

PromptBridge: Cross-Model Prompt Transfer for Large Language Models

Yaxuan Wang^{1,2}, Quan Liu², Zhenting Wang², Zichao Li², Wei Wei², Yang Liu¹, Yujia Bao²

¹University of California, Santa Cruz, ²Center for Advanced AI, Accenture

Large language models (LLMs) underpin applications in code generation, mathematical reasoning, and agent-based workflows. In practice, systems access LLMs via commercial APIs or open-source deployments, and the model landscape (e.g., GPT, Claude, Llama) evolves rapidly. This rapid evolution forces frequent model switches driven by capability, cost, deployment constraints, and privacy. Yet prompts are highly model-sensitive: reusing a prompt engineered for one model on another often yields substantially worse performance than a prompt optimized for the target model. We term this phenomenon *Model Drifting*. Through extensive empirical analysis across diverse LLM configurations, we show that model drifting is both common and severe. To address this challenge, we introduce *PromptBridge*, a training-free framework that preserves prompt effectiveness under model switches, enabling cross-model prompt transfer without costly per-task or per-model re-optimization. *PromptBridge* requires only a small set of alignment tasks for calibration. It first applies *Model-Adaptive Reflective Prompt Evolution (MAP-RPE)* to obtain task- and model-specific optimal prompts via iterative reflective refinement and quantitative evaluation. Using the resulting calibrated prompt pairs for the source and target models, *PromptBridge* learns a cross-model prompt mapping. At test time, i.e., for an unseen task, given a source-model prompt, this mapping directly produces an optimized prompt for the target model. Experiments in single-agent and multi-agent settings show that *PromptBridge* consistently improves downstream accuracy while reducing migration effort. For example, when transferring from the source model to the target model, such as o3, *PromptBridge* yields a 27.39% improvement on SWE-BENCH and a 39.44% improvement on TERMINAL-BENCH relative to direct transfer. These results establish cross-model prompt transferability as a critical requirement for sustainable LLM system development and demonstrate that *PromptBridge* is an effective, practical, training-free solution for maintaining performance as models evolve. The code will be available at [PromptBridge](#).

1. Introduction

Large Language Models (LLMs) have demonstrated remarkable capabilities across diverse tasks, including code generation (Jimenez et al., 2024; Jain et al., 2024; Dong et al., 2025b,a), mathematical reasoning (Guan et al., 2025; Xia et al., 2025), agent-based workflows (Huang et al., 2023; Wang et al., 2024; Yuan et al., 2024; Liu et al., 2025; Zhang et al., 2025a; Wei et al., 2025b) and other complex problem-solving domains. LLM applications now range from single-turn question answering to complex, multi-agent systems (Li et al., 2025) that write code, analyze data, make decision (Su et al., 2025) and call external tools (Yao et al., 2023; Wu et al., 2024; Chen et al., 2023). In all cases, system behavior is governed by prompts that define the model’s role, guide its reasoning, and specify constraints and workflows. As applications mature, these systems accumulate carefully engineered prompt templates that encapsulate expert knowledge and are repeatedly reused across tasks.

New model releases (e.g., GPT (Achiam et al., 2023), Gemini (Comanici et al., 2025), Claude (Anthropic, 2024), Llama (Grattafiori et al., 2024)) routinely promise improvements in accuracy, latency, deployment, licensing, and privacy (Regulation, 2018; Neel & Chang, 2023). This churn raises a practical question:

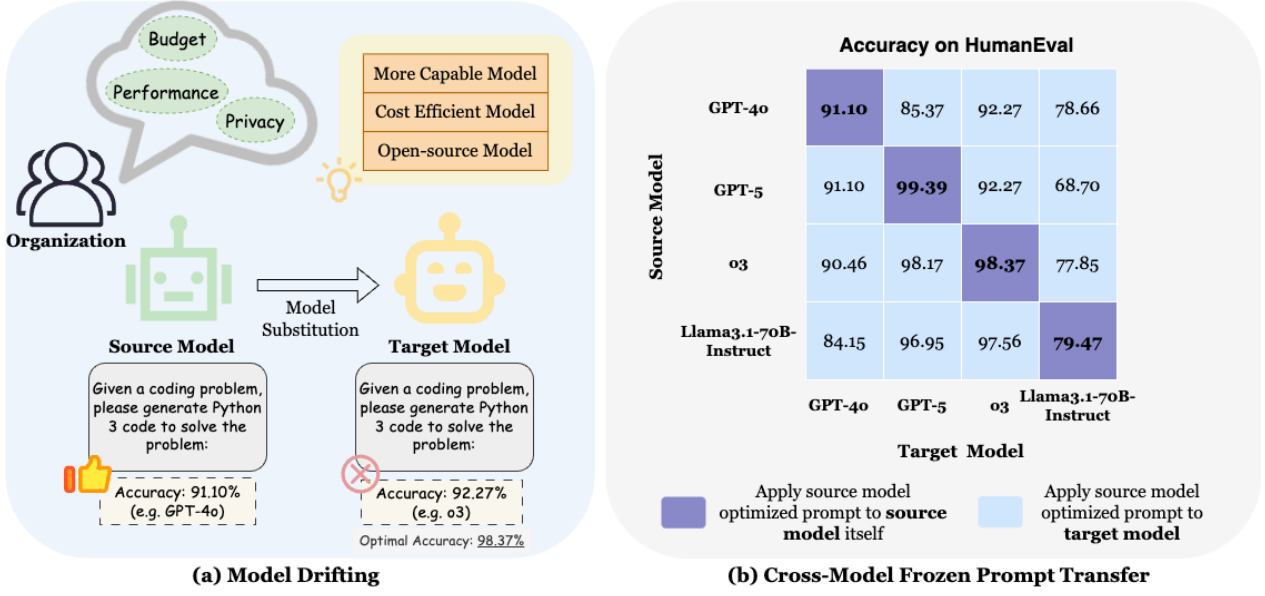


Figure 1 | **Illustration of Model Drifting Problem.** (a) Organizations may replace the source model with a more capable, cost-efficient, or open-source alternative. However, directly transferring the source-engineered prompt (e.g. GPT-4o) to a stronger target model (e.g. o3) yields only 92.27% accuracy on the HumanEval dataset, far below the target model’s achievable 98.37%, revealing that prompts engineered or optimized for one model do not reliably generalize to another. (b) Cross-model evaluation on HumanEval demonstrates model drifting: prompts designed for a source model frequently degrade when applied to different target models.

Can prompts engineered for one LLM transfer optimally to another without re-tuning?

Empirically, the answer is often no. Here we use transfer-learning terminology: the *source model* is the model on which prompts are tuned, and the *target model* is the new model to which those prompts are applied. Prompts engineered or optimized for a source model frequently underperform when applied to a target model, a phenomenon we term **Model Drifting**. Figure 1 illustrates the pattern. Panel (a) shows the common scenario: an organization replaces its source model (GPT-4o (OpenAI et al., 2024)) with a more capable alternative (o3 (OpenAI, 2025b)), yet the source-optimal prompt no longer remains optimal for the target. Panel (b) quantifies drift on HumanEval (Chen et al., 2021): the best prompt for GPT-5 (OpenAI, 2025a) reaches 99.39% on GPT-5 but yields 68.70% when transferred to Llama-3.1-70B-Instruct (Grattafiori et al., 2024), while Llama-3.1-70B-Instruct’s own optimal prompt achieves 79.47% on its own. The drift from GPT-5 to Llama-3.1-70B-Instruct is therefore 68.70% vs. 79.47%. Conversely, the Llama-3.1-70B-Instruct optimal prompt attains 96.95% on GPT-5, below GPT-5’s optimum of 99.39% (a 2.44-point drift). Additional analysis about this phenomenon appears in Section A.1. Such transfer gaps occur across model families and across coding, agentic, and planning tasks, creating a barrier to model migration and A/B testing without re-engineering prompt libraries.

Why does model drift arise? Vendors train and align models with different corpora, tokenization schemes, role tags, and human-feedback criteria. For instance, Llama 3 (Grattafiori et al., 2024) introduces an `ipython` role for tool calls—an instruction absent from GPT models. Models trained by different companies (e.g., Qwen (Yang et al., 2025), DeepSeek (Guo et al., 2025) vs. Mistral (Jiang et al., 2023)) ingest different linguistic and domain distributions. These alignment and interface

mismatches make a prompt fit one model but mis-specify another, forcing organizations to rewrite and re-validate prompts when swapping LLMs.

This work tackles the prompt-transfer problem at two levels:

(1) **Quantifying Model Drifting.** We formalize drift as the transfer gap to the target model’s own optimal prompt. For a task T , source model M_s , target model M_t , and performance metric $A(\cdot)$, let $p_{M,T}^*$ denote the model-preferred prompt that maximizes $A(M, T, p)$. We define the transfer gap from M_s to M_t on T as $\Delta(M_s \rightarrow M_t, T) = A(M_t, T, p_{M_s,T}^*) - A(M_t, T, p_{M_t,T}^*)$, i.e., the shortfall of the transferred source-optimal prompt relative to the target model’s own optimum. We empirically characterize this gap across models and tasks, revealing systematic drift even within model families.

(2) **PromptBridge: training-free cross-model prompt adaptation.** Rather than re-optimizing prompts anew for every model and task, PromptBridge uses a small calibration suite of *alignment tasks* for which we obtain source- and target-optimal prompts. Given an unseen task and a source prompt, PromptBridge maps the source prompt to a target-style prompt by leveraging these calibration anchors, enabling zero-shot, in-context prompt adaptation without access to evaluation data for the unseen task. To obtain high-quality anchors, we introduce **Model-Adaptive Reflective Prompt Evolution (MAP-RPE)**, an evaluation-guided, reflection-driven, island-based prompt search that produces model-preferred, task-specific prompts. MAP-RPE supports both our drift measurement (by approximating $p_{M,T}^*$) and PromptBridge’s calibration.

We evaluate PromptBridge across seven LLMs and eight benchmarks under single-agent and multi-agent settings, with comprehensive ablations and preliminary analysis of model drifting. PromptBridge consistently reduces the transfer gap and maintains or improves performance while substantially lowering migration overhead across different models.

Contributions. Our main contributions are summarized as follows:

- We introduce and empirically quantify the problem of **Model Drifting**, which arises when the underlying LLM in a system is replaced or updated, leading to performance shifts under the same prompt. We further establish the task of prompt transfer under model drifting and provide a general framework for addressing it.
- We propose a **Model-Adaptive Reflective Prompt Evolution (MAP-RPE)** method that calibrates prompts for a given specific model through reflective, metric-driven optimization, producing both task- and model-specific optimized prompts and serving the drift measurement process.
- We develop **PromptBridge**, a training-free cross-model prompt transfer framework that learns transferable prompt transformation mappings from alignment tasks and enables zero-shot prompt adaptation to unseen tasks across different target models.
- Extensive experiments across single-agent and multi-agent settings demonstrate the effectiveness of our framework and underscore the importance of the model drifting problem.

2. Related Work

Prompt Optimization A closely related line of research investigates how to optimize and adapt prompts to better elicit desirable behaviors from LLMs. Foundational strategies such as Chain-of-Thought prompting (Wei et al., 2022) substantially enhance reasoning abilities across diverse domains. Subsequent works explore automated prompt optimization using advanced online LLMs to iteratively refine prompts (Zhou et al., 2022; Cheng et al., 2023; Yang et al., 2023; Agarwal et al., 2024). Fernando et al. (2023) introduce PromptBreeder, which evolves a population of task prompts through mutation and fitness evaluation on a validation set. MultiPrompter (Kim et al., 2023) formulates prompt optimization as a cooperative game among multiple prompters that sequentially refine the

same prompt. MARS (Zhang et al., 2025b) adopts a multi-agent architecture that generates and evaluates candidate prompts via diverse search strategies. Beyond direct LLM-based optimization, evolutionary algorithms have also been applied to improve LLM performance. EvoPrompt (Guo et al., 2023) connects LLMs with evolutionary search to progressively populate optimized prompts, though it lacks integration with feedback. AlphaEvolve (Novikov et al., 2025) and its open-source variant OpenEvolve (Sharma, 2025) demonstrate that evolutionary rewriting within prompts can yield substantial performance gains in code generation tasks. GEPA (Agrawal et al., 2025) employs reflective feedback and Pareto-aware optimization to balance competing objectives when evolving prompt candidates. Despite the effectiveness of previous prompt optimization method, these methods often rely on large, proprietary LLMs and show limited transferability to smaller models (Wang et al., 2023a). To address this, Zhu et al. (2025) propose several design principles for constructing more generalizable prompts. The importance of prompt design that generalizes across model families remains underexplored. Rakotonirina et al. (2023) partially address this by inducing prompts via model mixing during training, producing prompts that generalize better across architectures.

Prompt Sensitivity Previous study has demonstrated that LLM is sensitive to prompt (Wei et al., 2025a; Zhu et al., 2025). Minor deviations between user-provided instructions and the training instructions can result in significant performance degradation (Wei et al., 2025a). The existing research primarily emphasizes the importance of adapting prompts to specific tasks (Zheng et al., 2024; Wang et al., 2023b), or specific modality (Zhang et al., 2025c), rather than specific LLMs. MAPO (Chen et al., 2024) first demonstrate that different prompts should be adapted to different LLMs to enhance their capabilities across various downstream tasks. The LLM agent might also suffer from prompt sensitivity (Verma et al., 2024), where simple modifications in the prompt can already exert significant but unexpected degradation of performance (Zhou et al., 2023; Liu et al., 2024). However, these works typically assume a fixed model. In contrast, the model drifting problem studied in this paper concerns how the optimal prompt itself evolves as the model architecture, or alignment procedure change over time or across model families. This phenomenon extends beyond prompt sensitivity: rather than a local perturbation effect, model drifting captures a systematic semantic shift in how different or updated LLMs interpret the same instruction, necessitating adaptive calibration and transfer methods to maintain performance consistency.

Optimizing AI systems and Agents Optimizing AI systems composed of multiple interacting LLM modules remains a fundamental challenge due to the complex dependencies among components. DSPy (Khatab et al., 2024) automates exemplar and instruction design to enhance the performance of multi-module LLM pipelines. TextGrad (Yuksekgonul et al., 2025) introduces a differentiable framework that backpropagates textual feedback across modules, allowing downstream performance signals to refine upstream prompts. MIPROv2 (Opsahl-Ong et al., 2024) jointly aligns instructions and few-shot examples via Bayesian optimization. Optimas (Wu et al., 2025) further improves multi-module systems by optimizing locally aligned rewards that are globally coordinated across the entire architecture. In multi-agent settings, MASS (Zhou et al., 2025) serves as a plug-and-play prompt and workflow optimizer that mitigates prompt sensitivity arising from inter-agent compounding effects. Together, these studies highlight the growing interest in system-level optimization methods that move beyond individual prompts or models to holistically improve multi-component AI systems.

Prompt Transfer Research on prompt transfer investigates how to generalize and reuse effective prompts across tasks, modalities, or models. Existing work focuses mainly on cross-task transfer (Zheng et al., 2024; Wang et al., 2023b), while Zhang et al. (2025c) explore cross-modality transfer, demonstrating that prompts designed for one modality can partially transfer to another. Su

et al. (2022) investigate the cross-model transferability of soft prompts by training a prompt projector through prompt tuning on a specific task. However, their findings are limited to soft prompt and pretrained language models such as RoBERTa, without exploring broader model families or more recent LLM architectures. Chen et al. (2024) introduce MAPO, a model-adaptive prompt optimizer that learns to tailor prompts for specific LLMs, showing that prompt effectiveness strongly depends on the underlying model family. Wang et al. (2025) propose a privacy-preserving soft prompt transfer framework, which trains soft prompts on smaller models and securely transfers them to larger ones, thereby improving transferability while maintaining privacy guarantees. Despite these advances, cross model prompt transfer remains underexplored for scenarios involving model drifting or complex LLM systems, where prompt adaptability and consistency are critical yet challenging to maintain.

3. Problem Formulation

We focus on prompt template transfer when switching from one LLM to another, which fall within the task- and model-adaptive prompt optimization area. The definition of Model Drifting is given below:

Model Drifting. *Given an LLM-based system $\Phi_{M,p}$, where M denotes the underlying LLM and p represents the prompt template that guide model behavior, we may wish to replace the source model M_s with a target model M_t in practice due to budget constraints, deployment requirements, or special demands. Such a replacement inevitably alters the system’s behavior and performance characteristics. We refer to the performance drift induced by the model substitution as model drifting.*

Formally, let $A(\cdot)$ denote a performance metric of an LLM-based system. When the underlying model is replaced, the same prompt that was previously optimal for the source model often yields suboptimal performance on the target model for task T :

$$\underbrace{A(M_t, T, p_{M_s, T}^*)}_{\text{after model switch}} < \underbrace{\max_{p_{M_t, T}} A(M_t, T, p_{M_t, T})}_{\text{target-optimal prompt performance}} \quad (1)$$

The core challenge is that the prompt $p_{M_s, T}^*$, which was optimal for M_s , is no longer optimal for M_t . Therefore, effective prompt recalibration or transfer is required to restore alignment between the prompt and the target model’s behavioral distribution.

In this work, we study model drifting under two system configurations: the single-agent setting and the multi-agent setting, with our primary focus on the single-agent case. The single-agent setting also fully encompasses traditional single-LLM tasks: one model is responsible for the entire problem, and model drifting simply means swapping the underlying LLM while keeping the overall task and agent structure fixed. In the multi-agent setting, model drifting can occur at two levels: Local model drifting refers to updating the source LLM of a specific agent to a new target model, while other agents remain unchanged; Global model drifting refers to simultaneously replacing the source models of all agents with their corresponding target models.

