

Meaning-Typed Programming: Language Abstraction and Runtime for Model-Integrated Applications

JAYANAKA L. DANTANARAYANA, University of Michigan, United States

YIPING KANG, University of Michigan, United States

KUGESAN SIVASOTHYNATHAN, University of Moratuwa

Jaseci Labs, Sri Lanka

CHRISTOPHER CLARKE, University of Michigan, United States

BAICHUAN LI, University of Michigan, United States

SAVINI KASHMIRA, University of Michigan, United States

KRISZTIAN FLAUTNER, University of Michigan, United States

LINGJIA TANG, University of Michigan, United States

JASON MARS, University of Michigan, United States

Software development is shifting from traditional logical programming to *model-integrated* applications that leverage generative AI and large language models (LLMs) during runtime. However, integrating LLMs remains complex, requiring developers to manually craft prompts and process outputs. Existing tools attempt to assist with prompt engineering, but often introduce additional complexity.

This paper presents **Meaning-Typed Programming (MTP)** model, a novel paradigm that abstracts LLM integration through intuitive language-level constructs. By leveraging the inherent semantic richness of code, MTP automates prompt generation and response handling without additional developer effort. We introduce the **by** operator for seamless LLM invocation, **MT-IR**, a meaning-based intermediate representation for semantic extraction, and **MT-Runtime**, an automated system for managing LLM interactions. We implement MTP in **Jac**, a Python superset language and find that MTP significantly reduces coding complexity while maintaining accuracy and efficiency. Our evaluation across diverse benchmarks and user studies demonstrates that MTP outperforms existing frameworks such as DSPy and LMQL by reducing lines of code by factors of 2.3-7.5 \times and 1.3-10.7 \times respectively. For math problems from the GSM8k dataset, MTP achieves accuracy rates approaching 90%, while reducing token usage in 10 out of 13 benchmarks. This leads to cost savings up to 4.5 \times and runtime speedups as high as 4.75 \times . Additionally, MTP demonstrates resilience even when 50% of naming conventions are suboptimal, establishing it as a practical, efficient solution for streamlining model-integrated application development.

1 INTRODUCTION

With the advent of generative AI and Large Language Models (LLMs), the way we develop software is rapidly evolving. While LLMs have been extensively used for code generation [37], another growing trend is their integration as core subcomponents within software to realize critical path features [4, 21, 34, 39]. These *model-integrated applications* leverage generative AI at runtime to perform critical program functionality, merging conventional programming with AI-driven capabilities. This shift fundamentally transforms software development, redefining how applications are designed and built.

Authors' Contact Information: Jayanaka L. Dantanarayana, University of Michigan, Ann Arbor, United States, jayanaka@umich.edu; Yiping Kang, University of Michigan, Ann Arbor, United States, ypkang@umich.edu; Kugesan Sivasothynathan, University of Moratuwa

Jaseci Labs, Moratuwa, Sri Lanka, kugesan.sivasothynathan@jaseci.org; Christopher Clarke, University of Michigan, Ann Arbor, United States, csclarke@umich.edu; Baichuan Li, University of Michigan, Ann Arbor, United States, patrli@umich.edu; Savini Kashmira, University of Michigan, Ann Arbor, United States, savinik@umich.edu; Krisztian Flautner, University of Michigan, Ann Arbor, United States, manowar@umich.edu; Lingjia Tang, University of Michigan, Ann Arbor, United States, lingjia@umich.edu; Jason Mars, University of Michigan, Ann Arbor, United States, profmars@umich.edu.

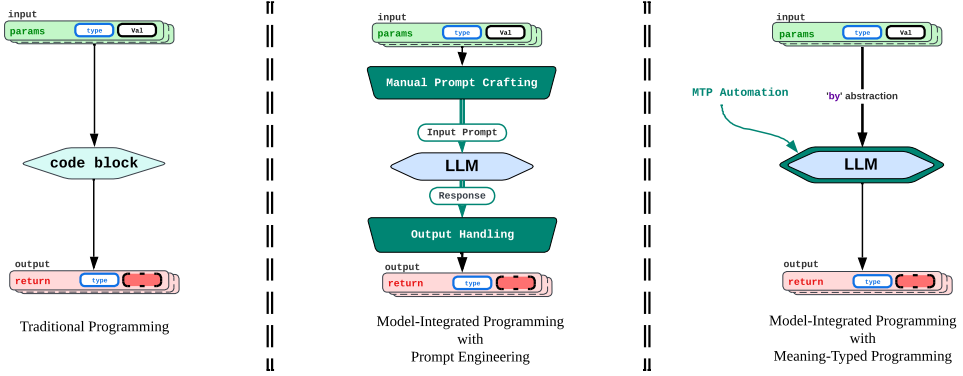


Fig. 1. High level view of what traditional program functionality consists of (Left) vs. model integrated programming using prompt engineering (Middle) and Meaning-Typed Programming(MTP) paradigm (Right).

Building model-integrated applications remains a complex challenge due to the fundamental differences between operating assumptions. In conventional programming, code describes precise operations performed on explicitly defined variables such as objects (Figure 1 left). In contrast, LLMs process natural language text as input and produce text-based outputs. To integrate LLMs into programs, developers must manually construct textual input, a process known as *prompt engineering*.

Prompt engineering is the primary method of realizing model-integrated programming and can be tedious, complex, and imprecise. Developers must craft descriptive instructions, manage textual context, and handle the inherent variability in LLM output. This process includes a significant amount of trial-and-error, making it time-consuming and difficult to maintain.

Several efforts have been made to address these challenges by developing frameworks and tools to facilitate prompt engineering. Frameworks such as LangChain [19], Language Model Query Language (LMQL) [4], DSPy [15], and SGLang [46] have been introduced to assist developers in generating and managing prompts. However, while these frameworks provide more sophisticated tooling to help bridge the gap between conventional programming and the incorporation of AI models, in practice, they introduce additional layers of complexity for developers:

- (1) *Prompt design complexity*: Developers remain responsible for the manual construction of prompts, including determining the appropriate prompt language and selecting relevant information to include.
- (2) *Steep learning curve*: These systems require developers to familiarize themselves with new query languages, frameworks, or specialized syntax, complicating the integration process.
- (3) *Input/output conversion complexity*: Parsing LLM output and converting them back to objects remains a challenge, especially given the variability in output across different LLMs.

In this paper, we introduce *meaning-typed programming* (MTP), a novel approach to simplify the creation of model-integrated applications by (i) introducing new high level language abstractions that hide the complexities of *model-integration*, (ii) a novel technique for automating prompt generation using inherent code semantics, and (iii) a runtime system that automates LLM inference and output handling.

The **key insights** of this work are: (1) *In practice, programs are written to be readable by various developers to comprehend the intentions and semantics conveyed by the code. This semantic richness can be leveraged to automatically translate intent embodied in the code to prompts for the LLM.* (2) *As state-of-the-art LLMs continue to advance, AI becomes increasingly feasible to infer code intentions*

without requiring explicit descriptions or prompt construction by developers. This trend is not only supported by our empirical findings, but we surprisingly find less can perhaps be more in prompt engineering (see § 5).

There are four primary **challenges** that must be overcome in the design and implementation of MTP. Firstly, the introduced language abstraction should be simple, intuitive, and flexible. Secondly, developing an automated approach that provides LLMs with sufficient semantic information for accurate task performance is challenging. Thirdly, since model-integrated applications rely on dynamic values, prompts must be generated at runtime, requiring a lightweight system to handle prompt generation and LLM inference efficiently. Lastly, because LLMs generate non-deterministic outputs, the outputs must be carefully processed to ensure correct type conversion. To address these challenges, we introduce three key components to enable the seamless creation of model-integrated programs. These include

- (1) **The ‘by’ operator**, that enables the seamless integration of LLM functionality into code. It acts as a bridge between traditional and LLM operations, allowing developers to replace function bodies with a simple **by model_name** clause.

```
1 def calc_age(birth: date, today: date) -> int by gpt4():
```

As shown later in our evaluation, meticulously extracted relevant semantic information is sufficient for the LLM to produce the accurate result, significantly simplifying LLM integration into existing codebases.

- (2) **MT-IR**, a meaning-based intermediate representation at the compiler level that collects and organizes semantically rich information, preserving the meaning behind variable names, function signatures, and types, during runtime.
- (3) **MT-Runtime**, an automated runtime engine within the language’s virtual machine, triggered at by call sites. It binds runtime values to MT-IR dynamically, manages LLM interactions, generates context-aware prompts, and handles responses and errors. Through these three components, MTP allows developers to focus on application logic while automatically managing LLM integration, ensuring seamless and contextualized operations.

We implement meaning-typed programming (MTP) in a production-grade Python superset language called Jac [18, 24]. The **by** operator was introduced as a language primitive in Jac, and the MT-IR construction is implemented as a compilation pass in Jac’s ahead-of-time compiler. The MT-Runtime system is implemented as a language plugin which can be separately installed as a PyPi package. This end-to-end framework has been used to develop real-world applications and benchmarks utilized in our extensive evaluation.

Specifically, we make the following contributions:

- We introduce **meaning-typed programming**, a novel paradigm that leverages the inherent semantic richness of code to automate LLM integration into applications without requiring explicit prompt engineering.
- We introduce the **‘by’ operator**, a simple yet powerful language abstraction that provides both intuitive and flexible integration of LLM functionality into code. It enables developers to seamlessly incorporate language models for multiple code constructs including functions, methods, and object initialization while maintaining familiar programming patterns.
- We develop supporting compiler and runtime infrastructure: **MT-IR**, a semantically rich intermediate representation that extracts and organizes code meaning, and **MT-Runtime**, an automated execution engine that dynamically converts semantic information into prompts, manages LLM interactions, and handles type-safe response processing.
- We **implement MTP in Jac**, a production-grade Python superset language, conduct extensive evaluation and show that MTP significantly reduces development complexity, lines of

code, and costs while improving run-time performance and maintaining or exceeding the accuracy of existing approaches.

- We present multiple **case studies** demonstrating the application of MTP in complex problem domains, analyzing its impact on code structure, development workflow, and complexity reduction compared to traditional LLM integration methods.

We perform extensive evaluation to investigate the implications of MTP on code complexity, model-integration accuracy, inference token usage, runtime performance, and cost compared to prior work. We also assess how sensitive MTP is to poor coding practices, as our system is built on leveraging semantically rich code. Our evaluations show that MTP reduces lines of code by factors of 2.3-7.5 \times compared to LMQL and 1.3-10.7 \times compared to DSPy, demonstrating significant code complexity reduction. For accuracy on math problems from the GSM8k [8] dataset, MTP outperformed other frameworks with newer OpenAI and Llama models, achieving accuracies approaching 90%. Regarding inference costs, MTP demonstrated savings over DSPy on most benchmarks, with reductions up to 4.5 \times . Moreover, MTP enables program speedups of up to 4.75 \times over DSPy, indicating minimal runtime overhead. Our sensitivity studies confirm that MTP maintains effectiveness even when 50% of naming conventions are suboptimal. These results establish MTP as a practical, efficient solution for streamlining model-integrated application development.

2 MOTIVATING EXAMPLE

This section presents a motivating example to illustrate the concept of *model-integrated* programming and to highlight the complexities involved in the current implementation technique, i.e., prompt engineering. The example application used here is a video game implementation, where dynamically generating game levels is achieved through model integration with an LLM.

The section is organized as follows. We first introduce the example application and the code sketch of traditional implementation without an LLM (Figure 3). We then demonstrate model-integrated implementation using the industry-standard approach of prompt engineering (Figure 4), highlighting the significant overhead it imposes on developers. Lastly, we present implementation using our proposed MTP approach (Figure 5), demonstrating MTP’s capabilities in removing complexities associated with prompt engineering.

