules from the core prompt content. POML offers two primary mechanisms for applying styles: *inline attributes* and *external stylesheets*.

POML provides control over various presentation aspects known to influence LLM behavior, such as layout adjustments (e.g., chat versus block formatting), list styles (`<list listStyle="decimal">`), caption presentation (`<hint caption="Important Note" captionStyle="header">`), and overall verbosity [48, 76]. Beyond simple formatting, POML provides a crucial **syntax attribute**. This attribute allows developers to explicitly control the rendering format for specific components or the entire prompt (e.g., `<table syntax="html"/>` or `<poml syntax="json">`). This control is vital because different LLMs can exhibit varying sensitivities or capabilities when parsing structured formats like Markdown, JSON, or XML [7] embedded within the prompt text. The styling attributes can be nested, allowing, for example, a subsection of a predominantly Markdown-formatted prompt to be rendered as JSON, providing fine-grained control over the final string sent to the model. One way to apply these styling adjustments is via **inline attributes** on specific POML components — a convenient method for localized formatting familiar to web developers. Figure 4 ((b) compared to



**Figure 4: Demonstrating POML styling capabilities (§ 4.3).** (a) Default rendering of `<example>` components. (b) Inline chat and introducer attributes on the parent `<examples>` element modify its presentation (from chat messages to plain text with `**Input**`/`**Output**` captions). (c) A `<stylesheet>` applies global rules to POML in (a), controlling layout (chat=false), captions/prefixes (caption="Q:"/"A: "), and styles (captionStyle="plain"), resulting in customized output.

(a)) demonstrates how inline attributes can alter the presentation of `<example>` components.

For managing styles more systematically across multiple components or prompts, POML supports **external stylesheets**, typically defined in separate JSON files (e.g., `stylesheet.json` as shown in Figure 4 (c)). These files contain global formatting rules defined using a concise, JSON-like syntax. Style rules within the stylesheet target specific POML component types (e.g., `hint`, `table`) or user-defined classes (see § 7.2 for examples). Alternatively, these style rules can also be embedded directly within a POML document using the `<stylesheet>` tag, as shown in Figure 1 (c). Overall, stylesheets provide a mechanism to define styles that apply globally or to specific component types, offering a way to batch-manage or customize styles that might otherwise be set inline individually. They offer several advantages: (1) it provides centralized control, establishing a single source of truth for formatting; (2) it keeps the primary prompt logic cleaner by avoiding repetitive presentation attributes; (3) it simplifies style modifications by eliminating multi-point edits, enhancing maintainability. As a result, this approach directly facilitates systematic experimentation with different presentation formats to address LLM sensitivities or task requirements (**DG3**) without altering the core prompt content.

#### 4.4 Templating Engine

POML integrates a built-in templating engine to facilitate the creation of dynamic, data-driven prompts without execution environment dependencies (**DG2**). Inspired by established web templating systems [29, 34, 40, 69], this engine allows runtime customization without relying solely on external scripting. Key features include

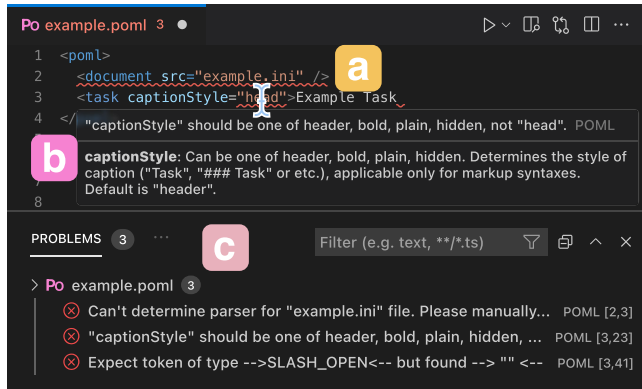
```
1 <poml>
2 <task>Review a collection of text files.
3   Generate a summary of their key themes.</task>
4 <hint>Large files (>10KB) can be skipped.</hint>
5 <let name="files" src="files.json" />
6 <div for="file in files">
7   <p>File name: {{ file.name }}</p>
8   <document if="file.size < 10*1024" src="{{ file.path }}" />
9   <p else>File size is too large ({{ file.size / 1024 }} KB).</p>
10 </div>
11 </poml>
```

**Figure 5: Example of POML’s templating engine (§ 4.4):** using `<let>` to load data (from `files.json`), for attribute to iterate over items, `{{ ... }}` for variable substitution (e.g., `{{ file.name }}`), and `if/else` attributes for conditional component rendering based on data values (embedding a `<document>` only if `file.size` is below a threshold)

variable substitution using `{{variable}}` syntax, iteration over data collections via the `for` attribute, conditional rendering using `if/else`-like attributes, and inline variable definition with the `<let>` tag. Figure 5 demonstrates these capabilities, showing how data loaded via `<let>` is iterated over with `for`, variables like `{{ file.name }}` are substituted, and component rendering is controlled conditionally based on `file.size`.

This integrated approach offers advantages over manual string manipulation, which is often error-prone and hard to debug [15]. It reduces redundancy by enabling reusable prompt structures populated with varying data. The declarative nature enhances readability, providing a clearer view of the prompt structure. Furthermore, the engine is independent from other programming languages, differentiating POML from other solutions relying on existing templating





**Figure 6: Real-time diagnostic feedback from the POML toolkit in the IDE.** (a) Errors indicated directly in the editor via inline highlighting (e.g., missing attribute). (b) Hovering over an error shows detailed validation messages and documentation. (c) A dedicated panel lists all detected issues (syntax errors, validation problems, file processing issues).

engines (e.g., [3]), contributing to a cohesive framework where structure, presentation, data integration, and dynamic logic are managed within a single language. Details on the templating engine’s syntax and capabilities are available in Appendix B.

## 5 Development Toolkit

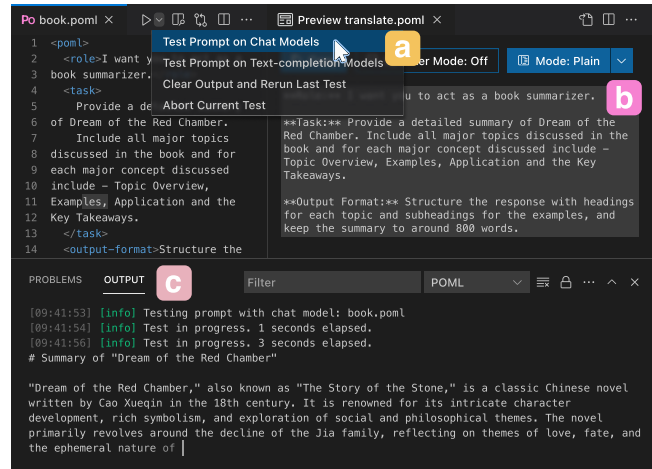
Complementing the core POML language specification, a comprehensive development toolkit is provided to enhance the efficiency, reliability, and collaborative nature of the prompt engineering process, directly addressing (DG4). While POML’s architecture allows potential integration with specialized prompt engineering environments like PromptIDE [79] or ChainForge [9], our initial implementation focuses on providing a feature-rich experience within Visual Studio Code, a widely adopted and familiar IDE. This toolkit encompasses two key aspects: (1) integrated POML IntelliSense for rapid issue detection and code authoring assistance, and (2) Software Development Kits (SDKs) enabling integration of POML into established software development ecosystems.

### 5.1 POML IntelliSense

POML IntelliSense refers to a suite of integrated features designed to simplify the prompt authoring process within the IDE, providing real-time feedback and assistance.

**Syntax Highlighting** Syntax highlighting, adapted from standard HTML conventions due to POML’s similar structure, visually differentiates component tags, attributes, and content. This improves code readability and aids in identifying structural elements.

**Hover and Auto Completion** Informative hover tooltips and context-aware auto-completion enhance productivity. Hovering over POML elements displays dynamically retrieved documentation, definitions, and examples (Figure 6 (b); Figure 1 (d)), ensuring up-to-date information. As the developer types, the system provides context-aware auto-completion suggestions (Figure 1 (f)). It intelligently suggests valid POML component names based on the



**Figure 7: Integrated interactive prompt testing within the POML development environment.** (a) Users initiate or abort tests against specific model types (e.g., chat or text-completion) via a context menu directly from the editor. (b) The live preview panel displays the prompt being tested. (c) The output panel shows test progress logs and streams the LLM’s response in real-time.

current nesting context, available attributes for a given component, and permissible values for certain attributes (e.g., predefined style options), reducing errors and accelerating development by facilitating syntax discovery.

**Inline Diagnostics** Real-time inline diagnostics provide immediate feedback within the editor (Figure 6; Figure 1 (e)). Potential syntax errors, invalid attributes, or structural issues are automatically detected and highlighted with descriptive messages attached to the problematic lines for rapid identification and correction. The system employs error tolerance, attempting to recover from detected issues and reporting multiple issues simultaneously (Figure 6 (c)), allowing developers to address several problems efficiently.

**Live Preview** Working in concert with inline error reporting is a dynamic live preview feature (Figure 7 (b); Figure 1 (g)). This feature renders a visual representation of the final prompt structure as the POML code is being written, updating automatically with each change. This immediate visual feedback helps verify component hierarchy, content, styling, and templating, reducing the cognitive load [8] of mentally translating markup. The preview mechanism supports multiple viewing modes tailored to different needs: a “speaker mode” for chat-based LLMs, a “text mode” for text completion models, a visually enhanced “render mode” for readability, and a “plain text mode” showing the exact raw string for debugging subtle spacing issues.

**Interactive Testing** Integrated interactive testing allows developers to initiate LLM tests directly within the editor (Figure 7). Users can select model types (e.g., chat, text-completion) via a context menu (Figure 7 (a)) and test the current prompt without context switching. The output panel (Figure 7 (c)) displays test progress

```

1 import { poml } from "poml";
2 const data = readPopulationData();
3 const prompt = <Text>
4   <Task captionStyle="bold">
5     Forecast the population in 2050.
6   </Task>
7   <Table syntax="csv" records={data} />
8 </Text>;
9 const messages = await poml(prompt);
10 const response = await invokeLm(messages);

```

(a) Node.js SDK Example Usage

```

1 import poml
2 data = open("globalwarming.pdf", "rb").read()
3 prompt = poml.Prompt()
4 with prompt:
5   with prompt.task():
6     prompt.text("Will 2025 be the warmest year?")
7     prompt.document(buffer=data)
8 messages = prompt.render()
9 response = openai.chat.completions.create(messages)

```

(b) Python SDK Example Usage

**Figure 8: Using POML Software Development Kits (SDKs) to integrate into programming workflows. (a) JavaScript/TypeScript SDK example using JSX-like tagged template literals. (b) Python SDK example using a context manager approach.**

logs and streams the LLM’s response in real-time, enabling immediate observation and early abortion of the generation in case it starts to deviate from the desired outcome (Figure 7 (a)). This tight feedback loop facilitates iterative refinement and assessment of prompt compatibility across LLMs, addressing model sensitivities (DG3) and accelerating the write-debug-validate cycle.

**LSP Server** The IntelliSense features described above – including hover information, auto-completion, and inline diagnostics – are powered by a dedicated implementation of the Language Server Protocol (LSP) [52]. Leveraging the standardized LSP allows these rich language-specific capabilities to be provided consistently. Caching and throttling are implemented at the side of LSP server, ensuring that the client (e.g., VSCode) remains responsive and does not experience performance degradation during heavy usage. This architectural choice not only enhances the development experience currently within VSCode but also creates the potential for integrating POML IntelliSense into other LSP-compatible editors (e.g., Emacs, PyCharm) in the future.

## 5.2 Node.js and Python SDKs

To facilitate integration into larger software applications, POML provides Software Development Kits (SDKs) for the Node.js (JavaScript / TypeScript) and Python environments, chosen for their prevalence in web development and AI respectively (Figure 8).