Prompt Transfer under Model Drifting. Let $S = \{S_1, S_2, \dots, S_m\}$ denote the set of *alignment tasks* used as model-adaptive prompt alignment data, and $T = \{T_1, T_2, \dots, T_q\}$ denote the set of *unseen tasks*. For each alignment task $S_i \in S$, let p_{M_s, S_i}^* and p_{M_t, S_i}^* represent the task-specific optimal prompts for the source model M_s and target model M_t , respectively. Similarly, for each unseen task $T_j \in T$, let p_{M_s, T_j}^* and p_{M_t, T_j}^* denote the corresponding optimal prompts. Given the unseen tasks T , prompt transfer aims to learn a mapping: $\mathcal{T}_{M_s \rightarrow M_t, S} : p_{M_s, T} \mapsto p_{M_t, T}^*$, where the mapping captures systematic relationships between model- and task-specific prompt structures. The transfer function \mathcal{T} is learned from the

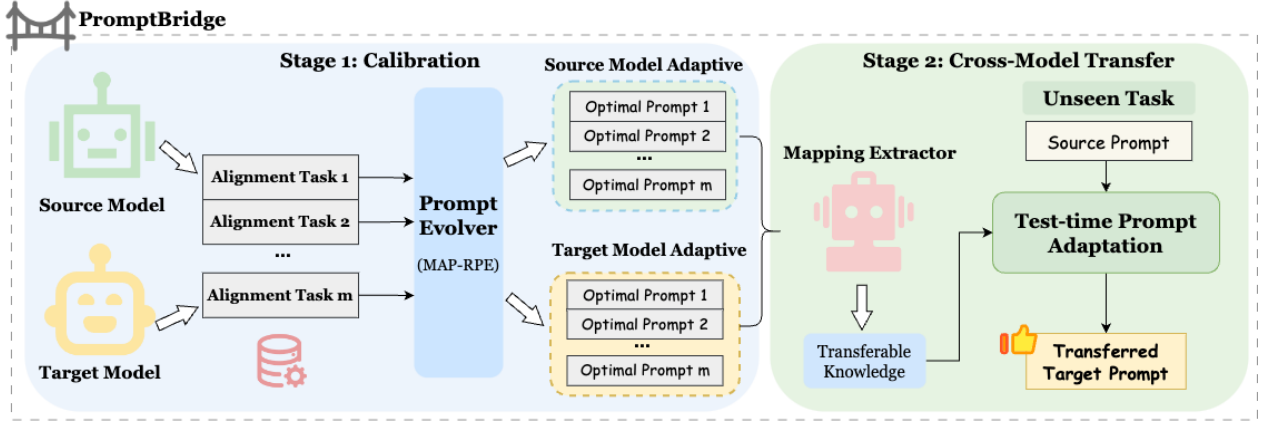


Figure 2 | **Overview of the proposed PromptBridge framework.** The framework operates in two stages: 1) Calibration, where Model-Adaptive Reflective Prompt Evolution (MAP-RPE) optimizes task-specific prompts for each model using quantitative feedback; and 2) Cross-Model Transfer, where pairs of calibrated source–target prompts are used to learn a transferable prompt mapping function \mathcal{T} . This learned mapping acts as a bridge that enables the system to synthesize target-compatible prompts for unseen tasks, supporting zero-shot prompt transfer without additional training.

alignment tasks S , for which both the source-model prompts and the corresponding target-model prompts are either known or can be obtained through calibration methods. Once learned, this transformation is applied to unseen tasks, where only the source-model prompt $p_{M_s, T}$ is available. Applying \mathcal{T} produces an adapted target-model prompt $\hat{p}_{M_t, T}$ that aims to (i) match or surpass the performance of the source-model prompt $A(M_s, T, p_{M_s, T})$, and (ii) approximate as closely as possible the performance of the true task-optimal target-model prompt $A(M_t, T, p_{M_t, T}^*)$.

4. Method

In this section, we first introduce **PromptBridge**, a model-adaptive prompt transfer framework that leverages calibrated prompts on alignment tasks to learn transferable prompt knowledge across different LLMs. We then present our Model-Adaptive Reflective Prompt Evolution (MAP-RPE) method, which optimizes prompts for a specific model in a reflective and adaptive manner. MAP-RPE is used as the calibration mechanism to obtain optimal prompts on the alignment tasks, which subsequently support the model-adaptive transfer process. Figure 2 provides an overview of the proposed approach.

4.1. Cross-Model Prompt Transfer

Transfer Mapping Learning. Let $S = \{S_1, S_2, \dots, S_m\}$ denote the set of alignment tasks used for learning model-adaptive prompt alignment. For each task S_i , we first obtain the task-specific optimal prompts using MAP-RPE, denoted as $\text{Evolver}(\cdot)$:

$$p_{M_s, S_i}^* = \text{Evolver}(p_{M_s, S_i}), \quad p_{M_t, S_i}^* = \text{Evolver}(p_{M_t, S_i})$$

The objective of the transfer module is to learn the mapping:

$$\mathcal{T}_{M_s \rightarrow M_t; S} : p_{M_s, S_i}^* \mapsto p_{M_t, S_i}^* \quad S_i \in S. \quad (2)$$

which captures how prompts should be systematically reformulated when moving from the source model to the target model. To extract the transferable knowledge, we employ a *Mapping Extractor*

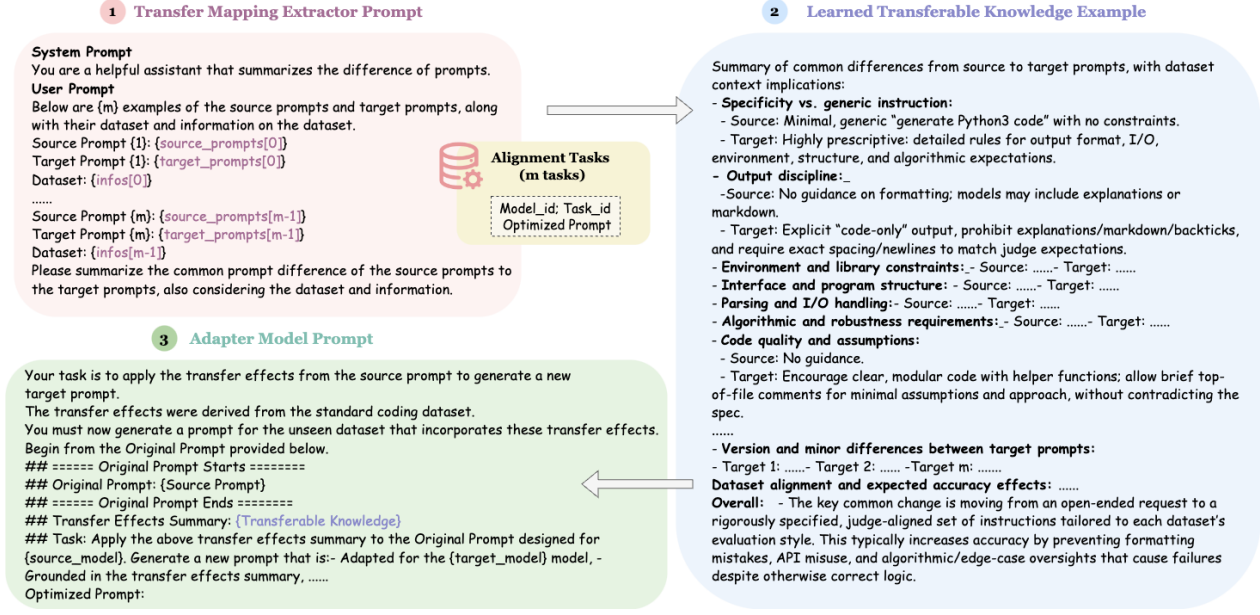


Figure 3 | **A demonstration of PromptBridge.** 1) shows the prompt template used by the Mapping Extractor for learning transferable mappings between source and target prompts. 2) presents an example of the learned transferable knowledge, which may contain multiple sections analyzing cross-model prompt differences. 3) illustrates the prompt template used by the Adapter Model, which applies the learned mapping to adapt prompts for unseen tasks.

(e.g., a high-capability LLM) to analyze each pair $(p_{M_s, S_i}^*, p_{M_t, S_i}^*)$ and generate a concise description of the model-specific transformation pattern. By leveraging the LLM’s ability to summarize structural and stylistic differences, the mapping-model distills the consistent prompt-level adjustments required when transferring prompts from the source model to the target model. This distilled summary constitutes the transferable knowledge, describing how task-optimal prompts for the source model should be reformulated to align with the behavioral, stylistic, and reasoning characteristics of the target model.

Test-Time Prompt Adaptation. For an unseen task $T_j \in T$ with its source prompt p_{M_s, T_j} , PromptBridge applies the learned mapping in context to synthesize a target-model-compatible prompt:

$$\hat{p}_{M_t, T_j} = \mathcal{T}_{M_s \rightarrow M_t}(p_{M_s, T_j}) \approx p_{M_t, T_j}^*, \quad T_j \in T. \quad (3)$$

where \mathcal{T} denotes the transformation extracted from the alignment tasks. At test time, this transformation is implemented through an *adapter model* that conditions both the source prompt and the learned structural mapping. The adapter predicts the structural edits, formatting adjustments, and semantic refinements required for compatibility with the target model, and then generates a fully adapted prompt \hat{p}_{M_t, T_j} . This design enables efficient, zero-shot prompt adaptation to new tasks and target models without additional training or calibration, allowing our approach to generalize beyond the alignment tasks used to learn \mathcal{T} .

PromptBridge introduces a new paradigm for cross-model prompt generalization. It 1) formulates prompt transfer under model drifting as a learnable transformation problem rather than performing independent prompt re-optimization or manually human design, allowing an optimized prompt from one model to be systematically transformed into an effective prompt for another. Details on the feasibility of cross-model prompt transfer are presented in [Section A.2](#) in the Appendix; 2)

extracts model-level differences through linguistic summarization instead of relying on heavy gradient-based signals; and 3) enables zero-shot adaptation to different unseen tasks utilizing the mapping learned from alignment tasks without re-optimization. Together, these features make PromptBridge a scalable, effective, and low-cost bridge for aligning prompt behavior across continuously evolving LLM architectures. Figure 3 illustrates how PromptBridge learns transferable mappings with the Mapping Extractor and applies them through the Adapter Model to adapt prompts for unseen tasks. The complete prompt template used for this process is provided in Section E.1 in the Appendix.

4.2. Calibration: Model-Adaptive Reflective Prompt Evolution

Inspired by AlphaEvolve (Novikov et al., 2025), we propose **Model-Adaptive Reflective Prompt Evolution (MAP-RPE)**, a calibration method that enables a specific model to iteratively evolve and align its prompt through reflective feedback, behavior-aware scoring, and island-based population evolution. Unlike AlphaEvolve, which directly evolves program outputs through population-based code mutation and selection, MAP-RPE focuses on evolving the prompt template itself. This shift allows the model to adaptively improve its own prompt without modifying internal parameters or accessing gradients. The framework leverages a *reflection model* to analyze past generations and propose refined prompts, while the *specific model* executes these prompts to generate responses that provide behavioral and performance feedback.

Reflective Prompt Optimization Process. Instead of typical generate-and-select paradigm (Wang et al., 2023a), our method begins with the source prompt p_{M,S_i} and iteratively refines it for the specific model. For each iteration, a batch of randomly selected inputs from the current task S_i is used to query the target model, producing responses that are then evaluated using designed quantitative metrics, such as performance score and behavior score, to capture both task accuracy and response quality. The reflection model aggregates these results, identifies weaknesses or inconsistencies, and proposes an improved prompt. Through repeated reflection and feedback, the prompt gradually converges to an optimized version p_{M,S_i}^* , where M denotes the specific model.

Island-Based Evolution for Diversity. To prevent premature convergence to local optima, MAP-RPE maintains an island-based population of diverse prompts. Each island stores variants that capture different behavioral patterns or task-solving strategies. In each iteration, the current best-performing prompt serves as the parent to generate new candidate prompts inspired by accumulated feedback. These candidates are re-evaluated and reflected, allowing the system to continuously balance exploration and exploitation in the prompt space.

Compared with GEPA Unlike GEPA (Agrawal et al., 2025), which relies primarily on textual feedback and Pareto-aware optimization over prompt candidates, our method adopts a model-adaptive and metric-driven optimization paradigm. It introduces explicit quantitative metrics, including performance, behavioral, and task-specific scores to evaluate how well each prompt aligns with the target model’s performance characteristics. Guided by these evaluations, the optimizer iteratively refines the source prompt into an optimal variant p_{M,S_i}^* tailored to the specific target model. By maintaining diverse islands of prompt candidates and addressing the prompt–model mismatch problem common in cross-model transfer, MAP-RPE achieves more structured, efficient, and adaptive prompt calibration than GEPA’s general reflection-based framework.

Overall, MAP-RPE provides a gradient-free, model-adaptive calibration mechanism that evolves prompts through structured reflection and population diversity, enabling each specific model to

Table 1 | Pass@1 Accuracy on several code generation datasets when switching the original model GPT-4o to the target models, including o3, o4-mini, Llama3.1-70B-Instruct. Direct Transfer refers to apply the source prompt to the target model directly. The higher the accuracy the better. The best transfer results are highlighted in **Bold**.

Method	HumanEval	MBPP	APPS	xCodeEval	CodeContests
Source Model: GPT-4o	91.10	79.80	12.00	37.03	5.45
Target Model: o3					
Direct Transfer	92.27	77.92	32.67	66.04	48.61
GPT-5 Optimizer	92.68	79.93	36.44	72.64	31.88
MIPROv2 (Opsahl-Ong et al., 2024)	93.70	63.81	32.67	65.09	45.65
GEPA (Agrawal et al., 2025)	96.95	80.01	32.89	69.49	46.19
PromptBridge (Ours)	97.15	80.44	36.44	74.84	56.36
Target Model: o4-mini					
Direct Transfer	96.54	79.09	37.78	72.01	40.12
GPT-5 Optimizer	98.17	70.03	38.44	74.53	57.70
MIPROv2 (Opsahl-Ong et al., 2024)	96.95	77.92	35.33	74.84	48.28
GEPA (Agrawal et al., 2025)	97.56	76.49	37.33	81.13	55.15
PromptBridge (Ours)	98.37	80.60	38.00	77.67	58.79
Target Model: Llama3.1-70B-Instruct					
Direct Transfer	68.70	65.57	10.00	21.38	20.20
GPT-5 Optimizer	78.66	73.38	7.78	23.58	19.80
MIPROv2 (Opsahl-Ong et al., 2024)	78.25	68.26	11.33	18.87	21.82
GEPA (Agrawal et al., 2025)	35.98	25.86	10.22	22.01	18.99
PromptBridge (Ours)	79.88	73.64	11.11	24.53	23.84

achieve its optimal task-specific prompt. Algorithm 1 shows the reflective loop of MAP-RPE.

5. Experiments

We conduct extensive experiments to demonstrate the effectiveness of the proposed prompt transfer method (§ 5.3) and prompt calibration methods (§ 5.4) in single agent and multi-agent settings. Our experiments cover all three real-world application cases: transferring to a more advanced model, a more cost-efficient model, and an open-source model, demonstrating the effectiveness and generalization of the proposed cross-model transfer method.

5.1. Alignment and Unseen Tasks.

Alignment Tasks. Alignment tasks refer to the set of training tasks used to calibrate and align the prompts of the source and target models. These tasks provide the foundation for learning the transferable prompt mapping in our method. We use synthetic-code-generations (Wei et al., 2023) and the train split of CodeContests (Li et al., 2022) as the alignment tasks, which includes basic code generation and complex competitive coding problems.

Unseen Tasks. For extensive evaluation, we use five coding benchmark datasets: two from basic programming and three from complex competitive programming domains. Two simple coding generation benchmarks are: HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021). Due to the absence of sample I/O in MBPP, the evaluation involves randomly removing one test-case from

MBPP-ET (Dong et al., 2025a) for each problem and provide this test-case as a sample I/O for the problem following (Islam et al., 2024). Three complex competitive benchmarks are: Automated Programming Progress Standard (APPS) (Hendrycks et al., 2021), xCodeEval (Khan et al., 2023), and CodeContests (Li et al., 2022).

Moreover, we use SWE-bench Verified (Jimenez et al., 2024) that evaluates how well LLM agents automatically solve the coding task. Given a code repository and a task instruction, the agent is expected to make changes to the repository in order to fulfill the task. The Verified subset is a human-filtered subset of 500 instances. Terminal-Bench (Team, 2025) measures the model’s ability to accomplish complex tasks in a terminal environment. We also evaluate on the validation split of TravelPlanner (Xie et al., 2024). In the two-stage mode, it uses the ReAct framework for information collection and it allows to assess how different LLMs perform under a uniform tool-use framework. The agents are required to give the plan directly based on the information collected by themselves, without employing any other planning strategies. In the sole-planning mode, it aims to assess if the strategies proven effective in other planning benchmarks maintain their efficacy in TravelPlanner.

5.2. Experimental Setup

Baselines. We evaluate three training-free prompt transfer baselines: 1) GPT-5 Optimizer, a powerful prompt optimizer that improving existing prompts and migrating prompts for GPT-5 and other OpenAI models; 2) MIPROv2 (Opsahl-Ong et al., 2024), a widely used prompt optimizer and has been integrated into the DSPy (Khatab et al., 2024) framework. All MIPROv2 optimization runs are performer with the *auto = heavy* setting. 3) GEPA (Agrawal et al., 2025), a prompt optimizer that thoroughly incorporates natural language reflection to learn high-level rules from trial and error. Prompt optimizer with star (MIPROv2*, GEPA*) refer to directly optimized on the test benchmark itself with randomly select 50 questions from it. The detailed experiment setup introduction can be found in the Appendix B. The exact optimized prompt templates are provided in the Appendix E.

Implementation Details. We use GPT-4o (OpenAI et al., 2024) as the source model, o3 (OpenAI, 2025b), o4-mini (OpenAI, 2025b), Llama3.1-70B-Instruct (Grattafiori et al., 2024), Qwen3-32B (Yang et al., 2025) and Gemma3-27B-it (Team et al., 2025) as the target models for coding benchmarks. We utilize GPT-5 (OpenAI, 2025a) as the Mapping Extractor and Adapter Model to transfer prompt. For SWE-Bench, we utilize the framework of mini-SWE-agent (Yang et al., 2024). We use the Terminus framework for evaluation on Terminal-Bench following Zeng et al. (2025). We utilize MapCoder (Islam et al., 2024) as our multi agent system framework. It consists of four LLM agents specifically designed to handle coding problems: recalling relevant examples (Retrieval Agent), planing (Planning Agent), code generation (Coding Agent) and debugging (Debugging Agent).

Evaluation Metric and Setup. We report average Pass@1 accuracy for all coding benchmarks, considering a model successful if its first generated solution is correct. When calling API, we set the temperature as 0.0 to reduce the randomness of each experiments. For o3 and o4-mini, the Azure deployed models do not allow setting temperature manually, so we use the default setting to evaluate on benchmarks, and to reduce the randomness, we run each experiment several times and report the average result. For Qwen3-32B model, on APPS and CodeConetest, we set the max tokens to be 10000. For all datasets, we use temperature = 0.7. For each experiment, we run three times and report the average accuracy to ensure stability and reproducibility of the results. For the SWE-Bench Verified benchmark, we instead report the resolved rate. For Terminal-Bench, we report the accuracy.

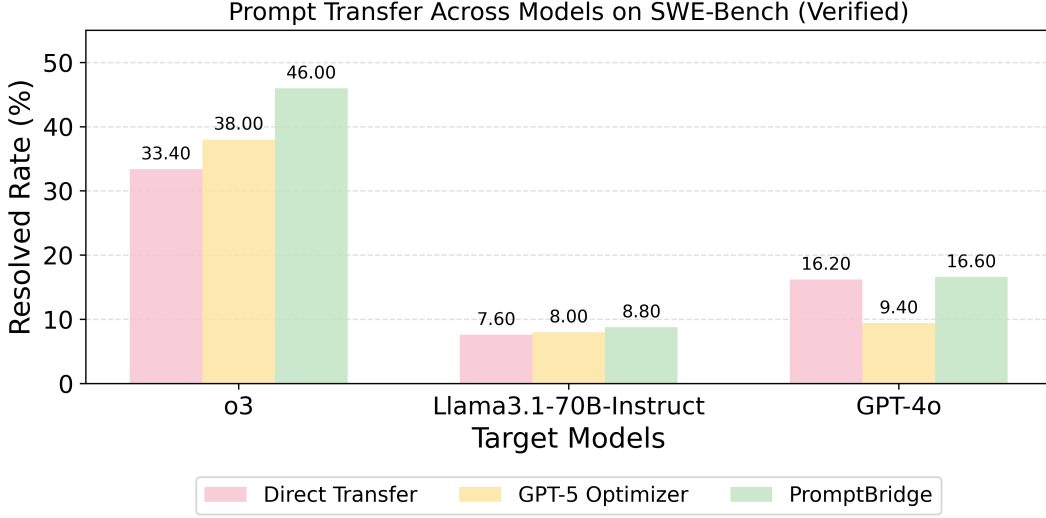


Figure 4 | **Results on SWE-Bench Verified.** We use o4-mini as the source model and evaluate transferability on o3, Llama-3.1-70B-Instruct, and GPT-4o as target models. On the source, the resolved rate with the default prompt is 38.60%, which increases to 42.20% after calibration with the optimized prompt. Direct Transfer refers to directly applying the optimized prompt from o4-mini to the other models without further adaptation.