2.1 Traditional Programming

The video game implementation used here is a real program developed based on the pygame [25] tutorials. In this video game, the player must complete the current game level to proceed to the next. An example of the video game level progression is shown in § 5 (Figure 12), where each level includes a set of configurations (maps, enemies, difficulties, etc). Figure 2 presents the basic data structures used in the video game implementation, including the Level object (lines 14–21), which encapsulates the configurations and a nested Map object (lines 9–12).

```

1 class Position:
2     x: int
3     y: int
4
5 class Wall:
6     start_pos: Position
7     end_pos: Position
8
9 class Map:
10     walls: list[Wall]
11     enemies: list[Position]
12     player_pos: Position
13
14 class Level:
15     level_id: int
16     difficulty: int
17     width: int
18     height: int
19     num_wall: int
20     num_enemies: int
21     map: Map

```

Fig. 2. A typical implementation for a video game level development usecase. Code shows how the data-structure definitions for the video game application.

```

1 def get_next_level(prev_level: Level) -> Level:
2     ...
3     ...
4     ## code for generating/fetching the next
5     level
6     ...
7     return next_level
8
9 new_level = get_next_level(curr_level)

```

Fig. 3. Traditional implementation of game levels for the non model-integrated usecase.

```

1  import openai
2  import json
3
4  def get_next_level(prev_level: Level) -> Level:
5      prev_level_json = json.dumps(prev_level.__dict__, default=
        =lambda o: o.__dict__)
6
7      prompt = f"""
8      You are an expert game level generator. Given the current
9      level details, generate the next level with
10     increased difficulty.
11
12     """
13     """ Current Level:
14     {prev_level_json}
15
16     """
17     """ Rules for Next Level:
18     - Increase difficulty
19     - Increase width and height slightly (max 10% each)
20     - Increase the number of walls and enemies based on
21       difficulty progression
22     - Generate a valid map with walls and enemy positions
23     - Ensure 'player_pos' is within bounds and not colliding
24       with walls
25
26     """
27     """ Output Format:
28     Provide the next level in JSON format:
29     {{
30     'level_id': <int>,
31     'difficulty': <int>,
32     'width': <int>,
33     'height': <int>,
34     'num_wall': <int>,
35     'num_enemies': <int>,
36     'map': {{
37       'walls': {{{'start_pos': {{'x': <int>, 'y': <int>}}}},
38       'end_pos': {{'x': <int>, 'y': <int>}}}},
39       'enemies': {{{'x': <int>, 'y': <int>}}}},
40       'player_pos': {{'x': <int>, 'y': <int>}}
41     }}
42     """
43
44     response = openai.ChatCompletion.create(
45         model='gpt-4',
46         messages=[{"role": "system", "content": "You are an
47           expert in procedural game level generation."},
48           {"role": "user", "content": prompt}],
49         temperature=0.7,
50         max_tokens=500
51     )
52
53     level_dict = json.loads(response['choices'][0]['message
54       'content'])
55
56     def parse_position(data: dict) -> Position:
57         return Position(x=data['x'], y=data['y'])
58
59     def parse_wall(data: dict) -> Wall:
60         return Wall(
61             start_pos=parse_position(data["start_pos"]),
62             end_pos=parse_position(data["end_pos"])
63         )
64
65     def parse_map(data: dict) -> Map:
66         walls = [parse_wall(w) for w in data.get("walls", [])]
67         enemies = [parse_position(e) for e in data.get("
68           enemies", [])]
69         player_pos = parse_position(data["player_pos"])
70         return Map(walls=walls, enemies=enemies, player_pos=
71           player_pos)
72
73     next_level = Level(
74         level_id=level_dict["level_id"],
75         difficulty=level_dict["difficulty"],
76         width=level_dict["width"],
77         height=level_dict["height"],
78         num_wall=level_dict["num_wall"],
79         num_enemies=level_dict["num_enemies"],
80         map=parse_map(level_dict["map"])
81     )
82
83     return next_level

```

Fig. 4. Implementation of the video game level generation using prompt-engineering

manually-craft prompt to an LLM, and the output of the LLM (in a text format) is used to construct the next level. The LLM innovation is in lines 37-42, which eliminates the need to manually implement the level generation algorithms. However, this example illustrates the significant complexity surrounding the LLM invocation/integration.

get_next_level The level progression logic is achieved through `get_next_level` function, which uses the previous level as input and returns the next level. Figure 3 shows the code sketch, including a placeholder (lines 2-6) for the function body. In traditional programming, developers must implement complex algorithms and logic to generate the next level based on the current one, often leveraging procedural generation techniques, which can be both challenging and time-consuming. However, by integrating an LLM, a level can be generated automatically without developers' explicit coding.

2.2 State-of-the-Art Model-Integration Approach

Code synthesis vs. Model integration. Before presenting the current approach to model integration, it is important to highlight the difference between model-integration and the use of LLMs for code synthesis [37]. LLM-enabled code synthesis uses LLMs to generate code that is executed at runtime. In contrast, *model-integrated* code [39] leverages an LLM to dynamically perform tasks within a program. During runtime, the LLM inference generates output values based on given inputs to the LLM. In this example, *model-integrated* approach uses the input (previous level object) to query an LLM at runtime. The output of the LLM inference will be used as the next level configuration.

Current approaches to model-integration primarily rely on prompt engineering as the industry standard, which involves crafting input prompts to guide LLMs in generating desired responses. It involves structuring inputs, optimizing the wording and context of the prompt to maximize the accuracy and relevance of the model's output. Figure 4 presents the prompt engineering-based implementation of `get_next_level` functionality, where the input of the previous level object is described in a

Manual Prompt Crafting. The example includes a manually crafted prompt template (lines 7-35). As shown here, to effectively improve LLM accuracy, the prompt consists of several key components: a problem description (line 8), an input object (including a description of the input and the `prev_level` object in JSON format), and a set of rules (lines 14-18) to guide the LLM’s generation. In addition, developers need to explicitly specify the output format (lines 20-34), which in this case follows a JSON schema for the `Level` object. Since the `Level` object contains deeply nested objects, structuring the output format is tedious and labor-intensive. This prompt engineering process is time-consuming and often requires an iterative trial-and-error process.

Output Type Conversion. In addition to prompt crafting, developers must also convert the LLM’s output into a return-typed object compatible with traditional programming. In this example, LLM’s output follows the JSON object format, as specified in the prompt, which then must be converted into a `Level` object that is understandable by the rest of the application (lines 47-70). This includes mapping the JSON output to the corresponding attributes of the `Level` class (lines 62-70). To handle deeply nested objects, three additional helper functions are defined on lines 47, 50, and 56. This added complexity makes model-integration with prompt engineering more cumbersome and less intuitive.

Error Handling and Retries. Due to the non-deterministic nature of LLMs, the generated output may not always be consistent with the expected format, leading to program runtime errors that developers must anticipate and handle. In addition to error-handling, a retry mechanism may be needed to reattempt the LLM invocation process until a correctly formatted output is produced. Although the example code snippet does not include error handling and retries due to space constraints, incorporating such mechanisms adds additional complexity and requires development effort to ensure robustness and reliability.

In conclusion, this example demonstrates that while model-integration can leverage LLMs to improve traditional programming, the complexities associated with prompt engineering may introduce significant cognitive and implementation burden. Several frameworks and tools have been developed to simplify the interaction with LLMs. While these systems offer valuable capabilities for prompt engineering and LLM integration, they usually incur additional complexity. LMQL [4] requires developers to learn a new query language, while DSPy [15] requires the developers to annotate all essential information in code, replacing prompt-engineering complexities with annotation overhead. The comparison to prior work is presented later in § 5.

2.3 MTP Vision for Model-Integrated Programming

To address the challenges in the current model-integration approach, our research aims to *design a language abstraction and system to streamline the process of integrating LLMs into traditional programming by abstracting and automating prompt generation, output type conversions, and error handling*.

Figure 5 illustrates our envisioned approach, where adding only a few lines of code enables LLM integration, eliminating the need for prompt engineering and its associated complexities. Moreover, this approach is significantly more intuitive than prior works. The code snippet provides a preview of what we introduce in this paper as *Meaning-Typed Programming*. Line 5 of the code snippet shows how model-integration is achieved. MTP hides prompt-engineering using a language abstraction, automating the prompt generation, LLM inference and output handling. We will describe details on our language abstraction in § 4.2.

When hiding prompt engineering, the information required for accurate LLM inference should be available to generate prompts. As shown in Figure 5, the intent of the `get_next_level()` function on line 5 is evident from the code’s semantics. This observation leads to the key insight of our

paper: *programs are written to be readable by various developers to comprehend the intentions and semantics conveyed by the code. This semantic richness can be leveraged to automatically translate intent embodied in the code to prompts for the LLM.*

Figure 6 presents three additional illustrative examples showing how meaning can be derived from various code elements:

- (1) The meaning of the variable `dob` can be inferred from its name, type, and value. Here `dob` is a commonly used abbreviation for date of birth which still retains adequate semantic meaning.
- (2) The meaning of the `french2english()` function can be inferred from its name, input parameter, and return type. These semantics indicate the function takes French text as input and returns an English translation of that text.
- (3) The semantics of the `description()` method in the `book` class can be understood by analyzing the class, variable names, and an example function call.

With the help of a compiler and runtime techniques, LLMs can infer the intended meaning from a program to perform model-integration tasks as shown in our evaluations, later in the §5.

3 PROBLEM, OBJECTIVES AND CHALLENGES

Our goal is to leverage semantically rich, human-readable code expressions to automate prompt-engineering. Instead of requiring developers to manually craft detailed prompts or annotations, we harness the inherent structure and semantics of the code itself to generate meaningful LLM inputs.

To achieve this goal, we define **three key objectives**:

- O₁** - Design simple, intuitive and flexible language abstractions for developers while hiding prompt engineering and other complexities.
- O₂** - Design techniques to extract relevant and sufficient semantic information from code for automatic prompt synthesis.
- O₃** - Design a system that combines semantic information and dynamic information to manage prompt synthesis and output interpretation.

When designing a solution that achieves these objectives, **four unique challenges arise**:

- C₁** - How do we design a language abstraction that supports multiple model integration methods while remaining easy to learn and use, minimizing the developer learning curve?
- C₂** - What semantic information must be extracted to generate an LLM prompt that ensures program accuracy without incurring cost and runtime overhead?
- C₃** - What type of runtime system design is appropriate to combine static semantic information with dynamic information to manage LLM input?
- C₄** - How do we design a system that handle output type conversion and error handling without increasing token cost and runtime cost?

```

1 from mtp.core.llms import OpenAI
2
3 llm = OpenAI(model_name="gpt-4o")
4
5 def get_next_level(prev_levels: list[Level]) -> Level by llm
6     ()
7
8 new_level = get_next_level(prev_levels)

```

Fig. 5. Implementation of model-integration as envisioned in this paper, with Meaning-Typed programming.

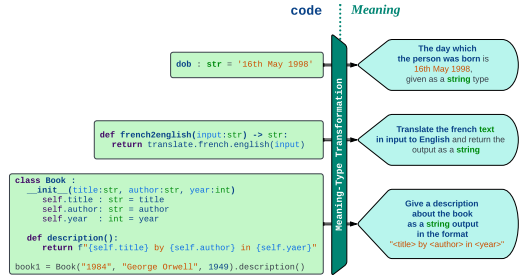


Fig. 6. Examples of *meaning* embedded in code semantics.

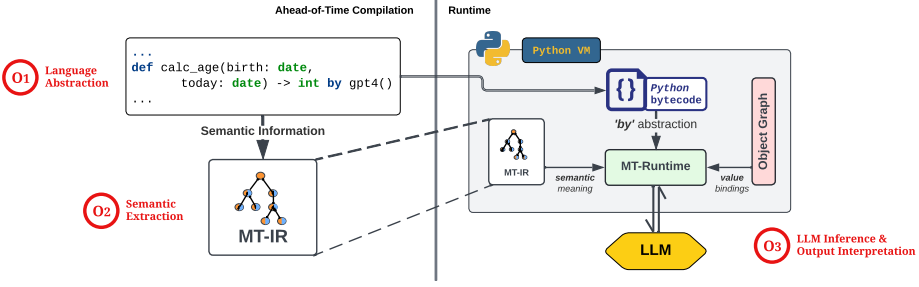


Fig. 7. Meaning-Typed Programming system overview.

