

SAGE: A FRAMEWORK FOR SEMANTIC-ALIGNMENT-GUIDED ENGINEERING OF PROMPTS AND FINE-TUNING IN INDUSTRIAL CONTROL TASKS

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models show great potential for code generation tasks, but automatic code generation for industrial control systems still faces challenges such as inaccurate semantic understanding, a lack of alignment evaluation, and a shortage of domain-specific fine-tuning models. Given the stringent requirements for real-time performance, security, logical rigor, and correct execution of industrial control code, existing general-purpose methods struggle to meet these demands. Therefore, this paper proposes a semantic alignment-guided prompt engineering approach for industrial control tasks. The approach consists of three core components: first, a dataset of function prompt formats covering five structured prompt patterns and a selection of 1,500 prompt examples for industrial control tasks is constructed; second, a semantic alignment analysis metric is designed to evaluate the semantic correctness and task consistency of code generated by different models; and third, an alignment-guided fine-tuning strategy is proposed, leveraging prompt-output-intent triples to enhance the model’s generation capabilities for industrial control tasks. Experiments are conducted on five mainstream 7B models: DeepSeek-7B, Qwen2.5-7B, InternLM2-7B, Mistral-7B, and Gemma-7B. Results show that after fine-tuning, the executable performance of Mistral-7B and DeepSeek-7B increased from 0.719 to 0.886 and from 0.676 to 0.837, respectively, and the BLEU scores increased from 3.79 to 7.45 and from 3.45 to 6.62, respectively. All models maintained intent consistency (Intent = 1.000). Gemma-7B and Qwen2.5-7B showed decreases in executable performance, success rate, and BLEU, suggesting possible overfitting or distribution mismatch issues. The method proposed in this paper significantly improves the code executable performance and semantic alignment of some models in industrial control scenarios. It also reveals the sensitivity of model architecture to fine-tuning strategies, providing an important reference for subsequent architecture-aware alignment optimization.

1 INTRODUCTION

In recent years, generative artificial intelligence has developed rapidly and demonstrated its potential in many fields. Among them, code generation, as one of the important application areas of generative artificial intelligence, has received increasing attention. Code generation technology aims to reduce the workload of developers and improve programming efficiency by automatically generating source code that meets specific requirements.

The research on automatic code generation has a long history. In the early days, it mainly relied on compilers to directly map high-level languages into machine code. Although the conversion from languages such as Fortran and COBOL to assembly was achieved, it lacked optimization and portability and was difficult to adapt to complex application scenarios (Aho et al., 2006). Since 1990, modern compilers represented by GCC and LLVM have introduced cross-process optimization and multi-target backends, allowing the same optimization logic to be applied to multiple languages and platforms (Lattner & Adve, 2004). Around 2000, the model-driven development (MDD) method emerged, generating code through abstract models such as UML and SysML, and combining tools such as Simulink to achieve rapid prototyping of control systems (Schmidt, 2006). After 2010, it

054 has become possible to automatically generate interlocking, control logic, and other codes based
055 on requirement rules or industry standards (such as IEC 61131-3). Although the degree of automa-
056 tion has increased, the generated results are often rigid and difficult to handle complex semantics
057 and unstructured inputs (De Smet et al., 2008). With the emergence of large language models such
058 as GPT-4, code generation has begun to move from “rule-driven” to “intention-driven.” By under-
059 standing natural language instructions, the model can directly generate complete functions, control
060 logic, and even cross-domain hybrid codes (Chen et al., 2021), effectively alleviating the problem
061 that traditional modeling methods are difficult to cover complex semantics.

062 Unlike general code generation tasks, industrial control code generation requires not only functional
063 correctness but also real-time performance, security, and standard compliance (Krüger et al., 2012).
064 Its tasks involve ladder diagram instruction sets, PID parameter self-tuning, and industry standards
065 such as IEC 61131-3 and ISO26262. In this context, although large language models have the
066 ability to generate syntactically correct code, they still have obvious deficiencies in understanding
067 the domain semantics of industrial control. For example, the emergency stop (ESD) logic of the
068 PLC must strictly follow the timing constraint of “cutting off the output first and then triggering the
069 alarm”. If the generated logic is reversed, it may cause serious safety risks. At the same time, the
070 “hallucination” problem common in large models is even more fatal in industrial control scenarios
071 (Ouchani & Fakh, 2024). Fictional functions or communication rules have limited impact in gen-
072 eral applications, but in control systems, they may cause motor overload and burnout or equipment
073 communication failure, thereby inducing system-level accidents.

074 In response to the above challenges, some studies in recent years have explored the application of
075 large language models to industrial control code generation tasks. For example, the AlphaCodium
076 system proposed by Zhao et al. uses control flow tracking and rapid engineering mechanisms to
077 enhance the model’s ability to model temporal dependencies and function structures, alleviating the
078 hallucination problem to a certain extent (Ouchani & Fakh, 2024). The LLM4PLC framework pro-
079 posed by Fakh et al. combines semantic templates with formal verification mechanisms to generate
080 interlocking and control logic programs that comply with the IEC61131-3 standard, effectively im-
081 proving the reliability and standardization of the generated code (Fakh & Ouchani, 2024a). In addi-
082 tion, there are also studies that attempt to further improve the generation stability and interpretability
083 of large models in industrial environments through constraint injection, template guidance, and other
084 methods, and promote their application in actual engineering scenarios.

085 However, most of these methods remain at the experimental verification stage and lack a unified
086 prompt design and evaluation system, making them difficult to meet actual industrial needs. Cur-
087 rently, key challenges facing the industrial control field include: first, the lack of a suitable, high-
088 quality, and adaptable set of industrial prompts; second, the lack of methods to detect semantic
089 alignment between these prompts and generated results; and third, the lack of fine-tuned industrial
090 control models for evaluation and comparison. Therefore, there is an urgent need to build an in-
091 dustrial control system capable of prompt regulation, alignment detection, and model fine-tuning to
092 support performance evaluation and optimization of basic models using metrics such as BLEU and
093 Pass@k.

