

From Knowledge Concepts to Process-Based Skills: A Feedback-Driven Ability Induction Framework for Student Modeling

Anonymous ACL submission

Abstract

Ability representations are central to student modeling because they connect item requirements to students’ observed responses. Most existing approaches use expert-defined knowledge concepts as ability units, which describe *what* content an item involves but often miss *how* it is solved when solving requires multi-step procedures. We study a process-based alternative by inducing ability units from solution procedures. Specifically, we treat large language model (LLM) reasoning traces as noisy process observations, and propose a closed-loop framework that extracts operation-level signals from structured LLM reasoning and consolidates them into stable, reusable process-based abilities. To reduce noise and keep induced abilities aligned with the target items, we combine semantic consolidation with a global feasibility signal from downstream student modeling, and use a TextGrad-style discrete controller to refine generation constraints and consolidation hyperparameters. Experiments on two benchmark datasets show that replacing expert concepts with induced process-based abilities improves student modeling performance and yields coherent ability structures with informative granularity and coverage.¹

1 Introduction

A central goal in intelligent education is to infer learners’ ability states from their observed learning interactions, so as to support the assessment of learning progress and personalized instructional interventions (Wang et al., 2020; Abdelrahman et al., 2023; Dong et al., 2025). To relate problem requirements to student performance, most student modeling pipelines rely on expert-defined knowledge concepts as their basic ability units. These concepts are typically curriculum-aligned labels and are widely used in cognitive diagnosis (CD) (Liu

¹Anonymous code repository: <https://anonymous.4open.science/status/ProcessAbilityInduction>

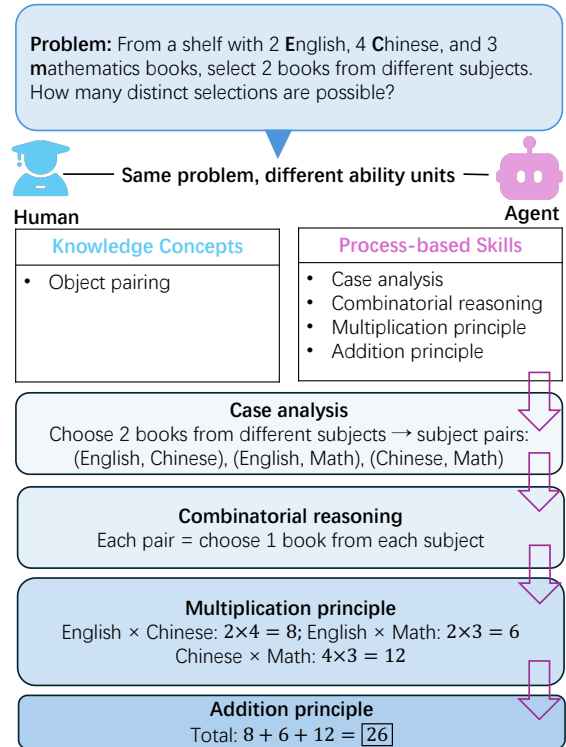


Figure 1: Comparison between knowledge-based and process-based ability representations on a multi-step combinatorial counting problem.

et al., 2023a) and knowledge tracing (KT) (Song et al., 2022; Scarlatos et al., 2025). However, such concept-based units characterize *what* content a problem assesses rather than *how* it is solved.

This limitation becomes particularly evident when problems require multi-step reasoning. As shown in Figure 1, a single problem can be decomposed into multiple reusable procedural operations, such as case analysis, combinatorial reasoning, and the application of the multiplication and addition principles, which are not explicitly captured by a single concept label. From the perspective of student modeling, representing all such procedures with a single concept-level ability hides important differences in how learners approach and solve the

056 problem. This observation motivates a complemen- 107
057 tary representation question: 108

058 *Can we automatically induce reusable,* 109
059 *process-based ability units that replace* 110
060 *expert-defined concepts while remaining* 111
061 *compatible with existing student model-* 112
062 *ing frameworks?* 113

063 Recent advances in Large Language Models 114
064 (LLMs) make process-based ability induction fea- 115
065 sible by explicitly externalizing intermediate rea- 116
066 soning steps in text, for example, through Chain-of- 117
067 Thought prompting (Wei et al., 2022; Liang et al., 118
068 2023). By exposing sequences of operations, in- 119
069 termediate states, and reasoning order, such traces 120
070 provide a rare opportunity to observe otherwise im- 121
071 plicit problem-solving processes. However, these 122
072 traces often exhibit substantial stylistic variability, 123
073 redundancy, and noise (Turpin et al., 2023; Tutek 124
074 et al., 2025). As a result, directly treating raw rea- 125
075 soning text as ability representations is unreliable. 126
076 Nevertheless, when aggregated across many prob- 127
077 lems, LLM-generated traces exhibit recurring struc- 128
078 tural regularities, such as repeated operation types 129
079 and reasoning patterns. We therefore treat reason- 130
080 ing traces as noisy but information-rich process- 131
081 level observations, which require explicit abstrac- 132
082 tion and consolidation to yield stable and reusable 133
083 procedural abilities. 134

084 Building on this perspective, we propose a frame- 135
085 work for process-based ability induction that explic- 136
086 itly separates process observation and semantic ab- 137
087 straction. Rather than operating on individual items 138
088 in isolation, our framework performs induction at 139
089 the dataset level. At each iteration, a frozen LLM 140
090 generates Structured Chain-of-Thought traces for 141
091 each item under explicit generation constraints, en- 142
092 suring that reasoning steps are annotated with nor- 143
093 malized operation labels. The operation-level sig- 144
094 nals extracted from all items are then aggregated 145
095 and consolidated into canonical procedural abili- 146
096 ties through semantic clustering, after which each 147
097 cluster is assigned a concise and reusable canon- 148
098 ical name, yielding a global item–skill mapping 149
099 that captures recurring procedural patterns shared 150
100 across the dataset. 151

101 Textual coherence alone does not guarantee that 152
102 induced skills are useful for student modeling. In 153
103 particular, assigning skills independently to indi- 154
104 vidual items can be misleading, as the effectiveness 155
105 of a skill representation depends on how skill as- 156
106 signments across all items jointly support student

107 modeling. To address this, we introduce a student 108
109 modeling-based feasibility evaluation that assesses 109
110 the quality of the induced item–skill mapping us- 110
111 ing validation performance on downstream student 111
112 modeling tasks (Liu et al., 2023a; Scarlatos et al., 112
113 2025). This feasibility evaluation produces a single 113
114 global signal based on validation performance of 114
115 downstream student models. The signal reflects 115
116 the quality of the induced ability representation at 116
117 the dataset level, since ability assignments across 117
118 different items must be considered jointly rather 118
119 than evaluated item by item. 119