5.3. Prompt Transfer

Single-Agent Setting. Table 1 presents the generalization performance of the proposed PromptBridge across unseen coding tasks and target LLMs, demonstrating its effectiveness in transferring optimized prompts between models. The first row reports results on the test datasets using the source model, GPT-4o. The target models, o3, o4-mini, and Llama-3.1-70B-Instruct, represent three practical deployment scenarios: a more capable model, a cost-efficient model, and an open-source alternative, respectively. For relatively simple code generation tasks (HumanEval, MBPP), PromptBridge achieves the highest accuracy. On o3, GEPA performs comparably on MBPP but does not consistently match PromptBridge’s effectiveness. Across more challenging problem-solving benchmarks such as APPS, xCodeEval, and CodeContests, PromptBridge demonstrates substantial improvements over direct transfer, achieving relative gains of 11.5%, 13.33%, and 15.92% for o3; 0.58%, 7.86%, and 46.54% for o4-mini; and 11.1%, 14.73%, and 18.02% for Llama-3.1-70B-Instruct, respectively. Since the baseline methods (GPT-5 Optimizer and MIPROv2) are not model-adaptive, their performance remains inconsistent, occasionally competitive but lacking robustness across tasks. GEPA, which requires a task-specific LLM, can implicitly encode model preferences. However, as GEPA directly applies the optimized prompt learned from the training dataset to the test benchmark, it risks overfitting to training-specific features, resulting in reduced generalization. Overall, the results across diverse models and benchmarks highlight the robustness and generalization capability of the proposed transfer method. More results can be found in Table 7 in the Appendix.

Figure 4 presents the transfer results on SWE-Bench Verified. PromptBridge consistently improves performance over direct transfer, demonstrating its effectiveness in adapting optimized prompts across models. The optimized prompt for o4-mini achieves a 9% gain compared with the default prompt. PromptBridge achieves substantial improvements over direct transfer, with relative gains of 27.39% for o3, 15.79% for Llama-3.1-70B-Instruct, and 2.5% for GPT-4o. However, the GPT-5 Optimizer, despite being a strong baseline, fails to deliver stable improvements. On GPT-4o, its performance even falls below direct transfer, with only 9.40%. Figure 5 presents the transfer results on Terminal-

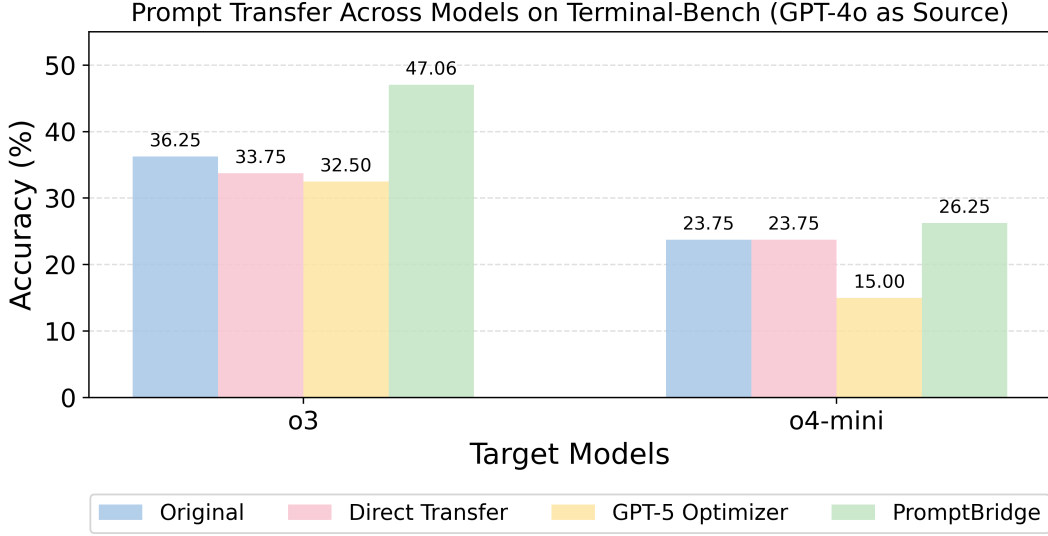


Figure 5 | **Results on Terminal-Bench.** We use GPT-4o as the source model and evaluate transferability on o3 and o4-mini as target models. On the source, the accuracy with the default prompt is 15%, which increases to 18.75% after calibration with the optimized prompt. Original denotes the baseline performance using the default prompt template from Terminus. Direct Transfer refers to directly applying the optimized prompt from GPT-4o to the other models without further adaptation.

Bench. When transferred to different target models, PromptBridge achieves relative improvements of 39.44% on o3 and 10.53% on o4-mini compared with the direct transfer baseline. In contrast, the GPT-5 Optimizer performs poorly when handling cross-model transfer without leveraging transferable knowledge, even yielding worse performance than direct transfer. On o3, the performance of direct transfer does not improve over original prompt. This finding suggests that a prompt optimized for one model or system may not remain optimal across different LLMs, underscoring the importance of discovering model-adaptive prompts for robust cross-model performance. These results further demonstrate the effectiveness and necessity of our framework in mitigating performance degradation across models. More analysis on SWE-Bench and Terminal-Bench can be found in [Section C.2](#) and [Section C.3](#) in the Appendix, respectively.

Multi-Agent System. [Table 2](#) reports the Pass@1 accuracy on the HumanEval dataset in a multi-agent setting under both local and global model drifting scenarios. For global model drifting, PromptBridge consistently improves performance compared with direct transfer, whereas the GPT-5 Optimizer does not yield consistent gains across different target LLMs. For local model drifting, where only specific agents (e.g., coding or planning agents) are switched, PromptBridge is not always the top performer; however, for the debugging agent, it consistently achieves the best results. This observation aligns with the observation that the debugging agent plays a central role in the overall multi-agent workflow as indicated in [Islam et al. \(2024\)](#). In the case of global model drifting, PromptBridge jointly updates all four agent prompts, which may enhance coherence among agents and lead to superior performance. Although the overall performance improvement is not large, our emphasis is on maintaining performance stability across different LLMs after transfer, initially demonstrating that PromptBridge mitigates degradation when model drift occurs in the multi-agent system.

Non-Coding Domain. We also conduct experiments in the non-coding domain. [Table 3](#) presents the results under the two-stage mode. The source model is GPT-4o, and the target model is o3. Using

Table 2 | Pass@1 Accuracy on HumanEval using different calibration methods in multi-agent setting (MapCoder) when model drifting occurs local or global. Source model is GPT-4o, target model is Llama3.1-70B-Instruct. Global Agents denotes that all agents in this multi-agent system switch from source model to target model; Local Coding Agent denotes that the coding agent switches from source model to target model.

Method	Global Model Drifting		Local Model Drifting	
Dataset	All agents	Coding Agent	Planning Agent	Debugging Agent
Source Model	92.68	92.68	92.68	92.68
Target Model: Llama3.1-70B-Instruct				
Direct Transfer	87.59	90.24	92.07	88.41
GPT-5 Optimizer	87.80	87.80	90.24	85.96
PromptBridge (Ours)	87.80	89.63	87.20	89.02
Target Model: o3				
Direct Transfer	94.51	94.51	95.73	93.29
GPT-5 Optimizer	93.29	94.51	96.34	94.51
PromptBridge (Ours)	95.73	93.90	95.12	95.12

Table 3 | Results on TravelPlanner validation set under Two Stage Mode. It involves two agent: React Agent and Planner Agent. Global Model Drifting denotes that all agents switch from the original model to the target model; Local Model Drifting denotes that the planner agent switches from the source model to the target model.

Metric	Final Pass Rate	CC Pass Rate		HC Pass Rate		Delivery Rate
Method		Micro	Macro	Micro	Macro	
Source Model: GPT-4o (Original)	0.0	50.01	2.22	7.62	2.78	66.67
Target Model: o3 (Original)	1.67	41.04	2.22	12.14	10.56	58.89
Global Model Drifting						
Source Model: GPT-4o (Optimized)	2.78	61.25	5.56	19.52	15.56	83.33
Frozen Direct Transfer	0.56	63.96	1.11	28.81	26.67	82.78
PromptBridge - optimized	2.22	61.46	4.44	18.81	15.56	85.56
Local Model Drifting (Planner)						
Source Model: GPT-4o (Optimized)	3.33	60.56	4.44	11.67	13.33	85.56
Frozen Direct Transfer	1.11	67.64	3.89	28.81	23.33	87.78
PromptBridge - optimized	2.78	67.85	4.44	38.33	30.00	86.67

Table 4 | Pass@1 accuracy on multiple code generation datasets using different calibration methods, evaluated on `o3` and `o4-mini`. Frozen Prompt denotes directly applying the prompt of `GPT-4o` to the models without adaptation. It can also be viewed as a default prompt. Average denotes the mean accuracy across the five code benchmarks, and higher values indicate better performance. Higher values indicate better performance. The best results are highlighted in **bold**.

Method	HumanEval	MBPP	APPS	xCodeEval	CodeContests	Average
Model: o3						
Frozen Prompt	92.27	77.92	32.67	66.04	48.61	63.50
MIPROv2*	97.56	47.52	36.00	71.38	40.00	58.49
GEPA*	98.17	83.79	33.11	69.18	46.87	66.22
MAP-RPE	98.37	80.86	36.67	74.53	58.79	69.84
Model: o4-mini						
Frozen Prompt	96.54	79.09	37.78	72.01	40.12	65.11
MIPROv2*	98.78	15.95	38.66	78.30	32.73	52.88
GEPA*	98.78	79.26	34.89	76.73	56.23	69.18
MAP-RPE	96.95	79.60	39.33	76.42	47.88	68.04

the default prompt leads to pass rates of 0.0% and 1.67%, respectively. When the optimized prompts are applied to both the React and Planner agents, the pass rates increase: Global reaches 2.78%, and Local achieves 3.33%. In this case, optimizing the React prompt slightly decreases performance. PromptBridge achieves consistent improvements over direct transfer and outperforms the default prompt configuration. Notably, even in this non-coding setting, PromptBridge enhances performance despite being trained exclusively on coding-domain datasets. The results of TravelPlanner under the sole-planning mode are provided in the [Section C.4](#) in the Appendix.

5.4. Prompt Calibration

Prompt Calibration aims to generate model-adaptive and task-specific optimal prompts. To evaluate the effectiveness of different optimization strategies, we compare our proposed MAP-RPE with several existing methods, including GEPA and MIPROv2, using five code generation benchmarks. As shown in [Table 4](#), we evaluate optimized prompts on `o3` and `o4-mini` as target specific models. For the `o3`, MAP-RPE consistently yields high accuracy across all datasets and has the highest average accuracy, achieving the best performance on xCodeEval (74.53) and CodeContests (58.79), indicating that our method effectively tailors prompts to model-specific behavior. On `o4-mini`, MAP-RPE maintains balanced performance, achieving the top results on MBPP and APPS, while avoiding the instability seen in MIPROv2. It achieves average accuracy comparable to that of GEPA. Across all settings, MAP-RPE consistently delivers strong and stable performance, outperforming frozen prompt and achieving competitive or superior results compared with prior SOTA prompt optimization frameworks. These results demonstrate that reflective, metric-guided evolution enables robust calibration of prompts to model-specific and task-specific characteristics, forming the foundation for reliable prompt adaptation within the proposed PromptBridge framework.

5.5. Ablation Study

Ablation Study: Understanding Cross-Model Prompt Transferability To further evaluate the effectiveness of PromptBridge, we conduct an ablation study comparing three transfer strategies: 1)

Table 5 | Pass@1 Accuracy on several code generation datasets when switching the original model GPT-4o to the target models, including o3, o4-mini, Llama3.1-70B-Instruct. Direct Transfer refers to apply the source prompt to the target model directly. The higher the accuracy the better. The best transfer results are highlighted in **Bold**.

Method	HumanEval	MBPP	APPS	xCodeEval	CodeContests
Source Model: GPT-4o	91.10	79.80	12.00	37.03	5.45
Target Model: o3					
Direct Transfer	92.27	77.92	32.67	66.04	48.61
One-shot ICL	89.63	79.43	33.11	71.66	52.32
Few-shot ICL	87.40	77.70	36.89	71.38	51.63
PromptBridge (Ours)	97.15	80.44	36.44	74.84	56.36
Target Model: o4-mini					
Direct Transfer	96.54	79.09	37.78	72.01	40.12
One-shot ICL	97.56	79.18	35.33	76.10	57.78
Few-shot ICL	97.56	80.52	37.33	77.36	57.17
PromptBridge (Ours)	98.37	80.60	38.00	77.67	58.79
Target Model: Llama3.1-70B-Instruct					
Direct Transfer	68.70	65.57	10.00	21.38	20.20
One-shot ICL	41.26	36.78	10.22	19.18	22.83
Few-shot ICL	78.46	72.04	5.11	21.70	3.64
PromptBridge (Ours)	79.88	73.64	11.11	24.53	23.84

One-shot in-context learning (ICL), which uses the single highest-scoring example source-target pair as the demonstration; 2) *Few-shot ICL*, which uses the top five pairs as demonstrations; and 3) our *PromptBridge*, which first distills the systematic prompt differences across alignment tasks and then applies this mapping to unseen tasks. Table 5 shows that one-shot ICL and few-shot ICL provide only limited improvements and are often unstable due to their high sensitivity to the specific example pairs used during in-context conditioning. This instability likely arises because the model imitates superficial, task-specific patterns in the provided examples rather than learning a generalizable transformation between models. In contrast, PromptBridge consistently achieves the highest transfer performance. By summarizing cross-model prompt differences into an explicit transformation rule, it captures stable and reusable adaptation patterns that generalize beyond the observed examples. These results demonstrate that learning a high-level, model-level transformation is substantially more robust and effective than relying solely on direct ICL examples.

Ablation study on the calibration question number. Figure 6 reports the Pass@1 accuracy on HumanEval under different calibration strategies (Generation-Selection, see Section B.1 in the Appendix) and varying numbers of calibration questions. Across all models, MAP-RPE consistently achieves the highest accuracy, demonstrating its effectiveness in adaptively aligning prompt semantics to model-specific behaviors. For instance, on Llama3.1-70B-Instruct, MAP-RPE achieves up to 79.5% Pass@1, significantly outperforming the default prompt baseline (68.7%); on o3, MAP-RPE reaches 98.4%. The ablation on the number of calibration questions ($n \in 5, 20, 50$) further illustrates that performance improves steadily as n increases, since a larger calibration set exposes the optimizer to a more diverse range of reasoning and formatting patterns. Notably, MAP-RPE is sample-efficient: even with only five calibration questions, it already achieves substantial gains over the default prompt

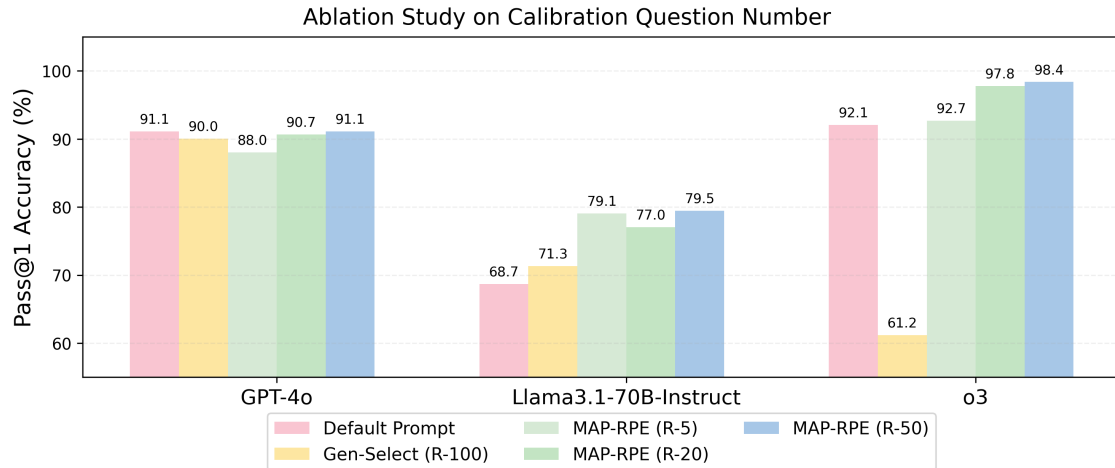


Figure 6 | **Ablation study on the number of calibration questions.** Pass@1 accuracy on HumanEval using different calibration methods to derive the optimal prompt for each model. Default Prompt denotes the original designed prompt for solving coding tasks. Gen-Select refers to generating multiple prompt candidates and selecting the best-performing one.

on Llama3.1-70B-Instruct and o3. Compared with random Generation-Selection, which often exhibits unstable performance and may occasionally find a locally optimal candidate, MAP-RPE remains consistently strong, highlighting the advantage of reflective, metric-guided calibration and systematic adaptivity over heuristic search.

6. Conclusions

We investigated the overlooked challenge of *prompt transfer under model drifting*: the degradation that occurs when prompts tuned for one LLM are reused on another. We formalized this problem by empirical analysis and introduced a practical framework, PromptBridge, that (1) evolves task- and model-specific prompts and (2) learns cross-model transfer mapping to translate prompts for unseen tasks when switching models. Across single-agent and multi-agent settings, our experiments show consistent gains in performance, indicating that efficient prompt transfer is important for sustainable LLM development. Treating prompts as migratable artifacts rather than static instructions offers a practical path to maintain performance as LLM ecosystems evolve. Our framework provides a new paradigm to make model switching more reliable and efficient.

Future Work. In future, we plan to invest in research across the following directions.