The **Node.js SDK** (Figure 8 (a)) allows programmatic POML definition using familiar JavaScript paradigms like JSX syntax or React-style functional components [33, 90]. This approach supports creating typed, composable prompt components that can be versioned and shared using standard tools like npm. It also enables developers to define custom, extensible POML components using

TypeScript or JavaScript, allowing the core language to be augmented within specific application contexts.

Similarly, the **Python SDK** (Figure 8 (b)) offers an idiomatic interface using context managers and a fluent API builder pattern, inspired by libraries like `yattag` [43]. This design facilitates integration into Python-based AI/ML workflows, for instance, dynamically creating prompts within data processing pipelines or interacting with frameworks like `LangChain` [3]. Importantly, the Python tooling also supports a standalone mode, allowing direct processing of .poml files via scripts or command-line tools, ideal for rapid testing as demonstrated in our TableQA case study (§ 7.2).

Crucially, both SDKs support programmatic generation and manipulation of stylesheets. This capability enables advanced use cases, such as the automated exploration and optimization of prompt styles. They are also compatible with popular LLM client libraries, simplifying the integration of POML-generated prompts into the LLM communication layer. Furthermore, the language-agnostic design of the core POML specification also allows for the potential development of SDKs in other languages (e.g., C#, Java, Go) in the future.

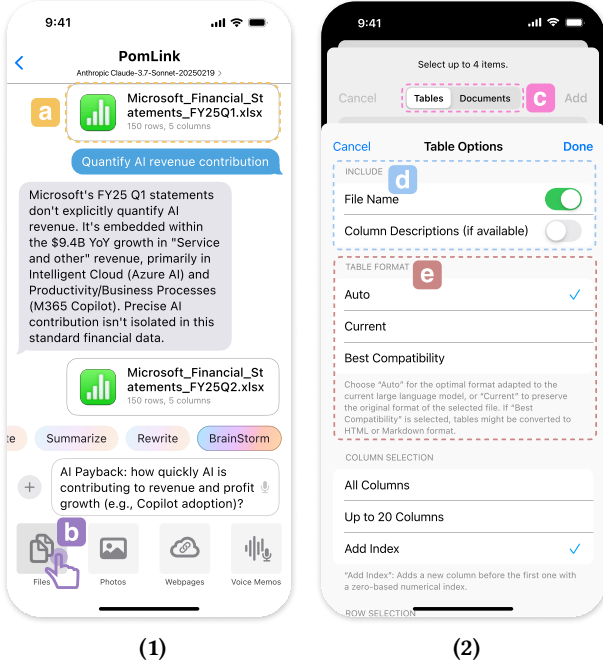
## 6 Implementation

**Core Technology Stack** The core POML engine is developed using TypeScript, leveraging the React framework combined with Server-Side Rendering (SSR) [71]. React’s component model naturally supports POML’s design philosophy centered on separable and reusable prompt components. SSR facilitates the efficient handling of asynchronous operations, such as processing external data sources like images or documents referenced within prompts.

**Codebase Overview** The POML codebase currently comprises approximately 14.8k lines of TypeScript code. This implementation realizes the POML language specification, encompassing 37 distinct components (categorized into 14 basic structural, 7 data, 11 intention, and 5 utility components) and a total of 283 attributes. The codebase is divided into approximately 11.3k lines for the core engine functionality, 3.4k lines for the Visual Studio Code extension features (§ 5.1), and 96 lines for a lightweight Python SDK wrapper (§ 5.2). Comprehensive testing validates the implementation through 10 test suites containing 115 distinct test cases. These tests cover diverse scenarios, including basic tag parsing, attribute validation, complex multi-modal data integration, nested component structures, and the correct application of styling rules, verifying robustness of core POML features.

**Rendering Architecture** The POML implementation employs a three-pass rendering architecture to enhance modularity and extensibility. First, a **Parser** pass validates the input POML markup, executes templating logic (loops, conditionals), resolves variables, applies styles, and transforms the source code into React JSX components. Second, **React** processes these components, and generates a comprehensive **Intermediate Representation (IR)**. This IR is a structured tree containing the resolved content, computed styles, and associated metadata. Third, a dedicated **Writer** pass traverses the IR and serializes it into the final target output format (e.g., Markdown, JSON, plain text). This separation of concerns significantly enhances flexibility; for instance, it enables support for new





**Figure 9: PomLink iOS interface, powered by POML. (1) Main chat screen showing LLM interaction: (a) linked Excel data within the conversation; (b) buttons for adding various file types. (2) Table configuration panel, derived from POML’s <table> component features, where users can select (c) documents/tables, configure (d) inclusion options (e.g., file name), and choose (e) table formatting preferences (including an “Auto” option)**

output formats by implementing additional Writers without modifying the core parsing or IR generation stages. Furthermore, this architecture potentially allows support for alternative input syntaxes beyond XML, such as Markdown extensions or YAML, by developing corresponding Parser modules. It also facilitates performance optimizations like IR caching. A detailed explanation of the rendering architecture is available in Appendix C.

## 7 Case Studies

### 7.1 PomLink – Prototype of iOS Agent

To demonstrate POML’s utility in developing complex, data-intensive applications, we created PomLink, an iOS agent prototype. This case study highlights POML’s capability to streamline the integration of diverse data sources within a functional LLM-powered mobile application. PomLink functions as an LLM agent that utilizes context derived from “linked” files, which users can upload via dedicated interface buttons supporting various types including documents, tables, images, webpages, and voice memos (Figure 9 (1)), all of which can be tailored to personalized settings. For example, tapping the “Files” button (Figure 9 (b)) initiates a pop-up window, asking the user to select multiple documents or tables from the file system (Figure 9 (c)), with an optionally adjusted inclusion style, which are then linked to the agent’s session. Once files are linked, users

can have a conversation with the agent (Figure 1 (2); Figure 9 (1)). This interface enables both customized “ask anything” queries and the use of predefined task shortcuts like “Translate”, “Summarize”, “Rewrite”, and “Brainstorm”. These shortcuts operate on the context provided by *both* the linked files and the rich multi-modal chat history. The **primary objectives** of this case study were twofold: first, to validate POML’s effectiveness by integrating it into a functional mobile application prototype; and second, to assess how POML’s features aided the development process, particularly in handling diverse data inputs and managing prompt complexity within this application.

PomLink’s core agent functionality relies extensively on POML for structuring the prompts used to interact with LLMs. This structure incorporated several key POML features, visible in the example prompt within the development environment shown in Figure 1 (b). Specifically, we used:

- An <include> component to modularly incorporate a common system prompt stored in a separate file, enhancing reusability and maintainability.
- The <conversation> component to present the conversational history between the user and the agent systematically, crucial for maintaining context in interactive sessions.
- Components like <role>, <task>, and <output-format> to provide a clear structure based on the RTF prompting framework [104], organizing instructions for both the LLM and human developers.
- Various data components (§ 4.2) to integrate content from the linked files selected by the user.
- The styling system (§ 4.3) to define formatting rules centrally and create distinct style profiles (e.g., compact vs. verbose), allowing adaptation to specific tasks or backend LLMs.

Integrating varied data sources into the prompt context is a key challenge in developing applications like PomLink. POML addresses this by providing components to reference diverse data types structurally and the flexibility to specify their processing and presentation to the LLM. For instance, the <document> component was used to embed content from user-selected text document files (e.g., PDF documents). Attributes within this component allowed control over details like page selection or whether to preserve original formatting and embedded images. Similarly, the <table> component ingested data from table files (e.g., Excel or CSV). This component allowed specifying input formats and ensured consistent rendering of table data (e.g., as Markdown or CSV) within the prompt. Figure 9 (2) displays the user-facing options derived directly from POML’s <table> component capabilities for customizing data processing, such as toggling file name inclusion (Figure 9 (d)), incorporating column descriptions, selecting rendering formats, and choosing between full or partial table rendering. Notably, the “Auto” table format option (Figure 9 (e)) leverages findings from our TableQA study (§ 7.2) to automatically select the empirically determined best format for the target LLM. Representing other data types like voice memos (<audio>), web pages (<webpage>), images (<img>), or managing conversation history (<conversation>) was achieved straightforwardly using POML’s unified framework for context data aggregation.

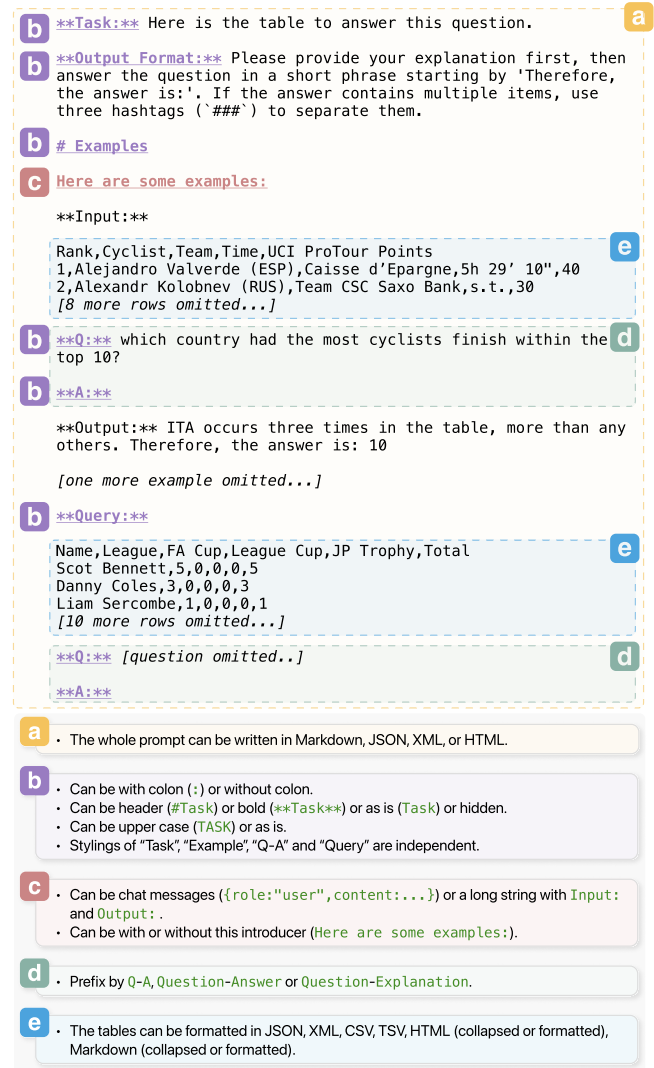
PomLink was implemented using the React Native framework, a popular choice for cross-platform mobile application development. The development process for PomLink highlighted POML’s value for rapid prototyping. A functional prototype was completed by *one of the POML development team in just 2 days*. Around 90% of this time was dedicated to iOS environment configuration, UI development, and simulator deployment, rather than core prompt engineering tasks. The application’s core logic required *only 6 POML prompts*: 1 common system prompt, 4 task-specific prompts, and 1 prompt processing “ask-anything” queries, averaging a concise *35 lines of code* each. This efficiency stems from POML’s ability to abstract complex data handling. The built-in data components managed file parsing and formatting, eliminating the need for custom data processing code typically required in traditional approaches. This allowed the developer to focus primarily on the application’s logic and user interface. Furthermore, prompt styling was managed effectively through POML’s internal styling system; presentation adjustments (e.g., captions, spacing, syntax) were made efficiently by modifying the central stylesheet JSON (Figure 1 (c)), without altering the main prompt logic.