120 Finally, to allow this external feedback to guide 120
121 induction, we adopt a TextGrad-style (Yuksekgonul 121
122 et al., 2025) discrete controller. The controller 122
123 translates feasibility feedback into updates of gen- 123
124 eration constraints and consolidation hyperparam- 124
125 eters, enabling the induced skill space to progres- 125
126 sively stabilize over multiple iterations. Through- 126
127 out this process, the language model remains frozen 127
128 and all improvements arise from feedback-driven 128
129 refinement of generation constraints and induction 129
130 hyperparameters. 130

131 In summary, our contributions are: 130

- 131 • We propose a feedback-driven framework that 131
132 automatically induces fine-grained, process- 132
133 based ability units from reasoning traces with- 133
134 out requiring manual skill design, leveraging 134
135 large language models to infer reusable skills, 135
136 and providing a more expressive alternative 136
137 to expert-defined knowledge concepts for stu- 137
138 dent modeling. 138
- 139 • We design a structured process observation 139
140 and induction pipeline that converts noisy lan- 140
141 guage model reasoning traces into reusable 141
142 process-based skills by enforcing explicit op- 142
143 eration labels and consolidating semantically 143
144 similar operations. 144
- 145 • We introduce a global feasibility signal from 145
146 student modeling as external feedback, and 146
147 develop a TextGrad-style discrete controller 147
148 to iteratively refine generation constraints and 148
149 induction hyperparameters. 149
- 150 • We conduct comprehensive evaluations on 150
151 two benchmark datasets, demonstrating con- 151
152 sistent improvements over human-defined 152
153 knowledge concepts on cognitive diagnosis 153
154 and knowledge tracing tasks, as well as more 154
155 informative and coherent skill-level ability 155
156 structures. 156

2 Related Work

2.1 Ability Representation in Student Modeling

Cognitive Diagnosis (CD) and Knowledge Tracing (KT) are two core paradigms in student modeling. Most existing approaches rely on a predefined and fixed set of expert-defined knowledge concepts as ability units to link problem requirements with students’ observed responses (Gao et al., 2022; Liu et al., 2019). Such concept-based representations are primarily content-oriented, describing *what* knowledge a problem involves rather than *how* it is solved.

A large body of prior work has focused on improving student modeling performance through model-level advancements, including more expressive architectures, inference mechanisms, and temporal dynamics (Li et al., 2025; Pandey and Karypis, 2019; Minn et al., 2018; Cui et al., 2023). These approaches typically assume the underlying ability units to be fixed and sufficient, leaving differences in problem-solving procedures to be implicitly absorbed by the model during training.

In contrast, our work shifts the focus from model optimization to ability representation. We investigate whether student abilities can be represented at a finer, process-based level by inducing ability units from language-model-generated reasoning traces. Rather than modifying CD or KT models, we use their predictive performance as evaluation signals to assess the quality of the induced ability representations.

2.2 Chain-of-Thought as Process-Level Signals

Chain-of-Thought (CoT) prompting (Wei et al., 2022) enables large language models to produce intermediate reasoning steps in natural language, making problem-solving processes observable. Although CoT traces are not faithful explanations of internal model computation and can exhibit noise and stylistic variability (Turpin et al., 2023), aggregated traces often reveal recurring structural patterns, such as repeated operations and reasoning motifs.

Most existing work leverages CoT to improve reasoning performance, provide post-hoc explanations, or supervise model training (Plaat et al., 2025). In contrast, we treat CoT as a source of process-level observations rather than a final representation. Rather than analyzing individual CoT

traces in isolation, we generate structured CoT traces for all items in the dataset and analyze them collectively at the system level. By aggregating and consolidating CoT-derived operations across multiple items, we aim to induce stable and reusable procedural abstractions that can serve as generalizable skill labels.

3 Methodology

3.1 Problem Setup and Overview

We consider an educational dataset consisting of assessment items (e.g., math word problems or programming exercises), denoted by \mathcal{I} . Our goal is to induce a reusable, process-based *skill space* that captures *how* items are solved, rather than only *what* content they assess. Formally, each item $i \in \mathcal{I}$ is associated with an induced skill set $\mathcal{S}_i \subseteq V$, where V is a global vocabulary of process-based skills.

A key challenge is that raw LLM reasoning traces are noisy and stylistically variable: the same procedure may appear under different surface forms, while overly aggressive merging can collapse distinct procedures. To address this, we employ an iterative closed-loop induction framework (Figure 2), in which the skill space is repeatedly refined through structured process observation, semantic consolidation, student modeling-based feasibility evaluation, and TextGrad-style discrete feedback updates.

System state. At induction round r , the system maintains the following components:

- a skill vocabulary $V^{(r)}$, defining the currently available operation labels;
- generation-side constraints $\mathcal{C}^{(r)}$, which comprise prompt instructions and structured schemas that govern all LLM-based generation steps in the induction process, including (i) structured process observation for producing step-wise reasoning traces, and (ii) semantic consolidation for assigning canonical names to operation clusters; these constraints encourage reuse of labels in $V^{(r)}$ while allowing NEW: labels only when necessary;
- consolidation hyperparameters $\Theta^{(r)}$, which control how extracted operation phrases are clustered and promoted into canonical skills.

We denote the complete system configuration at round r as $\text{State}^{(r)} = (V^{(r)}, \mathcal{C}^{(r)}, \Theta^{(r)})$.