093 To address these issues, this paper proposes a semantic alignment guidance method for industrial
094 control code generation. The contributions of this paper can be summarized as follows: It proposes
095 a unified prompt design method, a semantic alignment evaluation mechanism, and an alignment-
096 driven model fine-tuning strategy. Through techniques such as structured prompt construction, se-
097 mantic label recognition, and alignment-guided training, it systematically addresses key issues such
098 as a lack of prompt sets, difficulty in intent detection, and model mismatch, resulting in significantly
099 improved generation results and reliable alignment of industrial semantics.

100 First, we constructed the Functional Framework Prompt (FFP) dataset. This dataset contains 1,500
101 prompt words across 10 categories, addressing the current lack of a suitable, high-quality set of
102 prompt words. These prompt words cover various aspects of structured text content. This dataset
103 provides essential information for subsequent evaluations such as alignment and fine-tuning of large
104 language models, making it suitable for comparing and improving code generation performance.

105 Secondly, a Semantic Alignment Analysis (SAS) method was proposed to evaluate the alignment
106 of code generation results with industrial control semantics. This method constructs a key semantic
107 tag system to annotate and match the generated code to ensure that functions such as sensor reading,

108 actuator control, and interlocking logic are correctly implemented. This method quantifies the con-
109 sistency of model output with the intended control intent. This method is based on a cross-language
110 semantic taxonomy that covers key control intents such as flow control, temperature control, and
111 interlocking logic. SAS quantifies the alignment of generated code with the intended functional
112 categories, providing a unified evaluation framework for different programming languages.

113 Third, we propose an alignment-guided LoRa fine-tuning (AGFT) strategy. This strategy leverages
114 semantic alignment scores to select and prioritize high-quality prompt-output pairs. Based on this,
115 LoRa is used for efficient fine-tuning, enabling large language models (LLMs) to better adapt to in-
116 dustrial control tasks while maintaining generation safety and domain compliance. By constructing
117 prompt-output-intent (POI) triples, AGFT enhances the model’s ability to perceive control objectives
118 and improves metrics such as code generation feasibility, task consistency, and success rate.

119 The remainder of this paper is organized as follows: "Related Work" outlines the basic preparations.
120 "Proposed Method" presents the experimental workflow. "Function Prompt Format (FFP)" intro-
121 duces our prompt word set and sample standard code set; "Semantic Alignment Analysis (SAS)"
122 introduces an evaluation system for aligning prompts and answers in industrial control systems;
123 and "Alignment-Guided Fine-Tuning (AGFT)" compares the performance of the LoRA fine-tuned
124 model with the original model. "Experiments" presents the experimental platform and experimental
125 results, followed by a discussion and analysis. "Conclusion" summarizes this paper and identifies
126 any shortcomings and future work.

127 128 2 RELATED WORK

129 130 2.1 PROMPT ENGINEERING

131
132 Prompt engineering has become a key technology for adapting large language models (LLMs) to
133 downstream tasks such as code generation. Zhou et al. (2024b) first proposed a standard for evaluat-
134 ing the reliability of LLM output from the perspective of trust and consistency, providing a reference
135 for subsequent prompt design. Subsequently, Shin et al. (2020) compared the performance differ-
136 ences between prompt-based and fine-tuned models in code-related tasks, while Lester et al. (2021)
137 proposed a parameter-efficient prompt fine-tuning method, providing a practical path for actual sys-
138 tems. In order to improve prompt efficiency, Zhou et al. (2024a) proposed a cost-effective and
139 efficient prompt generation method based on search. At the theoretical and methodological level,
140 Liu et al. (2023) systematically reviewed the prompt framework and proposed a classification of
141 methods and applications; Mohanty et al. (2025) further extended it to adaptive prompts, proposed a
142 framework driven by learnability, and emphasized risk control in multimodal environments. At the
143 same time, Kojima et al. (2022) proposed a zero-shot reasoning prompt method, providing new ideas
144 for practitioners and researchers. Furthermore, Gu et al. (2024) explored the role of prompt tags in
145 enhancing instruction tuning performance, highlighting the importance of tag-level prompt design.
146 Subsequently, Springer et al. (2024) compared various prompt engineering techniques through em-
147 pirical research, revealing the effectiveness of prompt design in different task scenarios. Finally,
148 Wang et al. (2024c) focused on the relationship between target code complexity and prompt strate-
149 gies, emphasizing the necessity of complexity-aware prompts.

150 2.2 ALIGNMENT/EVALUATION METRICS

151
152 Although prompt engineering has improved the adaptability of LLM in code generation, ensuring
153 that the generated results are aligned with human needs remains a key issue, which has led to the
154 research on alignment and evaluation metrics. Aligning large language models (LLMs) with human
155 values, preferences, and safety requirements has become a key research direction. At the evaluation
156 level, Chiang et al. (2024) proposed a semantic relevance metric designed specifically for prompt
157 evaluation to capture a deeper level of alignment between prompt representations and generated
158 outputs, providing a useful supplement to traditional reward models or BLEU-type metrics. Sub-
159 sequently, Shen et al. (2023) conducted a fundamental review of LLM alignment methods, focusing
160 on supervised fine-tuning and reinforcement learning with human feedback (RLHF). On this basis,
161 Wang et al. (2024b) proposed a more comprehensive taxonomy covering RLAI and PPO-based
methods and pointed out that there are still challenges in building a scalable and generalizable align-
ment framework. For multimodal scenarios, Yu et al. (2025) explored the adaptability of alignment

162 methods when combining language models with visual or audio modalities, emphasizing the com-
163 plexity of preference alignment in cross-modal environments. Finally, Zhang et al. (2024b) proposed
164 a roadmap for a scalable automatic alignment process, combining automated feedback collection and
165 tuning loops, focusing on the scalability and stability of the evaluation.