C_1 involves designing a simple and intuitive language abstraction, minimizing the learning curve and cognitive load. At the same time, it must be flexible enough to seamlessly integrate across different code locations. Designing a system that is both intuitive and flexible is challenging, as the abstraction must be simple enough to reduce the learning curve while also being usable orthogonally at multiple model-integration points.

C_2 involves extracting only the relevant semantics for model-integration points. A naive approach would be to include the entire source code in the prompt, but this would lead to high token costs, or even context window overflow. To optimize token usage, we must selectively extract only the necessary semantics. But this is particularly challenging because these semantics are often scattered across different code locations, requiring sophisticated codebase-wide analysis to identify and extract only what is necessary for accurate LLM prompting.

C_3 involves combining extracted code semantics with dynamic values to create effective prompts. This requires accessing variable data at runtime, which can vary in complexity depending on where the integration happens. For functions, it is straightforward since values are passed as parameters. However, for object methods, it becomes more complex because the system must also access relevant object attributes at runtime. To handle this, a sophisticated runtime system is needed to capture the execution context and combine dynamic values into the prompts.

C_4 involves correctly interpreting LLM-generated responses and converting them to appropriate types with robust error handling. This is particularly difficult because LLMs generate outputs non-deterministically and often deviate from expected formats. As demonstrated in the prompt-engineering example in §2.2, extracting correct results requires careful parsing, but inconsistencies in output structure make this process complex. Implementing a retry mechanism to correct such inconsistencies requires additional LLM inferences, increasing token usage and cost. Therefore, the system design must balance effective error handling with minimal impact on token consumption and runtime cost.

4 MEANING TYPED PROGRAMMING

To address the problem of automating and streamlining model-integrated application development discussed in §3, we introduce *Meaning Typed Programming (MTP)*. The MTP paradigm (Figure 7) comprises three key components that achieves each objective defined in §3, enabling seamless model-integration into programming languages:

- (1) The **by** operator as the language level abstraction. (To achieve O_1)
- (2) Meaning-Typed Intermediate Representation (MT-IR), generated during compilation. (To achieve O_2)
- (3) MT-Runtime, an automated runtime engine that manages the LLM integration. (To achieve O_3)

This section presents an overview of MTP as well as an in depth view into each component introduced in this paper with respect to the challenges mentioned in §3.

4.1 MTP Overview

The **by** operator acts as a bridge between traditional code and LLM operations, allowing developers to seamlessly integrate LLM functionality into their code by simply adding the **by** keyword. As shown in Figure 7, developers can replace a function’s implementation, such as `calc_age()`, with a straightforward **by** `<model_name>` clause. At runtime, the **by** operator invokes MT-Runtime to execute the desired functionality using the specified LLM, achieving O_1 . In this example, the system leverages GPT-4 to calculate age based on a given birthday and the current date. § 4.2 provides further details on the **by** operator.

MT-IR is an intermediate representation designed to capture and preserve the meaning of the program to be used at runtime upon invocation of the **by** keyword. During compilation, MT-IR is generated by analyzing and organizing semantically rich information in the code (see Figure 7 (left)). It preserves key elements such as variable names, function signatures, and object schema relevant to a given **by** call-site to provide access to the code’s original meaning by the LLM at run-time (see Figure 7 (right)) achieving O_2 . (see § 4.3 for more details).

MT-Runtime manages the complexities of LLM interaction, through prompt synthesis using MT-IR. It enhances execution by intelligent response parsing, and dynamic error handling, achieving O_3 . It operates within the language’s virtual machine (e.g., Python’s VM) and is specifically triggered at **by** call sites. MT-Runtime combines both static information from MT-IR and dynamic values from the language’s object graph. By seamlessly integrating with the language’s execution environment, MT-Runtime allows developers to focus on their application logic while automating LLM interaction (see § 4.4 for more details).

4.2 Meaning-Type Language Construct: ‘by’ operator

Hiding prompt engineering within *a language abstraction that is simple, intuitive and flexible is challenging* as discussed in §3 as C_1 . To address this we introduce the **by** operator as the core language abstraction of MTP. The base syntax of the **by** operator is,

`<code construct> by llm_ref(model_hyperparameters)`

Here, the workload defined by the code construct on the left-hand side of the **by** operator is model-integrated using the LLM specified on the right. Hence, this syntax is simple and intuitive, making it easy to learn. Additionally, this syntax provides flexibility by enabling developers to integrate LLMs at multiple points within standard code constructs. Furthermore, additional hyperparameters can be provided as arguments to the LLM object for further customization.

In this paper, we introduce three code construct locations to use **by** call sites, allowing developers to use MTP flexibly while maintaining simplicity and minimal overhead. These call sites include *function calls*, *object initialization*, and *member method calls*.

Function Calls - A function call includes a name, input arguments, and type annotations. As shown in Figure 8a, the function `calculate_age` determines a person’s age from their `dob` and `cur_year`, with an integer output indicating the number of years.

Object Initialization - In Figure 8b, the **by** call uses the partially initialized object (name being “Einstein”) to predict and fill in the remaining attributes (`dob`). MTP uses the class name, attribute names, types, and the provided value to generate a prompt for the LLM and converts the results to the correct type using MT-Runtime.

Member Method Call - Unlike standalone function calls, class methods access instance attributes, which are essential for semantic extraction. In Figure 8c, the `calculate_age` method of `Person` derives `dob` from the `Einstein` object instead of an argument. Integrating **by** operator invokes MT-runtime, passing method semantics, class name, attributes, and types via MT-IR.

4.2.1 Formal Semantics of the "by" Operator.

To precisely define the behavior of the **by** operator, we introduce a formal semantics that specifies its operational characteristics across different integration points. This formalization provides a rigorous foundation for understanding how the **by** operator transforms traditional code constructs into model-integrated operations.

Notation and Preliminaries. Let us define the following notation:

- \mathcal{M} represents the set of available language models.
- \mathcal{C} represents the set of code constructs eligible for model integration.
- $\llbracket \cdot \rrbracket$ denotes the semantic evaluation function.
- Γ represents the typing environment.
- σ represents the runtime state (variable bindings and object graph).

Function Integration Semantics. For a function definition with signature $f(p_1 : T_1, \dots, p_n : T_n) \rightarrow T_r$ **by** $m(\theta)$ where $m \in \mathcal{M}$ and θ represents model hyperparameters:

$$\frac{\Gamma \vdash f : (T_1, \dots, T_n) \rightarrow T_r \quad m \in \mathcal{M} \quad \theta \text{ valid for } m}{\llbracket \text{def } f(p_1 : T_1, \dots, p_n : T_n) \rightarrow T_r \text{ by } m(\theta) \rrbracket_{\Gamma} = \Gamma'} \quad (1)$$

where Γ' extends Γ with the binding $f \mapsto \langle (T_1, \dots, T_n) \rightarrow T_r, m, \theta \rangle$ to record that f is implemented via the **by** operator.

When the function is called:

$$\frac{\Gamma(f) = \langle (T_1, \dots, T_n) \rightarrow T_r, m, \theta \rangle \quad \Gamma \vdash v_i : T_i \text{ for } i \in \{1, \dots, n\}}{\llbracket f(v_1, \dots, v_n) \rrbracket_{\sigma} = \text{MT-Runtime}(f, [v_1, \dots, v_n], [(T_1, \dots, T_n), T_r], m, \theta, \sigma)} \quad (2)$$

The MT-Runtime is invoked only at call time, not at definition time, and receives both the parameter values and their types.

Method Integration Semantics. For a method defined in class C as $meth(p_1 : T_1, \dots, p_n : T_n) \rightarrow T_r$ **by** $m(\theta)$ where $m \in \mathcal{M}$:

$$\frac{\Gamma \vdash C \text{ is a class} \quad m \in \mathcal{M} \quad \theta \text{ valid for } m}{\llbracket \text{def } meth(p_1 : T_1, \dots, p_n : T_n) \rightarrow T_r \text{ by } m(\theta) \text{ in class } C \rrbracket_{\Gamma} = \Gamma'} \quad (3)$$

where Γ' extends Γ with the binding $C.meth \mapsto \langle (T_1, \dots, T_n) \rightarrow T_r, m, \theta \rangle$.

When the method is called on an object instance:

$$\frac{\Gamma \vdash obj : C \quad \Gamma(C.meth) = \langle (T_1, \dots, T_n) \rightarrow T_r, m, \theta \rangle \quad \Gamma \vdash v_i : T_i \text{ for } i \in \{1, \dots, n\}}{\llbracket obj.meth(v_1, \dots, v_n) \rrbracket_{\sigma} = \text{MT-Runtime}(C.meth, obj, [v_1, \dots, v_n], [C, (T_1, \dots, T_n), T_r], m, \theta, \sigma)} \quad (4)$$

```
1 def calculate_age (cur_year: int, dob: str) -> int by llm
2
3 age = calculate_age(2024, "March 14, 1879")
```

(a) Function Call

```
1 class Person :
2     name : str
3     dob : str
4
5 Einstein = Person("Einstein") by llm
```

(b) Object Initialization

```
1 class Person :
2     name : str
3     dob : str
4
5     def calculate_age (cur_year: int) -> int by llm(
6         temperature=0.7)
7
8 Einstein = Person("Einstein", "March 14, 1879")
9 Einstein.calculate_age(2024)
```

(c) Member Method Call

Fig. 8. Meaning-Typed Programming supports three model-integration points.

Here, the MT-Runtime has access to the object’s state through the *obj* parameter, in addition to the method arguments. The type parameters include the class type C , parameter types, and return type.

Object Initialization Semantics. For an object initialization expression $C(a_1, \dots, a_k) \text{ by } m(\theta)$ where C is a class with attributes $(attr_1 : T_1, \dots, attr_n : T_n)$ and $k < n$:

$$\frac{\Gamma \vdash C : \text{class} \quad \Gamma \vdash a_i : T_i \text{ for } i \in \{1, \dots, k\} \quad m \in \mathcal{M} \quad \theta \text{ valid for } m \quad k < n}{\llbracket C(v_1, \dots, v_k) \text{ by } m(\theta) \rrbracket_\sigma = \text{obj}} \quad (5)$$

where:

$$\text{obj} = \text{MT-Runtime}(C, [v_1, \dots, v_k], [(T_1, \dots, T_k), (T_{k+1}, \dots, T_n)], m, \theta, \sigma) \quad (6)$$

Unlike functions and methods where MT-Runtime is invoked at call time after definition, for object initialization the **by** operation and MT-Runtime invocation happen at the same time—during object instantiation. The MT-Runtime fills in the remaining attributes $(attr_{k+1}, \dots, attr_n)$ based on the semantic meaning of the class and the provided values, returning a fully initialized object instance.

Type Safety and Guarantees. The **by** operator maintains type safety by enforcing the following rules:

- (1) *Input Type Preservation:* For all inputs v_i , the type T_i is preserved when passed to the MT-Runtime.
- (2) *Output Type Enforcement:* The MT-Runtime ensures that the output from the language model is converted to the expected return type T_r .
- (3) *Error Handling:* If the language model produces output that cannot be converted to T_r , the MT-Runtime raises a type error.

Temporal Behavior. The **by** operator exhibits different temporal behaviors depending on the integration point:

- For functions and methods, the **by** operator’s effect is split across two times:
 - At definition time: recording the function/method signature, model, and hyperparameters
 - At call time: invoking MT-Runtime with arguments to perform the actual computation
- For object initialization, the **by** operator’s effect occurs in a single step:
 - At instantiation time: invoking MT-Runtime to complete the partially specified object

This formalization ensures that model integration through the **by** operator maintains the type safety guarantees of the host programming language while enabling seamless integration of language model capabilities. The formal semantics also reinforces our approach to achieving \mathbf{O}_1 by providing a mathematically rigorous foundation for the intuitive language abstraction.