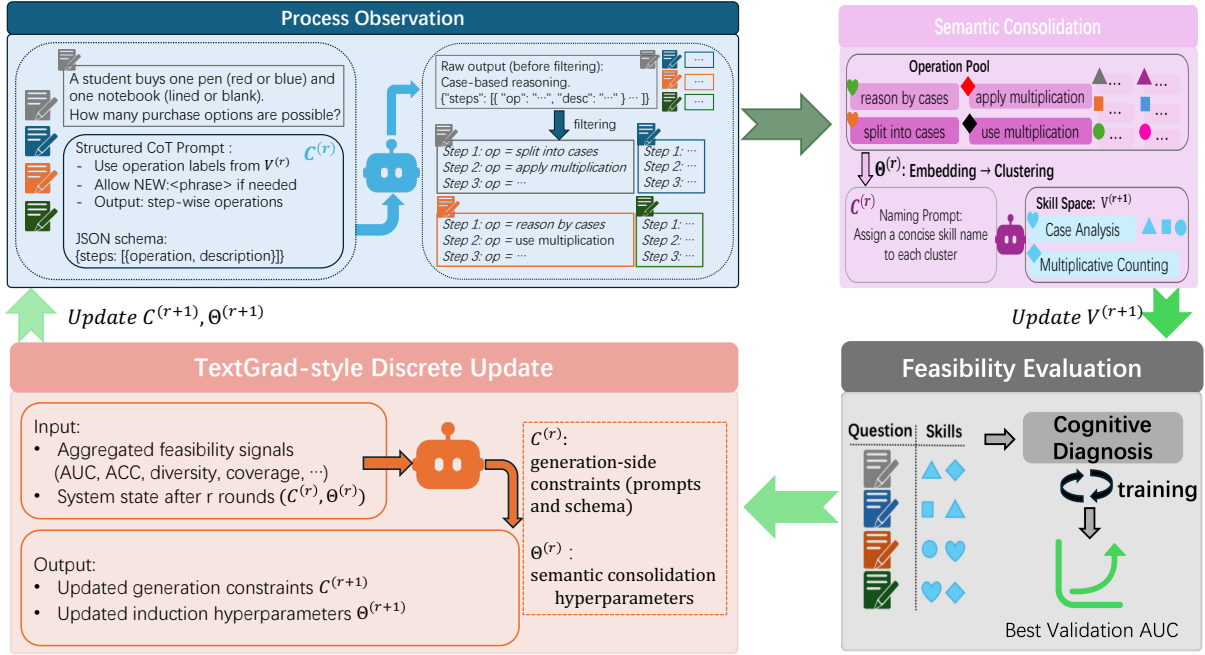


Figure 2: Overview of the proposed feedback-driven closed-loop framework for process-based ability induction.

Induction loop. Each induction round proceeds as follows. First, the LLM generates Structured Chain-of-Thought traces under $\mathcal{C}^{(r)}$ and $V^{(r)}$. Second, extracted operation phrases are semantically consolidated into canonical skills, yielding updated item-skill mappings and forming a new skill vocabulary $V^{(r+1)}$. Third, the induced skill space is evaluated as a whole using downstream student modeling performance as a feasibility signal. Finally, a TextGrad-style discrete controller invokes a frozen LLM to interpret the global feasibility feedback and to explicitly generate updated generation constraints $\mathcal{C}^{(r)}$ and consolidation hyperparameters $\Theta^{(r)}$, without any differentiable optimization.

3.2 Process Observation via Structured Chain-of-Thought

For each item $i \in \mathcal{I}$, we generate a *Structured Chain-of-Thought* (Structured CoT) using a frozen LLM under generation-side constraints $\mathcal{C}^{(r)}$:

$$\mathbf{c}_i^{(r)} = \text{LLM}(i; V^{(r)}, \mathcal{C}^{(r)}), \quad (1)$$

where the initial prompt and JSON schema for structured process observation are shown in Figure 3.

Structured CoT follows a fixed JSON schema that enforces step-wise reasoning with explicit operation annotations. Each step is labeled with an operation selected from the current vocabulary $V^{(r)}$, and a reserved *NEW:phrase* prefix is allowed when

no suitable label exists. This yields a machine-readable sequence of reasoning steps, enabling a consistent process observation across items.

From each trace, we extract a set of operation labels:

$$\mathcal{O}_i^{(r)} = g(\mathbf{c}_i^{(r)}), \quad (2)$$

where $g(\cdot)$ retains only step-level operation annotations. We treat $\mathcal{O}_i^{(r)}$ as a redundant and non-canonical process observation, rather than a finalized skill representation.

3.3 Semantic Consolidation into Canonical Skills

As illustrated in Figure 2, operation labels extracted from all items are aggregated into a global *operation pool*. Due to linguistic variability (e.g., “reason by cases” vs. “split into cases”), directly reusing these non-canonical phrases would lead to a fragmented skill space.

We perform semantic consolidation through three stages. First, each operation phrase is encoded into a semantic embedding using the pre-trained sentence encoder BGE-M3 (Chen et al., 2024). Second, operation embeddings are grouped via similarity-based clustering, governed by the consolidation hyperparameters

$$\Theta^{(r)} = \{\tau, m, u\}, \quad (3)$$

where τ denotes the similarity threshold, m the minimum cluster size, and u the minimum item-

level support required for a cluster to be promoted as a canonical skill. Third, for each stable cluster, we invoke a frozen LLM with a naming prompt to assign a concise and reusable canonical skill name. The initial naming prompt is provided in Appendix 4 and is later updated automatically as part of $\mathcal{C}^{(r)}$ through TextGrad-style discrete feedback. This step consolidates linguistically diverse operation phrases into a single canonical skill (e.g., *Case Analysis*).

Based on the resulting canonical skills, each item $i \in \mathcal{I}$ is assigned an induced skill set $\mathcal{S}_i^{(r)}$. The updated skill vocabulary is then constructed by aggregating induced skills across all items:

$$V^{(r+1)} = \bigcup_{i \in \mathcal{I}} \mathcal{S}_i^{(r)}. \quad (4)$$

3.4 Student Modeling-Based Feasibility Evaluation

To assess the dataset-level feasibility of the induced skill space as an ability representation, we introduce a student modeling-based feasibility evaluation, using student response prediction performance to estimate this global feasibility.

Concretely, we instantiate the evaluation using a representative cognitive diagnosis model, KaNCD (Wang et al., 2022), and quantify feasibility via validation AUC. This choice is not tied to KaNCD specifically, but reflects the general principle that a feasible ability representation should enable accurate prediction of student responses under standard student modeling frameworks. Accordingly, we define the global feasibility score at round r as:

$$\text{Score}^{(r)} = \text{AUC}_{\text{val}}. \quad (5)$$

A higher score indicates that the induced skill space, when considered as a whole, better supports modeling student response behavior at the dataset level.

We do not use knowledge tracing performance as a feedback signal during induction, as KT models are sequence-dependent and less sensitive to local item-skill structure. Using CD as the feasibility signal isolates the effect of skill representations from temporal modeling choices.

After completing all induction rounds, we select the round with the highest validation AUC and retain the corresponding skill space $V^{(r+1)}$ together with the item-skill assignments $\mathcal{S}^{(r)}$. The retained item-skill mapping is then fixed and used as the ability representation for downstream cognitive diagnosis and knowledge tracing models.

Table 1: Dataset statistics for the processed splits used in our experiments.

Statistic	XES3G5M	BEPKT
#Students	100	905
#Items	1055	551
#Interactions	5003	25853
#Concepts	285	100

3.5 TextGrad-Style Discrete Feedback Updates

To incorporate external feasibility feedback without relying on differentiable optimization, we adopt a TextGrad-style (Yuksekgonul et al., 2025) discrete update mechanism driven by a large language model.

Inputs. At iteration r , the update module takes the following as input: (i) feasibility signals aggregated across previous multiple update rounds (three rounds in practice), including validation AUC and auxiliary statistics (e.g., skill coverage and skill diversity), and (ii) the current system configuration ($\mathcal{C}^{(r)}, \Theta^{(r)}$).

Outputs. Given these inputs, the language model produces discrete proposals for: (i) updated generation-side constraints $\mathcal{C}^{(r+1)}$ (e.g., prompt-level instructions and schema rules), and (ii) updated semantic consolidation hyperparameters $\Theta^{(r+1)}$.