167 2.3 COMMAND ADJUSTMENT/FINE-TUNING

169 As alignment methods continue to evolve, instruction tuning, as a core means to improve the per-
170 formance of models in specific tasks, has gradually become a research focus. It has been proven
171 that instruction tuning is an important strategy for adapting large language models (LLMs) to spe-
172 cific domain tasks, especially in code generation. Wei et al. (2021) first proposed the basic idea of
173 instruction tuning and verified its effectiveness in code generation. Subsequently, He et al. (2024)
174 proposed a secure code generation pipeline that uses an instruction tuning model to enforce security
175 and compliance constraints. Zhong et al. (2024) presented an empirical method for transitioning
176 from prompt engineering to full fine-tuning, with special attention to actual deployment scenarios
177 and customized frameworks. As a supplement, Ma et al. (2024) proposed LLaMoCo, which op-
178 timizes code generation through instruction tuning to improve computational efficiency. In terms
179 of comparing different methods, Wang et al. (2023) conducted an empirical comparison between
180 instruction tuning and prompt engineering, while Sun et al. (2024) conceptually explained the core
181 principles of instruction tuning and distinguished it from other tuning paradigms. Liu et al. (2024b)
182 further explored a hybrid optimization strategy that combines fine-tuning with prompt engineering.
183 Finally, Xu et al. (2023) conducted a comprehensive analysis of code-specific instruction tuning
184 models and evaluated their performance on a variety of code-related tasks.

185 2.4 INDUSTRIAL CONTROL SYSTEMS

187 The above methods have made progress in general tasks, but in the high-reliability scenario of in-
188 dustrial control systems, they still need to be explored in combination with domain characteristics.
189 Large language models (LLMs) have recently been widely explored in the field of industrial au-
190 tomation and control systems. Zhang et al. (2024a) pointed out that AI-generated code may bring
191 critical supply chain risks in industrial environments, while Li et al. (2024) analyzed the internal
192 mechanism of LLM in detail and systematically explained the operation and response behavior
193 of the converter. Furthermore, Chen et al. (2025) explored the programming workflow based on
194 LLM from a practitioner’s perspective and emphasized potential problems in application scenarios.
195 In terms of open source research, Wang et al. (2024a) compared retrieval-augmented generation
196 (RAG), fine-tuning, and prompt engineering methods, and provided a benchmark for collaborative
197 evaluation. In terms of framework exploration, Fakih & Ouchani (2024b) proposed LLM4PLC, a
198 verifiable PLC programming framework based on LLM, which aims to address the correctness and
199 safety constraints of control logic. Liu et al. (2024a) proposed an optimization framework that in-
200 tegrates prompt engineering, RAG, and fine-tuning, while Klein & Smith (2024) proposed a formal
201 framework for “generative prompt engineering” to strengthen the theoretical foundation and mea-
202 surement indicators. At the same time, Wei et al. (2022) proposed a Chain-of-Thought prompting
203 framework and experimentally demonstrated the key role of prompt design in complex reasoning
204 tasks. In terms of specific method research, Singh & Patel (2024) analyzed the interaction between
205 fine-tuning and prompt design in LLM-based code generation and proposed a hybrid strategy suit-
206 able for engineering tasks; Zhao et al. (2024) proposed the AlphaCodium system, which combines
207 prompt engineering with control flow tracing to improve code generation performance.

208 3 METHODOLOGY

209
210 Figure 1 shows the framework of the method proposed in this paper. The Functional Framework
211 Prompt (FFP) dataset is proposed, which contains industrial control prompts and five prompt types.
212 Semantic Alignment Analysis (SAS) is designed to test the quality of FFP using semantic labels
213 and semantic alignment scores. Alignment-Guided Fine-Tuning (AGFT) is proposed to fine-tune
214 models such as DeepSeek-7B and Qwen2.5-7B using the Prompt-Output-Intent (POI) pool. Finally,
215 the fine-tuned models are systematically compared with the original models using metrics such as
execution rate, intent alignment, success rate, and BLEU.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

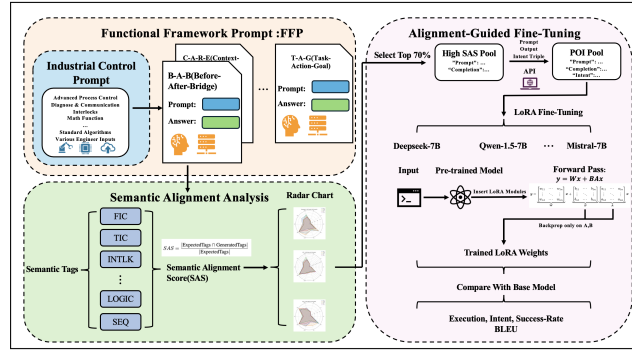


Figure 1: SAGE: semantic alignment analysis and alignment-guided fine-tuning with functional framework prompts

3.1 FFP: FUNCTIONAL FRAMEWORK PROMPT

We design five structured prompt formats to guide industrial control code generation. Each emphasizes a different reasoning aspect while remaining concise and interpretable. Due to space limitations, a detailed description is included in the Appendix.

3.1.1 BEFORE-AFTER-BRIDGE (B-A-B)

Emphasizes transformation from a baseline to an improved solution. **Example:** compare traditional PID control with a model predictive control implementation.

3.1.2 CONTEXT-ACTION-RESULT-EXAMPLE (C-A-R-E)

Links environment, action, and outcome to encourage descriptive reasoning. **Example:** reactor temperature control with an ANN replacing PID.

3.1.3 ROLE-INPUT-STEPS-EXPECTATION (R-I-S-E)

Highlights procedural reasoning and stepwise logic. **Example:** PLC startup sequence for a machine with heating, cooling, and feeders.

3.1.4 ROLE-TASK-FORMAT (R-T-F)

Maps task requests directly to a fixed-format output. **Example:** generate a PID controller as an IEC 61131-3 function block.

3.1.5 TASK-ACTION-GOAL (T-A-G)

Supports backward reasoning from goal to action, useful for diagnostics and optimization. **Example:** monitor PLC network health with watchdog logic.

3.2 SAS: SEMANTIC ALIGNMENT SCORE

3.2.1 CROSS-LINGUAL CONTROL SEMANTIC TAXONOMY

To ensure effective evaluation of functional correctness across industrial control code written in multiple languages—including Python, C++, MATLAB, and PLC Structured Text—we propose a language-independent control semantics taxonomy. This taxonomy covers key control intents commonly found in automation programs:

- **FIC: Flow Control** – Includes flow setpoint adjustment, flowmeter signal processing, and valve control to achieve stable flow.