4.3 Meaning-Typed Compilation and IR (MT-IR)

When a **by** call site is encountered, we perform static analysis to extract all relevant code semantics for automated prompt generation for LLM inference at runtime. Relevant semantics are identified using the **by** call-site, using definitions in equations 2, 4 and 5. As discussed in §3 \mathbf{C}_2 , this task is challenging due to distributed semantics. Figure 9a illustrates this challenge using the video game example from §2. In this example, the `Level` object types for the **by** call site in `game.py` are actually defined in a separate file, `level.py`. Furthermore, the dependent nested types used by `Level` reside in yet another file, `primitives.py`. This modular structure is common for improving code maintainability.

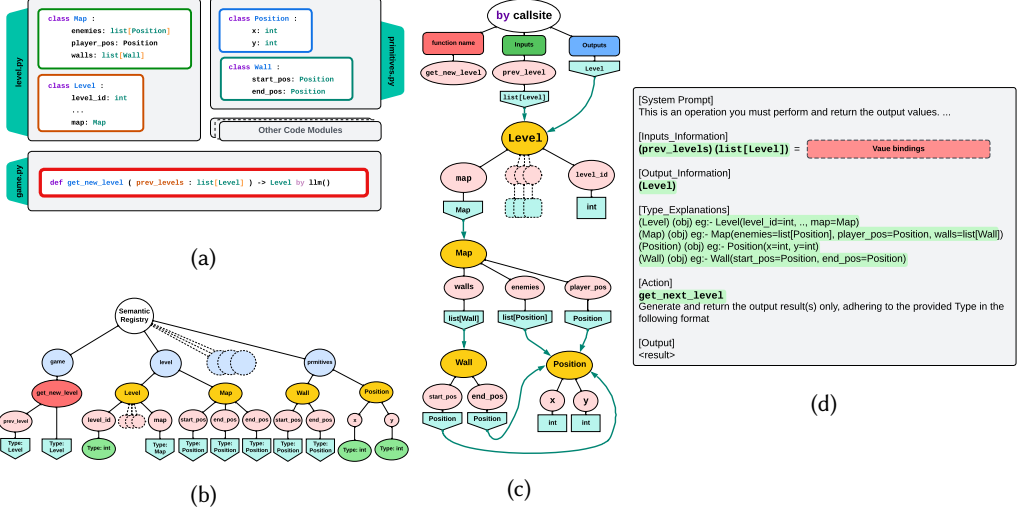


Fig. 9. Representation of how the MTIR is generated and used for the video game example in §2. (a) Code layout in modules. (b) Corresponding semantic registry containing all semantic information in the codebase. (c) MTIR relevant for the `get_next_level` `by` call-site. (d) Generated prompt using the MTIR (semantic information extracted using MTIR are highlighted in green).

Generating the MT-IR for each `by` call-site involves several key steps:

1. Semantic Registry Generation.

During the compilation phase of the language, we construct a *semantic registry* using the Abstract Syntax Tree (AST) generated for each code module. This registry serves as a comprehensive repository of all semantic information present in the code, including variable names, function names, and type information. As seen in Figure 9a, semantic information from multiple modules (`game.py`, `level.py`, and `primitives.py`) is integrated into a unified data structure. The resulting semantic registry, shown in Figure 9b, is a super set of the symbol-table of the language as it contains the definitions as well as the usages of objects, linking usages with definitions.

2. `by` Call-site Semantic Extraction.

The created semantic registry is used to traverse and find semantics relevant to each `by` call-site. The process begins

Algorithm 1 MT-IR Constructor Algorithm

Require: Code Base C , Semantic Registry S

Ensure: MT-IR Map M

```

1: Initialize MT-IR  $M \leftarrow \emptyset$ 
2: for each byi call-site in  $C$  do
3:   if byi is a function-call then
4:     Retrieve function name  $\sigma_i$  and signature from  $S$ 
5:     Retrieve formal parameters  $(p_1 : T_1, \dots, p_n : T_n)$  and return type  $T_r$  from  $S$ 
6:     Initialize  $M[\text{by}_i] \leftarrow \langle \sigma_i, (T_1, \dots, T_n) \rightarrow T_r, m, \theta \rangle$  {From equation 1}
7:   else if byi is a method-call then
8:     Retrieve method name  $\sigma_i$  and signature from  $S$ 
9:     Retrieve class  $C$ , formal parameters  $(p_1 : T_1, \dots, p_n : T_n)$  and return type  $T_r$  from  $S$ 
10:    Initialize  $M[\text{by}_i] \leftarrow \langle \sigma_i, C, (T_1, \dots, T_n) \rightarrow T_r, m, \theta \rangle$  {From equation 3}
11:   else if byi is an object-initialization then
12:     Set  $\sigma_i$  to 'Complete the object attributes'
13:     Retrieve class  $C$  with attributes  $(attr_1 : T_1, \dots, attr_n : T_n)$  from  $S$ 
14:     Retrieve provided attribute values  $(v_1, \dots, v_k)$  where  $k < n$ 
15:     Initialize  $M[\text{by}_i] \leftarrow \langle C, (T_1, \dots, T_k), (T_{k+1}, \dots, T_n), m, \theta \rangle$  {From equation 5}
16:   end if
17:   Extract Parameter Type Semantics:
18:   for each type  $T_j$  in input types of  $M[\text{by}_i]$  do
19:     if  $T_j$  is a primitive type then
20:       Add  $T_j$  directly to  $M[\text{by}_i]$ 
21:     else
22:        $type\_def_j \leftarrow \text{EXTRACTTYPEDEFINITION}(T_j, S)$ 
23:       Add  $type\_def_j$  to  $M[\text{by}_i]$ 
24:     end if
25:   end for
26:   Extract Return Type Semantics:
27:   for each type  $T_j$  in output/return types of  $M[\text{by}_i]$  do
28:     if  $T_j$  is a primitive type then
29:       Add  $T_j$  directly to  $M[\text{by}_i]$ 
30:     else
31:        $type\_def_j \leftarrow \text{EXTRACTTYPEDEFINITION}(T_j, S)$ 
32:       Add  $type\_def_j$  to  $M[\text{by}_i]$ 
33:     end if
34:   end for
35: end for
36: return  $M$ 

```

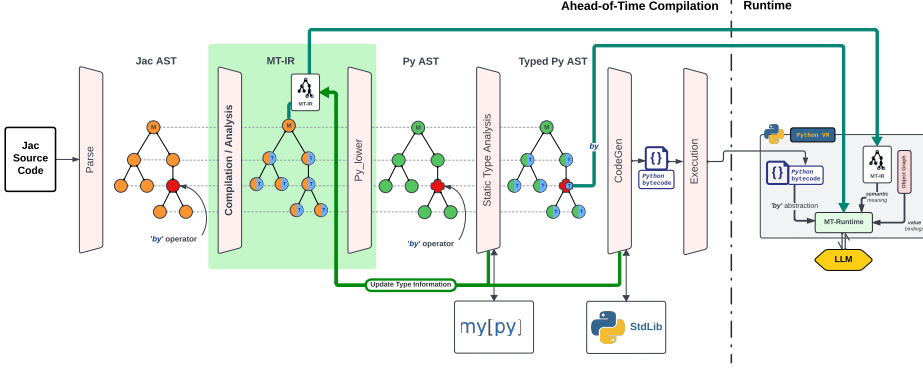


Fig. 10. The Jac compiler workflow on which Meaning-Typed compilation is implementation in the Jac programming language.

by extracting semantics from **by** call-site. For example, consider the **by** call-site in the `game.py` module, shown in Figure 9a for a the function call. As defined in equation 2, the **by** operator calls the MT-`Runtime` using $(f, [v_1, \dots, v_n], [(T_1, \dots, T_n), T_r], m, \theta, \sigma)$ as input information. These are the critical set of semantics relevant to the **by** call-site, which needs to be extracted first from the semantic registry as shown in Figure 9c. Next, the input and output types are checked to determine whether they are primitive types (e.g., integers, floats). If the types are non-primitive, we need to extract their definitions by recursively traversing the semantic registry, as outlined in lines 19 and 28 of the MT-`IR` Constructor Algorithm (Algorithm 1). Using this traversal, we need to extract the required semantics of the `Level` object, as illustrated in Figure 9c, which depends on `Map`. In turn, `Map` relies on `Wall` and `Position`. The use-definition links within the semantic registry enable this traversal, ensuring that all types are resolved down to primitive types. This approach effectively captures type information for nested types and addressing (C_2).

The extracted MT-`IR` can be used to fill in the semantic information in the prompt as shown in Figure 9d, where the highlighted text of the prompt represents all the semantics available in the MT-`IR` for the given **by** call-site.

Implementation of MT-`IR`. We implement the **by** language abstraction, MT-`IR`, and our compiler pass through `Jac`, a production-grade Python superset language [18, 24]. Typically, when a Python program is executed, the source code is compiled into Python bytecode via ahead-of-time compilation. The bytecode is then interpreted or compiled in a just-in-time fashion at runtime. `Jac` extends Python in the form of a PyPI package [17], providing a CLI command, `jac` that replaces the standard Python compiler by inserting a customized `jac` compilation phase as in Figure 10. The resulted python bytecode will then be interpreted using standard `python3`.

The compilation process for MTP in `Jac` involves six key steps through the `Jac` compiler as shown in Figure 10. Initially, an additional compiler pass generates MT-`IR` by extracting semantic elements from the AST. This information is then refined through multiple analysis passes. Next, the `Jac` AST is lowered to a Python AST, where specific code constructs are replaced with MT-`Runtime` calls (depicted as red nodes). Next, `MyPy` type checking is performed on the Python AST with resulting type information fed back into the MT-`IR` to enhance its semantic content. This enriched MT-`IR` is then registered into the MT-`Runtime` library, which remains available throughout the program execution. Finally, standard Python bytecode is generated from the modified Python AST,

preserving crucial semantic elements in the runtime that would otherwise be lost in conventional compilation in the MT-IR.

4.4 Meaning-Typed Runtime System

For model-integrated programs, LLM inference is done during the runtime as the prompt requires the semantic information as well as the dynamic values available when the `by` is called. This demands a dynamic system which can combine the MT-IR and dynamic value bindings to generate the final prompt for LLM inference. To this end, we introduce MT-Runtime, a novel automated runtime engine integrated into the language’s VM to realize O_3 . It is specifically triggered at the `by` call site abstractions, automatically managing the LLM invocation at runtime. MT-Runtime engine effectively addresses the final two challenges mentioned in §3 (C_3 and C_4).

In this Section, we describe the design of MT-Runtime, which functions through the following conceptual stages.

1. Prompt Synthesis. Upon the triggering of a `by` call, MT-Runtime first retrieves the MT-IR which is relevant to that particular `by` call-site and uses its information to fill a generalized prompt template. Based on the location of the `by` call-site, MT-Runtime constructs a prompt that incorporates static semantics, dynamic variable values relevant to the call-site, and the expected output type. As the MT-Runtime resides on the language VM as shown in Figure 10, it has access to the variable values at runtime. In Figure 11, this can be seen where the semantic information is fetched from the MT-IR while the variable values are fetched from the languages object graph. The synthesized prompt can be seen in Figure 9d. Here, the function name determines the ‘Action’ in the prompt, representing the intended task, while type information and nested types are included under ‘Type_Explanations’ in the prompt. Additionally, the inclusion of the output type and other schema information in the prompt forces the LLM to generate the output as per the defined schema, making it easier to convert in the type conversion stage. This stage address the challenge of generating meaningful but concise prompts discussed in 3 as C_3 .

2. Conversion from LLM Output to Target Output Variable. MT-Runtime queries the LLM with the dynamically synthesized input prompt and receives an output. The next step is to parse the LLM’s textual output and convert it into the target output variable addressing C_4 .