Formally, the update step is written as:

$$(\mathcal{C}^{(r+1)}, \Theta^{(r+1)}) = \Phi(\mathcal{C}^{(r)}, \Theta^{(r)}, \text{Score}^{(r)}), \quad (6)$$

where Φ denotes a language-model-driven discrete update operator. In practice, Φ is implemented by invoking a frozen LLM with TextGrad-style prompts to update generation constraints and induction hyperparameters, rather than performing differentiable optimization (see Appendix E).

4 Experiments

4.1 Experimental Setup

Datasets. We evaluate our framework on two educational datasets with expert-defined knowledge concepts: **XES3G5M** (Liu et al., 2023b), an online mathematics learning dataset, and **BEPKT** (Zhu et al., 2022), which contains programming learning trajectories. These datasets differ in domain, scale, and concept granularity, enabling evaluation

Table 2: Downstream performance on cognitive diagnosis and knowledge tracing using different skill representations. “Human Knowledge” denotes expert-provided knowledge concepts.

Model	XES3G5M				BEPKT			
	CD-AUC	CD-ACC	KT-AUC	KT-ACC	CD-AUC	CD-ACC	KT-AUC	KT-ACC
GPT Knowledge	0.6288	0.7360	0.6895	0.7409	0.6659	0.6558	0.6689	0.6335
Human Knowledge	0.6502	0.7240	0.6855	0.7490	0.6630	0.6551	0.6779	0.6395
Qwen3-0.6B	0.6768	0.7670	0.7013	0.7521	0.6641	0.6479	0.6617	0.6275
Qwen3-1.7B	0.6877	0.7671	0.6856	0.7382	0.6739	0.6487	0.6598	0.6220
Qwen3-4B	0.7008	0.7600	0.6998	0.7510	0.6723	0.6519	0.6906	0.6462
Qwen3-8B	0.6838	0.7561	0.7076	0.7531	0.6703	0.6602	0.6849	0.6458
Qwen3-14B	0.6883	0.7580	0.7041	0.7592	0.6793	0.6549	0.6782	0.6406
Qwen3-32B	0.7010	0.7622	0.7033	0.7520	0.6773	0.6623	0.6866	0.6423
Llama3-3B	0.6545	0.7610	0.6875	0.7452	0.6644	0.6514	0.6769	0.6343
Llama3-8B	0.6818	0.7690	0.6943	0.7516	0.6749	0.6591	0.6856	0.6442

across heterogeneous learning contexts. Both provide natural language problem descriptions, allowing large language models to generate reasoning traces. Dataset statistics are summarized in Table 1.

Skill representations and baselines. We consider two concept-based baselines: human expert-annotated knowledge concepts provided with each dataset, which represent the standard ability specification in cognitive diagnosis and knowledge tracing, and LLM-generated concept labels produced by GPT-4o from problem descriptions.

To evaluate the proposed closed-loop framework for process-driven skill induction, we instantiate the framework with frozen LLM backbones from two model families, Qwen3 (Yang et al., 2025) and Llama3 (Dubey et al., 2024), spanning multiple model sizes. Implementation details are provided in Appendix B.

Evaluation tasks and metrics. We evaluate each item-skill mapping from three complementary perspectives: task-level performance, structural properties, and expert-based qualitative assessment.

Task-level evaluation. We evaluate induced item-skill mappings on two downstream student modeling tasks: (i) Cognitive Diagnosis (CD) using KaNCD, and (ii) Knowledge Tracing (KT) using the classical DKT (Piech et al., 2015) model.

During induction, KaNCD is used as a feasibility evaluator, producing a validation score for each candidate item-skill mapping. This score is used solely to select the mapping that best supports student modeling, rather than to tune or compare different student models. After the induction process is completed, the selected item-skill mapping

is fixed. We then evaluate its downstream performance on the held-out test set using KaNCD for cognitive diagnosis and DKT for knowledge tracing, and report CD-AUC, CD-ACC, KT-AUC, and KT-ACC.

Structure-level evaluation. To characterize induced ability spaces beyond predictive performance, we report a set of structural statistics, including: (i) the total number of induced skills (#Skills); (ii) skill granularity measured by the average and 90th-percentile number of skills per item (Avg./Item and P90); (iii) coverage concentration measured by the Gini coefficient over skill frequencies (Gini) and the fraction of assignments covered by the 20 most frequent skills (Top-20). Formal definitions of the structure-level metrics are provided in Appendix A.

Expert-based evaluation. As a complementary evaluation, we conduct a questionnaire study with domain experts to qualitatively assess the induced skill representations.

Training protocol and reproducibility. All downstream models are trained using fixed architectures and identical data splits across skill representations, with the same optimization and hyperparameter selection procedure to ensure fair comparison. For reproducibility, the full inference settings for LLM generation, as well as the complete training configurations for KaNCD and DKT, are provided in Appendix B.

4.2 Main Results on Cognitive Diagnosis and Knowledge Tracing

Cognitive diagnosis performance. As shown in Table 2, process-induced skills consistently outper-

Table 3: Skill structure statistics on XES3G5M and BEPKT.

Model	XES3G5M					BEPKT				
	#Skills	Avg./Item	P90	Gini	Top-20	#Skills	Avg./Item	P90	Gini	Top-20
GPT Knowledge	291	1.28	2	0.64	0.48	511	2.8	4	0.55	0.32
Human Knowledge	285	1.10	1	0.54	0.35	100	1.76	3	0.70	0.78
Qwen3-0.6B	72	1.70	3	0.75	0.86	58	1.63	3	0.80	0.92
Qwen3-1.7B	157	2.60	4	0.79	0.75	86	4.56	6	0.82	0.92
Qwen3-4B	175	2.60	4	0.70	0.60	200	3.82	6	0.66	0.54
Qwen3-8B	213	2.70	4	0.72	0.60	144	4.79	7	0.74	0.72
Qwen3-14B	167	3.00	5	0.70	0.60	182	5.15	7	0.70	0.57
Qwen3-32B	223	3.50	5	0.67	0.50	161	5.59	8	0.70	0.64
Llama3-3B	149	2.30	4	0.69	0.59	130	2.82	4	0.69	0.70
Llama3-8B	164	4.30	7	0.71	0.64	204	4.71	7	0.71	0.59

form expert-defined knowledge concepts on cognitive diagnosis across both datasets. On XES3G5M, CD-AUC improves from 0.6502 using expert concepts to around 0.70 for several Qwen3 variants. Similar improvements are observed on BEPKT, where multiple medium-to-large models achieve the strongest CD performance. These gains are achieved without increasing the complexity of the model, suggesting that the improvement stems from more informative skill representations rather than model capacity.

Knowledge tracing performance. As shown in Table 2, process-induced skills achieve performance that is competitive with or better than the knowledge-based baseline on knowledge tracing across both datasets. On XES3G5M, KT-AUC increases from 0.6855 to over 0.70 with skills induced by several Qwen3 models, while performance on BEPKT remains stable or improves slightly. Notably, despite introducing higher-granularity skills, performance does not degrade on knowledge tracing. Importantly, even though the skills are induced solely under a cognitive diagnosis-based feedback signal, they yield consistent gains on knowledge tracing, demonstrating that the induced representations capture transferable procedural structure beyond the specific optimization objective.