- 270 • **TIC: Temperature Control** – Involves temperature measurement, PID regulation, heat-
271 ing/cooling commands, and thermal process optimization.
- 272 • **PIC: Pressure Control** – Covers pressure sensing, setpoint control loops, pressure relief
273 valve logic, and pressure safety regulation.
- 274 • **INTLK: Interlock Logic** – Monitoring of safety-critical conditions, trip logic implemen-
275 tation, emergency stop, and permissive checking.
- 276 • **ALM: Alarm/Fault Handling** – Detection of abnormal conditions, alarm triggering, op-
277 erator notification, and fault logging.
- 278 • **ACT: Actuator Commands** – Control signals for pumps, valves, motors, and other final
279 control elements.
- 280 • **SEN: Sensor Acquisition** – Reading and conditioning process signals, scaling anal-
281 og/digital inputs, and verifying sensor status.
- 282 • **LOGIC: Logical Structures** – IF/CASE decisions, Boolean operations, bitwise logic, and
283 conditional control flow.
- 284 • **SEQ: Sequential Control** – Process step management, state machines, sequential function
285 charts, and batch process sequencing.
- 286 • **RPT: Reporting and Logging** – Automatic report generation, event logging, trending, and
287 historical data processing.

290 The taxonomy supports consistent semantic evaluation across PLC structured text, Python, and
291 C++ platforms, facilitating the calculation of unified metrics such as the Semantic Alignment Score
292 (SAS).

294 3.2.2 STRUCTURE RADAR CHART

295 We analyze how different prompt structures in the FFP framework affect the code generation cor-
296 responding to the aforementioned semantic labels. For each hint structure (e.g., R-I-S-E, C-A-R-E,
297 T-A-G), we generate a set of code examples in multiple languages and label their control intent using
298 a rule-based parser.

300 3.2.3 SEMANTIC ALIGNMENT SCORE

302 To quantify how well a cue achieves its target functionality, we define a semantic alignment score
303 (SAS):

$$304 \quad SAS = \frac{|\text{ExpectedTags} \cap \text{GeneratedTags}|}{|\text{ExpectedTags}|}$$

306 Expected tags are determined based on prompt intention (e.g., a prompt requesting "interlock logic"
307 expects INTLK and ACT), while generated tags are extracted from code using multi-language se-
308 mantic rules. Cross-language semantic alignment provides a robust foundation for evaluating the
309 intent fulfillment of prompt-based code generation. By combining a unified control taxonomy with
310 the SAS metric, we create a reusable benchmark for aligning LLM output with domain-specific
311 control objectives.

313 3.3 AGFT: ALIGNMENT-GUIDED FINE-TUNING

314 To further improve the reliability and domain adaptability of large language models (LLMs) for
315 industrial control code generation, we propose an alignment-guided fine-tuning (AGFT) method.
316 This method leverages the alignment score between generated outputs and ground-truth structured
317 text to guide sample selection and optimization during fine-tuning. The overall pipeline consists
318 of three key components: selection of high-quality prompt-output-intent triples based on alignment
319 metrics, lightweight fine-tuning using LoRA, and evaluation metric feedback.

321 3.3.1 FILTER AND ENHANCE USING INDUSTRIAL CONTROL SAMPLES

322 The proposed approach begins with an industrial control benchmark framework, which contains
323 1,500 diverse prompt-response pairs covering diverse industrial control domains. Each generated

code example is evaluated using a SAS, which quantifies how well the model-generated response matches the expected structure, logic, and format specifications. Examples with low SAS scores are filtered out, while those with high SAS scores are retained and further enhanced using data augmentation strategies (e.g., explaining the prompt, changing parameter constraints, or slightly adjusting the target condition) to enrich the fine-tuning framework.

3.3.2 LIGHTWEIGHT PARAMETER EFFICIENT FINE-TUNING

To reduce computational overhead, we employ Low Rank Adaptation (LoRA) as a fine-tuning method for the backbone network. LoRA inserts trainable rank-factorized weight matrices into pre-trained model layers, enabling efficient adaptation while freezing most of the original parameters. By applying LoRA to the decoder layers and attention heads, we are able to guide the generative behavior of the LLM towards better syntactic and semantic alignment without retraining the entire model.

3.3.3 PROMPT-OUTPUT-INTENT ENHANCEMENT

We introduced a three-way alignment strategy, constructing "prompt-output-intent" (POI) triplets. The "intent" element is a refined natural language summary automatically extracted from the original prompt (e.g., "Generate PID control code to maintain the reactor temperature at 180°C"). During fine-tuning, these POI triplets enable the model to understand both the prompt and the underlying goal. This additional context helps the model maintain fidelity to the control intent, especially in safety-critical or multi-stage processes.

3.3.4 EXPERIMENTAL PLAN

To implement AGFT in our experiments, we performed the following steps:

- **Step 1:** Compute the SAS score on an existing industrial control dataset using either a rule-based or learning-based alignment scorer.
- **Step 2:** Select the top 70% of high-SAS examples and augment them with smaller hint variants, tripling the sample pool.
- **Step 3:** Form point-of-interest triplets by automatically generating intents and injecting the intent tags as soft hints.
- **Step 4:** Apply LoRA-based fine-tuning on an open-source LLM (e.g., DeepSeek-7B or Qwen-1.5-7B) using HuggingFace's `peft` library.
- **Step 5:** Evaluate the fine-tuned model on the industrial control dataset using BLEU and Pass@k, and compare it to the baseline model.

This alignment-guided approach allows the model to selectively learn from highly aligned demonstrations, achieving improved performance while ensuring safety and compliance in industrial scenarios.