The development workflow for PomLink was significantly accelerated by POML’s toolkit (§ 5), particularly the VSCode extension. Within the PomLink project, the extension’s live preview (Figure 1 (g)) provided immediate visual feedback on prompt rendering, including embedded document files (Figure 1 (h)), while integrated diagnostics (Figure 1 (e)) instantly flagged errors, streamlining debugging. Hover documentation (Figure 1 (d)) and auto-completion (Figure 1 (f)) reduced the need to consult external references. The interactive testing feature (Figure 7) allowed rapid iteration by enabling direct LLM calls and response viewing within the editor, considerably shortening the debug cycle for PomLink’s prompts. Integration into the PomLink application’s React Native backend was achieved using the POML Node.js SDK (Figure 8 (a)), which handled loading POML files, injecting dynamic data, and rendering the final prompt string. The text-based, modular nature of POML prompts also integrated seamlessly with Git for version control within the project (Figure 1 (a)).

In summary, the PomLink case study validated POML’s effectiveness as a practical tool for real-world LLM application development. It demonstrated how POML’s features streamline development: its structural markup (DG1), its data components (DG2), styling system (DG3), and development toolkit (DG4) collectively enabled rapid prototyping and simplified the management of complex, multi-modal prompts. The ability to create a versatile agent application prototype in just two days highlights POML’s potential to accelerate innovation in developing sophisticated LLM-powered software by managing prompt structure and integration complexities effectively.

## 7.2 TableQA – Deep Exploration of Prompt Styling

Prior research indicates that the specific structure and format used to present information within a prompt can significantly influence LLM performance [73, 76, 89]. However, systematically exploring this styling aspect is challenging due to the combinatorial complexity of generating prompt variants. This case study investigates



**Figure 10: Visualization of the prompt styling search space explored in our TableQA experiment. The diagram illustrates the dimensions of prompt variation investigated, including overall syntax structure (a), section formatting options (b, d), example presentation styles (c), and table representation alternatives (e). Systematically combining these choices resulted in an extensive search space of 74k unique prompt styles.**

the impact of prompt styling on LLM performance for the task of table-based question answering (TableQA). Our **primary goals** were twofold: first, to quantify the sensitivity of different LLMs to stylistic variations in prompts, and second, to demonstrate how POML’s features, particularly its styling system, facilitate systematic management and exploration of these prompt styles.

To conduct this investigation, we leveraged POML’s separation of content (markup) and presentation (stylesheets). We authored a single base POML prompt defining the TableQA task structure,

**Table 1: Optimal prompt styling configurations that yielded the highest accuracy for each LLM on the WikiTQ subset. Shows the specific combination of syntax and formatting choices (e.g., Overall Syntax (Figure 10 (a)), Example Body format (Figure 10 (c)), Table Syntax (Figure 10 (e))) for the best-performing style among 100 sampled styles.**

Model	Overall Syntax	Table Syntax	Instruction Header	Example Caption	Example Body	QA Caption	QA Body	User Input Caption
Claude 3 Haiku	XML	HTML (Ugly)	-	-	-	-	-	-
DeepSeek V3	Markdown	HTML	Plain-Upper	Hidden	Hidden	Bold-Colon	Question-Answer	Hidden
Gemini 2.0 Flash	HTML	HTML	Header-Upper	Hidden	Introducer	Bold-Colon	Q-A	Header-Upper
GPT-3.5 Turbo	Markdown	TSV	Plain-Upper-Colon	Header	Hidden	Hidden	Hidden	Hidden
GPT-4o Mini	Markdown	Markdown	Plain-Upper	Header	Hidden	Bold-Colon	Question-Answer	Plain-Upper
LLaMA 3 70B	Markdown	Markdown	Header-Colon	Header-Colon	Chat w. Introducer	Header-Colon	Question-Explanation	Header-Colon
Mistral AI 8x7B	Markdown	XML	Bold-Upper-Colon	Bold-Upper-Colon	Chat	Bold-Colon	Question-Answer	Bold-Upper-Colon
Phi-3 Medium	JSON	CSV	-	-	-	-	-	-

**Table 2: Impact of prompt styling variations on LLM accuracy for TableQA (WikiTQ subset). Shows minimum/maximum accuracy across 100 sampled styles, performance difference (Diff), relative improvement ((Max-Min)/Min), and style ranking stability (Self-correlation: mean Spearman correlation of style rankings across 1000 random data splits, see Figure 14).**

Model	Acc Min	Acc Max	Diff	Self-corr.
Claude 3 Haiku	0.138	0.555	0.417 (+303%)	0.866
DeepSeek V3	0.682	0.823	0.141 (+21%)	0.348
Gemini 2.0 Flash	0.717	0.830	0.113 (+16%)	0.314
GPT-3.5 Turbo	0.060	0.618	0.558 (+929%)	0.767
GPT-4o Mini	0.622	0.753	0.131 (+21%)	0.473
LLaMA 3 70B	0.152	0.657	0.505 (+333%)	0.606
Mistral AI 8x7B	0.177	0.481	0.304 (+172%)	0.849
Phi-3 Medium	0.007	0.322	0.314 (+4450%)	0.931

including intention components, few-shot examples, and placeholders for the question and table data (provided in Appendix E). We then defined a set of stylistic variations within POML stylesheets (example also available in Appendix E). These stylesheets controlled elements like the overall prompt syntax (e.g., Markdown, JSON, XML), the rendering format of embedded table data (e.g., CSV, Markdown, HTML), the presentation style of instructions and examples (e.g., caption styles, visibility, chat formatting), and the structure of the question-answer section, as visualized in Figure 10.

From this extensive space, we randomly sampled 100 styles without replacement for evaluation. These 100 distinct prompt styles were tested on 283 samples (10% of the data) from the WikiTableQuestions (WikiTQ) validation dataset [64]. Accuracy was measured by comparing the LLM’s generated answer and the ground truth answer with evaluation tools provided by WikiTQ. The evaluation included 8 low-cost LLMs (priced under \$0.3 per million input tokens as of February 2025): Claude-3-Haiku [6], DeepSeek-V3 [22], Gemini-2.0-flash [21], GPT-3.5-turbo-0125 [61], GPT-4o-mini [63], LLaMA3-70B [84], Mistral-AI-8x7b [38], and Phi3-medium [55].

**Results** The experimental results confirmed a significant dependency between prompt styling and TableQA accuracy, as detailed in Table 2. The degree of sensitivity to styling, however, varied considerably across models. Some models were extremely sensitive: GPT-3.5-Turbo’s accuracy ranged from 6% to 61.8%, a relative improvement of 929%, while Phi-3 Medium improved by 4450% (from 0.7% to 32.2% accuracy) between its worst and best styles. Other

**Table 3: Table format preferences across LLMs for the TableQA task. For each model, the table shows the top two preferred table syntax formats (e.g., CSV, Markdown, HTML) based on the highest average exact match accuracy achieved by prompt styles utilizing that specific format within our sample of 100 styles. The corresponding mean accuracy for each format is shown in parentheses.**

Model	Best Format	Second Best Format
Claude 3 Haiku	CSV (0.449)	Markdown (0.424)
DeepSeek V3	HTML (0.801)	JSON (0.800)
Gemini 2.0 Flash	XML (0.791)	HTML (0.788)
GPT-3.5 Turbo	Markdown (Collapse) (0.524)	CSV (0.510)
GPT-4o Mini	JSON (0.702)	Markdown (0.694)
LLaMA 3 70B	Markdown (0.601)	HTML (0.601)
Mistral AI 8x7B	XML (0.365)	HTML (Ugly) (0.348)
Phi-3 Medium	CSV (0.169)	JSON (0.116)

models, such as DeepSeek V3 and Gemini 2.0 Flash, were less sensitive, with performance varying by 21% and 16% respectively, highlighting model-dependent sensitivity. The self-correlation scores in Table 2 (especially 0.931 for Phi-3 Medium, 0.866 for Claude 3 Haiku) suggest that the relative performance ranking of different styles is stable for many models across different data subsets. Our further analysis revealed that the optimal prompt style is model-specific. Table 1 details the combination of styling features (e.g., overall syntax, table syntax, example formatting) that yielded the best performance for each model. This model-specific optimal styling highlights the challenge of prompt portability and the need for adaptable styling mechanisms. Overall, this TableQA case study underscores the critical role of prompt styling in achieving optimal LLM performance, which POML addresses.

A critical styling dimension for TableQA is the format used to represent the table data itself. Our results, summarized in Table 3, reveal diversity in the optimal table formats across different models, consistent with previous findings [36, 103]. While some models performed best with simple formats like CSV, others preferred structured representations like HTML, XML, JSON, or Markdown variants. This variation emphasizes the importance of tailoring table representations to the target LLM for optimal performance. These findings directly informed the development of our PomLink application (§ 7.1). The “Auto” table format option within PomLink uses these results, automatically selecting the table syntax identified as optimal for the user’s chosen backend LLM based on this study.

On the other hand, the experiment demonstrates the value of POML’s design, especially the separation of content (markup) from presentation (stylesheets) (DG3). This study required only one single, concise base POML prompt (30 lines of code, see Appendix E). By programmatically combining stylesheet variations, we generated a combinatorial search space encompassing 73,926 unique prompt style configurations derived from the single base POML file. Using the POML Python SDK (§ 5.2), we loaded the base “.poml” file and rendered final prompts by combining it with the generated stylesheets and the specific table/question data for each test case. By enabling systematic variation and optimization of styling parameters independently of the core prompt logic, POML provides a mechanism for adapting prompts to different LLM characteristics and maximizing performance through experimentation, reducing the engineering effort compared to manually managing numerous prompt variations.

## 8 User Study

We conducted a formal user study to evaluate the usability, effectiveness, and developer experience of POML and its toolkit in practical scenarios. The study aimed to assess POML’s utility for representative prompt engineering tasks and gather qualitative feedback to identify its strengths and limitations.

**Participants** We recruited seven participants with diverse technical backgrounds, including software engineers, researchers, and students. Participants had no prior exposure to POML or its use cases before the study. Participant backgrounds, self-reported prompt engineering experience levels, and prior use of LLMs in application development are detailed in Table 4. This diversity was sought to gather feedback reflecting different potential user experiences and needs, especially regarding prior experience developing LLM-related applications.

**Tasks** We designed five distinct tasks of increasing complexity to probe different facets of POML’s capabilities.

- **Task 1 (T1)** involved rewriting an existing plain text prompt into POML and subsequently restyling its output presentation to YAML format using POML’s styling features.
- **Task 2 (T2)** focused on utilizing POML’s document handling features (<document>) to process TODO items embedded within a Microsoft Word document, showcasing its ability to integrate and present document content effectively.
- **Task 3 (T3)** required participants to analyze stock market data (alongside its visualization) provided in an Excel spreadsheet (<table>) and prompt an LLM for investing recommendations, testing table integration and image presentation.
- **Task 4 (T4)** explored meta-prompting concepts by asking participants to use an LLM (initially GPT-4o) to generate POML code based on a list of POML examples; they were then asked to adapt the prompt to target a smaller code-completion model, CodeGemma-2B [82], evaluating POML’s role in managing prompts for different models and the perceived difficulty of this adaptation.
- **Task 5 (T5)**, the most complex task, involved translating content from two subtitle files (in TSV format, one from a 22-minute animation, the other from a 45-minute news report)

while preserving specific formatting conventions and timestamps, testing POML’s ability to read, process, and present tabular data subsets in specific formats under constraints.

The code and data associated with these tasks are available in the supplementary material.

**Procedure** Each participant engaged in a single session lasting approximately 90 minutes. To manage session time effectively, participants were assigned a randomly selected subset of the five tasks. Participants used their own laptops, requiring only VSCode, without additional runtime environment requirements. Sessions began with participants watching a 7-minute introductory video tutorial explaining POML’s core concepts and usage. Subsequently, they were provided with two quickstart code examples and a concise Markdown manual detailing POML syntax and features. They were provided with a VSCode workspace containing the task descriptions, data, and preconfigured access to the GPT-4o and CodeGemma-2B LLMs. Participants were instructed to use a think-aloud protocol, verbalizing their thought processes, challenges, and insights while completing the assigned tasks within the provided VSCode environment equipped with the POML language extension. Following the task completion phase, participants engaged in a semi-structured verbal interview guided by a predefined list of 11 key questions. These questions explored their task experience (completion, time, difficulties), prior prompt usage, perceived usefulness of POML (scenarios, valuable/difficult features, learnability), potential workflow integration, feedback on the VSCode tooling, suggestions for improvement, and any encountered bugs (the full list is provided in Appendix F).