During the prompt synthesis phase, the prompt instructs the model to generate output following the predefined Python object schema. MT-Runtime leverages Python’s built-in `ast.literal_eval()` to evaluate the output text and construct the corresponding variable instance. For example, as in Figure 8b, the variable `Einstein` is of the custom type `Person` with attributes `name: str` and `dob: str`. In this case, MT-Runtime instructs the LLM to generate an output as `Person(name="Albert Einstein", dob="03/14/1879")`, which can be directly evaluated into a valid `Person`-typed object.

When MT-Runtime encounters an output that cannot be converted into the target variable type, it constructs a revised prompt that attempts to rectify the error as illustrated in Figure 11. This revised prompt includes details of the current incorrect output and the expected type, guiding the model toward the correct format. MT-Runtime queries the LLM with this revised input and

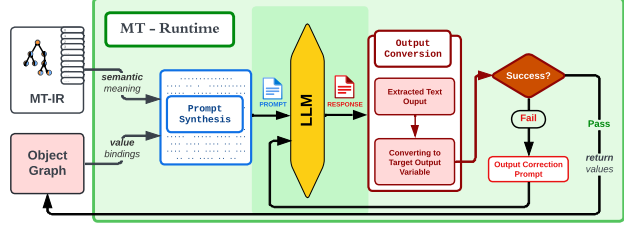


Fig. 11. Meaning-typed Runtime System (MT-Runtime).

Table 1. Benchmark Suite Descriptions.

Benchmarks	Task Description
Math Problem	Given a grade school math question an integer answer is returned. [15]
Translation	Translates a given text from English to French.[4]
Essay Reviewer	Reviews an essay and assigns a letter grade based on the quality of the essay. [46]
Joke Generator	Generates jokes based on a given topic, e.g., jokes with punchline. [4]
Expert Answer	Given a question, predicts the expert profession best suited to answer and generates an answer. [46]
Odd Word Out	Identifies the odd word out in a list of words along with a reason for why it is the odd word out. [4]
MCQ Reasoning	Answers multiple-choice questions by providing a reason for the chosen answer.[4]
Personality Finder	Given a public figure’s name, predicts their personality type.
Template	Given a template object and a set of variables, fills in the template with the variables.[4]
Text to Type	Given a text input and a structured data type, extracts the structured data into the structured data type.[4]
Taskman	Given a list of tasks, assigns priority ranks and estimates completion times for each task.
Level Generator	Generates new video game-levels based on previous levels.
Wikipedia	Given a question and an external tool (e.g., Wikipedia), retrieves the relevant answer using ReAct methodology. [4]

attempts to parse its output again. These retry attempts continue until either the type conversion succeeds or the maximum number of retries specified by the developer is reached. If the LLM fails to generate a correct output at the end, the MT-Runtime raises an exception with a type error. While errors are handled elegantly through this approach, the revised prompt is mostly short and concise since it only aims to correct the type mismatches. This offers an efficient solution for C_4 , handling errors gracefully.

5 EVALUATION

In this section, we present our evaluation of MTP in comparison to prior work (LMQL [4] and DSPy [15]) aimed at reducing the complexities associated with prompt engineering. With reduced developer effort, we strive to ensure that MTP maintains accuracy without increasing token usage with respect to prior works. Additionally, our lightweight runtime system is designed to minimize both cost and runtime overhead. Furthermore, we assess the robustness of MTP in scenarios where poor coding practices diminish the semantic richness of the code. Our evaluation aims to answer the following research questions.

RQ₁ - *How effectively does MTP reduce development complexity for model-integration?*

RQ₂ - *How does the accuracy of the MTP compare to prior work?*

RQ₃ - *What are the trade-offs for token usage, cost, and runtime in MTP?*

RQ₄ - *How resilient is MTP to sub-optimal coding practices?*

To evaluate these research questions, we introduce, (i) two case-studies, (ii) a user study and (iii) a suite of benchmarks.

- (1.) **Case Studies (RQ₁).** We first qualitatively evaluate the reduction of complexities for model-integration using MTP (§ 5.1).
- (2.) **Lines of code (LOC) (RQ₁).** In addition to the user study, we evaluate MTP for its effectiveness in reducing programming complexity through a quantitative study measuring lines of code across a comprehensive suite of benchmarks (§ 5.2).
- (3.) **User study (RQ₁).** We conducted a user study to evaluate MTP’s effectiveness in reducing complexity and increasing developer’s productivity (§ 5.3).
- (4.) **Accuracy and Trend across LLM evolution (RQ₂).** We evaluate the program accuracy when using MTP versus prior work across 10 various LLMs (§ 5.4).
- (5.) **Token Usage, Cost and Runtime (RQ₃).** We measured these performance metrics for MTP compared to DSPy (§ 5.5).
- (6.) **Sensitivity To Coding Practices (RQ₄).** Here we analyze how the MTP paradigm can be affected through bad coding practices using the game level generation benchmark. (§ 5.6)

Experimental Setup. We use OpenAI APIs for experiments with GPT models. We host llama models on a local server with Nvidia 3090 GPU with 24GB VRAM and 64GB RAM. We instrument the source code to profile token usage. Python cProfile[32] is used to measure program runtime. Our benchmark suite consists of problems representative of applications with model-integration. To construct the suite, we use the entire set of benchmarks and examples from prior work [4], supplemented with all relevant benchmarks in other work [15, 46] (Table 1).

5.1 A Case Study : Dynamic Level Generation in Video Game Development

§2 presented an example use-case of model-integration for automatic level generation of a video game. This video game is a real program developed using an on-line tutorial [35] using the pygame [25] library. Figure 12 represents three consecutive levels generated by an LLM through model-integration.

First we examine how this can be achieved using prior work, LMQL [4] and DSPy [15] which aims to simplify prompt engineering. Figure 13 includes the two code snippets for implementation from each framework.

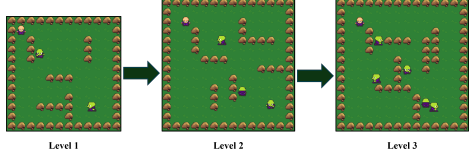


Fig. 12. A sequence of three game levels dynamically generated by an LLM at runtime.

LMQL [4]. Language Model Query Language, introduces a SQL-like query language designed specifically for LLM prompting, with a focus on constraining model outputs. This framework enables developers to enforce structured responses through predefined rules and output sachems. As shown in Figure 13a, LMQL provides a structured approach to our video game level progression use case. While the core prompt (line 8-34) contains similar instructions to traditional prompt engineering methods seen in Figure 4, LMQL’s key advantage lies in its explicit output format definition (lines 22-33). These constraints specify the expected data types and formats for each attribute, guiding the LLM to generate properly structured responses. During inference, LMQL operates through an iterative process. It generates one output field at a time, using uppercase text placeholders to mark generation points. Each successfully generated value is then injected back into the prompt, and the process repeats until all required values are obtained. While LMQL improves control over LLM-generated outputs, it does not fully abstract away prompt-engineering. Developers still need to construct base prompts manually requiring developer effort. This limits its effectiveness in simplifying model-integrated application development.

DSPy [15]. This is a developer-friendly framework for model integration, offering a more structured approach than LMQL. For simple tasks like question answering, developers can use `dspy.Signature` to define tasks with structured inputs and outputs. However, for complex tasks such as game level generation, which requires typed outputs, DSPy still necessitates manual code modifications. Figure 13b illustrates our DSPy implementation of game level transitioning. Here, lines 30–34 define a `dspy.Signature` with a task description and explicit input/output fields. A key challenge is DSPy’s inability to infer what a `Level` object represents without parsing the full codebase. As a workaround, developers must explicitly annotate class definitions, ensuring they inherit from Pydantic’s `BaseModel` (lines 7, 11, 15, and 21), with attributes defined as `Field` objects containing descriptions. While DSPy improves usability over generic prompt engineering and LMQL, its reliance on Pydantic introduces a steep learning curve and forces developers to modify existing code accordingly.

In contrast, our MTP implementation introduces a simple abstraction which reduces the complexities for developer in terms of learning-curve and developer effort as a few lines of code

```

1 import json
2 import lmql
3
4 @lmql.query
5 def generate_next_level(prev_level_json: str):
6     '''lmql
7     argmax
8     "You are an expert game level generator. Given the
9     current level details, generate the next level
10    with increased difficulty.
11
12    """
13    """ Current Level:
14    (prev_level_json)
15
16    """ Rules for Next Level:
17    - Increase difficulty by at least 1.
18    - Increase width and height slightly (max 10% each).
19    - Increase the number of walls and enemies based on
20    difficulty progression.
21    - Generate a valid map with walls and enemy positions.
22    - Ensure 'player_pos' is within bounds and not
23    colliding with walls.
24
25    """ Output Format:
26    Provide the next level in JSON format:
27    {
28      'level_id': [LEVEL_ID:int],
29      'difficulty': [DIFFICULTY:int],
30      'width': [WIDTH:int if WIDTH > current_width and
31      WIDTH <= current_width * 1.1],
32      'height': [HEIGHT:int if HEIGHT > current_height and
33      HEIGHT <= current_height * 1.1],
34      'num_wall': [NUM_WALL:int if NUM_WALL >=
35      current_num_wall],
36      'num_enemies': [NUM_ENEMIES:int if NUM_ENEMIES >=
37      current_num_enemies],
38      'map': {
39        'walls': [WALLS:list of {'start_pos': {'x': [X1:
40        int], 'y': [Y1:int]}, 'end_pos': {'x': [X2:
41        int], 'y': [Y2:int]}]],
42        'enemies': [ENEMIES:list of {'x': [EX:int], 'y':
43        [EY:int]}],
44        'player_pos': {'x': [PLAYER_X:int if PLAYER_X in
45        (0, WIDTH-1)], 'y': [PLAYER_Y:int if
46        PLAYER_Y in (0, HEIGHT-1)]}
47      }
48    }
49
50 def get_next_level(prev_level: Level -> Level:
51 prev_level_json = json.dumps(prev_level.__dict__,
52 default=lambda o: o.__dict__)
53 level_dict = json.loads(generate_next_level(
54 prev_level_json=prev_level_json))
55
56 def parse_position(data: dict) -> Position:
57 ...
58
59 def parse_wall(data: dict) -> Wall:
60 ...
61
62 def parse_map(data: dict) -> Map:
63 ...
64
65 return Level(
66 level_id=level_dict["level_id"],
67 difficulty=level_dict["difficulty"],
68 width=level_dict["width"],
69 height=level_dict["height"],
70 num_wall=level_dict["num_wall"],
71 num_enemies=level_dict["num_enemies"],
72 map=parse_map(level_dict["map"])
73 )

```

(a) Model-Integrated Game Levels : LMQL

```

1 import dspy
2 from pydantic import BaseModel, Field
3
4 llm = dspy.OpenAI(model="gpt-4o", max_tokens=1024)
5 dspy.settings.configure(llm=llm)
6
7 class Position(BaseModel):
8     x: int = Field(description="X Coordinate")
9     y: int = Field(description="Y Coordinate")
10
11 class Wall(BaseModel):
12     start_pos: Position = Field(description="Start Position
13     of the Wall")
14     end_pos: Position = Field(description="End Position of
15     the Wall")
16
17 class Map(BaseModel):
18     walls: list[Wall] = Field(description="Walls in the Map")
19     small_obstacles: list[Position] = Field(description="
20     Obstacles")
21     enemies: list[Position] = Field(description="Enemies in
22     the Map")
23     player_pos: Position = Field(description="Player
24     Position in the Map")
25
26 class Level(BaseModel):
27     level_id: int = Field(description="Level No")
28     difficulty: int = Field(description="Difficulty of the
29     Level")
30     width: int = Field(description="Width of the Map")
31     height: int = Field(description="Height of the Map")
32     num_wall: int = Field(description="Number of Walls in
33     the Map")
34     num_enemies: int = Field(description="Number of Enemies
35     in the Map")
36     map: Map = Field(description="Layout of the Map")
37
38 class get_next_level(dspy.Signature):
39     """Create Next Level"""
40     prev_levels: list[Level] = dspy.InputField(desc="Last
41     Played Levels")
42     level: Level = dspy.OutputField(desc="Level")
43
44 new_level_map = dspy.TypedPredictor(get_next_level)(level=
45 prev_levels).level

```

(b) Model-Integrated Game Levels : DSPY

```

1 from mtp.core.llms import OpenAI
2
3 llm = OpenAI(model_name="gpt-4o")
4
5 def get_next_level(prev_levels: list[Level]) -> Level by
6 llm()
7
8 new_level = get_next_level(prev_levels)

```

(c) Model-Integrated Game Levels : MTP

Fig. 13. Implementation of the video game usecase using (a) LMQL and (b) DSPY frameworks. (c) This gives a preview how the same performance can be achieved using the MTP paradigm introduced in the paper.