- **Broader Model Families and Alignment-aware Variants.** Extending the study to additional model families, such as Mistral and DeepSeek, and to alignment-aware variants (e.g., SFT-, DPO-, or GRPO-aligned models) would provide a more complete characterization of how alignment methods reshape prompt patterns.
- **Behavioral and Stylistic Drift.** Beyond performance, large models often exhibit noticeable shifts in style, verbosity, structure, or safety behavior when moving across models (e.g., GPT-4o vs. GPT-5). We plan to conduct a systematic study of how stylistic and behavioral patterns drift, and how such drift affects prompt robustness.
- **Community Benchmarks for Prompt Portability.** We plan to develop shared benchmarks and standardized corpora specifically designed to evaluate cross-model prompt portability across families, scales, and alignment settings.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Eshaan Agarwal, Joykirat Singh, Vivek Dani, Raghav Magazine, Tanuja Ganu, and Akshay Nambi. Promptwizard: Task-aware prompt optimization framework. *arXiv preprint arXiv:2405.18369*, 2024.
- Lakshya A Agrawal, Shangyin Tan, Dilara Soylu, Noah Ziemis, Rishi Khare, Krista Opsahl-Ong, Arnav Singhvi, Herumb Shandilya, Michael J Ryan, Meng Jiang, et al. Gepa: Reflective prompt evolution can outperform reinforcement learning. *arXiv preprint arXiv:2507.19457*, 2025.
- AI Anthropic. Introducing the next generation of claude, 2024.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chen Qian, Chi-Min Chan, Yujia Qin, Yaxi Lu, Ruobing Xie, et al. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors in agents. *arXiv preprint arXiv:2308.10848*, 2(4):6, 2023.
- Yuyan Chen, Zhihao Wen, Ge Fan, Zhengyu Chen, Wei Wu, Dayiheng Liu, Zhixu Li, Bang Liu, and Yanghua Xiao. Mapo: Boosting large language model performance with model-adaptive prompt optimization. *arXiv preprint arXiv:2407.04118*, 2024.
- Jiale Cheng, Xiao Liu, Kehan Zheng, Pei Ke, Hongning Wang, Yuxiao Dong, Jie Tang, and Minlie Huang. Black-box prompt optimization: Aligning large language models without model training. *arXiv preprint arXiv:2311.04155*, 2023.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025.
- Yihong Dong, Jiazheng Ding, Xue Jiang, Ge Li, Zhuo Li, and Zhi Jin. Codescore: Evaluating code generation by learning code execution. *ACM Transactions on Software Engineering and Methodology*, 34(3):1–22, 2025a.
- Yihong Dong, Xue Jiang, Jiaru Qian, Tian Wang, Kechi Zhang, Zhi Jin, and Ge Li. A survey on code generation with llm-based agents. *arXiv preprint arXiv:2508.00083*, 2025b.
- Chrisantha Fernando, Dylan Banarse, Henryk Michalewski, Simon Osindero, and Tim Rocktäschel. Promptbreeder: Self-referential self-improvement via prompt evolution. *arXiv preprint arXiv:2309.16797*, 2023.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

- Xinyu Guan, Li Lyna Zhang, Yifei Liu, Ning Shang, Youran Sun, Yi Zhu, Fan Yang, and Mao Yang. rstar-math: Small llms can master math reasoning with self-evolved deep thinking. *arXiv preprint arXiv:2501.04519*, 2025.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Qingyan Guo, Rui Wang, Junliang Guo, Bei Li, Kaitao Song, Xu Tan, Guoqing Liu, Jiang Bian, and Yujiu Yang. Connecting large language models with evolutionary algorithms yields powerful prompt optimizers. *arXiv preprint arXiv:2309.08532*, 2023.
- Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, et al. Measuring coding challenge competence with apps. *arXiv preprint arXiv:2105.09938*, 2021.
- Dong Huang, Jie M Zhang, Michael Luck, Qingwen Bu, Yuhao Qing, and Heming Cui. Agentcoder: Multi-agent-based code generation with iterative testing and optimisation. *arXiv preprint arXiv:2312.13010*, 2023.
- Md Ashraful Islam, Mohammed Eunus Ali, and Md Rizwan Parvez. Mapcoder: Multi-agent code generation for competitive problem solving. *arXiv preprint arXiv:2405.11403*, 2024.
- Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. Mistral 7b, 2023. URL <https://arxiv.org/abs/2310.06825>.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=VTF8yNQm66>.
- Mohammad Abdullah Matin Khan, M Saiful Bari, Xuan Long Do, Weishi Wang, Md Rizwan Parvez, and Shafiq Joty. xcodeeval: A large scale multilingual multitask benchmark for code understanding, generation, translation and retrieval. *arXiv preprint arXiv:2303.03004*, 2023.
- Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Saiful Haq, Ashutosh Sharma, Thomas T Joshi, Hanna Moazam, Heather Miller, et al. Dspy: Compiling declarative language model calls into state-of-the-art pipelines. In *The Twelfth International Conference on Learning Representations*, 2024.
- Dong-Ki Kim, Sungryull Sohn, Lajanugen Logeswaran, Dongsu Shim, and Honglak Lee. Multiprompter: Cooperative prompt optimization with multi-agent reinforcement learning. *arXiv preprint arXiv:2310.16730*, 2023.
- Yubo Li, Xiaobin Shen, Xinyu Yao, Xueying Ding, Yidi Miao, Ramayya Krishnan, and Rema Padman. Beyond single-turn: A survey on multi-turn interactions with large language models. *arXiv preprint arXiv:2504.04717*, 2025.

- Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien de Masson d’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal, Alexey Cherepanov, James Molloy, Daniel Mankowitz, Esme Sutherland Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-level code generation with alphacode. *arXiv preprint arXiv:2203.07814*, 2022.
- Fengyuan Liu, Nouar AlDahoul, Gregory Eady, Yasir Zaki, and Talal Rahwan. Self-reflection makes large language models safer, less biased, and ideologically neutral. *arXiv preprint arXiv:2406.10400*, 2024.
- Siwei Liu, Jinyuan Fang, Han Zhou, Yingxu Wang, and Zaiqiao Meng. Sew: Self-evolving agentic workflows for automated code generation. *arXiv preprint arXiv:2505.18646*, 2025.
- Seth Neel and Peter Chang. Privacy issues in large language models: A survey. *arXiv preprint arXiv:2312.06717*, 2023.
- Alexander Novikov, Ngân Vu, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco JR Ruiz, Abbas Mehrabian, et al. Alphaevolve: A coding agent for scientific and algorithmic discovery. *Google DeepMind*, 2025.
- OpenAI. Gpt-5 system card. Technical report, OpenAI, August 2025a. URL <https://cdn.openai.com/gpt-5-system-card.pdf>.
- OpenAI. Openai o3 and o4-mini system card. Technical report, OpenAI, April 2025b. URL <https://cdn.openai.com/pdf/2221c875-02dc-4789-800b-e7758f3722c1/o3-and-o4-mini-system-card.pdf>.
- OpenAI, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and . . . Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. URL <https://arxiv.org/abs/2410.21276>. arXiv:2410.21276.
- Krista Opsahl-Ong, Michael J Ryan, Josh Purtell, David Broman, Christopher Potts, Matei Zaharia, and Omar Khattab. Optimizing instructions and demonstrations for multi-stage language model programs. *arXiv preprint arXiv:2406.11695*, 2024.
- Nathanaël Carraz Rakotonirina, Roberto Dessi, Fabio Petroni, Sebastian Riedel, and Marco Baroni. Can discrete information extraction prompts generalize across language models? *arXiv preprint arXiv:2302.09865*, 2023.
- Protection Regulation. General data protection regulation. *Intouch*, 25:1–5, 2018.
- Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*, 2019.
- Asankhaya Sharma. Openevolve: an open-source evolutionary coding agent, 2025. URL <https://github.com/codelion/openevolve>.
- Hongjin Su, Ruoxi Sun, Jinsung Yoon, Pengcheng Yin, Tao Yu, and Sercan Ö Arik. Learn-by-interact: A data-centric framework for self-adaptive agents in realistic environments. *arXiv preprint arXiv:2501.10893*, 2025.
- Yusheng Su, Xiaozhi Wang, Yujia Qin, Chi-Min Chan, Yankai Lin, Huadong Wang, Kaiyue Wen, Zhiyuan Liu, Peng Li, Juanzi Li, et al. On transferability of prompt tuning for natural language processing. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 3949–3969, 2022.

- Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*, 2025.
- The Terminal-Bench Team. Terminal-bench: A benchmark for ai agents in terminal environments, Apr 2025. URL <https://github.com/laude-institute/terminal-bench>.
- Mudit Verma, Siddhant Bhambri, and Subbarao Kambhampati. On the brittle foundations of react prompting for agentic large language models. *arXiv preprint arXiv:2405.13966*, 2024.
- Junlin Wang, Jue Wang, Ben Athiwaratkun, Ce Zhang, and James Zou. Mixture-of-agents enhances large language model capabilities. *arXiv preprint arXiv:2406.04692*, 2024.
- Xinyuan Wang, Chenxi Li, Zhen Wang, Fan Bai, Haotian Luo, Jiayou Zhang, Nebojsa Jojic, Eric P Xing, and Zhiting Hu. Promptagent: Strategic planning with language models enables expert-level prompt optimization. *arXiv preprint arXiv:2310.16427*, 2023a.
- Xun Wang, Jing Xu, Franziska Boenisch, Michael Backes, Christopher A Choquette-Choo, and Adam Dziedzić. Efficient and privacy-preserving soft prompt transfer for llms. *arXiv preprint arXiv:2506.16196*, 2025.
- Zhen Wang, Rameswar Panda, Leonid Karlinsky, Rogerio Feris, Huan Sun, and Yoon Kim. Multitask prompt tuning enables parameter-efficient transfer learning. *arXiv preprint arXiv:2303.02861*, 2023b.
- Chenxing Wei, Yao Shu, Mingwen Ou, Ying Tiffany He, and Fei Richard Yu. Paft: Prompt-agnostic fine-tuning. *arXiv preprint arXiv:2502.12859*, 2025a.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Jason Wei, Zhiqing Sun, Spencer Papay, Scott McKinney, Jeffrey Han, Isa Fulford, Hyung Won Chung, Alex Tachard Passos, William Fedus, and Amelia Glaese. Browsecomp: A simple yet challenging benchmark for browsing agents. *arXiv preprint arXiv:2504.12516*, 2025b.
- Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. Magicoder: Source code is all you need. *arXiv preprint arXiv:2312.02120*, 3, 2023.
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, et al. Autogen: Enabling next-gen llm applications via multi-agent conversations. In *First Conference on Language Modeling*, 2024.
- Shirley Wu, Parth Sarthi, Shiyu Zhao, Aaron Lee, Herumb Shandilya, Adrian Mladenovic Grobelnik, Nurendra Choudhary, Eddie Huang, Karthik Subbian, Linjun Zhang, et al. Optimas: Optimizing compound ai systems with globally aligned local rewards. *arXiv preprint arXiv:2507.03041*, 2025.
- Shijie Xia, Xuefeng Li, Yixin Liu, Tongshuang Wu, and Pengfei Liu. Evaluating mathematical reasoning beyond accuracy. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 27723–27730, 2025.
- Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and Yu Su. Travelplanner: A benchmark for real-world planning with language agents. *arXiv preprint arXiv:2402.01622*, 2024.

- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun Chen. Large language models as optimizers. In *The Twelfth International Conference on Learning Representations*, 2023.
- John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik R Narasimhan, and Ofir Press. SWE-agent: Agent-computer interfaces enable automated software engineering. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://arxiv.org/abs/2405.15793>.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023.
- Siyu Yuan, Kaitao Song, Jiangjie Chen, Xu Tan, Dongsheng Li, and Deqing Yang. Evoagent: Towards automatic multi-agent generation via evolutionary algorithms. *arXiv preprint arXiv:2406.14228*, 2024.
- Mert Yuksekgonul, Federico Bianchi, Joseph Boen, Sheng Liu, Pan Lu, Zhi Huang, Carlos Guestrin, and James Zou. Optimizing generative ai by backpropagating language model feedback. *Nature*, 639(8055):609–616, 2025.
- Aohan Zeng, Xin Lv, Qinkai Zheng, Zhenyu Hou, Bin Chen, Chengxing Xie, Cunxiang Wang, Da Yin, Hao Zeng, Jiajie Zhang, et al. Glm-4.5: Agentic, reasoning, and coding (arc) foundation models. *arXiv preprint arXiv:2508.06471*, 2025.
- Guibin Zhang, Kaijie Chen, Guancheng Wan, Heng Chang, Hong Cheng, Kun Wang, Shuyue Hu, and Lei Bai. Evoflow: Evolving diverse agentic workflows on the fly. *arXiv preprint arXiv:2502.07373*, 2025a.
- Jian Zhang, Zhangqi Wang, Haiping Zhu, Jun Liu, Qika Lin, and Erik Cambria. Mars: A multi-agent framework incorporating socratic guidance for automated prompt optimization. *arXiv preprint arXiv:2503.16874*, 2025b.
- Ningyuan Zhang, Jie Lu, Keqiuyin Li, Zhen Fang, and Guangquan Zhang. Release the powers of prompt tuning: Cross-modality prompt transfer. In *The Thirteenth International Conference on Learning Representations*, 2025c.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.
- Hongling Zheng, Li Shen, Yong Luo, Tongliang Liu, Jialie Shen, and Dacheng Tao. Decomposed prompt decision transformer for efficient unseen task generalization. *Advances in Neural Information Processing Systems*, 37:122984–123006, 2024.
- Han Zhou, Xingchen Wan, Lev Proleev, Diana Mincu, Jilin Chen, Katherine Heller, and Subhrajit Roy. Batch calibration: Rethinking calibration for in-context learning and prompt engineering. *arXiv preprint arXiv:2309.17249*, 2023.
- Han Zhou, Xingchen Wan, Ruoxi Sun, Hamid Palangi, Shariq Iqbal, Ivan Vulić, Anna Korhonen, and Serkan Ö Arik. Multi-agent design: Optimizing agents with better prompts and topologies. *arXiv preprint arXiv:2502.02533*, 2025.
-

Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. Large language models are human-level prompt engineers. In *The eleventh international conference on learning representations*, 2022.

Zixiao Zhu, Hanzhang Zhou, Zijian Feng, Tianjiao Li, Chua Jia Jim Deryl, Mak Lee Onn, Gee Wah Ng, and Kezhi Mao. Rethinking prompt optimizers: From prompt merits to optimization. *arXiv preprint arXiv:2505.09930*, 2025.

Appendix Contents

A	Preliminary Details	23
A.1	Necessity of Cross-Model Prompt Transfer	23
A.2	Feasibility of Cross-Model Prompt Transfer	25
B	Experimental Settings	26
B.1	Baseline Methods	26
B.2	Dataset Description	27
B.3	Calibration Details and Setup	28
B.4	MapCoder Setup	29
B.5	TravelPlanner Setup	29
C	More Experimental Results	30
C.1	Coding Benchmarks	30
C.2	SWE-Bench	31
C.3	Terminal-Bench	31
C.4	TravelPlanner	32
C.5	Consistency Analysis	32
C.6	Ablation Study	35
D	Discussions	35
E	Prompt Template	36
E.1	PromptBridge	36
E.2	Optimized Prompt for Coding Benchmarks	39
E.3	Optimized Prompt for SWE-Bench and Terminal-Bench	41
E.4	Optimized Prompt for TravelPlanner	44
E.5	Optimized Prompt for Coding Benchmarks using Baseline Methods	49

A. Preliminary Details

A.1. Necessity of Cross-Model Prompt Transfer

Detailed setup for Figure 1. For each model, we first apply our prompt evolution method to obtain the optimal prompt template, using a calibration number of 50. In other words, the results of MAP-RPE (Random-50) in Table 4 correspond to the diagonal entries in figure (b). The reported results represent the average accuracy over three runs. When transferring the optimized prompt

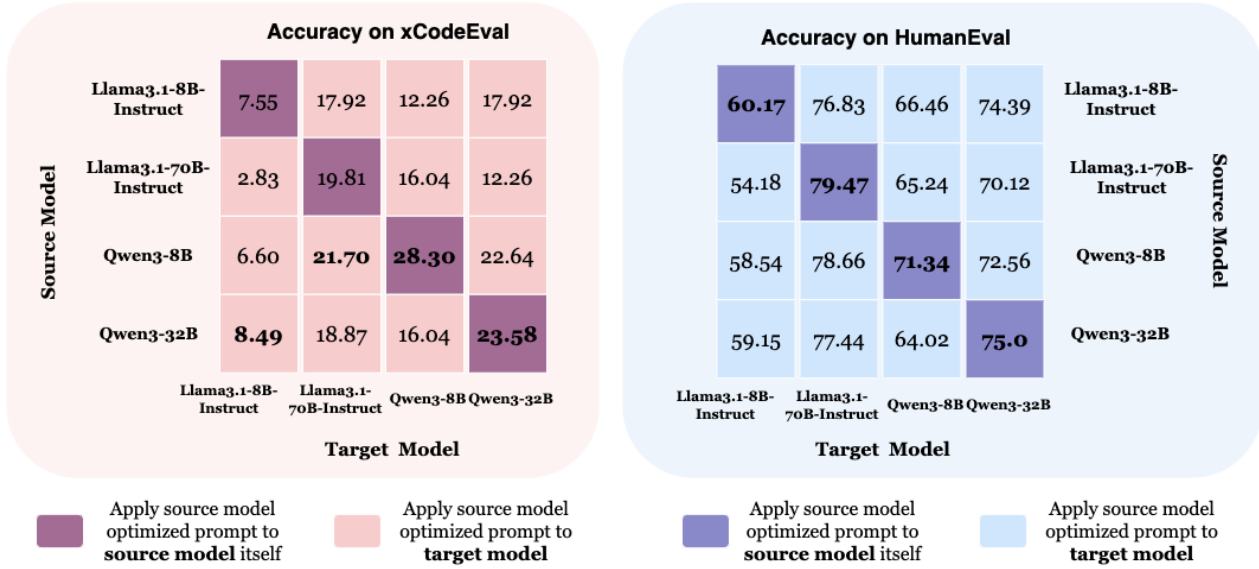


Figure 7 | **Cross-model frozen prompt transfer on xCodeEval (left) and HumanEval (right).** Each row shows the source model and each column the target model. Darker diagonal cells indicate the accuracy obtained when applying a model’s own optimized prompt, while lighter off-diagonal cells show performance when transferring that optimized prompt to a different model without further adaptation. The consistent diagonal–off-diagonal gap reveals substantial model drifting: prompts carefully optimized for one model fail to generalize to others, even within the same family.

from a source model to a target model, the transferred prompt does not necessarily yield the best performance. We also conduct experiments within the same GPT family.

Model Drifting Demonstration Beyond the results shown in Figure 1, we further evaluate cross-model transfer within the same model family but at different scales. As shown in Figure 7, both xCodeEval and HumanEval exhibit a consistent pattern of model drifting, even when source and target models share the same architecture. In nearly all cases, the diagonal entries, where a model uses its own optimized prompt, achieve the highest accuracy, while transferring that prompt to other scales within the same family leads to substantial degradation (e.g., Llama3.1-70B-Instruct → Llama3.1-8B-Instruct, Qwen3-32B → Qwen3-8B). The drop can be dramatic, reaching up to 50–70% relative loss on xCodeEval, indicating that the optimal prompt for a model does not smoothly scale across model sizes. Although HumanEval displays higher absolute accuracy, the same diagonal–off-diagonal gap persists, confirming that even closely related models do not share a stable prompt optimum. These results demonstrate that model drifting arises not only across different model families, but also within a single family across scales, further underscoring the need for model-adaptive prompt transfer mechanisms such as PromptBridge .