All user study sessions were screen and audio recorded for subsequent detailed analysis. Comprehensive telemetry data was automatically logged, including the time spent on each task, interactions with Language Server Protocol (LSP) features (such as hover help activations, code completion triggers and acceptances), and the complete version history of the POML code written by each participant. Facilitators maintained a minimal intervention policy, offering assistance only when a participant encountered a significant impediment preventing further progress. The evaluation involved a mixed-methods approach, combining qualitative analysis of the verbal feedback and facilitator observations with quantitative analysis of task completion metrics and telemetry data.

### 8.1 Task Completion and Component Usages

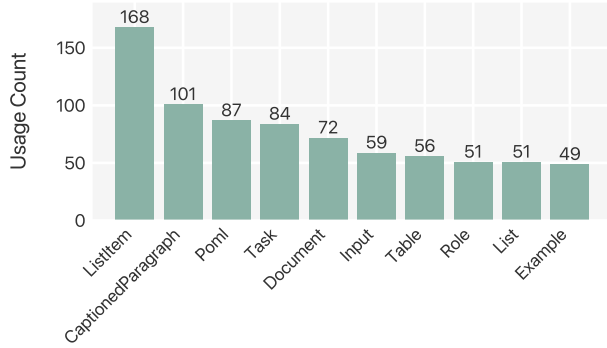
Participants demonstrated high overall task completion rates across the assigned activities. Table 4 details the completion status for each participant across the five tasks. Average task completion times varied considerably based on task complexity and participant differences, as also shown in Table 4. Simpler tasks, such as T1 (Rewriting/Restyling) and T2 (Document TODOs), were typically completed within approximately 10–15 minutes by assigned participants, suggesting effective basic usability.

Tasks evaluating fundamental POML capabilities were generally completed successfully. Participants generally completed T1 and T2 relatively quickly. T2 specifically demonstrated the perceived value of POML’s document handling features, although one participant (P7) noted occasional performance lag with larger documents. T3 (Stock Analysis) was also successfully completed by most assigned



**Table 4: User study results showing task completion and tool interaction metrics. Left Columns: Participant backgrounds. Middle Columns: Task completion status (✓ = completed, ⚡ = partially completed, — = not assigned) and approximate completion times. Right Columns: Interaction metrics (hover help usage, code completion attempts, items per suggestion, suggestion acceptance rate). Task complexity increased from T1 (prompt rewriting) to T5 (subtitle translation), with T5 requiring more time. Participants showed varied engagement with IDE assistance features.**

Participant	Prompt Experience		Task Completion					LSP Metrics			
	Rate	LLM-App-related	T1	T2	T3	T4	T5	Hover	Completion	Items/Comp	Accept (%)
P1	Novice	No	—	✓ 10min	✓ 10min	—	✓ 1h	109	908	2.38	21.21
P2	Advanced	Yes	✓ 20min	—	✓ 15min	⚡ 20min	—	68	286	2.18	49.37
P3	Intermediate	No	—	✓ 10min	—	—	✓ 1h	35	402	3.00	57.14
P4	Intermediate	No	—	✓ 10min	✓ 10min	⚡ 10min	—	24	226	2.50	50.00
P5	Novice	No	✓ 10min	✓ 5min	—	—	⚡ 40min	21	544	2.55	27.27
P6	Advanced	Yes	—	✓ 10min	✓ 15min	✓ 20min	—	120	1282	2.11	23.66
P7	Intermediate	Yes	—	✓ 10min	✓ 15min	✓ 15min	⚡ 20min	22	1734	2.16	50.82
<b>Average</b>	—	—	15min	9min	13min	16min	45min	57.0	768.9	2.41	39.92



**Figure 11: Frequency of POML component usage across user study sessions. The chart shows the total count of each component type used by participants. List items (`<item>`) were most frequently used, followed by captioned paragraphs (`<cp>`), and the root `<poml>` tag. Data components (`<document>`, `<table>`) and intention components (`<task>`, `<role>`) were also commonly used, showing engagement with POML’s data and intention components.**

participants, effectively demonstrating POML’s integration and presentation capabilities for tabular data. Some participants, however, encountered minor difficulties with the specific syntax for selecting table data subsets in T3. T4 (Meta-Prompting) presented greater challenges. While most participants successfully generated initial POML code using GPT-4o, adapting the prompt for the smaller CodeGemma-2B model proved difficult within the session’s time constraints, reflecting inherent challenges in model-specific prompt adaptation. T5 (Subtitle Translation) represented the most complex task, exhibiting the lowest completion rate and the longest average completion times among assigned participants, as detailed in Table 4. Common difficulties included managing precise TSV formatting, handling potential token limits with the longer input file, and preventing undesired line merging or timestamp inaccuracies. This task effectively probed the limits of managing complex format preservation under constraints using POML.

Analysis of component usage frequency, illustrated in Figure 11, reveals participants’ engagement with POML’s core structuring elements. Basic structural components like list items (`<item>`) and captioned paragraphs (`<cp>`) were the most frequently employed after the root `<poml>` tag. Intention components (e.g., `<task>`, `<input>`, `<role>`, `<example>`) and data integration components (`<document>`, `<table>`) also saw significant usage. This pattern indicates that participants actively used POML’s features for structured prompting and data handling throughout the study tasks.

Interaction metrics, presented in the rightmost columns of Table 4, indicate diverse patterns of reliance on the provided IDE tooling assistance. We observe that hover help usage varied widely among individuals, while code completion features were frequently triggered by most participants. However, the acceptance rates for completion suggestions (ranging from 21% to 57%) suggest mixed perceived utility or accuracy of the LSP-provided suggestions. Observations during the sessions confirmed differing programming habits: some participants (P3, P4, P7) used auto-completion extensively, while others (P6) frequently used the hover feature to explore documentation, and some relied more on manual typing or consulting the provided documentation.

## 8.2 User Experience and Qualitative Feedback

Participants’ experiences with POML were largely positive, shaped by their technical backgrounds and task complexity. Feedback consistently highlighted POML’s utility in structuring prompts and integrating data, while also pinpointing areas for potential refinement (§ 9).

**Feedback on Core Language Features** POML’s capability for integrating diverse data types emerged as a frequently praised strength. Participants consistently valued the built-in components for handling various file types (`<document>`, `<table>`, etc.), emphasizing the reduction in manual data preparation effort compared to traditional methods. For example, P7 noted the convenience of reading “various different files... PDFs, Word documents, and Excel files,” while P3 found features like “selected records” and “selected columns ... quite helpful.” P4 appreciated how POML simplified previously “troublesome” tasks like “reading documents... reading images,” enabling quick data loading via simple commands. This

positive reception aligns with the observed usage of data components in tasks T2, T3, and T5, as shown in Figure 11.

The styling system was also well-received. Participants recognized the value of format control (P1 found the format changing feature “quite cool,” a sentiment echoed by P7 regarding the value of format control). They also praised its suitability for structured tasks like format conversion. P5 described POML as “suitable for structured work... especially for file format conversion.” While tasks within the study can often be resolved with inline attributes only (e.g., in T1, T3, T5), P2 successfully used an external `<stylesheet>` in T1, demonstrating the feasibility of the approach. However, deeper engagement with complex stylesheets was limited by the study’s scope, leaving the full potential for systematic format management largely recognized conceptually.

**Development Experience** The integrated development toolkit received predominantly positive feedback for enhancing the prompt authoring workflow. The live preview feature was universally praised for providing immediate visual feedback on the rendered prompt structure, thereby reducing cognitive load and improving clarity (P1, P4). The integrated testing feature, particularly its real-time streaming of LLM outputs, was also valued, with P6 finding the streaming capability “surprisingly impressive.” Participants frequently adopted an iterative edit-preview-test cycle within the IDE, consistent with common prompt engineering practices [102]. Usage patterns for IDE assistance features varied, as discussed in § 8.1. Hover help was frequently used by some participants. For instance, P6 mentioned “frequently using the hover feature to explore available documentation”. Auto-completion saw extensive use by others (e.g., P3, P4, P7), indicating active engagement with these assistance tools despite differing individual habits.

**Learning Curve** Participants’ prior technical experience, particularly familiarity with HTML/XML markup, significantly influenced their perceived learning curve. Those with relevant backgrounds (P2 found it “elegant like React”; P7 was “familiar with markup languages”) generally found the syntax intuitive. Conversely, participants without such experience (P1 stated “I don’t know HTML very well ... don’t know attributes need to be quoted”) reported a steeper initial learning curve and relied more heavily on the live preview and provided documentation. P6 initially found the number of components somewhat “overwhelming.” The provided documentation (video, examples, manual) served as crucial initial learning resources, frequently consulted by participants (especially P1, P3, P6). Despite varying adaptation speeds, all participants appeared to gain comfort with basic POML usage within 10–15 minutes of starting their first task, aided by the IDE’s live preview and diagnostics.

**Workflow Integration** Participants described diverse existing strategies for managing prompts, ranging from informal methods like using note-taking applications (P1 stated “I currently manage my prompts through OneNote”) or simple text files (P5), to embedding prompts directly within application code (P2 mentioned “If I write in Python, I often write prompts messily [in Python string]”). Despite these varied habits, a majority (5 of 7) expressed potential for integrating POML into their workflows, particularly for two types of use cases. First, for managing frequently used or complex

*personal prompts*, participants valued the organization POML provides (P1 explained “if I put them in a [code] repo, I can find them,” contrasting with the disorganization of OneNote). Second, for *developing LLM applications*, participants with LLM-App development experience saw significant value in treating POML files like source code. This approach facilitates standard version control practices (P2 and P6 explicitly mentioned “Git integration”) and improves collaboration through shared, standardized prompt definitions (P6 reflected that “it might be more conducive to my subsequent prompt management”). Integration with Python, though not reflected in the given tasks, was frequently cited as critical for adoption in application development contexts (P1, P5, P6), as well as programmatic prompt generation and incorporation into existing AI/ML pipelines.

**Envisioned Values** Beyond the study tasks, participants envisioned several practical applications where POML’s structure and data handling would be beneficial. P1 suggested using POML to manage context for writing tasks by embedding related research papers as documents: “let it read several related works first... then let it follow some instructions to write [my article].” P3 similarly saw potential for academic writing support. P2 highlighted the convenience for table processing, stating that without POML, it “would require writing a lot of code and formatting... but with POML it would be convenient,” envisioning its use for “running benchmarks.”

**Unexpected Usage Patterns and Desires** Observations revealed potential use cases beyond POML’s primary design focus. Some participants (P1, P3, P6) appeared to use data components (e.g., `<document>`) partly as a convenient way to preview file contents within the IDE, suggesting secondary value as an integrated data viewer. P3 appreciated how POML could “import a file directly, then help me read it out, help me display it.” Occasional comments also hinted at a desire among some advanced users (P6, P7) for lower-level customization options regarding component rendering. Interest was also expressed in exploring agent-like capabilities by writing multiple distinct prompts within a single `.poml` file, with P5 suggesting “potential for extending the framework towards more complex task orchestration.”