(as shown in Figure 13c) offers the same functionality each of these frameworks offer, for 58 lines (LMQL) and 36 (DSPY) lines of code changes.

MTP implementation maintains the same function signature and uses the semantics of this signature and the data structures defined in Figure 2 to automate the prompt generation process. The MT-runtime engine described in §4.4 automates and removes the need to parse the LLM output and convert it to a custom-typed variable, saving developer effort.

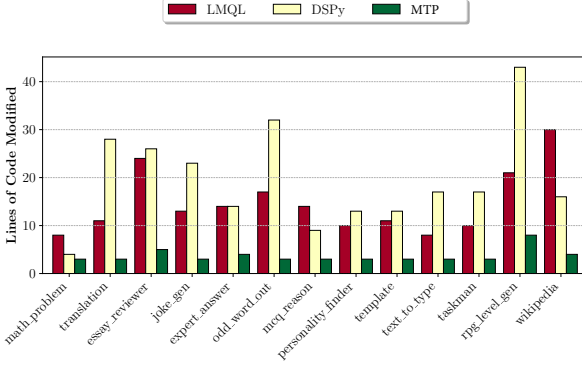


Fig. 14. Lines of Code modifications required for neuro-integration across benchmark programs (excluding manual code for response parsing).

Table 2. Lines of Code comparison between DSPy, LMQL, and MTP across multiple benchmark applications.

Problem	LoC for LLM Integration with MTP	
	vs. LMQL	vs. DSPy
Math Problem	↓ 2.7×	↓ 1.3×
Translation	↓ 3.7×	↓ 9.3×
Essay Reviewer	↓ 4.8×	↓ 5.2×
Joke Generator	↓ 4.3×	↓ 7.7×
Expert Answer	↓ 3.5×	↓ 3.5×
Odd Word Out	↓ 5.7×	↓ 10.7×
MCQ Reasoning	↓ 4.7×	↓ 3.0×
Personality Finder	↓ 3.3×	↓ 4.3×
Template	↓ 3.7×	↓ 4.3×
Text to Type	↓ 2.7×	↓ 5.7×
Taskman	↓ 3.3×	↓ 5.7×
RPG Level Generator	↓ 2.3×	↓ 4.8×
Wikipedia (ReAct)	↓ 7.5×	↓ 4.0×
Average	↓ 4.0×	↓ 5.3×

5.2 Lines of Code (LOC) Comparison

We evaluate MTP for its effectiveness in reducing programming complexity through a quantitative study measuring lines of code across our benchmark suite (Table 1). Each benchmark was implemented using all three frameworks, LMQL, DSPy and MTP. Figure 14 presents the lines of code modifications required for model-integration across benchmarks. LOC for LLM integration counts the amount of lines needed to be added or modified when a developer is integrating an LLM into traditional code. Figure 14 shows that minimum code changes required to integrate an LLM into traditional programming when using MTP, while DSPy and LMQL require much larger LOC across the board. This is due to the simple yet intuitive abstraction introduced in MTP.

Table 2 summarizes the LOC reduction achieved by MTP, on average, 5.3x and 4.0x reduction compared to DSPy and LMQL respectively. This large reduction is because DSPy requires tedious type annotations and LMQL requires manual programmatic prompt engineering. In contrast, MTP requires minimum code change for LLM integration by simply leveraging the `by` operator and hiding all complexity.

5.3 User Study Evaluating MTP against LMQL and DSPy

To quantify the degree of reduction in developer complexity and effort achieved through MTP, compared to prior work, we conducted a user study consisting of 20 software developers.

User Study Protocol. We issued a call soliciting developers with solid software engineering background to participate in the study. During the study, 20 Participants were first given one day to familiarize themselves with MTP, DSPy, and LMQL, through standard documentations, tutorials and available online resources. Participants were then tasked with implementing three progressively challenging programming tasks utilizing all three frameworks to integrate LLMs in traditional programming. These tasks were selected from our benchmark suite including: *Essay Evaluator (Easy)*, *Task Manager (Medium)* and *Game Level Generator (Hard)*. Each participant was allocated 90 minutes per task. After completing coding, participants filled out a questionnaire and provided feedback on the framework.

Table 3. Questions asked in the user study questionnaire, grouped under five criteria: Ease of Use, Clarity of Documentation, Learning Curve, Efficiency of Problem-Solving, and Overall Satisfaction.

Category	Question #	Statement
Ease of use	Q1	How easy was it to set up and start using the framework?
	Q2	How intuitive do you find the syntax and structure of the framework?
	Q3	How would you rate the ease of performing common tasks with the framework?
	Q4	How quickly will be able to integrate the framework into your existing projects?
Clarity of Documentation	Q5	How clear and understandable do you find the official documentation?
	Q6	How well is the documentation structured to find information quickly?
	Q7	How helpful are the provided examples and tutorials in understanding the framework?
	Q8	How complete and detailed are the code examples in the documentation?
Learning Curve	Q9	How long did it take you to feel comfortable with the basic features of the framework?
	Q10	How easy was it to learn and implement advanced features of the framework?
	Q11	How well does the framework support users in learning and utilizing advanced concepts?
Efficiency of Problem-Solving	Q12	How efficient is the framework in solving problems compared to others you have used?
	Q13	How easy was it to learn and implement advanced features of the framework?
Overall Satisfaction	Q14	Overall, how satisfied are you with the framework?
	Q15	Would you recommend this framework to others?
	Q16	How likely are you to continue using this framework in your future projects?
	Q17	What are the main reasons for your rating in the previous question?

Success Rate. Figure 15 presents the success rates across three programming tasks using DSPy, LMQL, and MTP. Success is defined as when the implementation a participate developed generates accurate outputs across test inputs. Overall, implementation utilizing MTP achieved the highest success rate in 2 out of the 3 tasks. It also performs the most consistently across all tasks. These results demonstrate MTP’s effectiveness in enabling programmers to leverage LLM capabilities intuitively. In contrast, we observe that participants found using LMQL to be the most challenging, with the lowest success rates. Qualitative feedback from participants indicated that this is because LMQL requires developers to manually craft prompts while DSPy and MTP automate the process. For the most challenging problem, *Game Level Generator (Hard)*, developers achieved zero success with LMQL, highlighting the difficulty to manually craft prompt and than manually interpret output to convert to a complex object type.

User Study Feedback. Our questionnaire consists of 17 questions grouped under five criteria: Ease of Use, Clarity of Documentation, Learning Curve, Efficiency of Problem-Solving, and Overall Satisfaction (Table 3). The average user scores for

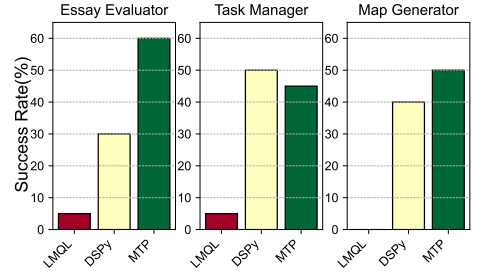


Fig. 15. User study success rates across LMQL, DSPy, and MTP.

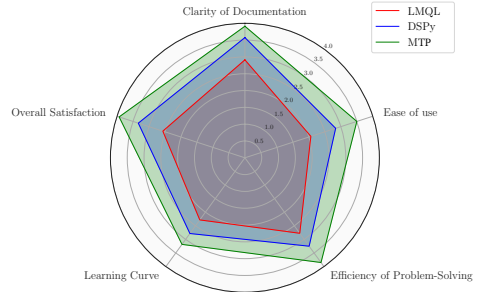


Fig. 16. User evaluation of LMQL, DSPy and MTP on five usability criteria

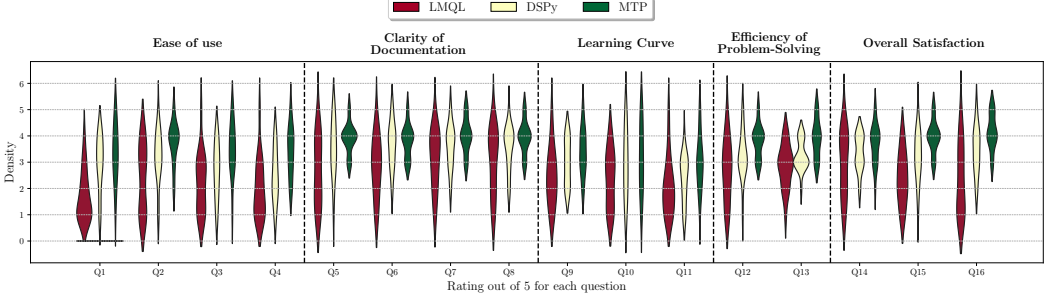


Fig. 17. Violin plot showing the distribution of scores given by participants for each framework across the five criteria: Ease of Use, Clarity of Documentation, Learning Curve, Efficiency of Problem-Solving, and Overall Satisfaction.

each framework across the five criteria are shown in Figure 16. MTP was consistently rated the better framework compared to LMQL and DSPy. DSPy generally scored in the mid-range, while LMQL scored the lowest across all criteria, suggesting it introduced the most complexity for programmers.

The Violin plot in Figure 17 shows the distribution of scores given by participants for each framework. Across all criteria, MTP exhibits higher density around the upper rating regions (higher scores) compared to LMQL and DSPy, indicating a generally more favorable rating. Notably, for categories such as Ease of use, efficiency of problem solving, MTP achieved higher density peaks near the higher scores, reflecting higher perceived developer productivity and lower complexity of MTP compared to the other frameworks.

Open-ended participant feedback further highlighted MTP’s strengths, with comments such as *"Learning LMQL was challenging, and it ultimately proved inadequate for accomplishing the required tasks."*, *"MTP provided a much simpler and more precise solution, making it the most efficient tool for the tasks at hand."* and *"MTP code is shorter than the other two frameworks... My favorite feature is the automatic filling of object attributes using by llm." conversions between data types can be done easily."*.

Using the results of the case study, user study and LoC reduction evaluation we can conclude that the abstraction offered through MTP reduces complexity for developers streamlining model-integration compared to other available solutions.

5.4 Program accuracy

In this section, we aim to evaluate if MTP’s complexity reduction comes at a cost of accuracy loss. In addition, we aim to investigate the impact of LLM evolution on program accuracy when using MTP compared to previous work.

Accuracy across benchmark suite. We first use the benchmark suite described in Table 1 for our evaluation using GPT-4o. Since 11 out of 12 benchmarks have text outputs with no reference output, we manually evaluate the output for its correctness. MTP achieved the highest accuracy by producing correct output on 12 out of the 12 benchmarks, while DSPy only achieved correctness 5 out of the 12. We observed that DSPy’s common issues are *"incorrectly formatted outputs"* and *"incorrectly typed outputs"*. These issues were not observed in MTP, demonstrating the robustness of MTP’s automatic prompt generation and type checking capabilities. LMQL is able to produce correct output on all 12 benchmarks. However, LMQL requires developers to manually craft prompts as opposed to the automated prompt generation of MTP and DSPy, resulting in a high complexity and lower developer productivity, as further demonstrated in our user study (§ 5.3).