Default prompt in Multi-Agent System Model drifting is not limited to optimized prompts tailored for specific source models; it can also occur with default, human-designed prompts that are intended to be universal across models but are not guaranteed to be optimal for any of them. Here, we evaluate the default performance using the original prompt templates for GPT-4o, Llama3.1-70B-Instruct, and other models. Then we use GPT-5 optimizer to refine the debugging agent’s prompt, producing an optimized version for comparison. Although default prompt templates are typically designed by humans to work broadly across models, they are seldom optimal for each specific model (see Table 2). Since model drifting is defined as the performance change caused by substituting the underlying

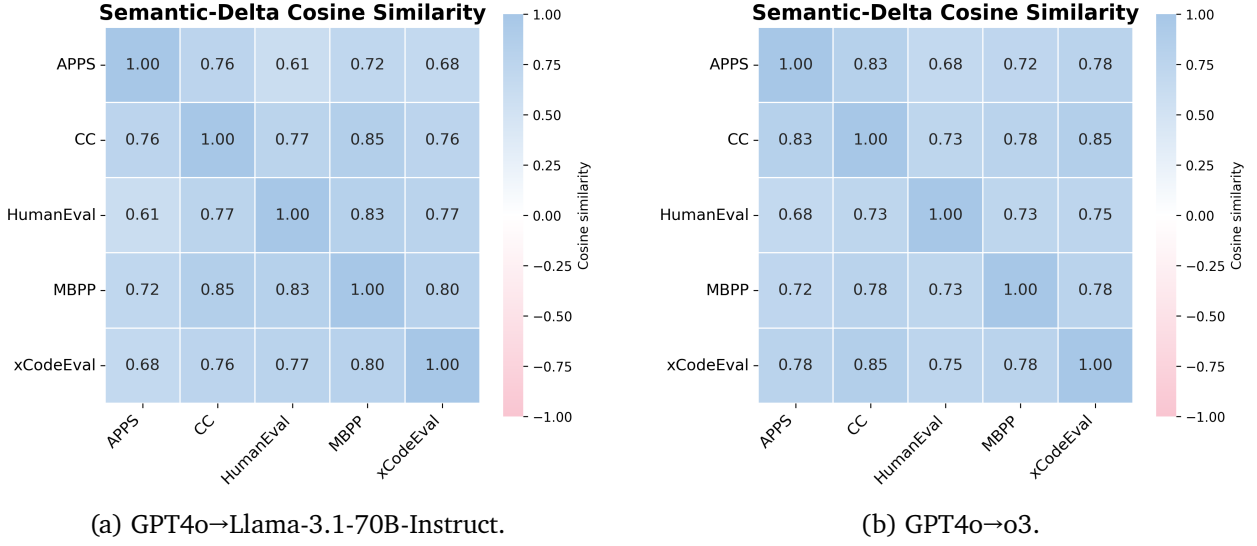


Figure 8 | **Semantic-delta cosine similarity across tasks for cross-model prompt transfer.** We compute the semantic-delta embeddings which captures how each task’s output distribution shifts under prompt transfer, and measure the cosine similarity between tasks. Higher similarity indicates shared transfer patterns. (a) compares GPT-4o → Llama-3.1-70B-Instruct; (b) compares GPT-4o → o3.

model, such universal designs without optimization can still exhibit model drifting when directly transferred. This observation underscores the importance of applying targeted prompt transfer or adaptation to existing prompt templates when deploying systems across different LLMs.

A.2. Feasibility of Cross-Model Prompt Transfer

We hypothesize that there exists a shared pattern of prompt change that generalizes across tasks when transferring between models. To examine this, we study whether the edit patterns required to transform an optimal source prompt into its target counterpart remain consistent across tasks. If such pattern exists, it would indicate that model adaptation follows a stable semantic trajectory, suggesting that prompts can be transferred between models through shared transformation principles rather than being re-optimized from scratch for every task. We evaluate this feasibility from semantic perspectives, analyzing how prompt deltas behave across models and tasks.

Experimental Setup We consider fixed source→target pairs (e.g., from GPT-4o to o3; from GPT-4o to Llama-3.1-70B-Instruct). For each pair we hold both models fixed during analysis. We use five code generation benchmarks, which are unseen tasks in the main text. Here, our goal is not to evaluate task performance itself but to analyze the pattern changes; thus, the specific tasks are not the focus. For each task t , we construct $\mathcal{P}_{\text{source},t}^*$, which refers to the source-optimal prompt found on the source model; $\mathcal{P}_{\text{target},t}^*$, which refers to the target-optimal prompt found on the target model. We compute the Semantic delta $v_t = e(\mathcal{P}_{\text{target},t}^*) - e(\mathcal{P}_{\text{source},t}^*)$, where $e(\cdot)$ is a fixed text embedding (here we use all-MiniLM-L6-v2 model¹) We compute pairwise cosine similarities $\cos(v_i, v_j)$ for all task pairs $i \neq j$, higher values indicate a shared direction of edits from source to target.

Semantic Delta Embedding Similarity Figure 8 presents the pairwise cosine similarity matrices of semantic deltas across five tasks for two source→target model pairs. Blue regions indicate higher

¹<https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>

alignment between task-specific prompt adaptations. The consistently positive similarities reveal that source-to-target prompt transformations exhibit shared latent structures, suggesting the existence of a transferable pattern across tasks and thus supporting the feasibility of cross-task prompt transfer. However, embedding similarity alone does not imply functional equivalence. To further validate feasibility, we assess performance changes when applying PromptBridge to unseen tasks, demonstrating that learned transfer effects can generalize beyond the observed task space. While our current work primarily focuses on establishing the model drifting problem and developing a general cross-model prompt transfer framework, these findings hint that incorporating task-level adaptation could further enhance transfer effectiveness. Exploring such cross-task prompt transfer represents an interesting direction for future work to make prompt transfer more powerful and context-aware.

B. Experimental Settings

B.1. Baseline Methods

In this section, we introduce all the baselines used in the experiments.

GPT5 Optimizer For the five coding benchmarks, we provide GPT-5 with the instruction shown in [Listing 1](#), prompting it to generate an optimized prompt based on the original one. Note that this process is essentially zero-shot, as the model is only given the name of the target model without any additional examples or training data.

```

You are an expert in optimizing prompts for code generation tasks - {dataset
}. Your task is to optimize the following prompt for generating Python code
that solves a specific problem. The original prompt is provided below.
Your goal is to make the prompt more concise, clear, and effective for the {
specific_model} model, while ensuring it remains understandable and retains
the original intent. You may add additional guidance or context that you
believe will help the model generate better responses.
## ===== Original Prompt Starts =====
## Original Prompt: {original_prompt}
## ===== Original Prompt Ends =====
Please provide an optimized version of the prompt that is suitable for the {
specific_model} model and can elicit better response.

Optimized Prompt:

```

Listing 1 | The prompt used in GPT5 Optimizer baseline for five coding datasets.

MIPROv2 (Opsahl-Ong et al., 2024) is a widely used prompt optimizer and has been integrated into the DSPy (Khattab et al., 2024) framework. It works by jointly optimizing both instructions and demonstrations using Bayesian optimization. For each problem module, it first bootstraps candidate sets of instructions and demonstrations, assigning uniform priors over their utilities. Candidate assignments are proposed with the Tree-Structured Parzen Estimator (TPE), and the Bayesian model is updated based on evaluation scores to favor high-performing candidates. The most probable sets of instructions and demonstrations are then selected and validated to obtain the final optimized program configuration. All MIPROv2 optimization runs are performed with the *auto = heavy* setting. For our training dataset, we use a textual similarity score between the prediction and the ground truth response, serving as a simple evaluation metric. We combine all the training dataset into a whole one to perform the optimizer.

Table 6 | Alignment Tasks and Unseen Tasks information.

Dataset Name	Domain	Data Number
Alignment Tasks		
synthetic-code-generations (Wei et al., 2023)	Coding	-
CodeContests (Train) (Li et al., 2022)	Coding	-
Unseen Tasks		
HumanEval (Chen et al., 2021)	Coding	164
MBPP (Austin et al., 2021)	Coding	397
APPS (Hendrycks et al., 2021)	Coding	150
xCodeEval (Khan et al., 2023)	Coding	106
CodeContests (Li et al., 2022)	Coding	156
SWE-bench Verified Jimenez et al. (2024)	Agent	500
Terminal Bench (Team, 2025)	Agent	80
TravelPlanner - Validation (Xie et al., 2024)	Planning	180

GEPA (Agrawal et al., 2025) is a prompt optimizer that thoroughly incorporates natural language reflection to learn high-level rules from trial and error. It works iteratively-proposing a new candidate in every iteration by improving some existing candidates using either reflective prompt mutation or system aware merge. To avoid local optimum, it introduces Pareto-based candidate sampling, which filters and samples from the list of best candidates per task.

MIPROv2* and GEPA* denote the variants where the optimizers are applied on the test benchmark. Specifically, we randomly select 50 questions from the test dataset and use them to derive the optimized user instruction. All other parameters are kept identical to those used in MIPROv2 and GEPA, respectively.

Calibration: Generation-Selection Wang et al. (2023a) Given n questions from a task T_i , generate the optimal prompt p_i for the task on target model. Generation-selection is a naive method, first we utilize GPT-5 to generate several candidate prompts, then we use the similarity between the generated response and the original response which is generated by the source model as a selection metric to select those candidates with highest similarity score.

B.2. Dataset Description

For the training dataset, we randomly sample 30 subtasks from the synthetic-code-generation dataset, with each subtask containing 100 question–answer pairs. For the CodeContests dataset, we select 24 tasks based on their difficulty levels. For each subtask or task, we apply the proposed MAP-RPE method to calibrate the optimal prompt and record its corresponding evaluation results. In total, this yields $30 + 24 = 54$ alignment tasks used for calibration. To evaluate the performance, we use five coding benchmark datasets: two from basic programming and three from complex competitive programming domains. Following Islam et al. (2024), the problem set size of HumanEval, MBPP, APPS, xCodeEval, and CodeContests are 164, 397, 150, 106, 156, respectively. We also evaluate on SWE-Bench Verified, Terminal-Bench, and TravelPlanner.

Algorithm 1: Model Adaptive Reflective Prompt Evolution (MAP-RPE) for Model Calibration

Input: Task dataset S_i , target model M_t , source model M_s (optional), prompt database \mathcal{B} with K islands, reflection model \mathcal{M}_{ref} , maximum global iterations G , local evolution steps L per calibration question, evaluation metrics \mathcal{E} (e.g., performance, behavioral score).

Output: Optimized prompt p_{M_t, S_i}^*

Initialization: Load dataset S_i and construct the initial template $p_0 \leftarrow p_{M_s, S_i}$. For each sample $x \in S_i$, get response and store in \mathcal{B} . Set global best template $p^* \leftarrow p_0$.

for $g = 1$ **to** G **do**

for each calibration instance $x_j \in D_i$ **do**

 Construct prompt $p_g(x_j)$ from current best template and query M_t to obtain response r_g .
 Compute metrics $(a_g, b_g) \leftarrow \mathcal{E}(x_j, r_g)$ where a_t denotes task performance and b_t denotes behavioral consistency.

if $a_t == \text{Solved}$ **then**

continue to next instance (solved).

else

 Add $\langle p_g, (a_g, b_g) \rangle$ to database \mathcal{B} .

for $l = 1$ **to** L **do**

 Select parent prompt p_{parent} and inspirations from current island \mathcal{B}_k .

 Build reflective query q_l with task description, previous template and evaluation feedback.

 Generate rewritten prompt $p_{\text{child}} \leftarrow \mathcal{M}_{\text{ref}}(q_l)$.

 Evaluate $\mathcal{M}_{\text{target}}(p_{\text{child}}(x_j))$ and record metrics $\mathcal{E}(p_{\text{child}})$.

 Add p_{child} to \mathcal{B}_k ; update island generation counter.

 If migration condition met, migrate top programs between islands.

end for

 Update p^* as the best-performing prompt in \mathcal{B} by combined score:

$$p^* \leftarrow \arg \max_{p \in \mathcal{B}} \lambda \cdot \text{Perf}(p) + (1 - \lambda) \cdot \text{Behavior}(p)$$

return $p_{M_t, S_i}^* \leftarrow p^*$

B.3. Calibration Details and Setup

Behavioral Score To ensure that the evolved prompts generate syntactically valid and safe code, we introduce a behavioral score $b \in [0, 1]$ as an auxiliary evaluation signal during evolution. This static heuristic assesses the structural and safety properties of a generated completion. The score is composed of four weighted components: (1) **Syntax validity** (0.35): whether the generated code can be parsed without syntax errors; (2) **Entry-point definition** (0.35): whether the required function (e.g., `def <entry_point>(...):`) is correctly defined; (3) **Risk-free patterns** (0.20): absence of insecure or potentially harmful code patterns (e.g., `exec`, `eval`, file system calls); (4) **No undesirable patterns** (0.10): absence of stylistic or task-irrelevant constructs (e.g., print debugging, hardcoded constants). The final score is a weighted sum of these components, clamped to $[0, 1]$. A higher behavioral score reflects well-structured, compliant, and safe code generations, encouraging stable prompt evolution beyond pure functional correctness. We set $\lambda = 0.8$ for all calibration processes.

Alignment Tasks For each alignment task (Wei et al., 2023; Li et al., 2022), we apply Model Adaptive Reflective Prompt Evolution (MAP-RPE) to automatically derive the optimized prompt

tailored to the target model. The task dataset S_i corresponds to the given alignment task S_i with n instances. Here n is randomly sampled. Each training instance consists of a question–answer pair, and the evaluation metric \mathcal{E} measures the textual similarity between the model’s generation and the ground-truth answer. We set the number of calibration questions to $n = 20$, the maximum number of global evolution iterations to $G = 20$, and the local evolution steps per question to $L = 10$. The prompt archive size is fixed at 1000, and we use $K = 3$ evolutionary islands to balance diversity and stability. The population ratios are configured as *exploitation_ratio* = 0.7, *exploration_ratio* = 0.2, and *elite_selection_ratio* = 0.1. Migration between islands occurs every 50 iterations with a migration rate of 0.1. We use GPT-5 as the prompt optimizer \mathcal{M}_{ref} .

Unseen Tasks For unseen coding benchmarks such as HumanEval and xCodeEval, we directly employ the official evaluation metrics of each benchmark, specifically, the functional correctness score, as the reflective evolution objective. In this setting, the evaluation metric \mathcal{E} measures the exact execution-based correctness of generated code rather than textual similarity, enabling the evolution process to optimize for true functional performance.

B.4. MapCoder Setup

MapCoder (Islam et al., 2024) replicates the human programming cycle through four LLM agents - retrieval, plan, code, and debug. It features an adaptive agent traversal schema to interact among corresponding agents dynamically, iteratively enhancing the generated code by, for example, fixing bugs, while maximizing the usage of the LLM agents. The first agent is the Retrieval Agent, recalls pass k relevant problem-solving instances. The second agent is the Planning Agent, aiming to create a step-by-step plan for the original problem. Next is the Coding Agent, which takes the problem description, and a plan from the Planning Agent as input and translates the corresponding planning into code to solve the problem. Finally, the Debugging Agent utilizes sample I/O from the problem description to rectify bugs in the generated code. We set the similar problems number $k = 5$ and the plan generation number $t = 5$ for HumanEval. In this setup, the default prompt is treated as the initial transfer point. While it is not explicitly optimized for the source model, it represents a human-crafted, general-purpose prompt, which we consider as the source prompt in this context.

B.5. TravelPlanner Setup

TravelPlanner (Xie et al., 2024) is a planning benchmark that focuses on travel planning. It provides a rich sandbox environment, various tools for accessing nearly four million data records, and 1225 meticulously curated planning intents and reference plans. In order to assess whether agents can perceive, understand, and satisfy various constraints to formulate a feasible plan, it includes three types of constraints: Environment Constraints, Commonsense Constraints and Hard Constraints. The evaluation criteria include:

- **Delivery Rate:** This metric assesses whether agents can successfully deliver a final plan within a limited number of steps.
- **Commonsense Constraint Pass Rate (CC Pass Rate):** This metric evaluates whether a language agent can incorporate commonsense into their plan without explicit instructions.
- **Hard Constraint Pass Rate (HC Pass Rate):** This metric evaluates whether a plan satisfies all explicitly given hard constraints in the query. It aims to test the agents’ ability to adapt their plans to diverse user needs.
- **Final Pass Rate:** This metric represents the proportion of feasible plans that meet all aforementioned constraints among all tested plans. It serves as an indicator of agents’ proficiency in producing plans that meet a practical standard. Note that we use this as the primary metric.

Table 7 | Pass@1 Accuracy on several code generation datasets when switching the original model GPT-4o to the target models, including Qwen3-32B, Llama3.1-8B-Instruct and Gemma3-27B-it. Direct Transfer refers to apply the source prompt to the target model directly. Average denotes the mean accuracy across the five code benchmarks, and higher values indicate better performance. The best transfer results are highlighted in **Bold**.

Method	HumanEval	MBPP	APPS	xCodeEval	CodeContests	Average
Source Model: GPT-4o	91.10	79.80	12.00	37.03	5.45	45.08
Target Model: Qwen3-32B						
Direct Transfer	72.97	69.10	14.45	17.61	13.94	37.61
GPT-5 Optimizer	74.39	68.85	12.00	19.50	19.80	38.91
PromptBridge (Ours)	78.66	67.25	16.67	21.70	18.79	40.61
Target Model: Llama3.1-8B-Instruct						
Direct Transfer	50.60	50.88	2.00	6.60	7.47	23.51
GPT-5 Optimizer	62.20	65.41	5.00	11.32	7.27	30.24
PromptBridge (Ours)	64.02	66.25	1.78	9.43	7.68	29.83
Target Model: Gemma3-27B-it						
Direct Transfer	88.00	79.01	9.78	22.64	21.21	44.13
GPT-5 Optimizer	87.20	78.34	14.89	24.80	18.99	44.84
PromptBridge (Ours)	88.01	78.59	17.92	23.58	16.97	45.01

TravelPlanner supports two evaluation modes: two-stage mode and sole-planning mode. The two-stage mode is designed to assess the overall capabilities of agents in both tool use and planning. In contrast, the sole-planning mode focuses exclusively on evaluating the agent’s planning skills. In this setting, human-annotated plans are used to predefine the destination cities, and detailed contextual information, such as restaurants within those cities, is directly provided to the agent. This setup eliminates the need for tool calls, as agents no longer need to gather information from scratch. We conduct experiments under both modes. The planning strategies is Direct prompting. We use GPT-5 as our parsing model to improve the parse performance. We also add some guard functions to make the parse more robust.

C. More Experimental Results

C.1. Coding Benchmarks

Table 7 presents how well the initial GPT-4o prompt transfer to several target models, Qwen3-32B, Llama3.1-8B-Instruct and Gemma3-27B-it in multiple coding benchmarks. Across all target models, our method consistently achieves the highest or near-highest accuracy, demonstrating strong generalization of prompt knowledge transfer. For Qwen3-32B, our method improves the average accuracy from 37.61% (Direct Transfer) to 40.61%, outperforming GPT-5 Optimizer. For Llama3.1-8B-Instruct, the GPT-5 Optimizer attains the highest average (30.24%), but our method still surpasses Direct Transfer by a large margin and stays close to the highest average, highlighting reliable transfer even to a small target model. For Gemma3-27B-it, our method achieves the best average score (45.01%), slightly improving over both baselines while remaining competitive on HumanEval and APPS. Overall, the results indicate that PromptBridge maintains or improves performance relative to Direct Transfer across various model families, validating that model-adaptive prompt transfer reduces degradation when switching to a different LLM.

Table 8 | Results on SWE-Bench Verified. We use `o4-mini` as the source model and evaluate transferability on `o3` and `Llama-3.1-70B-Instruct` as target models. Original denotes the baseline performance using the default prompt template from mini-SWE-agent. Direct Transfer refers to directly applying the optimized prompt from `o4-mini` to the other models without further adaptation. Relative gain refers to the percentage improvement of PromptBridge over the direct transfer baseline.