**Usage Overhead** The frequent use of intention and structural components (Figure 11) suggests that most participants appreciated the explicit structure POML imposes for managing prompt complexity in such scenarios. However, a trade-off regarding usage overhead was noted by several participants (P1, P4, P6). For simple, one-off prompts, the explicit component structure could introduce typing or cognitive overhead compared to plain text, described by P1 and echoed by P5 as potentially “using a sledgehammer to crack a nut.” While POML’s design attempts to minimize this for simple cases (e.g., the root `<poml>` tag is optional, as discussed in § 4.1), the need to create a file and the mental burden caused by `.poml` suffix still represents a higher initial cost than typing directly into a chat interface. This feedback suggests POML’s value proposition is strongest for complex, data-intensive prompts where structure, lifecycle maintainability, and reusability are paramount, while it may be overkill for simple, transient plain text prompts.



## 9 Discussion

The case studies and user study offered practical insights into POML’s application and adoption. Although POML’s strengths were confirmed, the evaluation also highlighted areas for refinement based on user feedback and task challenges. These findings inform this discussion on improvements, limitations, and future directions.

### 9.1 Suggestions for Improvement

**Usability and Developer Experience Improvements** Participants suggested several usability enhancements to improve the learning curve and development workflow. A recurring request was for more comprehensive documentation, including searchable references and practical examples covering advanced styling and templating (P2 asked for “documentation that’s easy to look up”). Clearer, more actionable error messages were desired to expedite debugging, as current messages were sometimes found opaque (P1 and P3 stated “Error messages are hard to understand”). Language or tooling simplifications, such as sensible default attributes (P2) or clearer distinctions between similar components (P6), were also proposed. Performance optimization, particularly for large document handling (P7) and code assistance responsiveness, remains crucial. Refining auto-completion accuracy (P7) and mitigating conflicts with tools like GitHub Copilot (P2) were noted as important for a smoother experience.

**Feature Requests** User requests focused on enhancing workflow integration and expanding POML’s functionality. Participants desired options to directly save LLM test outputs to files (P6), bypassing manual copying. Improved mechanisms for managing multi-turn conversations within the POML tool were commonly sought by almost all participants. Greater user control over LLM selection and API key management during testing was also requested (P3, P5). Support for additional document formats like LaTeX, particularly for academic use cases (P3), was suggested. Some novice users indicated a desire for integrated guidance on prompt best practices or model-specific formatting (P1), reflecting known challenges in prompt engineering interfaces [102].

### 9.2 Limitations

**Accessibility** The current POML VSCode extension (especially the live preview) presents significant accessibility barriers. It lacks optimization for users relying on assistive technologies like screen readers or requiring high-contrast modes. Key areas needing improvement include keyboard navigability, text equivalents for UI elements, sufficient color contrast, and font size scaling support. Addressing these gaps based on established guidelines [92] is essential to ensure broader usability.

**User Study** Our user study findings (§ 8) are subject to methodological limitations. The 90-minute sessions restricted deep exploration of advanced features. The small participant sample ( $N = 7$ ), though diverse, consisted mainly of technical users, potentially limiting generalizability. The controlled lab setting might not fully reflect real-world project complexities. Furthermore, most participants did not deeply engage with complex stylesheets during the tasks. To help mitigate these limitations, we are actively working

to open-source POML on GitHub and release the extension on the VSCode Extension Marketplace, which will allow us to gather feedback from a broader and more diverse audience using POML in real-world contexts.

### 9.3 Future Work

Based on user feedback, case study insights, and identified limitations, future work will prioritize enhancing POML’s usability, features, and ecosystem. We will keep improving the developer experience through better accessibility support, augmented documentation, refined error diagnostics, performance optimizations, and enhanced SDK usability. Key feature developments include more sophisticated multi-turn conversation management, flexible output handling (e.g., direct file saving), and an expanded template library, potentially incorporating model-specific styling insights (§ 7.2). The structured nature of POML also makes it suitable for exploring automated prompt engineering techniques [44, 48, 105], which is another promising direction for future research.

Beyond core enhancements, we aim to promote POML adoption and demonstrate its value across various domains. We envision POML as a general-purpose prompt markup language. Potential applications include serving as a meta-prompt language in agent systems [50], a structured approach to configure applications like Cursor [2], or a facilitator for industry-level LLM response evaluation workflows as seen in tools like the Azure AI SDK [1]. Further investigation into its utility in complex RAG pipelines, educational tools, domain-specific component libraries, and collaborative prompt engineering environments is also warranted.

## 10 Conclusion

This paper presented POML, a novel Prompt Markup Language designed to address critical challenges in contemporary prompt engineering. POML introduces a component-based structure for enhanced clarity and maintainability, specialized components for effective integration of diverse data types, and a decoupled styling system to systematically manage LLM format sensitivity. Complementing the language, an integrated development toolkit, featuring IDE support and multi-language SDKs, aims to streamline prompt authoring, testing, and management. The effectiveness and practical utility of POML were empirically validated through two distinct case studies alongside a formal user study assessing usability and developer experience. Collectively, these contributions establish POML as a structured, maintainable, and versatile paradigm that addresses prevalent prompt engineering difficulties, thereby enhancing developer workflows, prompt reusability, and collaborative efforts in building sophisticated LLM applications.

## References

- [1] Azure AI. 2024. Evaluate your Generative AI application locally with the Azure AI Evaluation SDK. <https://learn.microsoft.com/en-us/azure/ai-foundry/how-to/develop/evaluate-sdk>.
- [2] Cursor AI. 2025. Generate Cursor Project Rule (.mdc). <https://cursor.directory/generate>.
- [3] LangChain AI. 2022. LangChain. <https://www.langchain.com/>.
- [4] LangChain AI. 2023. LangSmith. <https://www.langchain.com/langsmith>.
- [5] Angular. 2010. Introduction to components and templates. <https://v17.angular.io/guide/architecture-components>.
- [6] Anthropic. 2024. The Claude 3 Model Family: Opus, Sonnet, Haiku. <https://assets.anthropic.com/m/61e7d27f8c8f5919/original/Claude-3-Model-Card.pdf>.

- [7] Anthropic. 2024. Use XML tags to structure your prompts. <https://docs.anthropic.com/en/docs/build-with-claude/prompt-engineering/use-xml-tags>.
- [8] anysphere. 2024. priompt. <https://github.com/anysphere/priompt>.
- [9] Ian Arawjo, Chelse Swoopes, Priyan Vaithilingam, Martin Wattenberg, and Elena L Glassman. 2024. ChainForge: A Visual Toolkit for Prompt Engineering and LLM Hypothesis Testing. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–18.
- [10] AutosseyAI. 2024. prxmpt. <https://github.com/AutosseyAI/prxmpt>.
- [11] Stephen H Bach, Victor Sanh, Zheng-Xin Yong, Albert Webson, Colin Raffel, Nihal V Nayak, Abheesht Sharma, Taewoon Kim, M Saiful Bari, Thibault Fevry, et al. 2022. Promptsource: An integrated development environment and repository for natural language prompts. *arXiv preprint arXiv:2202.01279* (2022).
- [12] Tim Berners-Lee. 1991. Re: status. Re: X11 BROWSER for WWW. <https://lists.w3.org/Archives/Public/www-talk/1991SepOct/0003.html>.
- [13] Bradybry. 2023. ChatXML: A proposal for a structured LLM prompt method. <https://github.com/Bradybry/chatXML>.
- [14] Tom B Brown. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165* (2020).
- [15] Fernando Miguel Carvalho, Luis Duarte, and Julien Gouesse. 2020. Text web templates considered harmful. In *Web Information Systems and Technologies: 15th International Conference, WEBIST 2019, Vienna, Austria, September 18–20, 2019, Revised Selected Papers 15*. Springer, 69–95.
- [16] ChatPDF. 2024. ChatPDF. <https://www.chatpdf.com/>.
- [17] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374* (2021).
- [18] Nan Chen, Yuge Zhang, Jiahang Xu, Kan Ren, and Yuqing Yang. 2024. VisEval: A Benchmark for Data Visualization in the Era of Large Language Models. *IEEE Transactions on Visualization and Computer Graphics* (2024).
- [19] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168* (2021).
- [20] Hai Dang, Lukas Mecke, Florian Lehmann, Sven Goller, and Daniel Buschek. 2022. How to prompt? Opportunities and challenges of zero- and few-shot learning for human-AI interaction in creative applications of generative models. *arXiv preprint arXiv:2209.01390* (2022).
- [21] Google Deepmind. 2024. Introducing Gemini 2.0: our new AI model for the agentic era. <https://blog.google/technology/google-deepmind/google-gemini-ai-update-december-2024/>.
- [22] DeepSeek-AI. 2024. DeepSeek-V3 Technical Report. *arXiv:2412.19437* [cs.CL].
- [23] Ant Design. 2015. Ant Design. <https://ant.design/>.
- [24] Michael Desmond and Michelle Brachman. 2024. Exploring Prompt Engineering Practices in the Enterprise. *arXiv preprint arXiv:2403.08950* (2024).
- [25] Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Jason E Weston. [n. d.]. Chain-of-Verification Reduces Hallucination in Large Language Models. In *ICLR 2024 Workshop on Reliable and Responsible Foundation Models*.
- [26] MDN Web Docs. 2025. CSS: Cascading Style Sheets. <https://developer.mozilla.org/en-US/docs/Web/CSS>.
- [27] MDN Web Docs. 2025. HTML: HyperText Markup Language. <https://developer.mozilla.org/en-US/docs/Web/HTML>.
- [28] MDN Web Docs. 2025. JavaScript. <https://developer.mozilla.org/en-US/docs/Web/JavaScript>.
- [29] Django documentation. 2024. Templates. <https://docs.djangoproject.com/en/5.1/topics/templates/>.
- [30] edspencer. 2024. mdx-prompt. <https://github.com/edspencer/mdx-prompt>.
- [31] Bahare Fatemi, Jonathan Halcrow, and Bryan Perozzi. 2024. Talk like a Graph: Encoding Graphs for Large Language Models. In *The Twelfth International Conference on Learning Representations*.
- [32] fixie ai. 2024. AIJSX — The AI Application Framework for Javascript. <https://github.com/fixie-ai/ai-jsx/>.
- [33] gensx inc. 2024. gensx. <https://github.com/gensx-inc/gensx>.
- [34] Handlebars. 2011. Handlebars. <https://handlebarsjs.com/>.
- [35] Jia He, Mukund Rungta, David Koleczek, Arshdeep Sekhon, Franklin X Wang, and Sadid Hasan. 2024. Does Prompt Formatting Have Any Impact on LLM Performance? *arXiv preprint arXiv:2411.10541* (2024).
- [36] Xinyi He, Yihao Liu, Mengyu Zhou, Yeye He, Haoyu Dong, Shi Han, Zejian Yuan, and Dongmei Zhang. 2025. TableLoRA: Low-rank Adaptation on Table Structure Understanding for Large Language Models. *arXiv:2503.04396* [cs.CL] <https://arxiv.org/abs/2503.04396>
- [37] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring Massive Multitask Language Understanding. In *International Conference on Learning Representations*.
- [38] Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Léo Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2024. Mixtral of Experts. *arXiv:2401.04088* [cs.LG] <https://arxiv.org/abs/2401.04088>
- [39] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. 2024. SWE-bench: Can Language Models Resolve Real-world Github Issues?. In *The Twelfth International Conference on Learning Representations*. <https://openreview.net/forum?id=VTF8yNQm66>
- [40] Jinja. 2022. Jinja – Jinja documentation. <https://jinja.palletsprojects.com/en/3.1.x/>.
- [41] latitude dev. 2024. latitude-llm. <https://github.com/latitude-dev/latitude-llm>.
- [42] Avraham Leff and James T Rayfield. 2001. Web-application development using the model/view/controller design pattern. In *Proceedings fifth ieee international enterprise distributed object computing conference*. IEEE, 118–127.
- [43] leforestier. 2016. yattag. <https://github.com/leforestier/yattag>.
- [44] Zekun Li, Baolin Peng, Pengcheng He, Michel Galley, Jianfeng Gao, and Xifeng Yan. 2024. Guiding large language models via directional stimulus prompting. *Advances in Neural Information Processing Systems* 36 (2024).
- [45] Demiao Lin. 2024. Revolutionizing Retrieval-Augmented Generation with Enhanced PDF Structure Recognition. *arXiv:2401.12599* [cs.AI] <https://arxiv.org/abs/2401.12599>
- [46] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual Instruction Tuning. *arXiv:2304.08485* [cs.CV] <https://arxiv.org/abs/2304.08485>
- [47] Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang. 2023. AgentBench: Evaluating LLMs as Agents. *arXiv preprint arXiv:2308.03688* (2023).
- [48] Yuanye Liu, Jiahang Xu, Li Lina Zhang, Qi Chen, Xuan Feng, Yang Chen, Zhongxin Guo, Yuqing Yang, and Peng Cheng. 2025. Beyond Prompt Content: Enhancing LLM Performance via Content-Format Integrated Prompt Optimization. *arXiv:2502.04295* [cs.CL] <https://arxiv.org/abs/2502.04295>
- [49] Arvid Lunnemark. 2023. Prompt Design. <https://arvid.xyz/posts/prompt-design/>.
- [50] mannaandpoem. 2025. Building OpenManus as a Service. <https://openmanus.org/>.
- [51] mattnigh. 2023. ChatGPT3-Free-Prompt-List: A free guide for learning to create ChatGPT3 Prompts. <https://github.com/mattnigh/ChatGPT3-Free-Prompt-List>.
- [52] Microsoft. 2016. Language Server Protocol. <https://microsoft.github.io/language-server-protocol/>.
- [53] Microsoft. 2023. Guidance. <https://github.com/guidance-ai/guidance>.
- [54] Microsoft. 2023. Semantic Kernel. <https://github.com/microsoft/semantic-kernel>.
- [55] Microsoft. 2024. Phi-3 Technical Report: A Highly Capable Language Model Locally on Your Phone. *arXiv:2404.14219* [cs.CL] <https://arxiv.org/abs/2404.14219>
- [56] Yasuhiko Minamide. 2005. Static approximation of dynamically generated web pages. In *Proceedings of the 14th international conference on World Wide Web*. 432–441.
- [57] Aditi Mishra, Utkarsh Soni, Anjana Arunkumar, Jinbin Huang, Bum Chul Kwon, and Chris Bryan. 2023. Promptaid: Prompt exploration, perturbation, testing and iteration using visual analytics for large language models. *arXiv preprint arXiv:2304.01964* (2023).
- [58] MUI. 2014. Material UI. <https://mui.com/>.
- [59] OpenAI. 2020. Playground. <https://platform.openai.com/playground>.
- [60] OpenAI. 2023. ChatGPT. <https://chatgpt.com/>.
- [61] OpenAI. 2024. GPT-3.5 Turbo. <https://platform.openai.com/docs/models/gpt-3.5-turbo>.
- [62] OpenAI. 2024. GPT-4 Technical Report. *arXiv:2303.08774* [cs.CL] <https://arxiv.org/abs/2303.08774>
- [63] OpenAI. 2024. GPT-4o System Card. <https://cdn.openai.com/gpt-4o-system-card.pdf>.
- [64] Panupong Paspatur and Percy Liang. 2015. Compositional Semantic Parsing on Semi-Structured Tables. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. 1470–1480.
- [65] pezzolabs. 2024. pezzo. <https://github.com/pezzolabs/pezzo>.
- [66] {Structured} Prompt. 2023. {Structured} Prompt. <https://structuredprompt.com/>.
- [67] promptfoo. 2023. promptfoo. <https://www.promptfoo.dev/>.
- [68] PromptML. 2023. PromptML (Prompt Markup Language). <https://www.promptml.org/>.
- [69] React. 2013. Introducing JSX. <https://legacy.reactjs.org/docs/introducing-jsx.html>.
- [70] React. 2013. React. <https://react.dev/>.
- [71] React. 2024. Server Components. <https://react.dev/reference/rsc/server-components>.