4.3 Skill Structure Analysis

Skill granularity. As shown in Table 3, human-defined knowledge concepts are relatively coarse-grained, with each item associated with few skills on average. In contrast, process-induced skills exhibit substantially higher granularity, as reflected by a larger number of unique skills and higher per-

item skill counts. For example, on BEPKT, the average number of skills per item increases from 1.76 under expert concepts to over 5 for several induced settings. This suggests that the induced representations decompose problems into multiple reusable procedural components, rather than collapsing them into monolithic content-level labels.

Coverage concentration and balance. As shown in Table 3, we further examine how skill usage is distributed across items using the Gini coefficient and Top-20 coverage. Highly concentrated distributions indicate that a small number of generic skills dominate many assignments. For instance, Qwen3-0.6B shows extremely high Top-20 coverage, suggesting over-generalized skills. In contrast, medium-capacity models yield more balanced distributions, with lower concentration and less dominance of frequent skills, indicating a more expressive and evenly utilized skill space.

4.4 Ablation Study

Role of multi-round feedback. As shown in Table 4, multi-round feedback consistently improves both task performance and skill structure compared to one-round induction. For example, with Qwen3-8B on XES3G5M, CD-AUC increases from 0.6678 to 0.6838 and KT-AUC increases from 0.7031 to 0.7076, while the average number of skills per item increases from 2.5 to 2.7. Similar trends are observed on BEPKT and with the Llama3-8B model. These results suggest that iterative feedback enables more stable and refined skill induction than single-pass approaches.

Role of feasibility evaluation (KaNCD). The setting “w/o KaNCD” removes the cognitive di-

Table 4: Ablation results on XES3G5M and BEPKT. “Multi rounds” denotes the full closed-loop induction, while “One round” disables iterative feedback updates.

Model	Setting	XES3G5M					BEPKT				
		CD-AUC	CD-ACC	KT-AUC	KT-ACC	Avg./Item	CD-AUC	CD-ACC	KT-AUC	KT-ACC	Avg./Item
Qwen3-8B	One round	0.6678	0.7403	0.7031	0.7403	2.5	0.6632	0.6581	0.6811	0.6387	3.1
	Multi rounds	0.6838	0.7561	0.7076	0.7531	2.7	0.6703	0.6602	0.6849	0.6458	4.79
	w/o KaNCD	0.6725	0.7431	0.6882	0.7392	2.6	0.6667	0.6541	0.6814	0.6425	3.6
	w/o BGE-M3	0.6698	0.7483	0.6917	0.7401	2.6	0.6690	0.6409	0.6813	0.6398	3.3
	w/o TextGrad	0.6794	0.7509	0.6875	0.7382	2.3	0.6641	0.6498	0.6809	0.6440	3.9
Llama3-8B	One round	0.6733	0.7620	0.6890	0.7461	3.5	0.6677	0.6495	0.6743	0.6293	3.4
	Multi rounds	0.6818	0.7690	0.6943	0.7516	4.3	0.6749	0.6591	0.6856	0.6382	4.71
	w/o KaNCD	0.6744	0.7620	0.6292	0.6959	3.3	0.6736	0.6453	0.6786	0.6332	3.6
	w/o BGE-M3	0.6723	0.7610	0.6827	0.7507	3.4	0.6722	0.6445	0.6753	0.6316	3.5
	w/o TextGrad	0.6704	0.7566	0.5602	0.6649	3.8	0.6748	0.6492	0.6747	0.6259	4.0

agnosis model from the induction loop, thereby disabling feasibility-based feedback during skill induction. As shown in Table 4, this leads to degraded performance on both datasets, indicating that without student modeling-based evaluation, the induction process is less effective at selecting and refining useful skill representations. These results highlight the importance of incorporating downstream feasibility signals to guide process-based ability induction.

Role of semantic consolidation. The setting “w/o BGE-M3” removes semantic clustering of generated operation phrases, preventing linguistically similar operations from being consolidated into canonical skills. As shown in Table 4, this leads to degraded performance and less stable skill structures. For example, with Qwen3-8B on XES3G5M, CD-AUC drops from 0.6838 to 0.6698 when BGE-M3 is removed. This indicates that semantic consolidation is crucial for reducing surface-level variability and enabling consistent reuse of procedural skills.

Role of discrete feedback control. The setting “w/o TextGrad” disables discrete feedback updates. As a result, the induction process cannot adapt to structural imbalances, such as over-fragmentation or over-concentration of skills. This leads to unstable skill spaces, reflected by lower performance and inconsistent average skill counts across settings. These results highlight the importance of feedback-driven hyperparameter updates for stabilizing the closed-loop induction process.

4.5 Expert Evaluation of Ability Representations

Ten education experts evaluated ten questions each across three evaluation aspects, resulting in 300 judgments comparing LLM-generated and expert-defined ability representations. Experts preferred the LLM-generated skills in 82.3% of cases, with 10.7% favoring expert-defined skills and 7% indicating comparable quality. Detailed analyses are provided in Appendix D.

5 Conclusion

In this paper, we propose a feedback-driven framework for inducing process-based ability representations from problem-solving traces, with the goal of moving beyond expert-defined knowledge concepts to discover reusable procedural skills that reflect how problems are solved. Specifically, we collect reasoning traces generated by large language models for all items in a dataset and treat them as process-level observations that exhibit substantial surface-level redundancy, where multiple expressions may correspond to the same underlying procedural operation. These traces are progressively refined through structured generation, semantic consolidation, and external feasibility feedback, resulting in stable and behaviorally meaningful dataset-level ability spaces. Experiments on two real-world educational datasets demonstrate consistent improvements in student modeling and reveal informative ability structures in terms of granularity and coverage, suggesting that process-based ability induction offers a practical and scalable alternative to traditional knowledge-based representations.

600 Limitations

601 A current limitation of this work is that our em-
602 pirical evaluation is restricted to mathematics and
603 programming domains, using datasets that primar-
604 ily reflect learning behavior within a single grade
605 level. While these settings allow us to validate the
606 effectiveness of process-based ability induction in
607 controlled educational scenarios, student abilities
608 often transfer across different knowledge areas and
609 evolve over multiple grade levels. As a result, the
610 extent to which the induced process-based skills
611 generalize across subjects and support longitudinal
612 modeling of student learning remains an open ques-
613 tion. In future work, we plan to extend our frame-
614 work to broader educational contexts by collecting
615 and incorporating datasets spanning multiple do-
616 mains and grade levels, in order to better assess
617 its applicability to long-term and cross-curricular
618 student modeling.