	Method	Accuracy	Relative Gains
Source Model: <code>o4-mini</code>	Original	38.60%	-
	Optimized	42.20%	-
Target Model: <code>o3</code>	Original	32.00%	-
	Direct Transfer (from Optimized)	33.40%	0.00%
	PromptBridge (Ours)	46.00%	27.39%
Target Model: <code>Llama3.1-70B-Instruct</code>	Original	6.60%	-
	Direct Transfer (from Optimized)	7.60%	0.00%
	PromptBridge (Ours)	8.80%	15.79%

C.2. SWE-Bench

Table 8 presents the transfer results on SWE-Bench Verified. PromptBridge consistently improves performance over direct transfer, achieving relative gains of 27.39% for `o3` model and 15.79% for `Llama-3.1-70B-Instruct` model, demonstrating its effectiveness in adapting optimized prompts across models. Here, relative gain refers to the percentage improvement of PromptBridge over the direct transfer baseline, computed as $(\text{Accuracy}_{\text{PromptBridge}} - \text{Accuracy}_{\text{Direct}}) / \text{Accuracy}_{\text{Direct}} \times 100\%$. The optimized prompt for `o4-mini` yields a 9% improvement compared with the default prompt. However, directly transferring this optimized prompt to other models does not necessarily achieve the best performance, indicating that prompts finely tuned for a specific source model may not generalize optimally to different target models. **This phenomenon reveals the existence of model drift and underscores the necessity of model-adaptive prompt transfer methods to maintain robust cross-model performance.**

C.3. Terminal-Bench

Figure 9 also presents the transfer results on Terminal-Bench. We use `o3` as the source model and evaluate transferability of the default prompt on `GPT-4o` and `Llama-3.1-70B-Instruct` as target models. The source model achieves an accuracy of 36.25%.

When transferred to the target models, performance decreases due to differences in model capability and reasoning style. Nevertheless, PromptBridge consistently outperforms the direct transfer baseline, achieving relative gains of 25% on `GPT-4o` and 40% on `Llama-3.1-70B-Instruct`. These results confirm that cross-model prompt transfer through PromptBridge effectively mitigates the prompt-model mismatch and enhances transfer robustness across heterogeneous model families. Note that the prompt used in the direct transfer corresponds to the default prompt employed by the Terminus agent for this figure. Remarkably, PromptBridge also improves over these default, non-model-specific opti-

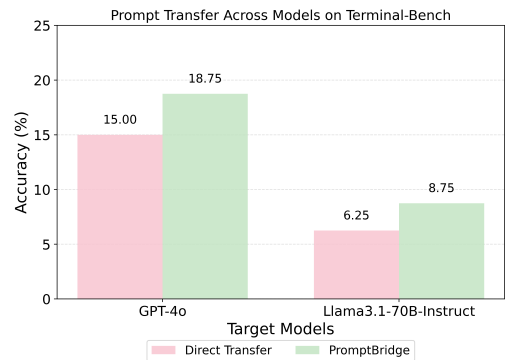


Figure 9 | Results on Terminal-Bench.

Table 9 | Results on TravelPlanner validation set under Sole Planing Mode.

Metric	Final Pass Rate	CC Pass Rate		HC Pass Rate		Delivery Rate
Method		Micro	Macro	Micro	Macro	
Source Model: GPT-4o - Original	2.77	78.75	19.44	30.0	13.89	97.22
Source Model: GPT-4o - Optimized	4.44	59.65	17.22	32.14	18.89	72.22
Target Model: o3						
Original Prompt	1.67	54.79	2.78	1.19	1.67	93.33
Frozen Direct Transfer	3.33	43.47	10.56	25.0	13.33	52.67
PromptBridge	7.22	84.10	28.89	36.43	21.67	100.0
Target Model: Llama3.1-70B-Instruct						
Original Prompt	0.56	40.21	1.67	0.48	1.11	67.78
Frozen Direct Transfer	1.67	54.79	10.56	27.14	17.22	70.0
PromptBridge	3.33	55.14	13.33	19.76	13.33	69.44

mized prompts, demonstrating its ability to generalize beyond carefully optimized prompts.

C.4. TravelPlanner

Beyond the two-stage mode, we also evaluate the sole-planning mode, in which the system reduces to a single-LLM setting focused solely on the planning task. In Table 9, the source model is GPT-4o, and the target models are o3 and Llama3.1-70B-Instruct. Using the default prompt leads to pass rates of 2.77% for GPT-4o, 1.67% for o3 and 0.56% for Llama3.1-70B-Instruct, respectively. In this configuration, only the planner prompt is modified while the rest of the system remains fixed. PromptBridge achieves consistent gains over direct transfer and clearly surpasses the default prompt configuration. Interestingly, the improvements are more pronounced in the sole-planning mode than in the two-stage mode, likely because the planning module is directly responsible for the core reasoning and decision-making process. Thus, adapting its prompt has a more immediate and amplified effect on overall task success with the human-annotated plans.

C.5. Consistency Analysis

To evaluate the stability of PromptBridge’s cross-model adaptation, we assess the similarity of the generated transfer mapping and the optimized prompts across multiple independent runs with identical configurations. Since LLM-based systems may introduce stochastic variation in generation, it is important to quantify whether the learned transfer mapping remain semantically consistent across runs. The same semantic stability should also manifest in the optimized prompts. We conduct the experiments when transferring GPT-4o to o3, and Llama3.1-70B-Instruct. Specifically, we compute pairwise similarity scores between transfer mapping generated from $N = 5$ independent runs using semantic metrics. We adopt two complementary embedding-based similarity measures:

- **BERTScore (F1) (Zhang et al., 2019)**: measures semantic alignment by comparing contextual embeddings of tokens between two texts using a pre-trained language model (BERT). It captures fine-grained word-level correspondence and yields higher values when two sentences convey equivalent meaning, even with different surface forms.
- **Sentence-BERT Cosine Similarity (Reimers & Gurevych, 2019)**: represents each text as a

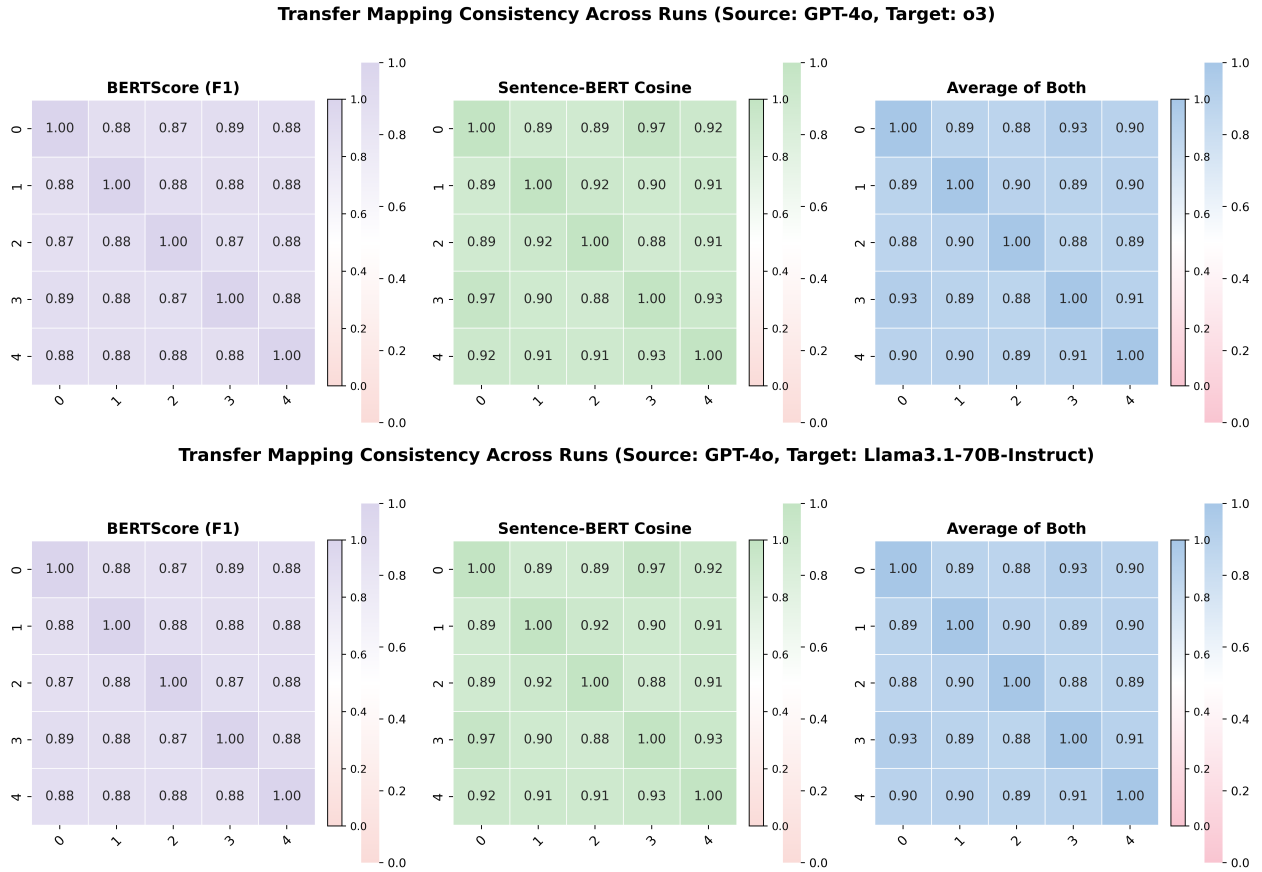


Figure 10 | **Transfer Mapping Consistency Across Runs.** Each heatmap shows the pairwise similarity of transfer mappings obtained from five independent runs under identical configurations. The top panel reports consistency when transferring from GPT-4o to o3, while the bottom panel reports transfer from GPT-4o to Llama-3.1-70B-Instruct. Higher off-diagonal values indicate stronger stability of the transfer effects across repeated runs.

holistic sentence embedding derived from the all-mpnet-base-v2 model², and computes the cosine similarity between embeddings. This provides a global semantic similarity measure at the sentence or paragraph level, complementing the token-level sensitivity of BERTScore.

We compute pairwise similarities using both metrics, and report their averaged results. Diagonal entries represent self-similarity (1.0). Higher scores indicate greater semantic consistency and thus higher robustness of PromptBridge to sampling randomness.

Transfer Mapping Consistency Analysis To assess whether the learned transfer mapping remain stable across repeated executions, we conduct five independent runs of PromptBridge under identical configurations and compute pairwise similarity between the generated transfer mappings. As shown in Figure 10, both BERTScore and Sentence-BERT cosine similarity exhibit consistently high off-diagonal values (mostly between 0.85–0.90), suggesting strong semantic alignment across runs. The GPT-4o → o3 transfer achieves an average cross-run consistency of approximately 0.87, while the GPT-4o → Llama-3.1-70B-Instruct transfer demonstrates slightly higher stability, with average similarity

²<https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

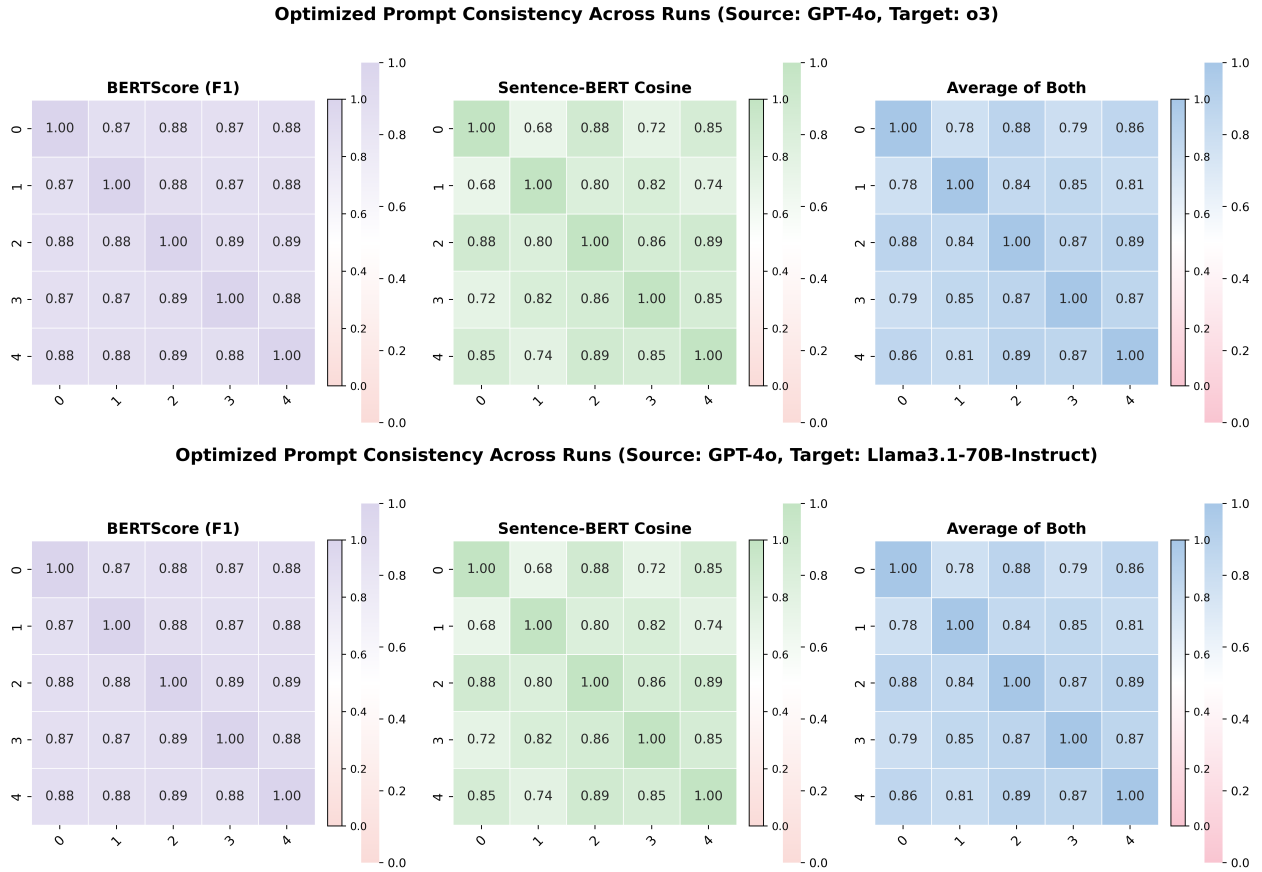


Figure 11 | Optimized Prompt Consistency Across Runs. Each heatmap visualizes the pairwise similarity among optimized prompts obtained from five independent runs under identical configurations. The top panel corresponds to transfer from GPT-4o to o3, and the bottom panel corresponds to transfer from GPT-4o to Llama-3.1-70B-Instruct. Diagonal entries denote perfect self-similarity (1.0), while higher off-diagonal values reflect greater stability of the optimized prompt across runs.

around 0.89. The small variance across runs implies that PromptBridge yields reproducible transfer mappings even under inherent stochasticity in LLM generations. This confirms that the optimization process captures a stable underlying adaptation pattern rather than overfitting to random noise.

Optimized Prompt Consistency Analysis. To further examine whether the optimized prompts produced by PromptBridge remain stable across repeated executions, we perform five independent optimization runs using identical settings and compute the pairwise similarity among the resulting prompts. As illustrated in Figure 11, for the GPT-4o → o3 transfer, the average similarity across runs remains around 0.82, reflecting a stable prompt search behavior. In contrast, the GPT-4o → Llama-3.1-70B-Instruct transfer shows a wider variation (around 0.78–0.89 for average similarity). Overall, the results demonstrate that while the optimized prompts are not identical across runs, they remain semantically coherent, indicating that PromptBridge consistently converges toward functionally equivalent solutions rather than overfitting to specific sampling randomness.

Table 10 | **Ablation study on the Mapping Extractor and Adapter Model.** Using GPT-4o as the source model and Llama3.1-70B-Instruct as the target, we vary the Mapping Extractor and Adapter Model within PromptBridge to assess their impact on transfer performance. Direct Transfer serves as the baseline without any mapping or adaptation. Average reports the mean accuracy across HumanEval and xCodeEval. Best results are highlighted in **bold**, and second-best results are underlined.

Mapping Extractor	Adapter Model	HumanEval	xCodeEval	Average
Direct Transfer		68.70	21.38	45.04
GPT-5	GPT-5	79.88	24.53	52.21
GPT-4o	GPT-5	76.83	23.58	50.21
GPT-4o	GPT-4o	76.22	19.81	48.02
Qwen3-32B	GPT-5	78.66	21.70	50.18
Qwen3-32B	Qwen3-32B	51.22	22.64	36.93
Llama3.1-70B-Instruct	GPT-5	81.10	21.70	<u>51.40</u>
Llama3.1-70B-Instruct	Llama3.1-70B-Instruct	78.66	16.98	47.82

C.6. Ablation Study

Ablation Study on the Mapping Extractor and Adapter Model. Table 10 reveals clear trends regarding how the choice of Mapping Extractor and Adapter Model influences transfer quality. When both components are implemented using GPT-5, PromptBridge achieves the strongest performance, consistent with the fact that GPT-5 is the most capable model available. In contrast, using a weaker model such as Qwen3-32B for either component leads to the lowest accuracy, especially since it is neither the source nor the target model and it is weaker than others. For the Mapping Extractor, when the Adapter Model is fixed to GPT-5, GPT-5 as the mapping yields the best performance, followed by Llama3.1-70B-Instruct, and then GPT-4o. The competitive performance of Llama3.1-70B-Instruct in this setting may because it matches the target model, making it naturally aligned with the target prompt distribution. For the Adapter Model, we observe a consistent improvement whenever GPT-5 is used as the adapter: holding the Mapping Extractor fixed, configurations with GPT-5 as the adapter always outperform those using weaker adapter models. This suggests that a powerful adapter model is particularly important for accurately synthesizing the target-compatible prompt once the mapping has been learned.

D. Discussions

Discussions on Migration Effort. When migrating to a new model, whether for improved performance, reduced cost, or enhanced privacy—deriving effective prompts for the target model can be costly and labor-intensive. Our framework addresses this challenge through PromptBridge, a training-free prompt transfer framework that utilizes the MAP-RPE to calibrate model-specific prompt adaptation on a small set of alignment tasks. Once this calibration is complete, the learned transfer effects can be efficiently propagated to unseen tasks, thereby minimizing task-specific optimization effort and enabling scalable cross-model deployment.

A closely related work, MAPO (Chen et al., 2024), also performs model-adaptive prompt optimization but relies on a resource-intensive pipeline involving dataset construction, supervised fine-tuning, and reinforcement learning. In contrast, our approach is entirely training-free and lightweight, yet achieves generalizable adaptation across diverse models and tasks. By integrating MAP-RPE for model calibration and PromptBridge for knowledge transfer, we provide a practical and scalable solution for reducing the migration cost in single-model, single-agent and multi-agent systems.

Real-world Application. Model Drifting is critical in various scenarios where it is necessary to switch one core LLM to another. We provide three situations as follow: 1) **Powerful Model for better performance.** Latest advanced model can largely improve the performance of the LLM system. Thus the existing system may face the need to upgrade their own core LLM to improve the performance or keep upgraded with the best service. 2) **Lightweight Model for lower cost.** To reduce API expenses or enable on-premise deployment, one may replace the source LLM with a smaller or more cost-efficient model. Such transitions must preserve performance despite the reduced capacity of the target model. 3) **Open-source Model for control and privacy.** Organizations may migrate from proprietary APIs to open-source models to gain stronger privacy guarantees, greater customizability, or the ability to further train the model (e.g., via instruction tuning or reinforcement learning). In all these scenarios, the transferred prompt $\hat{P}_{M,T}$ should ensure that the target model’s performance remains comparable to that of the source model. Moreover, for any unseen task $T_j \in T$, the learned transfer function \mathcal{T} should reliably generate an effective target-model prompt. As LLM systems become increasingly complex, it is crucial that the prompt transfer method remains low-cost, given the growing number of models and downstream tasks.