- [72] Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha. 2024. A systematic survey of prompt engineering in large language models: Techniques and applications. *arXiv preprint arXiv:2402.07927* (2024).
- [73] Abel Salinas and Fred Morstatter. 2024. The butterfly effect of altering prompts: How small changes and jailbreaks affect large language model performance. *arXiv preprint arXiv:2401.03729* (2024).
- [74] Tobias Schnabel and Jennifer Neville. 2024. Prompts As Programs: A Structure-Aware Approach to Efficient Compile-Time Prompt Optimization. *arXiv preprint arXiv:2404.02319* (2024).
- [75] Tobias Schnabel and Jennifer Neville. 2024. Symbolic Prompt Program Search: A Structure-Aware Approach to Efficient Compile-Time Prompt Optimization. *arXiv:2404.02319* [cs.CL]. <https://arxiv.org/abs/2404.02319>
- [76] Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2024. Quantifying Language Models’ Sensitivity to Spurious Features in Prompt Design or: How I learned to start worrying about prompt formatting. In *The Twelfth International Conference on Learning Representations*.
- [77] Azure AI Services. 2024. Chat Markup Language ChatML (Preview). <https://learn.microsoft.com/en-us/azure/ai-services/openai/how-to/chat-markup-language>.
- [78] Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. 2021. ALFWorld: Aligning Text and Embodied Environments for Interactive Learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*. <https://arxiv.org/abs/2010.03768>
- [79] Hendrik Strobelt, Albert Webson, Victor Sanh, Benjamin Hoover, Johanna Beyer, Hanspeter Pfister, and Alexander M Rush. 2022. Interactive and visual prompt engineering for ad-hoc task adaptation with large language models. *IEEE transactions on visualization and computer graphics* 29, 1 (2022), 1146–1156.
- [80] Yuan Sui, Mengyu Zhou, Mingjie Zhou, Shi Han, and Dongmei Zhang. 2024. Table Meets LLM: Can Large Language Models Understand Structured Table Data? A Benchmark and Empirical Study. *arXiv:2305.13062* [cs.CL]. <https://arxiv.org/abs/2305.13062>
- [81] Michiaki Tatsubori and Toyotaro Suzumura. 2009. HTML templates that fly: a template engine approach to automated offloading from server to client. In *Proceedings of the 18th international conference on World wide web*. 951–960.
- [82] CodeGemma Team. 2024. CodeGemma: Open Code Models Based on Gemma. *arXiv:2406.11409* [cs.CL]. <https://arxiv.org/abs/2406.11409>
- [83] Gemini Team. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv:2403.05530* [cs.CL]. <https://arxiv.org/abs/2403.05530>
- [84] Llama Team. 2024. The Llama 3 Herd of Models. *arXiv:2407.21783* [cs.AI]. <https://arxiv.org/abs/2407.21783>
- [85] Twig. 2009. Twig - The flexible, fast, and secure template engine for PHP. <https://twig.symfony.com/>.
- [86] uithub. 2024. uithub. <https://uithub.com>.
- [87] Juho Vepsäläinen, Arto Hellas, and Petri Vuorimaa. 2023. The State of Disappearing Frameworks in 2023. In *Proceedings of the 19th International Conference on Web Information Systems and Technologies, WEBIST 2023 (International Conference on Web Information Systems and Technologies, WEBIST - Proceedings)*, Francisco Garcia Penalvo and Massimo Marchiori (Eds.). SciTePress, Portugal, 232–241. doi:10.5220/0012174000003584 Publisher Copyright: Copyright © 2023 by SCITEPRESS - Science and Technology Publications, Lda. Under CC license (CC BY-NC-ND 4.0); International Conference on Web Information Systems and Technologies, WEBIST ; Conference date: 15-11-2023 Through 17-11-2023.
- [88] Juho Vepsäläinen, Arto Hellas, and Petri Vuorimaa. 2023. *The Rise of Disappearing Frameworks in Web Development*. Springer Nature Switzerland, 319–326. doi:10.1007/978-3-031-34444-2\_23
- [89] Anton Voronov, Lena Wolf, and Max Ryabinin. 2024. Mind your format: Towards consistent evaluation of in-context learning improvements. *arXiv preprint arXiv:2401.06766* (2024).
- [90] VSCode. 2024. Prompt-tsx. <https://github.com/microsoft/vscode-prompt-tsx>.
- [91] Vue.js. 2014. Components Basics. <https://vuejs.org/guide/essentials/component-basics.html>.
- [92] W3C. 2025. Accessibility Fundamentals Overview. <https://www.w3.org/WAI/fundamentals/>.
- [93] W3C. 2025. Accessibility Principles. <https://www.w3.org/WAI/fundamentals/accessibility-principles/>.
- [94] W3C. 2025. W3C. <https://www.w3.org/>.
- [95] Ming Wang, Yuanzhong Liu, Xiaoming Zhang, Songlian Li, Yijie Huang, Chi Zhang, Daling Wang, Shi Feng, and Jigang Li. 2024. LangGPT: Rethinking Structured Reusable Prompt Design Framework for LLMs from the Programming Language. *arXiv preprint arXiv:2402.16929* (2024).
- [96] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-Consistency Improves Chain of Thought Reasoning in Language Models. In *The Eleventh International Conference on Learning Representations*.
- [97] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.
- [98] Tongshuang Wu, Ellen Jiang, Aaron Donsbach, Jeff Gray, Alejandra Molina, Michael Terry, and Carrie J Cai. 2022. Promptchainer: Chaining large language model prompts through visual programming. In *CHI Conference on Human Factors in Computing Systems Extended Abstracts*. 1–10.
- [99] Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. 2024. Agentless: Demystifying llm-based software engineering agents. *arXiv preprint arXiv:2407.01489* (2024).
- [100] Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022. Webshop: Towards scalable real-world web interaction with grounded language agents. *Advances in Neural Information Processing Systems* 35 (2022), 20744–20757.
- [101] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2023. ReAct: Synergizing Reasoning and Acting in Language Models. In *The Eleventh International Conference on Learning Representations*.
- [102] JD Zamfirescu-Pereira, Richmond Y Wong, Bjoern Hartmann, and Qian Yang. 2023. Why Johnny can’t prompt: how non-AI experts try (and fail) to design LLM prompts. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–21.
- [103] Yuge Zhang, Qiyang Jiang, Xingyu Han, Nan Chen, Yuqing Yang, and Kan Ren. 2024. Benchmarking Data Science Agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, Bangkok, Thailand, 5677–5700. doi:10.18653/v1/2024.acl-long.308
- [104] Jeffrey Zheng. 2023. Role-Task-Format (RTF) framework for prompting. <https://x.com/thejeffreyzheng/status/166022396975511427>.
- [105] Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. [n. d.]. Large Language Models are Human-Level Prompt Engineers. In *The Eleventh International Conference on Learning Representations*.
- [106] Jingming Zhuo, Songyang Zhang, Xinyu Fang, Haodong Duan, Dahua Lin, and Kai Chen. 2024. ProSA: Assessing and understanding the prompt sensitivity of LLMs. *arXiv preprint arXiv:2410.12405* (2024).

## A POML vs. Other Prompt Markup Languages

There have been various markup languages and frameworks, each addressing specific needs, as summarized in Table 5.

JSX/TSX-based tools like AIJSX [32] and GenSX [33] primarily target agent and workflow orchestration using high-level components within the Node.js ecosystem. Others, such as VSCode Prompt TSX [90], Priompt [8], Prxmpt [10], and MDXPrompt [30], focus on prompt authoring using JSX/TSX or MDX syntax, often integrating closely with specific runtime environments (e.g., Node.js) or prioritizing features like context length management via priority systems. Specialized formats like ChatML [77] define role-based structures for conversational history, while PromptML [68] proposes a custom DSL focused on intention components but currently lacks integrated tooling. SAMMO [75] takes a distinct approach, using a Python API for programmatic manipulation and structure-aware optimization rather than direct markup authoring.