619 Ethical Considerations

620 This work studies automatic induction of process-
621 based ability representations from educational
622 problem-solving data. All experiments are con-
623 ducted on publicly available datasets that have
624 been anonymized and do not contain personally
625 identifiable information. We do not attempt to
626 infer sensitive attributes or re-identify individual
627 students, and our analysis is performed at the ag-
628 gregate level. The induced ability representations
629 are intended for research purposes, such as ana-
630 lyzing problem-solving processes and supporting
631 student modeling, rather than for high-stakes ed-
632 ucational decision making. As with other data-
633 driven educational models, the induced skills may
634 reflect biases present in the underlying datasets,
635 and their interpretations should therefore be treated
636 with caution, particularly when applied to new pop-
637 ulations or contexts. Large language models are
638 used only to generate intermediate reasoning traces
639 from problem descriptions. These traces are treated
640 as process-level observations and are not assumed
641 to be factually correct or pedagogically optimal.
642 The language models are not used for grading, de-
643 cision making, or generating the manuscript text it-
644 self. Finally, while process-based ability induction
645 has the potential to support more fine-grained anal-
646 ysis of learning behavior, we emphasize that any
647 real-world deployment should involve appropriate
648 human oversight and consideration of educational,
649 social, and ethical implications.

References

- Ghodai Abdelrahman, Qing Wang, and Bernardo Nunes. 2023. Knowledge tracing: A survey. *ACM Computing Surveys*, 55(11):1–37. 651–653
- Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2024. Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation. *arXiv preprint arXiv:2402.03216*. 654–658
- Jiajun Cui, Zeyuan Chen, Aimin Zhou, Jianyong Wang, and Wei Zhang. 2023. Fine-grained interaction modeling with multi-relational transformer for knowledge tracing. *ACM Transactions on Information Systems*, 41(4):1–26. 659–663
- Zhiang Dong, Jingyuan Chen, and Fei Wu. 2025. Knowledge is power: Harnessing large language models for enhanced cognitive diagnosis. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 164–172. 664–668
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*. 669–673
- Lina Gao, Zhongying Zhao, Chao Li, Jianli Zhao, and Qingtian Zeng. 2022. Deep cognitive diagnosis model for predicting students’ performance. *Future Generation Computer Systems*, 126:252–262. 674–677
- Xiaoyu Li, Shaoyang Guo, Jin Wu, and Chanjin Zheng. 2025. An interpretable polytomous cognitive diagnosis framework for predicting examinee performance. *Information Processing & Management*, 62(1):103913. 678–682
- Yuanyuan Liang, Jianing Wang, Hanlun Zhu, Lei Wang, Weining Qian, and Yunshi Lan. 2023. Prompting large language models with chain-of-thought for few-shot knowledge base question generation. *arXiv preprint arXiv:2310.08395*. 683–687
- Qi Liu, Zhenya Huang, Yu Yin, Enhong Chen, Hui Xiong, Yu Su, and Guoping Hu. 2019. Ekt: Exercise-aware knowledge tracing for student performance prediction. *IEEE Transactions on Knowledge and Data Engineering*, 33(1):100–115. 688–692
- Yingjie Liu, Tiancheng Zhang, Xuecen Wang, Ge Yu, and Tao Li. 2023a. New development of cognitive diagnosis models. *Frontiers of Computer Science*, 17(1):171604. 693–696
- Zitao Liu, Qiongqiong Liu, Teng Guo, Jiahao Chen, Shuyan Huang, Xiangyu Zhao, Jiliang Tang, Weiqi Luo, and Jian Weng. 2023b. Xes3g5m: A knowledge tracing benchmark dataset with auxiliary information. *Advances in Neural Information Processing Systems*, 36:32958–32970. 697–702

703	Sein Minn, Yi Yu, Michel C Desmarais, Feida Zhu, and Jill-Jenn Vie. 2018. Deep knowledge tracing and dynamic student classification for knowledge tracing. In <i>2018 IEEE International conference on data mining (ICDM)</i> , pages 1182–1187. IEEE.	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	756 757 758 759 760
708	Shalini Pandey and George Karypis. 2019. A self-attentive model for knowledge tracing. <i>arXiv preprint arXiv:1907.06837</i> .	Mert Yuksekogonul, Federico Bianchi, Joseph Boen, Sheng Liu, Pan Lu, Zhi Huang, Carlos Guestrin, and James Zou. 2025. Optimizing generative ai by backpropagating language model feedback. <i>Nature</i> , 639:609–616.	761 762 763 764 765
711	Chris Piech, Jonathan Bassen, Jonathan Huang, Surya Ganguli, Mehran Sahami, Leonidas J Guibas, and Jascha Sohl-Dickstein. 2015. Deep knowledge tracing. <i>Advances in neural information processing systems</i> , 28.	Renyu Zhu, Dongxiang Zhang, Chengcheng Han, Ming Gaol, Xuesong Lu, Weining Qian, and Aoying Zhou. 2022. Programming knowledge tracing: A comprehensive dataset and a new model. In <i>2022 IEEE International Conference on Data Mining Workshops (ICDMW)</i> , pages 298–307. IEEE.	766 767 768 769 770 771
716	Aske Plaat, Annie Wong, Suzan Verberne, Joost Broekens, and Niki Van Stein. 2025. Multi-step reasoning with large language models, a survey. <i>ACM Computing Surveys</i> .	A Additional Details of Structure-Level Metrics	772 773
720	Alexander Scarlato, Ryan S Baker, and Andrew Lan. 2025. Exploring knowledge tracing in tutor-student dialogues using llms. In <i>Proceedings of the 15th International Learning Analytics and Knowledge Conference</i> , pages 249–259.	This appendix provides formal definitions for the structure-level metrics reported in Table 3. Let \mathcal{I} denote the item set with $ \mathcal{I} = N$. For each item $i \in \mathcal{I}$, let \mathcal{S}_i be its induced skill set. Let V be the induced skill vocabulary and $ V = K$.	774 775 776 777 778
725	Xiangyu Song, Jianxin Li, Taotao Cai, Shuiqiao Yang, Tingting Yang, and Chengfei Liu. 2022. A survey on deep learning based knowledge tracing. <i>Knowledge-Based Systems</i> , 258:110036.	(i) Number of unique skills. We report the vocabulary size:	779 780
729	Miles Turpin, Julian Michael, Ethan Perez, and Samuel Bowman. 2023. Language models don’t always say what they think: Unfaithful explanations in chain-of-thought prompting. <i>Advances in Neural Information Processing Systems</i> , 36:74952–74965.	$\#\text{Skills} = V . \quad (7)$	781
734	Martin Tutek, Fateme Hashemi Chaleshtori, Ana Marasović, and Yonatan Belinkov. 2025. Measuring chain of thought faithfulness by unlearning reasoning steps. In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 9946–9971.	(ii) Granularity (Avg./Item and P90). Define the per-item skill count:	782 783
740	Fei Wang, Qi Liu, Enhong Chen, Zhenya Huang, Yuying Chen, Yu Yin, Zai Huang, and Shijin Wang. 2020. Neural cognitive diagnosis for intelligent education systems. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 34, pages 6153–6161.	$c_i = \mathcal{S}_i . \quad (8)$	784
745	Fei Wang, Qi Liu, Enhong Chen, Zhenya Huang, Yu Yin, Shijin Wang, and Yu Su. 2022. Neuralcd: a general framework for cognitive diagnosis. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 35(8):8312–8327.	Then the average number of skills per item is:	785
750	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	$\text{Avg./Item} = \frac{1}{N} \sum_{i \in \mathcal{I}} c_i. \quad (9)$	786
		Let $\{c_i\}_{i \in \mathcal{I}}$ sorted in non-decreasing order be $c_{(1)} \leq \dots \leq c_{(N)}$. The 90th percentile is:	787 788
		$\text{P90} = c_{(\lceil 0.9N \rceil)}. \quad (10)$	789
		(iii) Coverage concentration (Gini and Top-20 share). We define skill assignment frequency over the entire item set. For each skill $s \in V$, let	790 791 792 793
		$f_s = \sum_{i \in \mathcal{I}} \mathbf{1}\{s \in \mathcal{S}_i\}, \quad (11)$	794
		and let the total number of item-skill assignments be	795 796
		$F = \sum_{s \in V} f_s = \sum_{i \in \mathcal{I}} \mathcal{S}_i . \quad (12)$	797