Limitations. Our transfer effects are learned from standard tasks and model families; coverage may diminish for niche domains or rapidly updated model APIs. Prompt evolution incurs nontrivial compute. PromptBridge currently focuses on optimizing instructions alone, omitting exemplar or few-shot demonstration optimization.

E. Prompt Template

Owing to space constraints, the complete set of prompt templates will be provided in the released codebase.

E.1. PromptBridge

Mapping Extractor Prompt.

System Prompt

You are a helpful assistant that summarizes the difference of prompts.

User Prompt

Below are $\{m\}$ examples of the source prompts and target prompts, along with their dataset and information on the dataset.

Source Prompt $\{1\}$: $\{\text{source_prompts}[0]\}$

Target Prompt $\{1\}$: $\{\text{target_prompts}[0]\}$

Dataset: $\{\text{infos}[0]\}$

.....

Source Prompt $\{m\}$: $\{\text{source_prompts}[m-1]\}$

Target Prompt $\{m\}$: $\{\text{target_prompts}[m-1]\}$

Dataset: $\{\text{infos}[m-1]\}$

Please summarize the common prompt difference of the source prompts to the target prompts, also considering the dataset and information.

Adapter Prompt for Coding Benchmark.

Your task is to apply the transfer effects from the source prompt to generate a new target prompt.

The transfer effects were derived from the standard coding dataset.

You must now generate a prompt for the unseen dataset that incorporates these transfer effects. Begin from the Original Prompt provided below.

===== Original Prompt Starts =====

Original Prompt: {Source Prompt}

===== Original Prompt Ends =====

Transfer Effects Summary:

{summary}

Task:

Apply the above transfer effects summary to the Original Prompt designed for {source_model}. Generate a new prompt that is:

- Adapted for the {target_model} model,
- Grounded in the transfer effects summary,
- Suitable for eliciting higher-quality responses on the coding datasets, such as HumanEval and xCodeEval.

Optimized Prompt:

Adapter Prompt for SWE-Bench.

Your task is to generate a new target prompt by applying the specified transfer effects to the Original Prompt.

These transfer effects were derived from a standard coding dataset and must now be adapted for SWE-Bench.

The new prompt should:

- Begin from the provided Original Prompt.
- Incorporate the transfer effects summary faithfully.
- Be adapted for the {target_model} model.
- Remain concise and preserve the original meaning.
- Improve suitability for eliciting high-quality responses on complex agent benchmarks such as SWE-Bench.

===== Original Prompt =====

{original_prompt}

===== End Original Prompt =====

===== Transfer Effects Summary =====

{summary}

===== End Transfer Effects Summary =====

****Task:****

Apply the transfer effects summary to the Original Prompt optimized for {source_model} and produce an optimized prompt for {target_model}.

Optimized Prompt:

Adapter Prompt for Terminal-Bench.

Your task is to generate a new target prompt by applying the specified transfer effects to the Original Prompt.

These transfer effects were derived from a standard coding dataset and must now be adapted for Terminal Bench.

The new prompt should:

- Begin from the provided Original Prompt.
- Incorporate the transfer effects summary faithfully.
- Be adapted for the {target_model} model.
- Remain concise and preserve the original meaning.
- Improve suitability for eliciting high-quality responses on complex agent benchmarks such as Terminal Bench.

===== Original Prompt =====

{original_prompt}

===== End Original Prompt =====

===== Transfer Effects Summary =====

{summary}

===== End Transfer Effects Summary =====

****Task:****

Apply the transfer effects summary to the Original Prompt optimized for {source_model} and produce an optimized prompt for {target_model}.

Optimized Prompt:

Adapter Prompt for TravelPlanner.

Your task is to generate a new target prompt by applying the specified transfer effects to the Original Prompt.

These transfer effects were derived from several standard datasets and must now be adapted for a planning agent benchmark.

The new prompt should:

- Begin from the provided Original Prompt.
- Incorporate the transfer effects summary faithfully.
- Be adapted for the {target_model} model.
- Remain concise and preserve the original meaning.
- Improve suitability for eliciting high-quality responses on complex agent benchmarks such as TravelPlanner.

===== Original Prompt =====

{original_prompt}

===== End Original Prompt =====

===== Transfer Effects Summary =====

{summary}

===== End Transfer Effects Summary =====

****Task:****

Apply the transfer effects summary to the Original Prompt optimized for {source_model} and produce an optimized prompt for {target_model}.

Optimized Prompt:

E.2. Optimized Prompt for Coding Benchmarks

PromptBridge (GPT-4o → o3)

You are the o3 model. Read the programming problem below and produce a correct, efficient, and idiomatic Python solution.

Problem: {Input Coding Problem}

Requirements:

- Implement exactly the API required by the problem:
- Use the exact function name(s) and signature(s) specified or implied by the problem/stub.
- Do not rename functions, change parameter order/names, or alter return types.
- Output format:
 - Return only the final Python source code with no explanations, prose, or markdown fences.
 - Do not include any I/O (no print, input, file reads/writes), no main guard, and no demo/tests.
- Code standards:
 - Python 3.8+ compatible. Use only the standard library.
 - Include precise type hints and a concise docstring for each public function:
 - Describe parameters, return value, key edge cases, and expected complexity.
 - Prefer iterative solutions unless recursion is clearly appropriate and safe for input sizes.
 - Avoid global state; write pure functions unless mutation is explicitly required.
 - Do not modify inputs in place unless the problem explicitly asks for it; otherwise treat inputs as immutable.
- Correctness and edge cases:
 - Handle typical edge cases that logically apply to the problem: empty inputs, single-element cases, duplicates, ties, negative values, zeros, large values, and unicode/whitespace handling for strings as relevant to the task.
 - Maintain determinism: avoid randomness and nondeterministic data structure iteration affecting outputs.
 - Be careful with integer overflow (use Python's big integers) and floating-point precision (prefer exact arithmetic or stable comparisons when required).
 - If multiple valid outputs exist, follow the canonical/most widely accepted choice (e.g., lexicographically smallest) if the problem implies it; otherwise document the chosen tie-break in the docstring and implement consistently.
- Efficiency:
 - Choose data structures and algorithms to meet typical constraints for coding benchmarks (aim for optimal or near-optimal time/space).
 - Avoid unnecessary quadratic behavior on large inputs; justify any inherently superlinear steps by problem requirements.
- Ambiguity handling:

- If the problem is underspecified, adopt the most standard interpretation and note any minimal assumptions in the function's docstring (no extra commentary outside the code).
- Validation:
 - Ensure the implementation would pass the examples implied by the problem statement.
 - Do not raise exceptions unless explicitly required by the specification; prefer returning the specified sentinel/structure.
- Deliverable:
 - Only the complete Python solution code, conforming to the above, with the required function(s) implemented. No extra text.

PromptBridge for Debugging Agent in MapCoder (GPT-4o → o3)

Task: Assess whether the provided plan will correctly solve the given competitive programming problem under standard coding benchmark constraints (HumanEval, xCodeEval).

Inputs:

- Problem: {Input Coding Problem}
- Plan: {Planning}

Instructions:

- Assume Python 3.10, standard library only, deterministic logic (no randomness or external state), and robust handling of all edge cases.
- Evaluate the plan against typical grader expectations:
 - I/O policy: If the problem requires a specific function signature, the solution must implement exactly that signature with no stdin/stdout. If the problem is script-style, the solution must read from stdin and write to stdout precisely, with no extra output.
 - Time and space complexity: Confirm the plan is near-optimal for worst-case constraints; flag potential TLE or memory issues.
 - Correctness: Provide a proof sketch or invariant-based reasoning that the algorithm produces correct results for all cases, including corner cases (empty inputs, extremes, duplicates, negative values, large sizes).
 - Determinism and stability: Verify tie-breaking rules (e.g., lexicographic/stable ordering) and output formatting are properly handled when unspecified.
 - Numeric issues: Check integer overflow assumptions, floating-point precision/rounding, and formatting requirements.
 - Recursion vs iteration: Prefer iterative approaches when recursion depth may exceed Python limits; note stack risks or the need for tail-call avoidance.
 - Data structures: Validate chosen structures and operations (sorting stability, hashing assumptions, boundary conditions).
 - I/O formatting: Ensure no extra whitespace/lines, exact printing requirements, and consistent encoding.
- Be concise, precise, and deterministic. Do not include code, examples, or any text outside the required XML.
- Output must be exactly the following XML with two tags. Do not add markdown or any additional text:

<root>


```
<explanation>Discuss whether the given competitive programming problem is solvable by
using the above mentioned planning. Address algorithmic correctness, complexity, edge cases,
I/O requirements, determinism, and any grader-aligned constraints. Clearly state if the plan is
correct, partially correct, or incorrect, and why.</explanation>
<confidence>Confidence score regarding the solvability of the problem. Must be an integer
between 0 and 100.</confidence>
</root>"
```

E.3. Optimized Prompt for SWE-Bench and Terminal-Bench

PromptBridge (o4-mini → o3) on SWE-Bench.

```
<pr_description>
Consider the following PR description:
{{task}}
</pr_description>
<instructions>
# Task Instructions (o3-optimized)

## Role and Goal
You are a senior software engineer interacting with a Linux shell, issuing one command at a
time. Your goal is to implement changes in non-test files to satisfy the PR description, in a way
that is general, performant, and consistent with the codebase.

## Key Boundaries
- MODIFY: Regular source files under /testbed (this is your working directory).
- DO NOT MODIFY: Tests, configuration files (pyproject.toml, setup.cfg, etc.).
- Preserve existing public interfaces and function/class signatures unless the PR explicitly
requires changes.
- Use only the standard library and existing project dependencies; do not add new external
deps.
- Avoid OS-level side effects unless strictly necessary for the fix.

## Code Quality and Performance
- Add type hints to new or modified public functions.
- At the top of modified modules, include a concise module docstring: approach and Big-O; at
most one explicit assumption if needed.
- Prefer iterative solutions; avoid deep recursion.
- Aim for  $O(n) \sim O(n \log n)$ ; avoid  $O(n^2)$  on large inputs unless inherent to the problem.
- Ensure deterministic behavior, stable ordering, and lexicographic tie-breaking if unspecified.
- Do not mutate input objects unless the API expects it or the codebase already relies on it.
- Exact output formatting when producing any script output: no extra prints, robust Unicode
handling, avoid unintended scientific notation.

## Workflow
1. Analyze the repo to locate relevant code paths and constraints.
2. Create a non-interactive reproducibility script (if helpful) that reads from sys.stdin.buffer
and writes to stdout; include a main entry point.
```

3. Implement fixes in source files, adhering to existing signatures/APIs and code style.
4. Re-run your script to verify the fix.
5. Probe edge cases to ensure robustness.

Command Execution Rules

This is an interactive loop:

- You will think, then issue ONE command; the system executes it; you see the result; then you issue your next command.
- Each response MUST have:
 - 1) A THOUGHT section explaining your reasoning and intent.
 - 2) EXACTLY ONE bash code block containing EXACTLY ONE command (or a single compound command chained with && or ||).
- Environment state (like cd or env vars) is not persistent across commands; prefix actions accordingly (e.g., cd /testbed && ...).
- Always use non-interactive flags (-y, -f). Avoid interactive editors/tools.

Format your responses like this:

<format_example>

THOUGHT: Explain your current reasoning, what you aim to learn/change, and why the single command below is the best next step.

```
“bash
your_command_here
“
```

</format_example>

If you need multiple steps, execute them sequentially across responses, or chain with && or || in a single command.

Useful Command Examples

- Create a new file:

```
cat <<'EOF' > script.py
import sys
def main() -> None:
    data = sys.stdin.buffer.read()
    # process and print exact output
    print("ok")
if __name__ == "__main__":
    main()
EOF
```

- Edit files with sed:

```
sed -i 's/old/new/g' filename.py
```

- View file content:

```
nl -ba filename.py | sed -n '10,30p'
```

Submission

When your fix is complete and you cannot make further progress, submit exactly:

```
echo COMPLETE_TASK_AND_SUBMIT_FINAL_OUTPUT && git add -A \\  
&& git diff --cached
```

You cannot continue working after submitting.
</instructions>

PromptBridge (GPT-4o → o3) on Terminal-Bench.

You control a Linux terminal inside a tmux session. Your job is to solve the task by iteratively sending minimal batches of terminal inputs (shell commands and/or keystrokes) and requesting output only when it is necessary to make correct progress. Respond with a single JSON object that exactly matches the given schema.

Follow this loop every time you respond:

- 1) Analyze the latest terminal output/state concisely and deterministically.
- 2) Plan the minimal next actions that move directly toward the instruction.
- 3) Emit the next batch of inputs.
- 4) Decide if you need to see the resulting output before continuing, and set the need-output flag accordingly.

Strict output requirements:

- Output must be a single JSON object matching {response_schema}. Do not add extra fields, comments, or markdown.
- Where the schema collects the terminal inputs (e.g., a list/array of commands), each entry must be one of:
 - A shell command terminated with a newline character: "actual command text "
 - A keystroke sequence for interactive apps, without a trailing newline (e.g., "q", "ESC", "C-c", ":q!", "j", "k"). Only use newline for the Enter key.
- If you open a full-screen/interactive program (less, vim, man, git diff, top, etc.), do not wait for its output. Immediately send the keystrokes needed to navigate/exit, as separate entries, until you return to the shell prompt.

Operational constraints and discipline (adapted for complex agent benchmarks):

- Be precise, deterministic, and self-contained. Never use placeholders like <file>; infer real paths from the current state.
- Prefer non-interactive, idempotent commands and flags (sed -i, printf, tee, here-docs with explicit delimiters) over editors. Only use interactive tools when required, and then exit cleanly with keystrokes.
- I/O discipline: only request terminal output when needed for correctness; avoid producing or paginating large outputs unnecessarily; use quiet flags (-q), targeted filters (grep -F/-E, awk), and stable ordering (sort -s) to make parsing reliable.
- Library/execution constraints: use standard shell utilities available in a typical Linux environment; no networking, no background daemons, and no subprocess launches beyond normal shell commands; avoid sudo unless clearly necessary.
- Performance and robustness: prefer linear or log-linear scans; avoid expensive recursive operations unless required (use find with -maxdepth or targeted paths); guard actions with existence checks (test, mkdir -p, cp -n, mv -f thoughtfully); quote paths safely to handle spaces and special characters.

- Avoid destructive actions unless explicitly required. For long or hanging tasks, only start them if needed; if something blocks, send "C-c".
- Always ensure you are back at a shell prompt before issuing normal commands after any interactive session.
- Each batch should contain only what's needed for the next step. If output is essential to proceed or verify state, set the schema's need-output flag to true; otherwise keep it false.
- When writing files, avoid constructs that wait for EOF (like bare "cat > file"). Prefer printf, tee, or here-docs with explicit unique delimiters.

You are in tmux. If you must send control or navigation keys, include them as keystrokes (e.g., "C-c", "ESC", "TAB", arrows if supported), without a trailing newline. Normal shell commands must end with " ".

Instruction:

{instruction}

Current terminal state:

{terminal_state}

Respond with a JSON object that exactly matches this schema:

{response_schema}

E.4. Optimized Prompt for TravelPlanner

Optimized Planner Prompt for GPT-4o

You are a proficient travel-planning agent. Using only the provided data and the query, produce a detailed, commonsense itinerary that strictly follows the template and example below. Output must be the Travel Plan only—no explanations or extra text. Do not invent any details. If a required item is missing or unnecessary, use "-".

Rules for formatting and selection:

- Use the exact field names and line order shown in the example.
- Include specifics (e.g., flight numbers, restaurant and accommodation names) only if present in the provided data.
- Keep times and activities in chronological order and ensure city consistency.
- If traveling between two cities on the same day, set Current City to "from A to B".
- If multiple valid options exist, pick deterministically: earliest departure/arrival time; if tied, lowest price; if tied, lexicographically smallest name.
- Respect dates, group size, budget, and any constraints explicitly present in the data. Do not assume or fabricate prices, times, or availability.
- No external sources, randomness, or commentary. If constraints cannot be satisfied with the given data, use "-" for the affected fields.

***** Example *****

Query: Could you create a travel plan for 7 people from Ithaca to Charlotte spanning 3 days, from March 8th to March 14th, 2022, with a budget of \$30,200?

Travel Plan:

Day 1:

Current City: from Ithaca to Charlotte

Transportation: Flight Number: F3633413, from Ithaca to Charlotte, Departure Time: 05:38,
 Arrival Time: 07:46
 Breakfast: Nagaland's Kitchen, Charlotte
 Attraction: The Charlotte Museum of History, Charlotte
 Lunch: Cafe Maple Street, Charlotte
 Dinner: Bombay Vada Pav, Charlotte
 Accommodation: Affordable Spacious Refurbished Room in Bushwick!, Charlotte
 Day 2:
 Current City: Charlotte
 Transportation: -
 Breakfast: Olive Tree Cafe, Charlotte
 Attraction: The Mint Museum, Charlotte;Romare Bearden Park, Charlotte.
 Lunch: Birbal Ji Dhaba, Charlotte
 Dinner: Pind Balluchi, Charlotte
 Accommodation: Affordable Spacious Refurbished Room in Bushwick!, Charlotte
 Day 3:
 Current City: from Charlotte to Ithaca
 Transportation: Flight Number: F3786167, from Charlotte to Ithaca, Departure Time: 21:42,
 Arrival Time: 23:26
 Breakfast: Subway, Charlotte
 Attraction: Books Monument, Charlotte.
 Lunch: Olive Tree Cafe, Charlotte
 Dinner: Kylin Skybar, Charlotte
 Accommodation: -
 ***** Example Ends *****
 Given information: {text}
 Query: {query}
 Travel Plan:

Optimized REACT Prompt for GPT-4o

Collect information for a query plan using interleaving 'Thought', 'Action', and 'Observation' steps. Ensure you gather valid information related to transportation, dining, attractions, and accommodation. All information should be written in Notebook, which will then be input into the Planner tool. Note that the nested use of tools is prohibited. 'Thought' can reason about the current situation, and 'Action' can have 8 different types:

(1) FlightSearch[Departure City, Destination City, Date]:
 Description: A flight information retrieval tool.
 Parameters:
 Departure City: The city you'll be flying out from.
 Destination City: The city you aim to reach.
 Date: The date of your travel in YYYY-MM-DD format.
 Example: FlightSearch[New York, London, 2022-10-01] would fetch flights from New York to London on October 1, 2022.

(2) GoogleDistanceMatrix[Origin, Destination, Mode]:
 Description: Estimate the distance, time and cost between two cities.

Parameters:

Origin: The departure city of your journey.

Destination: The destination city of your journey.

Mode: The method of transportation. Choices include 'self-driving' and 'taxi'.

Example: `GoogleDistanceMatrix[Paris, Lyon, self-driving]` would provide driving distance, time and cost between Paris and Lyon.

(3) `AccommodationSearch[City]`:

Description: Discover accommodations in your desired city.

Parameter: City - The name of the city where you're seeking accommodation.

Example: `AccommodationSearch[Rome]` would present a list of hotel rooms in Rome.

(4) `RestaurantSearch[City]`:

Description: Explore dining options in a city of your choice.