4 EXPERIMENT

4.1 EXPERIMENTAL SOFTWARE AND HARDWARE SETTINGS

The hardware and software setup for this experiment is as follows:

Hardware: Personal laptop, Alibaba Cloud A10 servers (4, each with 24GB of video memory)

Software: Vscode (SSH) (for model training and running basic code), Github (for storing datasets and results), HuggingFace (for downloading the model to be fine-tuned)

4.2 DATASET CONSTRUCTION DETAILS

The dataset is constructed based on 100 commonly used prompt words for industrial control scenarios. Based on these 100 prompt words, five different prompt methods (B-A-B, C-A-R-E, R-I-S-E, R-T-F, and T-A-G) were used to expand these 100 prompt words. These prompt words were then

Table 1: SAS distributions across models

(a) DeepSeek-7B

Prompt	FIC	TIC	PIC	INTLK	ALM	ACT	SEN	LOGIC	SEQ	RPT
B-A-B	0.034	0.026	0.051	0.126	0.126	0.092	0.171	0.130	0.087	0.156
C-A-R-E	0.034	0.018	0.054	0.108	0.098	0.096	0.184	0.162	0.084	0.162
R-I-S-E	0.030	0.016	0.049	0.169	0.122	0.091	0.208	0.080	0.103	0.131
R-T-F	0.036	0.019	0.060	0.137	0.115	0.089	0.206	0.115	0.082	0.141
T-A-G	0.029	0.010	0.044	0.128	0.149	0.112	0.238	0.089	0.091	0.110

(b) GPT

Prompt	FIC	TIC	PIC	INTLK	ALM	ACT	SEN	LOGIC	SEQ	RPT
B-A-B	0.034	0.016	0.053	0.135	0.119	0.111	0.208	0.103	0.077	0.145
C-A-R-E	0.024	0.011	0.048	0.149	0.125	0.093	0.215	0.106	0.106	0.125
R-I-S-E	0.031	0.010	0.048	0.176	0.130	0.099	0.216	0.097	0.081	0.112
R-T-F	0.037	0.011	0.053	0.130	0.122	0.101	0.202	0.114	0.085	0.144
T-A-G	0.032	0.012	0.050	0.136	0.136	0.099	0.191	0.117	0.089	0.136

(c) Grok

Prompt	FIC	TIC	PIC	INTLK	ALM	ACT	SEN	LOGIC	SEQ	RPT
B-A-B	0.036	0.025	0.063	0.115	0.149	0.085	0.162	0.135	0.068	0.160
C-A-R-E	0.030	0.019	0.061	0.096	0.163	0.077	0.163	0.126	0.126	0.138
R-I-S-E	0.034	0.017	0.069	0.100	0.180	0.105	0.182	0.077	0.066	0.169
R-T-F	0.040	0.020	0.062	0.123	0.189	0.080	0.185	0.105	0.074	0.121
T-A-G	0.030	0.016	0.065	0.109	0.154	0.085	0.170	0.134	0.091	0.148

generated using three different mature large-scale model platforms (Grok, GPT, and Deepseek). Therefore, the dataset contains a total of 1,500 prompt words.

4.3 VISUALIZING SEMANTIC ALIGNMENT

Tables 1a, 1b, and 1c report the semantic alignment scores (SAS) distributions of different prompt structures on ten functional semantic labels for the Deepseek-7B, GPT, and Grok models, respectively. Overall, T-A-G consistently achieves the highest SAS on SEN labels across all models. For example, in Table 1a (Deepseek-7B), T-A-G scores 0.238 on SEN, significantly higher than B-A-B’s 0.171, demonstrating its optimal performance in generating sensor-related logic. The R-I-S-E prompt performs exceptionally well on interlock (INTLK) and report (RPT) labels. For example, in Table 1b (GPT), R-I-S-E achieves 0.176 on INTLK, surpassing other prompt structures and demonstrating its effectiveness in generating protection and diagnostic logic. In contrast, the scores of C-A-R-E and R-T-F are more evenly distributed. For example, in Table 1c (Grok), C-A-R-E’s scores across the ten labels show little variation, highlighting its generalization ability across different functional categories. Although the absolute scores vary between models, the relative performance trends of the prompt structures remain consistent. This result demonstrates the robustness and cross-model applicability of the Functional Prompt Format (FFP) dataset in guiding LLMs to achieve the target semantic intent. Furthermore, we use the radar chart in Figure 2 to visualize alignment trends, showing the SAS values for semantic labels for each prompt type across different models.

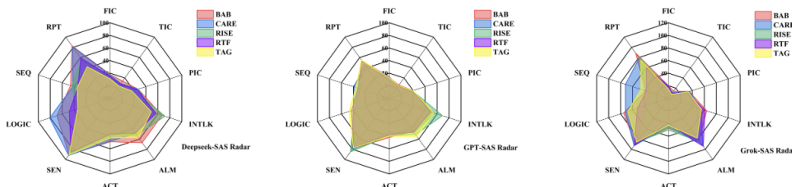


Figure 2: Radar chart of semantic alignment score

Table 2: Comparison of training losses and metric deltas across models

(a) Training losses over epochs

Epoch	DeepSeek-7B	Gemma-7B	InternLM2-7B	Mistral-7B	Qwen2.5-7B
0	1.49	2.45	2.22	0.91	1.92
1	0.99	2.14	2.03	0.62	1.08
2	0.94	2.41	1.74	0.51	0.72

(b) Pretrained vs. finetuned comparison on four metrics

Model	Exec _o	Exec _t	Intent _o	Intent _t	Succ _o	Succ _t	BLEU _o	BLEU _t
DeepSeek-7B	0.676	0.837	1.000	1.000	0.424	0.601	3.45	6.62
Gemma-7B	0.647	0.625	1.000	1.000	0.389	0.358	3.31	2.96
InternLM2-7B	0.702	0.641	1.000	1.000	0.418	0.385	3.33	3.09
Mistral-7B	0.719	0.886	1.000	1.000	0.473	0.585	3.79	7.45
Qwen2.5-7B	0.716	0.653	1.000	1.000	0.467	0.379	4.06	3.04

4.4 ALIGNMENT-GUIDED FINE-TUNING LOSS RESULTS

Table 2a shows the training loss curves for five different 7B-scale language models during alignment-guided fine-tuning over three epochs. The loss of all models shows a downward trend, demonstrating that the models are able to effectively learn from the alignment-augmented dataset. Mistral-7B and Qwen2.5-7B show the most significant decreases in training loss, indicating better convergence. In contrast, Gemma-7B exhibits an unusual upward trend after the first epoch, suggesting possible overfitting or instability during fine-tuning. DeepSeek-7B and InternLM2-7B achieve higher final loss values than Mistral and Qwen, indicating slower convergence, but their performance steadily improves. These results highlight the variability in how models respond to alignment-guided supervised learning and emphasize the importance of architecture-specific tuning strategies.