Gini coefficient. Let $\{f_s\}_{s \in V}$ sorted be $f_{(1)} \leq \dots \leq f_{(K)}$, and let $\mu = \frac{1}{K} \sum_{k=1}^K f_{(k)} = \frac{F}{K}$. We compute the Gini coefficient over skill frequencies as:

$$\text{Gini} = \frac{\sum_{k=1}^K \sum_{\ell=1}^K |f_{(k)} - f_{(\ell)}|}{2K^2\mu}. \quad (13)$$

A larger Gini indicates that a small number of skills dominate the assignments.

Top-20 share. Let $f_{(1)}^\downarrow \geq \dots \geq f_{(K)}^\downarrow$ denote frequencies sorted in descending order. We define:

$$\text{Top-20} = \frac{\sum_{k=1}^{\min(20, K)} f_{(k)}^\downarrow}{F}. \quad (14)$$

A larger Top-20 implies more concentrated coverage on a small set of frequent skills.

B Implementation Details and Experimental Configurations

LLM backbones and decoding settings. We instantiate the induction framework with instruction-tuned models from the Qwen3 and Llama3 families, covering a wide range of model scales. For Qwen3, we use models ranging from 0.6B to 32B parameters with a decoding temperature of 0.7, while for Llama3 we use 3B and 8B instruction-tuned variants with a decoding temperature of 0.3.

Datasets and splits. For each dataset, we split the interaction logs into training/validation/test sets with a ratio of 7:1:2.

Downstream student models. We select KaNCD as the representative cognitive diagnosis model and DKT as the representative knowledge tracing model. For KaNCD, we perform a grid search over the learning rate and embedding dimension, with the learning rate selected from 0.001, 0.002, 0.01, 0.02 and the embedding dimension from 8, 16, 32, 64, 128, 256, 512. For DKT, we conduct a grid search over the dropout rate and learning rate, where dropout is chosen from 0.0, 0.1, 0.2, 0.3 and the learning rate from $1e-4$, $3e-4$, $1e-3$, $3e-3$. All experiments are repeated with five different random seeds, and the reported results are averaged across runs.

Computational resources. All experiments are conducted on a server equipped with 30 CPU cores, 400GB system memory, and two NVIDIA H100 GPUs (80GB each). The Qwen3 and Llama3 family models are directly deployed on the server for inference during the induction process. The

downstream student models, including KaNCD and DKT, are lightweight and run with minimal GPU usage. For ChatGPT-4o, we access the model via the official API.

C Case Study: Process-Induced skills vs. Human Concepts

We present two illustrative examples to qualitatively compare process-induced skills with expert-defined knowledge concepts. Each example is shown as a boxed item for clarity.

Example 1: Outfit Combination

Problem. Xiaomei bought two tops and two skirts at a supermarket. How many different outfits can she make if she wears exactly one top and one skirt?

LLM-generated reasoning steps.

```
{
  "steps": [
    {
      "operation":
        "reason_by_cases",
      "description": "Decompose
        the selection into
        choosing a top and
        choosing a skirt"
    },
    {
      "operation": "enumeration",
      "description": "Count the
        number of available tops
        and skirts independently"
    },
    {
      "operation":
        "multiplication_rule",
      "description": "Combine
        independent choices by
        multiplying their counts"
    }
  ]
}
```

Induced (process-based) skills.

- Case-based decomposition
- Enumeration of possibilities
- Multiplication-based composition

Expert-provided concept.

- Item matching / pairing

Example 2: Factoring with a Common Factor

Problem. Compute:

$$178 \times 15 + 15 \times 22.$$

LLM-generated reasoning steps.

```
{
  "steps": [
    {
      "operation":
        "pattern_detection",
      "description": "Identify a
        shared multiplicative
        factor across terms"
    },
    {
      "operation":
        "factor_extraction",
      "description": "Rewrite the
        expression by factoring
        out the common term"
    },
    {
      "operation": "addition",
      "description": "Simplify
        the expression inside
        the parentheses"
    },
    {
      "operation":
        "multiplication",
      "description": "Compute the
        final result after
        simplification"
    }
  ]
}
```

Induced (process-based) skills.

- Detection of shared factors
- Algebraic transformation via factor extraction
- Basic arithmetic operations

Expert-provided concept.

- Integer construction and extraction

derived from curriculum-aligned routes. Experts were asked to provide judgments along four dimensions:

- **Q1 (Representational Adequacy):** Which ability representation more appropriately captures the key abilities required to solve the question?
- **Q2 (Process Fidelity):** Which representation better reflects learners’ actual problem-solving processes, rather than merely the involved content?
- **Q3 (Diagnostic Utility):** Which representation would be more informative for analyzing learner errors or ability differences?
- **Q4 (Issue Identification):** Does the LLM-generated ability representation exhibit any noticeable issues?

No predefined gold standard was assumed; all evaluations were based on expert judgment.

D.2 Quantitative Summary of Expert Judgments

For the preference-based questions (Q1–Q3), expert responses were aggregated as *LLM* (preference for LLM-generated abilities), *EXP* (preference for expert-defined abilities), or *SAME* (comparable quality or no clear preference).

In total, we collected 300 expert judgments (10 experts × 10 questions × 3 evaluation aspects). Table 5 summarizes the aggregated results.

Overall, experts preferred the LLM-generated process-based abilities in 82.3% of the judgments, compared to 10.7% favoring expert-defined representations, with 7.0% indicating comparable quality. This preference is consistent across representational adequacy, process fidelity, and diagnostic utility, and is strongest for diagnostic utility.