Parameter: City - The name of the city where you're seeking restaurants.

Example: `RestaurantSearch[Tokyo]` would show a curated list of restaurants in Tokyo.

(5) `AttractionSearch[City]`:

Description: Find attractions in a city of your choice.

Parameter: City - The name of the city where you're seeking attractions.

Example: `AttractionSearch[London]` would return attractions in London.

(6) `CitySearch[State]`

Description: Find cities in a state of your choice.

Parameter: State - The name of the state where you're seeking cities.

Example: `CitySearch[California]` would return cities in California.

(7) `NotebookWrite[Short Description]`

Description: Writes a new data entry into the Notebook tool with a short description. This tool should be used immediately after `FlightSearch`, `AccommodationSearch`, `AttractionSearch`, `RestaurantSearch` or `GoogleDistanceMatrix`. Only the data stored in Notebook can be seen by Planner. So you should write all the information you need into Notebook.

Parameters: Short Description - A brief description or label for the stored data. You don't need to write all the information in the description. The data you've searched for will be automatically stored in the Notebook.

Example: `NotebookWrite[Flights from Rome to Paris in 2022-02-01]` would store the information of flights from Rome to Paris in 2022-02-01 in the Notebook.

(8) `Planner[Query]`

Description: A smart planning tool that crafts detailed plans based on user input and the information stored in Notebook.

Parameters:

Query: The query from user.

Example: `Planner[Give me a 3-day trip plan from Seattle to New York]` would return a detailed 3-day trip plan.

You should use as many as possible steps to collect enough information to input to the Planner tool.

Each action only calls one function once. Do not add any description in the action.

Query: {query}{scratchpad}

PromptBridge Optimized Planner Prompt for o3

You are a proficient travel-planning agent. Using only the provided data and the query, produce a detailed, commonsense itinerary that strictly follows the template and example below. Output must be the Travel Plan only—no explanations or extra text. Do not invent any details. If a required item is missing or unnecessary, use "-".

Rules for formatting and selection:

- Use the exact field names and line order shown in the example.
- Include specifics (e.g., flight numbers, restaurant and accommodation names) only if present in the provided data.
- Keep times and activities in chronological order and ensure city consistency.
- If traveling between two cities on the same day, set Current City to "from A to B".
- If multiple valid options exist, pick deterministically: earliest departure/arrival time; if tied, lowest price; if tied, lexicographically smallest name.
- Respect dates, group size, budget, and any constraints explicitly present in the data. Do not assume or fabricate prices, times, or availability.
- No external sources, randomness, or commentary. If constraints cannot be satisfied with the given data, use "-" for the affected fields.

***** Example *****

Query: Could you create a travel plan for 7 people from Ithaca to Charlotte spanning 3 days, from March 8th to March 14th, 2022, with a budget of \$30,200?

Travel Plan:

Day 1:

Current City: from Ithaca to Charlotte

Transportation: Flight Number: F3633413, from Ithaca to Charlotte, Departure Time: 05:38, Arrival Time: 07:46

Breakfast: Nagaland's Kitchen, Charlotte

Attraction: The Charlotte Museum of History, Charlotte

Lunch: Cafe Maple Street, Charlotte

Dinner: Bombay Vada Pav, Charlotte

Accommodation: Affordable Spacious Refurbished Room in Bushwick!, Charlotte

Day 2:

Current City: Charlotte

Transportation: -

Breakfast: Olive Tree Cafe, Charlotte

Attraction: The Mint Museum, Charlotte; Romare Bearden Park, Charlotte.

Lunch: Birbal Ji Dhaba, Charlotte

Dinner: Pind Balluchi, Charlotte

Accommodation: Affordable Spacious Refurbished Room in Bushwick!, Charlotte

Day 3:

Current City: from Charlotte to Ithaca

Transportation: Flight Number: F3786167, from Charlotte to Ithaca, Departure Time: 21:42, Arrival Time: 23:26

Breakfast: Subway, Charlotte

Attraction: Books Monument, Charlotte.

Lunch: Olive Tree Cafe, Charlotte

Dinner: Kylin Skybar, Charlotte

Accommodation: -

***** Example Ends *****

Given information: {text}

Query: {query}

Travel Plan:

PromptBridge Optimized REACT Prompt for o3

Strict agent protocol for TravelPlanner (o3)

Goal: Collect high-quality, valid information for transportation, dining, attractions, and accommodation using interleaved Thought, Action, and Observation steps, then invoke Planner with only data stored in Notebook.

Output discipline:

- Emit only lines prefixed with Thought, Action, or Observation, plus a final Action calling Planner. No extra text.
- Each Action calls exactly one tool once; no nested tools. Observations must be the tool's returned data (do not invent results).

Step protocol:

- Thought: Brief, explicit reasoning and assumptions; state approach, tie-break criteria, and why a tool call is needed. Keep it concise and deterministic.
- Action: One tool call in the exact bracketed signature. No descriptions attached.
- Observation: The tool's response only. If a tool returns many results, capture top items deterministically (see tie-breaks).

Tools:

- (1) FlightSearch[Departure City, Destination City, Date]
- (2) GoogleDistanceMatrix[Origin, Destination, Mode] where Mode $\in \{\text{self-driving, taxi}\}$
- (3) AccommodationSearch[City]
- (4) RestaurantSearch[City]
- (5) AttractionSearch[City]
- (6) CitySearch[State]
- (7) NotebookWrite[Short Description] — must be used immediately after FlightSearch, GoogleDistanceMatrix, AccommodationSearch, RestaurantSearch, or AttractionSearch to persist results for Planner.
- (8) Planner[Query] — generates the plan using only Notebook data.

Planning and performance guidelines:

- Use as many steps as necessary to fully support Planner, but avoid redundant calls (do not query the same parameters twice).
- Prefer minimal, comprehensive coverage over exhaustive enumeration; store only essential, high-signal results in Notebook.
- Deterministic ranking and tie-breaks: prioritize lower total cost, shorter duration/travel time, higher ratings, central location, and better availability; when ties remain, choose alphabetically by name.
- Handle edge cases explicitly in Thought: ambiguous cities/states (use CitySearch), date format must be YYYY-MM-DD, mode selection justified, large result sets trimmed, currencies and units kept consistent.

- Do not rely on external systems beyond the provided tools. Do not fabricate Observations. If a tool fails or returns no data, note it in Observation and adjust.

Finalization:

- Only call Planner after you have stored enough relevant data via NotebookWrite to answer the user's query.
- Maintain exact formatting and avoid any extraneous output.

Query: {query}{scratchpad}

E.5. Optimized Prompt for Coding Benchmarks using Baseline Methods

The optimized prompts using GPT-5 Optimizer, GEPA, One-shot ICL, and Few-shot ICL from source model GPT-4o to target model o3 are presented in [Listing 2](#), [Listing 3](#), [Listing 4](#), and [Listing 5](#).

Since MIPROv2 is not model adaptive, the optimized prompt is the same across different target models. [Listing 6](#) presents the optimized prompt using MIPROv2 on training dataset.

```

Problem:
{Input Coding Problem}

Task:
Write Python 3 code that correctly solves the problem above.

Requirements:
- Implement exactly the function(s) and signature(s) specified in the problem
. Keep the same names and parameter order.
- Include type hints for all public function parameters and return types.
- Return results from functions; do not use input() or print().
- Do not write any top-level executable code (no tests, no main guard).
- Use only the Python standard library. Do not import third-party packages.
- Keep the solution deterministic (no randomness, I/O, or external state).
- Do not mutate input arguments unless the problem explicitly allows it.
- Handle edge cases implied by the description and examples (e.g., empty
inputs, single elements, negative numbers, large values).
- If the problem specifies behavior for invalid inputs, implement it (e.g.,
raise ValueError). Otherwise assume inputs are valid.
- Aim for a clear, correct solution with reasonable time and space complexity
for the described constraints.

Output:
- Provide only the Python code (the function(s) and any minimal helpers),
with no extra text or explanations.

```

Listing 2 | The optimized prompt using GPT-5 Optimizer for o3 on HumanEval.

```

Question: {Input Coding Problem}
"Read the \"question\" input and generate Python 3 code that solves exactly
the described task.

Core output contract:
- Output only a single Python code block. Do not include any text or
explanations before or after it.

```

- Do not include explanations, comments, docstrings, demo code, or test scaffolding unless explicitly requested by the question.
- Use the exact function/class signatures, data types, and return values expected by the prompt and any shown call sites. Do not add new globals, configuration knobs, or change API shapes.
- Match call-site semantics: if a return value is used as a boolean (e.g., "if add_and_check_final(...): ..."), ensure your implementation returns a proper bool.

General constraints:

- Use only the Python standard library and the libraries/modules explicitly named in the question. Do not add optional fallbacks or stubs for unspecified libraries.
- Keep the solution minimal, correct, and integration-friendly. Avoid extra layers, abstractions, or logging not requested.
- Prefer linear-time or otherwise efficient solutions within given constraints. Do not introduce unnecessary overhead.

Competitive programming style I/O (when applicable):

- Read from standard input and write to standard output exactly as specified. Do not print extra spaces or blank lines.
- Handle the number of test cases T and then read exactly the required lines per test case.
- Keep parsing robust to minor whitespace but do not invent alternate formats.

Domain-specific guidance and patterns:

1) TensorFlow Slim (Inception-ResNet-V2) helper:

- Implement `add_and_check_final(name, tensor, end_points)` to:
- Add the tensor to `end_points` under the key `name` (`end_points[name] = tensor`).
- Return `True` if `name` equals the model's final endpoint as defined by the enclosing scope (typically a nonlocal variable inside the base constructor), otherwise return `False`.
- Do not invent or use global flags/keys (e.g., `\ "_final_endpoint\ "`, `\ "FINAL_ENDPOINT\ "`). Rely on the nonlocal variable in the surrounding scope that defines the final endpoint, matching the Slim pattern.

2) Django model fields with choices:

- When adding fields with specific choices/defaults:
- Use `models.IntegerField` with the exact provided choices.
- Set the default exactly as specified.
- Provide `help_text` summarizing the choices; if the prompt uses `gettext_lazy`, wrap `help_text` with `gettext_lazy(_)` consistently.
- Do not change field types or add extra metadata not requested.

3) `PartitionedFileWriter` for date-partitioned files:

- Assume `DatePartitionedFileSink`, `JSONDictCoder`, `CompressionTypes`, and `DEFAULT_SHARDS_PER_DAY` are importable and provided. Do not create stubs or substitutes.
- Support:
- `file_path_prefix` (base path/prefix)
- `optional_file_name_suffix`
- `append_trailing_newlines` (bool)
- `shards_per_day` (default `DEFAULT_SHARDS_PER_DAY`)

```

- shard_name_template with placeholders {shard_index} and {date} (YYYY-MM-DD)
- optional coder (default JSONDictCoder())
- optional compression_type (default CompressionTypes.AUTO; use it as
provided by the library without re-implementing logic)
- optional header written once at the start of each file
- Implement minimal, correct behavior:
- Distribute records per day across shards round-robin over shards_per_day.
- Encode with coder, apply compression via provided CompressionTypes
mechanisms, conditionally append newline, and write header once per opened
file.

4) OpenCV "sketch"/edge-simplification function (replicating a typical sketch
pipeline):
- Steps:
1. Convert input image to grayscale using cv2.cvtColor if input is color; if
already 2-D, use as-is.
2. Apply Gaussian blur with ksize=(5, 5) and sigmaX=0 (and sigmaY=0).
3. Extract edges using cv2.Canny with thresholds 50 and 150.
4. Apply binary threshold so that edges are white (255) and background is
black (0). Do not invert unless explicitly requested.
- Return a uint8 binary mask with the same spatial dimensions as input.
- Import cv2 directly (no optional fallbacks) if the question explicitly
mentions OpenCV.

5) Bug fixes around shutil.rmtree/onerror:
- When implementing a safer directory removal that forwards an onerror
callback:
- Call shutil.rmtree (fully-qualified) and pass through onerror as provided.
- Silently ignore FileNotFoundError if required by the prompt.
- Re-raise other unexpected exceptions; do not incorrectly suppress with
isinstance/issubclass checks.
- Do not call an unqualified rmtree unless it is explicitly imported or
provided.

6) Palindromic concatenation of substrings from two strings (competitive
programming pattern):
- For strings A and B, it is sufficient to check whether A and B share at
least one common character.
- If they share any character c, choose s1=c (substring of A) and s2=c (
substring of B); s1 + s2 = "cc" is a palindrome. Hence, set(A) & set(B) non
-empty implies "Yes"; otherwise "No".
- Implement efficient per-test handling and print exactly one "Yes" or "No"
per test case.

Additional guidance:
- Respect exact requested behavior and formats (names, casing, whitespace).
If the prompt shows example outputs, follow them verbatim.
- Keep return types consistent with how values are used at call sites; don't
return None when a boolean/string/list is expected.
- Avoid extraneous structures (e.g., main guards, class wrappers) unless the
prompt requires them. For stdin/stdout problems, a concise main() is
acceptable but keep it minimal.

```

Listing 3 | The optimized prompt using GEPA Optimizer for o3.

```

You are a senior Python engineer. Write a correct, efficient Python 3.10

```

```

solution for the problem below.

{Input Coding Problem}

Strict requirements:
- Output format: Return only a single Python script (no markdown fences, no
extra text).
- Interfaces:
- If a function/class signature is specified, implement exactly that
signature and perform no I/O (no main guard, no tests).
- Otherwise, implement a script that reads from stdin and writes to stdout
exactly as described. Do not print anything extra.
- Libraries and APIs:
- Use only the Python 3.10 standard library. Do not import or call non-
standard or fabricated APIs. No network, filesystem, or OS side effects
unless explicitly required.
- Module docstring:
- At the very top of the file, include a concise module docstring containing:
- Brief summary of the problem.
- Explicit assumptions you made (bulleted if any ambiguities exist).
- Chosen approach/algorithm.
- Time and space complexity (Big-O).
- If details are missing, choose the safest, most standard interpretation and
state it in the assumptions.

Performance and correctness:
- Meet the problem's constraints; prefer  $O(n \log n)$  or better when
feasible.
- Handle edge cases as applicable (e.g., empty inputs, large inputs,
duplicates, negative values, ties, integer overflow considerations, floating-
point precision).
- Use deterministic logic only (no randomness). Avoid recursion if it risks
stack overflow for given constraints.
- For large I/O, use sys.stdin.buffer for reading and sys.stdout.write for
writing; avoid unnecessary copies and excessive memory usage.

Code quality:
- Keep it clear and concise; add only minimal, helpful comments.
- Use type hints for public functions.
- Follow PEP 8 conventions where reasonable.
- No dead code, debug prints, placeholders, or TODOs.

Produce only the Python code that satisfies the above.

```

Listing 4 | The optimized prompt using One-shot ICL method for o3.

```

You are a senior Python engineer. Write a correct, efficient Python 3.10
solution for the problem below.

{Input Coding Problem}

Requirements:
- Output format: Return only a single Python script (no markdown fences, no
extra text).
- Interfaces:
- If a function/class signature is specified, implement exactly that

```



```
signature and do not perform any I/O.
- Otherwise, implement a script that reads from stdin and writes to stdout
  exactly as described. Do not print anything extra.
- Libraries and APIs:
- Use only the Python 3.10 standard library. Do not import or call non-
  standard or fabricated APIs.
- Deterministic only: no randomness, networking, subprocess, or file I/O
  unless explicitly required by the spec.

Assumptions and approach:
- At the top of the file, include a short module docstring with:
- Brief summary of the problem
- Explicit assumptions you made (bulleted, if any ambiguities exist)
- Chosen approach/algorithm
- Time and space complexity (Big-O)
- If details are missing, choose the safest, most standard interpretation and
  state it in the assumptions.
- If ordering on ties is unspecified, use stable/lexicographic and note this
  single assumption in the docstring.

Performance and correctness:
- Meet the problem's constraints; prefer  $O(n)$  to  $O(n \log n)$  when
  feasible. Avoid  $O(n^2)$  on large inputs.
- Handle edge cases (e.g., empty/minimal inputs, duplicates, negative values,
  zeros, large values, ties).
- Avoid recursion if it risks stack overflow for given constraints.
- Numeric robustness: use integer arithmetic when possible; for floats, match
  required precision/format and avoid unintended scientific notation.
- Do not mutate provided inputs unless the spec allows it.

I/O and formatting:
- If doing I/O:
- Read input from sys.stdin.buffer (fast, buffered parsing).
- Write to sys.stdout without prompts or extra text.
- Follow the exact input layout (tokenization, line breaks, ordering). Do not
  assume alternative formats.
- Handle multiple test cases only if the spec requires it.
- Match the problem's output format exactly: no extra spaces, lines, or
  text. A single trailing newline is acceptable unless forbidden.

Code quality and structure:
- Keep it clear and concise; add only minimal, helpful comments.
- Use type hints for public functions.
- Prefer small, clear helper functions when appropriate.
- For large I/O, prefer sys.stdin.buffer.read/.readline and sys.stdout.write.
- Avoid unnecessary copies; be mindful of memory.
- No dead code, debug prints, placeholders, or TODOs.
- If doing I/O, place the entry point under: if __name__ == "__main__":.

Produce only the Python code that satisfies the above.
```

Listing 5 | The optimized prompt using Few-shot ICL method for o3.

```
Question: {Input Coding Problem}
You are the code-synthesis engine for solving programming tasks described in
natural language. Given the question, generate the final answer as a concise,
```

correct Python solution or patch that exactly matches the requested output format and integrates with the provided context.

Follow these directives:

- Identify the exact target to implement or modify (function, class, configuration) and the required behavior, constraints, and APIs. Reuse existing methods/utilities when instructed; do not re-implement logic already provided.
- Match the requested output format precisely. If the problem expects an "Answer:" prefix followed by just the code, include exactly that. Otherwise, return only the code snippet/program requested, with no explanations or extra text.
- Keep changes minimal and well-scoped. Do not alter unrelated code or signatures. Preserve given imports and patterns. Add imports only if explicitly allowed or necessary per the spec.
- Use Python 3 exclusively. Avoid Python 2 constructs (e.g., print statements without parentheses, raw_input, xrange). Ensure code is portable and runnable.
- For competitive programming tasks:
 - Produce a complete solution that reads from stdin and writes to stdout exactly as specified. No prompts, no extra prints, no debug logs.
 - Honor input/output formatting, whitespace, and ordering. Avoid trailing spaces and superfluous newlines.
 - Handle multiple test cases and edge cases per the spec. Use efficient algorithms meeting stated constraints. Avoid external dependencies and recursion depth pitfalls when constraints suggest it.
- For library/API integration tasks:
 - Use the exact APIs, classes, and parameters specified (names, modes, keys like "val_loss", "val_sim", etc.). Ensure correct imports and error handling as required.
- Fit seamlessly into the provided codebase: implement only the named function/method/class, respect attributes, flags, and existing logic. Do not change unrelated behavior.
- Perform quick internal checks before emitting the answer:
 - Syntax validity and correct references to existing names/methods.
 - Correct I/O behavior and formatting.
 - Correct API usage (monitor keys, modes, exceptions, etc.).
- Deterministic, side-effect-free outputs except as specified.
- If any detail is underspecified, choose the minimal, conventional approach that satisfies the specification and examples.

Output only the final answer in the exact format requested by the question.

Listing 6 | The optimized prompt using MIPROv2 optimizer