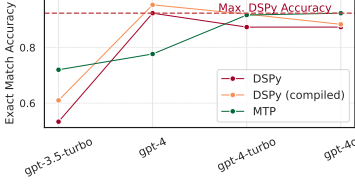


Fig. 18. Accuracy on GSM8K across OpenAI models.

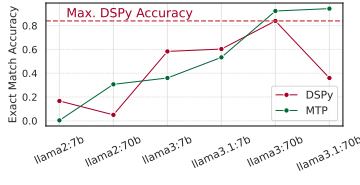


Fig. 19. Accuracy on GSM8K across Llama models.

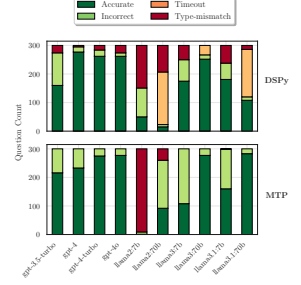


Fig. 20. Results breakdown for DSPy and MTP on GSM8K. (Timeout set at 120s).

Accuracy for GSM8k and Trends Across Model Evolution. In addition, we follow the same methodology as in previous work [15] for the evaluation of accuracy by focusing on the mathematical problem benchmark using a standard dataset, GSM8K [8], with expected answers. This dataset allows us to perform an objective evaluation of accuracy across a variety of inputs. We sample 300 question-answer pairs and evaluate the accuracy of the programs implemented using MTP and DSPy. Furthermore, we conduct the accuracy evaluation using 10 LLMs across generations throughout the model evolution to understand the accuracy trend as LLMs become more intelligent.

MTP vs. DSPy. As shown in Figure 18, for the latest GPT models (GPT-4-turbo and GPT-4o) DSPy and MTP perform similarly while MTP outperform slightly for the most recent GPT-4o model. This shows that although MTP hides all of the prompt engineering and type-conversion complexity, it achieves similar and sometimes even better performance accuracy through high quality automatic prompt synthesis and output interpretation.

Even more interestingly, we observe that *as the models continue to evolve and improve, the accuracy of programs implemented using MTP improves significantly*, while DSPy’s accuracy plateaued and also degraded a bit. *This trend indicates that as models get better, the need for complex prompt engineering actually is largely reduced.* LLM’s increasing capability to understand the traditional code, extract semantic meanings from the code and generate desirable output enabled MTP’s high-accurate automatic prompt synthesis and output interpretation.

MTP vs. DSPy (Compiled). It is also interesting to note that DSPy supports a compilation mode that requires additional training examples. It then searches across a set of prompt templates to optimize the prompt to improve accuracy, which could be time-consuming. We observe that MTP (with no training) outperforms DSPy (compiled with training) on all LLMs with the exception of GPT-4 (Figure 18). These results further highlight the effectiveness of MTP in automatically generating robust prompts using the semantic elements available in MT-IR.

We performed the same experiments with recent Llama models. Figure 19 shows a similar trend: MTP accuracy improves as the Llama model gets better, indicating the decreasing need for complex prompt engineering for more intelligent models.

Figure 20 presents the failure / success breakdowns between the two approaches. It shows that for DSPy, the inaccuracy when using GPT models is often due to “type mismatch”: e.g. the traditional code is expecting an integer output, yet the automatic type conversion in DSPy failed. For Llama, the degradation of DSPy’s accuracy at the latest model is mostly due to longer execution time of the LLMs (timeout set to be 120s).

These results demonstrate that MTP reduces complexity with little to no accuracy penalty and, in some cases.

5.5 Token Usage, Cost and Runtime

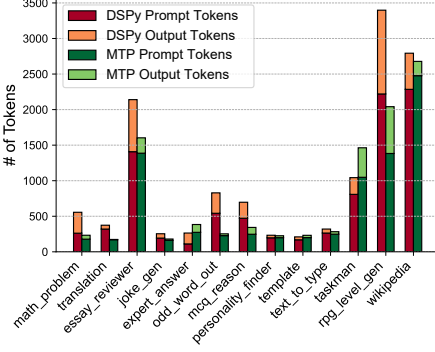


Fig. 21. Token usage comparison between MTP and DSPy.

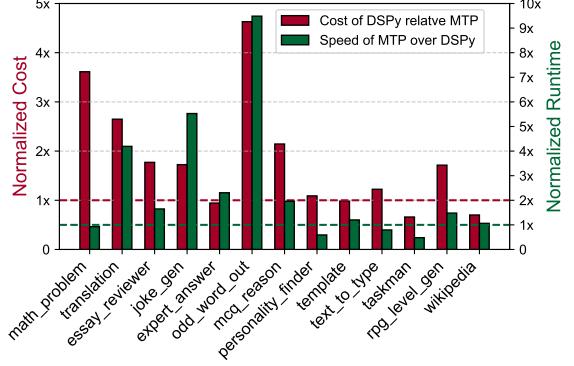


Fig. 22. Cost and runtime speed comparison between MTP and DSPy.

Token usage is an important metric since a higher token count leads to longer runtime (LLM inference time) and increased expenditure. Figure 21 presents the token usage comparison between DSPy (not-compiled) and MTP across the benchmark suite both for prompt (input) and completion (output) token usage. Across the benchmark suite, we observe that MTP consistently uses fewer tokens compared to DSPy, including the math problem using GSM8K dataset (first cluster of bars).

Figure 22 presents the runtime improvement and cost reduction achieved by MTP over DSPy. For each benchmark, the first bar shows the DSPy’s cost relative to MTP (left y-axis) and the second bar shows the MTP’s runtime improvement over DSPy (right-axis). The inference cost with the GPT-4o model is calculated using the OpenAI’s pricing formula [30]. The runtime is measured using cprofile. Figure 22 demonstrates that MTP achieves significant runtime improvement and cost reduction due to its more efficient token usage on majority of the benchmarks. As an example, for the math_problem benchmark, we can see more than 3.5x cost savings while maintaining same speed as DSPy. On the other hand, for the odd_word_out benchmark we can observe more than 4x cost reduction as well as runtime improvement over DSPy.

5.6 Sensitivity To Coding Practices

MTP leverages semantics, enabling automated prompt generation for model-integrated applications. Here, we evaluate how sensitive MTP is to poor coding practices, particularly when semantic richness is reduced.

To study this, we used the same video game level generation program from § 2. Starting with the original implementation, we systematically introduced poor coding practices and measured the probability of generating a correct output out of 100 consecutive runs. A correct output is a playable level with a player, walls, and enemies.

Figure 23 shows a gradual decline in correctness of output as semantic quality deteriorates. We first examined the effects of shortened but still meaningful names, a common practice for maintainability. This had no impact on correctness. Next, we evaluated progressively worse naming conventions by replacing meaningful identifiers with increasingly abbreviated or single letter alternatives. We tested four levels of degradation: renaming 25%, 50%, 75%, and 100% of identifiers.

At 25% and 50% renaming, MTP remained highly effective, with a correct output probability of 98%. However, as renaming increased to 70%, where nearly all data structures had poorly chosen names, correctness dropped to 70%. When renaming reached 100%, including function signatures, correctness declined sharply to 20%, highlighting the limits of MTP’s resilience.

Additionally, we tested cryptic naming conventions, where identifiers retained internal consistency but lost intuitive meaning (e.g., renaming "Level" to "Boss" and "Wall" to "Chasm"). As shown in Figure 23, this resulted in a 60% success rate, having a greater negative impact than 75% abbreviation but still performing better than complete semantic loss.

These results demonstrate MTP’s robustness, as the LLM can infer missing context, but excessive semantic degradation reduces accuracy.

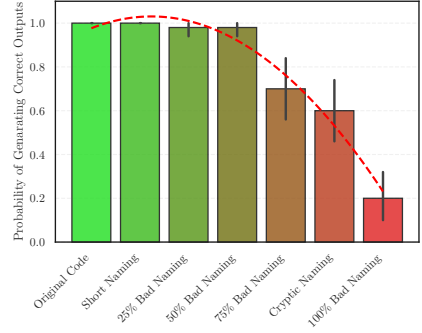


Fig. 23. Sensitivity of MTP to poor coding practices. The correct output probability declines slowly for first three experiments showing the resilience of MTP.

5.7 Using Multi-Modal LLMs to solve vision problems

There has been significant advancement recently in Multi-modal Generative AI models [36]. These models, such as LLaVa [22], QWEN [42], OpenAI’s GPT-4o, Anthropic’s Claude 3 and Google’s Gemini can process both text and images and perform visual tasks such as object detection and visual reasoning. In this example, we demonstrate how MTP generalizes leveraging such multi-modal generative models to build model-integrated programs using the same abstraction, reducing the developer effort at implementing such applications.

We select 5 vision benchmarks from related work Scallop [21]. This set of problems represents a wide range of multi-modal tasks such as visual reasoning and visual question-answering, as shown in Figure 24a. Using MTP, we are able to implement fully working solutions for these 5 vision problems using the same abstraction.

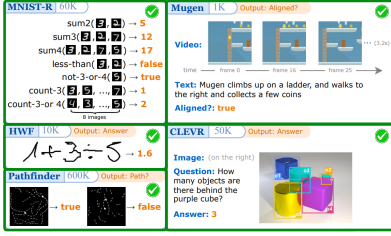
The code snippet Figure 24b shows our implementation of one of the benchmark problems (CLEVR), which involves reasoning about an input image and providing an answer to a user’s question.

As MTP treat every problem using semantics and types, using images as inputs just comes down to the the Image type used on line 6 on Figure 24. Hence, this case-study demonstrates the wide range of uses cases for MTP further reducing the need for the developer to learn frameworks for different model-integration problems, reducing developer complexities.

Our evaluation confirms that MTP simplifies LLM integration by reducing developer effort (\mathbf{RQ}_1) while maintaining accuracy (\mathbf{RQ}_2). By minimizing code complexity and improving usability, MTP outperforms DSPy and LMQL without increasing token usage and runtime costs (\mathbf{RQ}_3). Additionally, MTP shows a good degree of resilience to poor coding practices (\mathbf{RQ}_4).

6 RELATED WORK

There has been several recent efforts to help developers programming with LLMs. DSPy [15] and LMQL [4] and SGLang [46] introduce new libraries to help facilitating the integration of LLM into applications. While these frameworks provide additional tooling for LLM integration, they also introduce additional development complexity and new learning curves. LangChain [19] and LlamaIndex [23] focus on simplifying building LLM application while integrating with external data sources and can be leveraged together with our work. With this approach, the developer can



(a) Multi-modal solutions from Scallop [21].

```

1 # CLEVR benchmark problem implemented in MTP
2
3 from mtp.core.llms import OpenAI
4 from mtp.core.types import Image
5
6 llm = OpenAI(model_name="gpt-4o")
7
8 can_get_answer(img: Image, question: str) -> str by llm()
9
10 question:str = 'How many objects behind the purple cube?'
11 print(get_answer(Image('image.png'), question))

```

(b) Example CLEVR code snippet.

Fig. 24. Implemented multi-modal solutions for visual tasks from Scallop (Left) and CLEVR code snippet (Right).

retain the novel language-level abstraction that automates away much of complexity while benefit from the additional functionality introduced in these libraries.

There has been extensive studies in specialized compilation approach for AI and ML models [3, 6, 9, 12, 20, 31, 33, 40, 45]. In addition, recent works have also proposed novel runtime techniques and software system to accelerate training and inference of large scale models including large language models [1, 2, 5, 7, 11, 13, 14, 16, 26–29, 38, 41, 43, 44]. Furthermore, researchers have explored architectural support for ML/AI models [10, 47]. These work focus on accelerate training and inference of machine learning models. Our work focus on simplifying the development required to integrate LLM models in applications and can benefit from improved model performance.

7 CONCLUSION

Software is rapidly evolving from traditional logical code to model-integrated applications that leverage generative AI and large language models (LLMs) for application functionality. However, leveraging GenAI models in those applications is complicated and requires significant expertise and efforts. This paper presents meaning-typed programming (MTP), a novel approach to simplify the creation of model-integrated applications by introducing new language-level abstractions that hide the complexities of LLM integration. MTP automatically extracts the semantic meaning embedded in the traditional code, combined with dynamic information to dynamically synthesize prompts for LLMs and then convert the LLM output to the types that traditional code expects. We implemented MTP in a super-setted python language and demonstrate its effectiveness in reducing development complexity, improving runtime performance while maintaining high accuracy.