Table 2b presents a comparative evaluation of the pre-trained and fine-tuned 7B-scale language models on four key metrics: Executability (Exec), Intent Consistency (Intent), Task Success (Success), and BLEU score. Among all models, Mistral-7B and DeepSeek-7B show significant improvements after fine-tuning, particularly in Executability (from 0.719 to 0.886 and 0.676 to 0.837, respectively) and BLEU score (from 3.79 to 7.45 and 3.45 to 6.62, respectively), indicating improved syntactic and semantic quality of the generated code. Notably, all models maintain perfect intent consistency (Intent = 1.000), confirming that the functional objective of the hint is preserved. However, models such as Gemma-7B and Qwen2.5-7B show performance degradation in terms of executableness, success rate, and BLEU after fine-tuning, suggesting possible overfitting or misalignment with the training distribution. These results indicate that while alignment-guided fine-tuning can significantly improve code generation performance, its effectiveness is model-dependent, highlighting the importance of architecture-aware fine-tuning strategies.

5 CONCLUSION

In this study, we introduce a structured and scalable pipeline for industrial control code generation by proposing three key innovations. First, we design a Function Prompt Format (FFP) dataset that enhances the controllability and domain-specificity of prompts through five well-defined structural formats. Second, we develop a semantic alignment score (SAS) mechanism based on intent label matching to quantitatively evaluate the degree of alignment between the prompt structure and the generated code semantics. Third, we propose an alignment-guided fine-tuning (AGFT) method that leverages SAS to guide data selection and improves the adaptability of the LLM through lightweight techniques such as LoRA. Our evaluation on multiple LLMs demonstrates consistent improvements in execution success rate, BLEU score, and functional intent preservation. These contributions not only enhance the interpretability and controllability of code generation in industrial automation but also highlight the feasibility of deploying domain-adaptive LLMs in safety-critical environments. Future work includes automatic prompt structure recommendation based on task intent, integration with retrieval-augmented generation (RAG), and the construction of larger industrial datasets for ongoing alignment and benchmarking.

486 A APPENDIX

487 REFERENCES

- 488 Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman. *Compilers: Principles, Tech-*
489 *niques, and Tools (2nd Edition)*. Addison-Wesley, Boston, 2006.
- 490 Mark Chen, Jerry Tworek, Heewoo Jun, et al. Evaluating large language models trained on code,
491 2021. URL <https://arxiv.org/abs/2107.03374>.
- 492 Rui Chen, Hao Zhao, and Yi Sun. Programming with large language models: Practices and chal-
493 lenges. In *Proceedings of the ACM Symposium on Software Engineering*, New York, USA, 2025.
494 ACM.
- 495 Chih-Hui Chiang, Hsin Lee, and Yifan Zhou. Semantic relevance metrics for prompt evaluation in
496 large language models. In *International Conference on Learning Representations (ICLR)*, 2024.
497 URL <https://openreview.net/forum?id=abc123>.
- 498 Pieter De Smet, Bart De Decker, and Jan Vanthienen. Automatic code generation for industrial
499 automation using iec 61131-3. *Computers in Industry*, 59(8):858–866, 2008. doi: 10.1016/j.
500 compind.2008.06.003.
- 501 Rachid Fakh and Samir Ouchani. Llm4plc: A framework for verifiable code generation for indus-
502 trial control systems, 2024a. URL <https://arxiv.org/abs/2403.05561>.
- 503 Wissam Fakh and Samir Ouchani. Llm4plc: Harnessing large language models for verifiable pro-
504 gramming of plcs. In *Proceedings of the 2024 ACM/IEEE International Conference on Cyber-*
505 *Physical Systems (ICCP)*, New York, USA, 2024b. ACM/IEEE.
- 506 Jiawei Gu, Fang Wei, and Xinyun Chen. Prompt token design for enhancing instruction tuning. In
507 *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (ACL)*,
508 Stroudsburg, USA, 2024. ACL.
- 509 Jie He, Rui Zhang, and Fang Xu. Instruction tuning for secure code generation. In *Proceedings of*
510 *the IEEE Symposium on Security and Privacy (SP)*, Piscataway, USA, 2024. IEEE.
- 511 Felix Klein and John Smith. Generative prompt engineering: Models, methods, and metrics. In
512 *Lecture Notes in Computer Science (LNCS)*, pp. 56–70. Springer, Berlin, Germany, 2024.
- 513 Takeshi Kojima, Shixiang Gu, Machel Reid, et al. Large language models are zero-shot reasoners.
514 In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 35, pp. 22199–22213.
515 Curran, 2022. URL <https://arxiv.org/abs/2205.11916>.
- 516 Markus Krüger, Mario Trapp, Jan Kowalski, et al. Automatic code generation for safety-critical plc
517 systems. *IFAC Proceedings Volumes*, 45(2):119–124, 2012. doi: 10.3182/20120215-3-AT-3016.
518 00025.
- 519 Chris Lattner and Vikram Adve. Llvm: A compilation framework for lifelong program analysis &
520 transformation. In *Proceedings of the International Symposium on Code Generation and Opti-*
521 *mization (CGO)*, pp. 75–86, Palo Alto, USA, 2004. IEEE. doi: 10.1109/CGO.2004.1281665.
- 522 Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale: Parameter-efficient prompt
523 tuning for pre-trained language models. In *Proceedings of the 2021 Conference on Empirical*
524 *Methods in Natural Language Processing (EMNLP)*, pp. 3045–3059, Stroudsburg, USA, 2021.
525 ACL. doi: 10.18653/v1/2021.emnlp-main.243.
- 526 Qi Li, Zhao Wang, and Jun Xu. Inside transformers: A survey on architecture and interpretability
527 of large language models. *ACM Computing Surveys*, 2024.
- 528 Chen Liu, Xiang Zhang, and Jian He. A unified framework of prompt engineering, rag, and fine-
529 tuning for code generation. In *Proceedings of the International Joint Conference on Artificial*
530 *Intelligence (IJCAI)*, Palo Alto, USA, 2024a. IJCAI.