In contrast to these approaches, POML aims for a comprehensive, standalone solution distinguished by several key characteristics evident in the comparison. It provides a language-agnostic runtime, enhancing portability compared to tools tied to Node.js or Python. POML incorporates versatile components covering basic structure, intention expression, and data integration, offering arguably the most extensive built-in support among the listed tools for diverse data types like tables, documents, and webpages. POML also integrates a unique CSS-like styling system for decoupling presentation from content and provides a full VS Code IDE extension, seeking to unify structure, data handling, styling, and development tooling within a single cohesive framework more completely than other approaches listed in Table 5.

## B Templating Engine Design Details

To support the creation of dynamic and data-driven prompts, POML integrates a built-in templating engine. This capability is inspired by widely adopted templating systems in web development frameworks [5, 29, 34, 40, 69, 85, 91] and common practices in existing prompting frameworks [3, 53]. Providing an *integrated* engine enhances POML’s utility, particularly for scenarios where prompts need to be generated dynamically based on runtime data without external scripting languages (DG2), contributing to efficient development workflows (DG4).

The engine allows for the insertion of variable values directly into the prompt text using a familiar double curly brace syntax: `{{variable}}`. Its syntax for variables and control structures draws inspiration specifically from popular engines like Jinja and Handlebars [34, 40]. These variables can represent contextual placeholders that are populated at runtime from external data sources (e.g., documents, table records) or defined programmatically within the POML document itself using the `<let>` tag, enabling dynamic data binding and temporary variable assignment within the template logic.

Beyond simple variable substitution, the templating engine supports essential control structures for generating complex prompt logic dynamically. Iteration over data collections (like lists or arrays) is handled using a `for` attribute, allowing repetitive structures to be generated concisely (e.g., `<list><item for="file in files">{{file.description}}</item></list>` would create a

list item for each object in the `files` collection). Conditional rendering of prompt sections is achieved through `if/else`-like constructs, enabling parts of the prompt to be included or excluded based on the evaluation of variable values or logical expressions (e.g., conditionally including a detailed explanation section only if a `verbose` flag is set to `true`).

The inclusion of this integrated templating engine offers several advantages for prompt engineering. It significantly reduces redundancy by allowing developers to define reusable prompt structures that can be populated with different data inputs, rather than creating many similar static prompts. Compared to constructing prompts through manual string manipulation or concatenation in a general-purpose programming language, which can be often error-prone and less readable [15], POML’s templating approach makes the prompt structure more visually apparent and closer to a “What You See Is What You Get” (WYSIWYG) experience, simplifying the debugging process. The engine’s standalone capability provides dynamic prompting capabilities without mandating the use of an external programming language wrapper, differentiating POML from some other markup approaches that primarily serve as static data formats [13, 30, 68].

## C Three-Pass Rendering Framework

The POML implementation employs a three-pass rendering framework. This architecture processes the POML markup in three distinct stages: first, a “Parser” pass transforms the input into React JSX components; second, React processes these components to build a comprehensive Intermediate Representation (IR); and third, a “Writer” pass converts this IR into the desired final output format (e.g., Markdown, JSON). The primary motivation for this three-pass approach is the clear separation of concerns. The first pass (Parser) focuses exclusively on parsing the input markup, validating its structure, applying templating, stylesheets, and transforming it into React components. The second pass (React processing) focuses on generating a detailed IR with complete metadata. The third pass (Writer) is solely responsible for traversing the validated IR and serializing it into a specific target format. This separation minimizes interdependencies, ensuring that changes related to styling or data handling do not inadvertently break the final rendering logic.

This three-tier separation fosters flexibility and extensibility. Adding support for new output formats only requires implementing a new Writer module, without modifying the core parsing, templating, or styling logic. Feature additions are simplified as each pass addresses a distinct, well-defined set of tasks. From a performance perspective, the generated IR can be cached and reused to produce multiple output formats (e.g., generating both a plain text version and a chat-structured version simultaneously) without incurring the cost of re-parsing the original POML source, which is particularly beneficial for large prompts involving extensive data components.

The first pass, executed by the Parser module, begins by parsing the raw POML markup. It identifies all known POML tags (e.g., `<task>`, `<document>`, `<table>`) and their attributes. The parser validates the markup structure against the POML specification, flagging malformed tags or invalid attributes, often with a degree of error tolerance to handle minor issues without halting processing entirely

**Table 5: Comparison of POML with other Prompt Markup Languages**

Markup Language	Markup Syntax	Runtime Env.	Selling Point	Key Components	Data Support	Styling	Dev Tooling
<b>AI.JSX</b> [32]	JSX/TSX	Node.js	Agents & Workflows	High-level components (Chat completion, Tool use, Q&A)	Relies on external libraries	—	NextJS support, LangChainJS integration
<b>GenSX</b> [33]	JSX/TSX	Node.js	Agents & Workflows	High-level workflow/agent components	Text	—	TypeScript tooling
<b>VSCode Prompt TSX</b> [90]	TSX	Node.js (VS Code Ext)	Prompt (VS Code)	Role components (User/Asst), Tool (<ToolResult>), File (<FileLink>) integration	Chat history, VS Code API data (files), Tool results	Priority system (manages length)	VS Code Ext dev env (TSX/TS support)
<b>Prompt</b> [8]	JSX/TSX	Node.js	Writing Prompt	Message & Structural components (<scope>, <first>, etc.); Priority system (p/prel)	Image support, JSON data	Priority system (manages length)	Preview WebUI
<b>Prxmpt</b> [10]	JSX	Node.js	Writing Prompt	Utility elements (<kv>, etc.), Priority (<priority>), Space control, Data serialization (<json>, <yaml>)	Primitives (<num>), Dates (<datetime>), Objects (<json>, <yaml>)	Basic text format (<b>, <h1>, etc.); Space control	JSX/TS support, Standard Node.js tooling
<b>MDXPrompt</b> [30]	MDX (Mark-down + JSX)	Node.js	Writing Prompt	Mix structured/unstructured text, Default components (<Prompt>, <ChatHistory>)	Chat history, JSON data	—	MDX/React ecosystem
<b>ChatML (OpenAI)</b> [77]	Custom format (Roles)	Language-Agnostic	Writing Prompt	Role-based messages (system, user, assistant, tool)	Text, Images	Customizable chat templates	OpenAI SDKs
<b>PromptML</b> [68]	Custom DSL	Language-Agnostic	Writing Prompt	Intention components (@context, @objective, @instructions, etc.)	Text	—	Python Parser
<b>SAMMO</b> [75]	Python API + Mark-down	Python	Writing Prompt	Programmatic manipulation, Structure-aware optimization algos	Relies on Python data libs	—	Python dev env, LLM Runner integrations
<b>POML</b>	HTML-like DSL	<b>Language-agnostic</b>	Writing Prompt	<b>Versatile (basic structure + intention + data)</b>	<b>Comprehensive (tables, docs, webpages)</b>	<b>CSS-like styling system</b>	<b>Full VS Code IDE extension</b>

(§ 5.1). The template engine (§ 4.4) executes control flow attributes like for-loops and if-else conditions. Variable substitutions using the `{{...}}` syntax are performed, resolving references to internal variables or data loaded from external sources (e.g., via `<let src="...">`). Styling rules are applied from both the `<stylesheet>` tag and inline style attributes, resolving conflicts and computing final style properties for each element. The Parser then transforms these validated elements into React JSX components, preparing them for the second pass.

In the second pass, React processes the JSX components generated by the Parser. Using React’s virtual DOM concept, this pass

constructs a comprehensive Intermediate Representation (IR) - an in-memory tree of structured nodes corresponding to the POML elements, their resolved content, computed styles, and extensive metadata including original source positions and serialization hints. This IR contains 21 distinct node types (Appendix D) and captures all information needed for the final rendering stage, creating a clean, validated representation ready for the third pass.

The third pass is handled by a selected Writer module. Its purpose is to convert the structured Intermediate Representation (IR)



**Figure 12: The POML three-pass rendering architecture: (1)-(2) The Parser transforms POML markup into React JSX components, (2)-(3) React processes these components into a detailed Intermediate Representation (IR) that captures content structure, styling, and metadata, and (3)-(4) specialized Writers convert the IR into various output formats such as JSON or Markdown.**

generated by React into a final, usable output format, such as Markdown, plain text, or a specific JSON structure required by an LLM API.

The Writer traverses the IR tree, visiting each node (e.g., an IR node representing a `<list>`, `<img>`, or `<document>`). Based on the type of the IR node and the target output format, the Writer serializes the node’s content appropriately. For example, a list node might be rendered as a bulleted or numbered list in Markdown, while a table node could be formatted as CSV or a Markdown table depending on the Writer’s configuration and the style attributes captured in the IR. Furthermore, the Writer propagates and computes metadata associated with each node, such as the “speaker” attribute. This allows the Writer to structure the output appropriately for different contexts, for instance, by splitting the final rendered text into distinct sections based on the speaker (e.g., “human”, “system”) if required by the target format (like chat completion APIs).

Different Writer modules can be implemented and selected to produce various outputs from the exact same IR. This architecture

makes it straightforward to extend POML to support new target formats (e.g., LaTeX, specific API JSON schemas) by simply creating a new Writer implementation without modifying the earlier passes.

Writers are designed with error tolerance in mind. They can handle partially formed IR nodes (if errors occurred during earlier passes but were tolerated) or missing data references that could not be resolved. If configuration parameters specific to the writer (e.g., a required CSV delimiter) are missing or invalid, the writer typically alerts the developer but attempts to continue converting the parts of the IR that are unaffected.

This three-pass structure consistently maintains the separation of concerns between prompt logic/content and final presentation/layout. Updates to data sources (src attributes) or styling rules (`<stylesheet>`) captured during the earlier passes generally do not require any changes to the Writer logic. This reinforces POML’s commitment to a modular and robust approach to prompt engineering, aligning with design goals for Reusability (DG1) and Style Management (DG3).



## D Intermediate Representation (IR) Specifications

### D.1 Attributes Applicable to All Tags

The following attributes can be applied to any tag in the representation:

- **speaker** (*ai/human/system*): The speaker of the current content
- **original-start-index** (*integer*): The start offset of the element corresponding to the current one in the original document
- **original-end-index** (*integer*): The end offset of the element corresponding to the current one in the original document

### D.2 Tag Definitions

**any** Represents a generic container for arbitrary data values. Useful for storing dynamic or unstructured content.

- **type** (*string*): The data type of the value ('string', 'integer', 'float', 'boolean', 'array', 'object', 'buffer', 'null', or 'undefined').
- **name** (*string*): An optional identifier for the data.

**b** Represents text that should be displayed in boldface. Useful for highlighting important words or phrases.

**code** Represents a block or inline fragment of code. It can optionally include language and formatting attributes.

- **inline** (*boolean*): Indicates whether the code is inline (true) or a block element (false).
- **lang** (*string*): Specifies the programming language or syntax highlighting mode.
- **blank-line** (*boolean*): Inserts a blank line before and after the code block if inline = false.

**env** Represents a formatting environment or container to specify how nested content should be output.

- **presentation** (*string*): The output style or format mode ('markup', 'serialize', 'free', or 'multimedia').
- **markup-lang** (*string*): The specific markup language, required only if presentation = 'markup'.
- **serializer** (*string*): The name of the serializer, required only if presentation = 'serialize'.
- **writer-options** (*object*): Optional parameters passed to the writer constructor for customizing output.

**h** Represents a heading element.

- **level** (*integer*): Indicates the heading level. Typically ranges from 1 (highest level) to 6 (lowest level).

**i** Represents text that should be displayed in italics. Useful for emphasizing words or phrases.

**img** Represents an image element.

- **base64** (*string*): The base64-encoded image data.
- **alt** (*string*): Alternative text describing the image.
- **position** (*string*): The placement of the image relative to text, such as 'here', 'top', or 'bottom'.

- **type** (*string*): The image MIME type (e.g., 'image/jpeg', 'image/png').