References

- [1] Sohaib Ahmad, Hui Guan, Brian D Friedman, Thomas Williams, Ramesh K Sitaraman, and Thomas Woo. Proteus: A high-throughput inference-serving system with accuracy scaling. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 1*, pages 318–334, 2024.
- [2] Reza Yazdani Aminabadi, Samyam Rajbhandari, Ammar Ahmad Awan, Cheng Li, Du Li, Elton Zheng, Olatunji Ruwase, Shaden Smith, Minjia Zhang, Jeff Rasley, et al. Deepspeed-inference: enabling efficient inference of transformer models at unprecedented scale. In *SC22: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–15. IEEE, 2022.
- [3] Jason Ansel, Edward Yang, Horace He, Natalia Gimelshein, Animesh Jain, Michael Voznesensky, Bin Bao, Peter Bell, David Berard, Evgeni Burovski, et al. Pytorch 2: Faster machine learning through dynamic python bytecode transformation and graph compilation. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2*, pages 929–947, 2024.
- [4] Luca Beurer-Kellner, Marc Fischer, and Martin Vechev. Prompting is programming: A query language for large language models. *Proceedings of the ACM on Programming Languages*, 7(PLDI):1946–1969, 2023.

- [5] Renze Chen, Zijian Ding, Size Zheng, Chengrui Zhang, Jingwen Leng, Xuanzhe Liu, and Yun Liang. Magis: Memory optimization via coordinated graph transformation and scheduling for dnn. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3*, pages 607–621, 2024.
- [6] Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Meghan Cowan, Haichen Shen, Leyuan Wang, Yuwei Hu, Luis Ceze, Carlos Guestrin, and Arvind Krishnamurthy. Tvm: An automated end-to-end optimizing compiler for deep learning, 2018.
- [7] Seungbeom Choi, Sunho Lee, Yeonjae Kim, Jongse Park, Youngjin Kwon, and Jaehyuk Huh. Serving heterogeneous machine learning models on Multi-GPU servers with Spatio-Temporal sharing. In *2022 USENIX Annual Technical Conference (USENIX ATC 22)*, pages 199–216, Carlsbad, CA, July 2022. USENIX Association.
- [8] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- [9] Michael Davies, Ian McDougall, Selvaraj Anandaraj, Deep Machchhar, Rithik Jain, and Karthikeyan Sankaralingam. A journey of a 1,000 kernels begins with a single step: A retrospective of deep learning on gpus. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2*, pages 20–36, 2024.
- [10] Soroush Ghodrati, Sean Kinzer, Hanyang Xu, Rohan Mahapatra, Yoonsung Kim, Byung Hoon Ahn, Dong Kai Wang, Lavanya Karthikeyan, Amir Yazdanbakhsh, Jongse Park, et al. Tandem processor: Grappling with emerging operators in neural networks. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2*, pages 1165–1182, 2024.
- [11] Cong Guo, Jiaming Tang, Weiming Hu, Jingwen Leng, Chen Zhang, Fan Yang, Yunxin Liu, Minyi Guo, and Yuhao Zhu. Olive: Accelerating large language models via hardware-friendly outlier-victim pair quantization. In *Proceedings of the 50th Annual International Symposium on Computer Architecture*, pages 1–15, 2023.
- [12] Muyan Hu, Ashwin Venkatram, Shreyashri Biswas, Balamurugan Marimuthu, Bohan Hou, Gabriele Oliaro, Haojie Wang, Liyan Zheng, Xupeng Miao, Jidong Zhai, et al. Optimal kernel orchestration for tensor programs with korch. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3*, pages 755–769, 2024.
- [13] Qijing Huang, Minwoo Kang, Grace Dinh, Thomas Norell, Aravind Kalaiah, James Demmel, John Wawrzynek, and Yakun Sophia Shao. Cosa: Scheduling by constrained optimization for spatial accelerators. In *2021 ACM/IEEE 48th Annual International Symposium on Computer Architecture (ISCA)*, pages 554–566. IEEE, 2021.
- [14] Yimin Jiang, Yibo Zhu, Chang Lan, Bairen Yi, Yong Cui, and Chuanxiong Guo. A unified architecture for accelerating distributed DNN training in heterogeneous GPU/CPU clusters. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*, pages 463–479. USENIX Association, November 2020.
- [15] Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas T Joshi, Hanna Moazam, et al. Dspy: Compiling declarative language model calls into self-improving pipelines. *arXiv preprint arXiv:2310.03714*, 2023.
- [16] Seah Kim, Hyoukjun Kwon, Jinook Song, Jihyuck Jo, Yu-Hsin Chen, Liangzhen Lai, and Vikas Chandra. Dream: A dynamic scheduler for dynamic real-time multi-model ml workloads. In *Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 4*, pages 73–86, 2023.
- [17] Jaseci Labs. Jaclang pypi package, 2024. Accessed: 2024-10-18.
- [18] Jaseci Labs. Jaseci github repo, 2024. Accessed: 2024-10-18.
- [19] Langchain. Langchain, 2024. Accessed: 2024-10-18.
- [20] Chris Lattner, Mehdi Amini, Uday Bondhugula, Albert Cohen, Andy Davis, Jacques Pienaar, River Riddle, Tatiana Shpeisman, Nicolas Vasilache, and Oleksandr Zinenko. Mlir: Scaling compiler infrastructure for domain specific computation. In *2021 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*, pages 2–14. IEEE, 2021.
- [21] Ziyang Li, Jiani Huang, and Mayur Naik. Scallop: A language for neurosymbolic programming. *Proceedings of the ACM on Programming Languages*, 7(PLDI):1463–1487, 2023.
- [22] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023.
- [23] llama_index. llama_index, 2024. Accessed: 2024-10-18.
- [24] Jason Mars, Yiping Kang, Roland Daynauth, Baichuan Li, Ashish Mahendra, Krisztian Flautner, and Lingjia Tang. The jaseci programming paradigm and runtime stack: Building scale-out production applications easy and fast. *IEEE Computer Architecture Letters*, 22(2):101–104, 2023.
- [25] Will McGugan. *Beginning game development with Python and Pygame: from novice to professional*. Apress, 2007.
- [26] Xupeng Miao, Gabriele Oliaro, Zhihao Zhang, Xinhao Cheng, Zeyu Wang, Zhengxin Zhang, Rae Ying Yee Wong, Alan Zhu, Lijie Yang, Xiaoxiang Shi, et al. Specinfer: Accelerating large language model serving with tree-based speculative inference and verification. In *Proceedings of the 29th ACM International Conference on Architectural Support*

- for *Programming Languages and Operating Systems, Volume 3*, pages 932–949, 2024.
- [27] Xupeng Miao, Chunan Shi, Jiangfei Duan, Xiaoli Xi, Dahua Lin, Bin Cui, and Zhihao Jia. Spotsolve: Serving generative large language models on preemptible instances. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2*, pages 1112–1127, 2024.
 - [28] Wei Niu, Jiexiong Guan, Yanzhi Wang, Gagan Agrawal, and Bin Ren. Dnnfusion: accelerating deep neural networks execution with advanced operator fusion. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*, pages 883–898, 2021.
 - [29] Hyungjun Oh, Kihong Kim, Jaemin Kim, Sungkyun Kim, Junyeol Lee, Du-seong Chang, and Jiwon Seo. Exegpt: Constraint-aware resource scheduling for llm inference. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2*, pages 369–384, 2024.
 - [30] OpenAI. Openai api pricing, 2024. Accessed: 2024-10-18.
 - [31] Zaifeng Pan, Zhen Zheng, Feng Zhang, Ruofan Wu, Hao Liang, Dalin Wang, Xiafei Qiu, Junjie Bai, Wei Lin, and Xiaoyong Du. Recom: A compiler approach to accelerating recommendation model inference with massive embedding columns. In *Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 4*, pages 268–286, 2023.
 - [32] Python Software Foundation. *The Python Standard Library: cProfile — Profile Module*, 2024. Accessed: 2024-10-19.
 - [33] Nadav Rotem, Jordan Fix, Saleem Abdulrasool, Garret Catron, Summer Deng, Roman Dzhabarov, Nick Gibson, James Hegeman, Meghan Lele, Roman Levenstein, et al. Glow: Graph lowering compiler techniques for neural networks. *arXiv preprint arXiv:1805.00907*, 2018.
 - [34] Reshabh K Sharma, Jonathan De Halleux, Shraddha Barke, and Benjamin Zorn. Promptpex: Automatic test generation for language model prompts. *arXiv preprint arXiv:2503.05070*, 2025.
 - [35] ShawCode. Pygame rpg tutorial. YouTube, 2021. Accessed: 2024-09-24.
 - [36] Shakti N Wadekar, Abhishek Chaurasia, Aman Chadha, and Eugenio Culurciello. The evolution of multimodal model architectures. *arXiv preprint arXiv:2405.17927*, 2024.
 - [37] Jianxun Wang and Yixiang Chen. A review on code generation with llms: Application and evaluation. In *2023 IEEE International Conference on Medical Artificial Intelligence (MedAI)*, pages 284–289. IEEE, 2023.
 - [38] Zheng Wang, Yuke Wang, Jiaqi Deng, Da Zheng, Ang Li, and Yufei Ding. Rap: Resource-aware automated gpu sharing for multi-gpu recommendation model training and input preprocessing. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2*, pages 964–979, 2024.
 - [39] Irene Weber. Large language models as software components: A taxonomy for llm-integrated applications. *arXiv preprint arXiv:2406.10300*, 2024.
 - [40] Chunwei Xia, Jiacheng Zhao, Qianqi Sun, Zheng Wang, Yuan Wen, Teng Yu, Xiaobing Feng, and Huimin Cui. Optimizing deep learning inference via global analysis and tensor expressions. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 1*, pages 286–301, 2024.
 - [41] Daliang Xu, Mengwei Xu, Chiheng Lou, Li Zhang, Gang Huang, Xin Jin, and Xuanzhe Liu. Socflow: Efficient and scalable dnn training on soc-clustered edge servers. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 1*, pages 368–385, 2024.
 - [42] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yeqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. Qwen2 technical report, 2024.
 - [43] Gyeong-In Yu, Joo Seong Jeong, Geon-Woo Kim, Soojeong Kim, and Byung-Gon Chun. Orca: A distributed serving system for Transformer-Based generative models. In *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*, pages 521–538, Carlsbad, CA, July 2022. USENIX Association.
 - [44] Chenguang Zheng, Hongzhi Chen, Yuxuan Cheng, Zhezhen Song, Yifan Wu, Changji Li, James Cheng, Hao Yang, and Shuai Zhang. Bytegnn: efficient graph neural network training at large scale. *Proceedings of the VLDB Endowment*, 15(6):1228–1242, 2022.
 - [45] Lianmin Zheng, Chengfan Jia, Minmin Sun, Zhao Wu, Cody Hao Yu, Ameer Haj-Ali, Yida Wang, Jun Yang, Danyang Zhuo, Koushik Sen, Joseph E. Gonzalez, and Ion Stoica. Ansor: Generating High-Performance tensor programs for deep learning. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*, pages 863–879. USENIX Association, November 2020.
 - [46] Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Jeff Huang, Chuyue Sun, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E Gonzalez, et al. Efficiently programming large language models using sglang. *arXiv preprint*

arXiv:2312.07104, 2023.

- [47] Kai Zhong, Zhenhua Zhu, Guohao Dai, Hongyi Wang, Xinhao Yang, Haoyu Zhang, Jin Si, Qiuli Mao, Shulin Zeng, Ke Hong, et al. Feasta: A flexible and efficient accelerator for sparse tensor algebra in machine learning. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3*, pages 349–366, 2024.