- 540 Xiao Liu, Kai Zheng, and Wei Xu. Prompting frameworks for large language models: A survey.
541 *ACM Computing Surveys*, 2023.
- 542
- 543 Yan Liu, Jian He, and Lei Xu. Hybrid optimization of large language models: Combining fine-
544 tuning and prompt engineering. In *Proceedings of the 2024 Conference on Empirical Methods in*
545 *Natural Language Processing (EMNLP)*, Stroudsburg, USA, 2024b. ACL.
- 546
- 547 Rui Ma, Jian Li, and Lei Song. Llamoco: Instruction tuning of large language models for opti-
548 mized code generation. In *Proceedings of the 2024 IEEE International Conference on Software*
549 *Engineering (ICSE)*, Piscataway, USA, 2024. IEEE.
- 550
- 551 Abhishek Mohanty, Shreya Singh, and Soumen Chakrabarti. Adaptive prompting for multimodal
552 llms: Learnability and risk-controlled generalization. In *Proceedings of the 2025 Annual Meeting*
553 *of the Association for Computational Linguistics (ACL)*, Stroudsburg, USA, 2025. ACL.
- 554
- 555 Samir Ouchani and Rachid Fakih. A formal verification framework for llm-generated plc programs
556 in safety-critical applications, 2024. URL <https://arxiv.org/abs/2403.07756>.
- 557
- 558 Douglas C. Schmidt. Model-driven engineering. *IEEE Computer*, 39(2):25–31, 2006. doi: 10.1109/
559 MC.2006.58.
- 560
- 561 Yidong Shen, Yuntao Bai, Yixin Tang, et al. Large language model alignment: A survey, 2023. URL
562 <https://arxiv.org/abs/2307.12966>.
- 563
- 564 Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. Autoprompt:
565 Eliciting knowledge from language models with automatically generated prompts. In *Proceedings*
566 *of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp.
567 422–434, Stroudsburg, USA, 2020. ACL. doi: 10.18653/v1/2020.emnlp-main.346.
- 568
- 569 Arjun Singh and Kavya Patel. Fine-tuning and prompt engineering for llm-based code generation.
570 *Journal of Systems Architecture*, 145:103981, 2024. doi: 10.1016/j.sysarc.2024.103981.
- 571
- 572 Jonas Springer, Lukas Müller, and Anna Schmidt. A study of prompt engineering techniques for
573 code generation. In *Lecture Notes in Computer Science (LNCS)*, pp. 78–92. Springer, Berlin,
574 Germany, 2024.
- 575
- 576 Zhen Sun, Hao Liu, and Jun Zhao. Instruction tuning: Concepts, methods, and applications. *IEEE*
577 *Transactions on Neural Networks and Learning Systems*, 2024.
- 578
- 579 Peng Wang, Hao Liu, and Jun Zhang. Retrieval-augmented generation vs fine-tuning: An empirical
580 comparison. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Palo Alto, USA,
581 2024a.
- 582
- 583 Tao Wang, Han Zhao, and Yan Lin. A comprehensive survey of llm alignment techniques: Rlhf, rlaif,
584 ppo and beyond. In *Proceedings of the International Joint Conference on Artificial Intelligence*
585 *(IJCAI)*, Palo Alto, USA, 2024b. IJCAI.
- 586
- 587 Xin Wang, Ling Zhang, and Hao Sun. An empirical comparison of instruction tuning and prompt
588 engineering for code generation. In *Proceedings of the 2023 Conference on Empirical Methods*
589 *in Natural Language Processing (EMNLP)*, Stroudsburg, USA, 2023. ACL.
- 590
- 591 Yi Wang, Jing Liu, and Xin Chen. Selecting prompt engineering techniques for code generation
592 via complexity prediction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Palo
593 Alto, USA, 2024c.
- 594
- 595 Jason Wei, Maarten Bosma, Vincent Zhao, et al. Finetuned language models are zero-shot learners.
596 In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 34, pp. 19328–19341.
597 Curran, 2021.
- 598
- 599 Jason Wei, Xuezhi Wang, Dale Schuurmans, et al. Chain-of-thought prompting elicits reasoning in
600 large language models. In *Advances in Neural Information Processing Systems (NeurIPS)*, vol-
601 ume 35, pp. 24824–24837. Curran, 2022. URL <https://arxiv.org/abs/2201.11903>.

- 594 Jing Xu, Kai Wang, and Chen Li. Instruction-tuned models for code: A comprehensive study.
595 In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*
596 (*EMNLP*), Stroudsburg, USA, 2023. ACL.
- 597
- 598 Min Yu, Ling Zhang, and Ming Chen. Aligning multimodal llms with human preference: A survey.
599 In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*
600 (*EMNLP*), Stroudsburg, USA, 2025. ACL.
- 601
- 602 Kai Zhang, Xin Chen, and Yan Liu. Security risks of ai-generated code in industrial automation. In
603 *Proceedings of the IEEE International Conference on Industrial Informatics (INDIN)*, Piscataway,
604 USA, 2024a. IEEE.
- 605
- 606 Yifan Zhang, Zeyu Liu, and Wei Gao. Towards scalable automated alignment of llms: A sur-
607 vey. In *International Conference on Learning Representations (ICLR)*, 2024b. URL <https://openreview.net/forum?id=xyz456>.
- 608
- 609 Lin Zhao, Wei Ren, and Peng Huang. Code generation with alphacodium: From prompt engineering
610 to flow control. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 36.
611 Curran, 2024.
- 612
- 613 Zhen Zhong, Yu Lin, and Ming Zhou. From prompt engineering to full fine-tuning: Practical frame-
614 works for llm customization. In *Findings of the Association for Computational Linguistics: ACL*
615 *2024*, Stroudsburg, USA, 2024. ACL.
- 616
- 617 Han Zhou, Yi Jiang, and Qiang Liu. Cost-effective search-based prompt engineering for code gener-
618 ation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*
619 (*ACL*), Stroudsburg, USA, 2024a. ACL.
- 620
- 621 Ming Zhou, Jing Li, Hui Sun, et al. Trustworthy large language models: A survey and guideline for
622 evaluating alignment. *Journal of Artificial Intelligence Research*, 2024b.