**item** Represents a single item within a list. Typically used as a child element of "list".

**list** Represents an ordered or unordered list of items.

- **list-style** (*string*): The style of the list bullets or enumeration (e.g., 'star', 'dash', 'decimal').

**nl** Inserts newline characters.

- **count** (*integer*): Specifies how many newline characters to insert.

**obj** Represents a data object, typically stored in JSON format.

- **data** (*object*): A valid JSON object containing the structured data.

**p** Represents a paragraph of text. Useful for dividing content into readable blocks.

- **blank-line** (*boolean*): Inserts a blank line before and after the paragraph when true.

**s** Represents text that should be displayed with a strikethrough style.

**span** Represents an inline container for text without additional formatting. Useful for applying attributes without changing display structure.

**table** Represents a table structure containing rows and cells.

**tbody** Represents the body section of a table, containing the majority of data rows.

**tcell** Represents a single cell within a table row.

**text** Represents raw or unformatted text content.

**thead** Represents the header section of a table, typically containing column headings.

**trow** Represents a single row within a table, containing one or more cells.

**u** Represents text that should be displayed with an underline.

## E Details of TableQA Case Study

A detailed example of the prompt structure used in our TableQA case study is provided below, as discussed in the main text. Figure 13 illustrates the components involved.

Figure 13(a) presents the POML template crafted for this task. It defines the overall structure provided to the language model, encompassing the core task instruction, specific requirements for the output format (including explanation prefixes and separators), few-shot examples demonstrating the desired input-output behavior, and the placeholder for the actual query which includes the table data (columns and records) and the user question.

Figure 13(b) shows an example of the corresponding JSON stylesheet. This stylesheet controls the rendering and formatting aspects of the POML prompt, specifying details such as the syntax for different elements (e.g., markdown for POML, csv for tables), styling for instructional text and captions (e.g., bolding, headers, specific question/answer prefixes like "Q"- "A"), and the presentation of examples.

```

1 <poml>
2 <task className="instruction">
3   Here is the table to answer this question.
4 </task>
5 <output-format className="instruction">
6   Please provide your explanation first, then answer the question
7   in a short phrase starting by 'Therefore, the answer is:'.
8   If the answer contains multiple items, use three hashtags
9   (<code>###</code>) to separate them.
10 </output-format>
11 <examples>
12   <example for="example in examples">
13     <input>
14       <table columns="{{ example.table.columns }}"
15         records="{{ example.table.records }}" />
16       <qa>{{ example.question }}</qa>
17     </input>
18     <output>
19       {{ example.explanation }}
20       Therefore, the answer is: {{ example.answer }}
21     </output>
22   </example>
23 </examples>
24
25 <cp className="query" caption="Query" captionSerialized="query">
26   <table columns="{{ query.table.columns }}"
27     records="{{ query.table.records }}" />
28   <qa>{{ query.question }}</qa>
29 </cp>
30 </poml>

```

(a) Prompt (POML) for TableQA Case Study

```

1 {
2   "poml": {
3     "syntax": "markdown"
4   },
5   "table": {
6     "syntax": "csv"
7   },
8   ".instruction": {
9     "captionStyle": "bold",
10    "captionTextTransform": "none",
11    "captionColon": true
12  },
13  "examples": {
14    "captionStyle": "header",
15    "introducer": "Here are some examples:",
16    "chat": false
17  },
18  "example": {
19    "captionStyle": "hidden"
20  },
21  "qa": {
22    "captionStyle": "bold",
23    "questionCaption": "Q",
24    "answerCaption": "A"
25  },
26  ".query": {
27    "captionStyle": "bold"
28  }
29 }

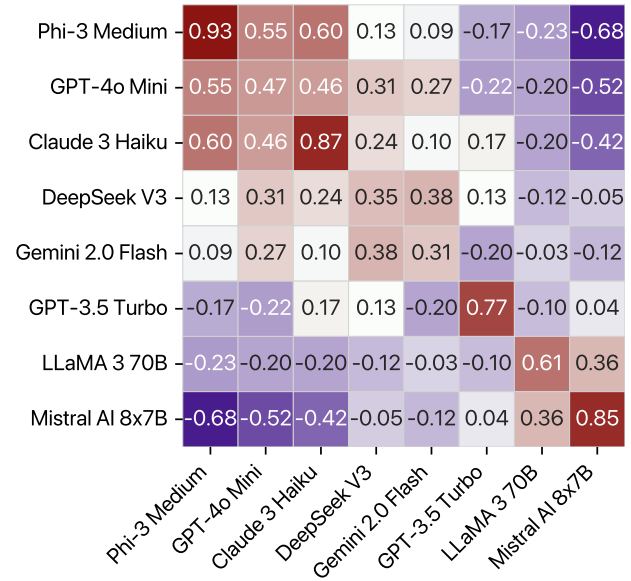
```

(b) Stylesheet Example for TableQA Case Study

**Figure 13: Example of POML usage for the TableQA case study (§ 7.2, Appendix E).** (a) The POML prompt template defining the task instructions, required output format, few-shot examples, and the query structure containing the table and question. (b) The corresponding JSON stylesheet example specifying formatting and presentation rules for elements within the POML prompt, such as overall syntax, caption styles, and introducer text.

This separation of content (POML) and presentation (stylesheet) allows for flexible prompt customization and management.

Analysis of style preferences across models, visualized in the correlation heatmap (Figure 14), revealed distinct model groupings. For instance, Phi-3 Medium, GPT-4o Mini, and Claude 3 Haiku formed a cluster exhibiting similar positive correlations in their style rankings, suggesting shared preferences. Conversely, Mistral



**Figure 14: Correlation matrix showing relationships between LLMs based on their performance rankings across 100 prompt styles in the TableQA task (§ 7.2, Appendix E).** Cells show Spearman correlation coefficients between model style rankings (Red: positive/similar preferences; Blue: negative/dissimilar preferences). The diagonal shows the self-correlation score for each model (see Table 2 for definition and values), indicating style ranking stability. The self-correlation score indicates the stability of the style performance ranking. It is computed by randomly splitting the 283 samples into two equal halves, calculating the accuracy of each of the 100 styles on both halves, and finding the Spearman correlation between the two resulting style rankings. This process is repeated 1000 times, and the mean correlation is reported; higher values indicate more stable style rankings across different data subsets.

AI 8x7B and LLaMA 3 70B formed another group, often showing negative correlations with the first cluster, indicating fundamentally different sensitivities to styling choices. Deeper analysis of individual feature impacts (Table 6) highlighted specific sensitivities with statistical significance. Using plain text for example bodies significantly benefited models like Claude 3 Haiku and Phi-3 Medium but harmed LLaMA 3 70B and Mistral AI 8x7B ( $p < 0.01$ ). Adopting a chat-like format for examples helped LLaMA 3 70B and Mistral AI 8x7B but was detrimental to several others ( $p < 0.01$ ). Hiding structural elements like QA captions negatively impacted most models but surprisingly benefited GPT-3.5 Turbo ( $p < 0.05$ ). Complex interaction effects were also observed, where the combined impact of certain style features differed from the sum of their individual effects, indicating non-linear relationships in how styling elements influence performance.

**Table 6: Impact analysis of specific prompt styling choices on TableQA accuracy across different LLMs (§ 7.2, Appendix E).** Each cell indicates whether a feature (row) significantly improved (✓), had no significant effect (−), or significantly harmed (✗) performance for a given LLM (column), based on statistical significance tests (Mann-Whitney U test comparing styles with vs. without the feature; significance levels indicated). The bottom rows show examples of interaction effects for specific feature combinations.

Style Preference	Claude	DeepSeek	Gemini	3.5-Turbo	4o-Mini	LLaMA-3	Mistral	Phi-3
Example Body = Text	✓ $p < 0.01$	✓ $p < 0.05$	−	−	✓ $p < 0.01$	✗ $p < 0.01$	✗ $p < 0.01$	✓ $p < 0.01$
Example Body = Introducer	✓ $p < 0.01$	−	✓ $p < 0.01$	✗ $p < 0.01$	✓ $p < 0.01$	−	✗ $p < 0.01$	✓ $p < 0.01$
Overall Syntax = HTML	−	−	✓ $p < 0.05$	✗ $p < 0.01$	✓ $p < 0.01$	−	✗ $p < 0.05$	✓ $p < 0.01$
Overall Syntax = XML	✓ $p < 0.05$	−	−	✓ $p < 0.05$	−	✗ $p < 0.05$	✗ $p < 0.01$	−
Table Syntax = CSV	✓ $p < 0.05$	✗ $p < 0.01$	✗ $p < 0.01$	−	−	−	−	✓ $p < 0.01$
Overall Syntax = Markdown	−	−	−	✓ $p < 0.01$	✗ $p < 0.01$	−	✓ $p < 0.01$	✗ $p < 0.01$
Example Body = Chat	✗ $p < 0.01$	−	−	−	✗ $p < 0.01$	✓ $p < 0.01$	✓ $p < 0.01$	✗ $p < 0.01$
Example Body = Chat w. Introducer	✗ $p < 0.01$	✗ $p < 0.05$	✗ $p < 0.05$	−	✗ $p < 0.01$	✓ $p < 0.05$	✓ $p < 0.01$	✗ $p < 0.01$
QA Caption = Hidden	✗ $p < 0.01$	−	✗ $p < 0.01$	✓ $p < 0.05$	✗ $p < 0.01$	−	−	✗ $p < 0.01$
Overall Syntax = Markdown, Example Body = Introducer	✓ $p < 0.01$	✓ $p < 0.05$	✓ $p < 0.05$	−	✓ $p < 0.05$	−	✗ $p < 0.01$	✓ $p < 0.01$
Instruction Header = Header, Example Body = Introducer	✓ $p < 0.05$	−	✓ $p < 0.05$	✗ $p < 0.01$	✓ $p < 0.05$	−	−	✓ $p < 0.05$
Example Body = Introducer, QA Caption = Bold-Upper	✓ $p < 0.05$	−	✓ $p < 0.05$	−	✓ $p < 0.05$	−	✗ $p < 0.05$	✓ $p < 0.01$
Example Body = Introducer, QA Body = Question-Explanation	✓ $p < 0.01$	−	✓ $p < 0.05$	✗ $p < 0.05$	✓ $p < 0.01$	−	−	✓ $p < 0.01$
Example Body = Text, QA Body = Question-Explanation	✓ $p < 0.01$	✓ $p < 0.05$	−	✓ $p < 0.05$	✓ $p < 0.01$	✗ $p < 0.01$	✗ $p < 0.01$	✓ $p < 0.05$
Overall Syntax = Markdown, Example Body = Text	✓ $p < 0.01$	✓ $p < 0.05$	−	✓ $p < 0.01$	✓ $p < 0.01$	✗ $p < 0.01$	✗ $p < 0.01$	✓ $p < 0.01$

## F Semi-structured Verbal Interview Questions

The following questions were used to guide the semi-structured interviews with users. The questions were designed to elicit detailed feedback on the user’s experience with POML, its features, and their suggestions for improvement.

- (1) What tasks do you find most challenging and why?
- (2) Do you have prior experience with using prompts? In what scenarios do you typically employ prompts?
- (3) In your opinion, for which scenarios is POML most useful?
- (4) In your opinion, for which scenarios is POML not useful?
- (5) Which parts or features of POML do you find most valuable, and why?
- (6) Which parts or features of POML do you find most difficult to use, and why?
- (7) Do you consider the POML syntax easy to learn and worthwhile to master?
- (8) How do you foresee integrating POML into your existing workflow or pipeline?
- (9) Do you like the VSCode integration of POML (preview, testing, auto-completion)? Why or why not?
- (10) What improvements or additional features do you suggest for POML or its VSCode integration?
- (11) Did you encounter any bugs or errors, and how did these affect your overall experience of POML?