D.3 Issue Analysis of LLM-Generated skills

Experts further assessed whether the LLM-generated skills exhibited noticeable issues (Q4). The distribution of identified issue types is reported in Table 6.

As shown in Table 6, in 70% of the evaluations, experts judged the LLM-generated skills to be reasonable and free of obvious issues. When issues were reported, they mainly concerned granularity mismatches, including skills perceived as overly

D Expert Evaluation of Ability Representations

D.1 Evaluation Protocol

Ten education experts participated in the evaluation. Each expert independently evaluated ten questions.

For each question, experts were presented with the problem text, a set of process-based skills induced from structured LLM reasoning traces, and a set of expert-defined knowledge-based abilities

coarse (9%) or overly fine-grained (9%). Other issues, such as the presence of irrelevant skills (8%) or semantic redundancy (4%), were relatively infrequent. Severe semantic errors were rarely observed, suggesting that most induced skills are meaningful and usable in practice.

Question	LLM	EXP	SAME	Total
Q1 (Adequacy)	78	13	9	100
Q2 (Process)	83	10	7	100
Q3 (Diagnostic)	86	9	5	100
All (Q1-Q3)	247	32	21	300

Table 5: Expert preferences between LLM-generated process-based skills and expert-defined knowledge-based representations across three evaluation aspects (Q1-Q3).

Issue Type	Count
Reasonable	70
Too coarse	9
Too fine	9
Irrelevant	8
Redundant	4

Table 6: Distribution of expert-identified issue types for LLM-generated skills (Q4).

E Prompt Templates for Process Observation and Feedback

To support reproducibility and clarify the interfaces used in our framework, we provide representative prompt templates for both process observation and feedback-driven updates. These prompts define structured, discrete interaction protocols with the language model, rather than trainable components. Throughout the induction process, all language model parameters remain fixed, and the prompts serve as the primary mechanism for controlling generation, abstraction, and feedback incorporation.

Specifically, Figure 3 shows the Structured Chain-of-Thought prompt used for process observation, which elicits step-wise reasoning with explicit operation annotations. This prompt may be further refined through feedback-driven updates (see Figure 6), while its core output interface remains unchanged. Figure 4 presents the prompt for symbolic naming of consolidated operation clusters, which affects only the canonical identifiers of induced skills and does not influence clustering or promotion decisions; this naming prompt

may likewise undergo minor refinements during feedback-driven updates. Figure 5 illustrates a representative TextGrad-style controller prompt used to update semantic consolidation hyperparameters based on global feasibility signals. In later update rounds, this controller may lead to refinements of generation-side constraints, such as modifying the Structured Chain-of-Thought prompt in Figure 3 to remove rules that allow the introduction of new operation labels (e.g., NEW: <short phrase>), thereby encouraging stricter reuse of an existing operation vocabulary. Finally, Figure 6 shows the prompt used to refine reasoning analysis and semantic consolidation templates in a feedback-driven manner. Through this mechanism, semantic consolidation hyperparameters may be adjusted across iterations; for example, the similarity threshold τ in $\Theta^{(r+1)}$ can change from an initial value of 0.72 to 0.68 in later stages to promote more stable consolidation under downstream feasibility feedback.

Together, these prompts specify the complete discrete control surface of the proposed framework, enabling reproduction of the induction process without modifying or fine-tuning the underlying language models and without manual rewriting of the prompt templates.

System role.

You are an expert problem-solving analyst.

Task.

Given a question, generate a structured reasoning trace in JSON format. The goal is to explicitly expose the problem-solving process using reusable operations.

Output format requirements.

- Output must be valid JSON.
- The JSON object must contain:
 - "steps": a list of reasoning steps
 - "reasoning_types": a list of operation labels used
 - "steps_count": the total number of steps

Step schema.

Each element in "steps" must include:

- "id": step index (starting from 1)
- "operation": operation used in this step
- "description": brief description of the step
- "evidence": supporting justification or intermediate result (optional)

Operation constraints.

- "operation" must be selected from a predefined vocabulary.
- If no suitable label exists, a new label may be introduced as NEW:<short phrase>.
- New labels should be avoided unless necessary; reuse is preferred.

Additional constraints.

- Reasoning should be correct, concise, and logically coherent.
- The output must contain JSON only, with no extra text.

Instruction.

Now solve the given question accordingly.

Figure 3: Structured Chain-of-Thought prompt used for process observation.

Role.

You are given a set of operation phrases that share similar procedural meaning.

Task.

Assign a concise canonical ability name that captures the shared procedure.

Input.

- A list of operation phrases representing a semantic cluster

Requirements.

- Provide a short, general name suitable as a reusable ability identifier.
- Do not consider frequency, clustering decisions, or downstream performance.

Output format.

The output must be a JSON object with the following fields:

- "name": canonical ability name
- "description": brief explanation of the ability

Figure 4: Prompt used for symbolic naming of induced skills. This step affects only identifiers and does not influence clustering or promotion decisions.

Role.

You are a controller responsible for adjusting discrete hyperparameters in the ability induction process based on external evaluation feedback.

Inputs (aggregated from recent rounds).

- Validation performance history (AUC, ACC)
- Induced ability count history
- Coverage statistics history
- Diversity statistics history

Task.

Analyze the current state of the induced ability space and determine whether it is:

- over-fragmented (too many overly fine-grained skills),
- over-concentrated (too few overly general skills), or
- reasonably balanced.

Based on this diagnosis, propose updated values for the predefined hyperparameters:

- tau: clustering similarity threshold
- m: minimum cluster size
- u: minimum item-level support

Constraints.

- Output only a JSON object with numeric values for tau, m, and u.
- All values must remain within predefined bounds.
- Do not introduce new parameters or modify the schema.

Output format.

- "tau": updated similarity threshold
- "m": updated minimum cluster size
- "u": updated minimum item-level support

Figure 5: A representative TextGrad-style controller prompt for discrete hyperparameter updates based on global feasibility signals.

Role.

You are a controller responsible for refining prompt templates used in the ability induction pipeline based on evaluation feedback.

Inputs.

- Current analyst prompt (system and user prefix)
- Current taxonomist prompt for cluster naming
- Aggregated evaluation feedback from recent rounds

Task.

Propose small, local refinements to the existing prompt templates to improve reasoning analysis and semantic consolidation. If no modification is necessary, return the original templates unchanged.

Constraints.

- Output only a JSON object.
- Preserve the original prompt structure.
- Do not introduce new prompt components.

Output format.

The output JSON must contain the following fields:

- "analyst.system": updated system prompt
- "taxonomist.name_cluster_prompt": updated cluster naming prompt

Figure 6: A representative TextGrad-style controller prompt for feedback-driven refinement of reasoning analysis and semantic consolidation prompts.