624 B DETAILED EXAMPLES OF FFP FORMATS

626 B.1 BEFORE-AFTER-BRIDGE (B-A-B)

628 The B-A-B format is inspired by storytelling and writing templates that are effective in persuasive
629 or comparative contexts. It includes:

- 630 • **Before:** A description of the traditional or current approach or challenge.
- 631 • **After:** A statement of the desired or improved future outcome.
- 632 • **Bridge:** A description of how the transformation will occur or how the model is expected to
633 operate.

634 This structure is particularly effective when comparing baseline control logic with optimized solu-
635 tions (e.g., traditional PID versus model predictive control). It guides LLMs in reasoning about the
636 change process, which is critical for industrial control improvements.

638 **Example:**

- 639 • **Before:** In gas turbines, precise temperature control is critical for performance, safety, and emis-
640 sions compliance. Turbine temperature fluctuations due to load or fuel variations can lead to
641 inefficiency, overheating, or damage. Manual or poorly adjusted controls often fail to maintain
642 stable turbine conditions.
- 643 • **After:** Develop an IEC 61131-3 structured text program that implements a PID loop to control
644 turbine temperature, adjust the inlet valve to maintain the setpoint, adhere to valve limits, and
645 ensure stable and responsive operation under varying conditions.
- 646 • **Bridge:** Develop IEC 61131-3 structured text code to calculate the error, update the PID term,
647 and adjust the inlet valve position within a safe range to maintain the turbine temperature setpoint
under varying loads.

648 B.2 CONTEXT-ACTION-RESULT-EXAMPLE (C-A-R-E)

649
650 The C-A-R-E format is based on descriptive reasoning. Its structure is as follows:

- 651 • **Context:** The operating setting or physical environment.
- 652 • **Action:** A specific control action or method.
- 653 • **Result:** The expected or desired outcome.
- 654 • **Example:** A specific scenario or example of a requirement.

655
656 This prompt is well-suited for tasks that require interpreting physical processes as code or rules,
657 particularly in batch control and safety logic scenarios.

658 **Example:**

- 659 • **Context:** In a chemical process, reactor temperature control is challenging due to its nonlinear
660 and time-varying dynamics.
- 661 • **Action:** Develop a Python-based ANN controller and train it using historical data to predict and
662 regulate reactor temperature.
- 663 • **Result:** Achieve stable temperature control with better adaptability and efficiency than a tradi-
664 tional PID controller.
- 665 • **Example:** Simulate a sudden change in feed temperature; demonstrate how the ANN can
666 smoothly adjust to the overshoot and oscillation of the PID controller.

668 B.3 ROLE-INPUT-STEPS-EXPECTATION (R-I-S-E)

669
670 The R-I-S-E format emphasizes procedural logic and helps guide the model through a defined work-
671 flow. Its components include:

- 672 • **Role:** The model’s assumed role or responsibilities (e.g., controls engineer).
- 673 • **Input:** Available data, sensor signals, or initial conditions.
- 674 • **Steps:** The sequential actions or algorithms to be followed.
- 675 • **Expectation:** The final outcome or performance goal.

676
677 This format is particularly well-suited for sequential logic, state transitions, or batch operations.

678 **Example:**

- 679 • **Role:** You are a PLC programmer responsible for developing the startup logic for a 3D bag-
680 making machine.
- 681 • **Input:** The machine contains eight heating stations, eight cooling stations, and dual feed rollers.
- 682 • **Steps:** Start the heating stations sequentially, enable the cooling system, and synchronize the
683 feeder and cutter timing.
- 684 • **Expectation:** Use structured text to ensure safe and synchronized startup of all machine compo-
685 nents.

687 B.4 ROLE-TASK-FORMAT (R-T-F)

688
689 The R-T-F format is the most concise and direct. It is best used to quickly scope code generation
690 tasks with fixed-format output. It includes:

- 691 • **Role:** The position or responsibility assumed by the model.
- 692 • **Task:** The specific action to be completed.
- 693 • **Format:** The desired output format (e.g., code, table, function block).

694
695 This format is often used in automated code generation platforms to enable rapid conversion from
696 human language to code.

697 **Example:**

- 698 • **Role:** Play the role of a control system engineer.
- 699 • **Task:** Implement a PID controller in structured text to maintain the liquid level in a storage tank.
- 700 • **Format:** Output a complete IEC 61131-3 function block, including initialization and control
701 logic.

702 B.5 TASK-ACTION-GOAL (T-A-G)
703

704 The T-A-G format is useful when users need the model to reason backward from a goal and define
705 its own path of action. Its structure is as follows:

- 706 • **Task:** General problem or high-level goal.
- 707 • **Action:** Proposed control strategy or approach.
- 708 • **Goal:** Operational outcome or process improvement.

709 T-A-G is well-suited for diagnostic code, optimization scenarios, and systems where the end result
710 guides logic design.
711

712 **Example:**

- 713 • **Task:** Monitor the communication health of a PLC network.
 - 714 • **Action:** Implement watchdog logic for OPC UA, Modbus, and Profinet.
 - 715 • **Goal:** Ensure continuous diagnostics and generate alerts on connection failures.
- 716

717 C LLM USAGE DISCLOSURE
718

719 In preparing this work, large language models (LLMs) were used only as auxiliary tools. Specifi-
720 cally, LLMs assisted in (i) refining the English writing for clarity and grammar, (ii) checking LaTeX
721 formatting and figure/table alignment, and (iii) generating small illustrative code snippets used as
722 examples in the appendix. All core contributions of this paper, including the design of the Functional
723 Framework Prompt (FFP) dataset, the development of the Semantic Alignment Score (SAS) metric,
724 the proposal of the Alignment-Guided Fine-Tuning (AGFT) strategy, as well as dataset construction,
725 experiments, and analysis, were entirely conceived, implemented, and validated by the authors. The
726 authors take full responsibility for the content of this paper.
727

728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755