# Agent Skill Framework: Perspectives on the Potential of Small Language Models in Industrial Environments

**Yangjie Xu[♠1], Lujun Li[♠1], Lama Sleem[1], Niccolo' Gentile[2],**
**Yewei Song[1], Yiqun Wang[1], Siming Ji[3], Wenbo Wu[4], Radu State[1],**
[1]University of Luxembourg, [2]Foyer S.A., [3]Princeton University, [4]Université Paris-Saclay

## Abstract

Agent Skill framework, now widely and officially supported by major players such as GitHub Copilot, LangChain, and OpenAI, performs especially well with proprietary models by improving context engineering, reducing hallucinations, and boosting task accuracy. Based on these observations, an investigation is conducted to determine whether the Agent Skill paradigm provides similar benefits to small language models (SLMs). This question matters in industrial scenarios where continuous reliance on public APIs is infeasible due to data-security and budget constraints requirements, and where SLMs often show limited generalization in highly customized scenarios. This work introduces a formal mathematical definition of the Agent Skill process, followed by a systematic evaluation of language models of varying sizes across multiple use cases. The evaluation encompasses two open-source tasks and a real-world insurance claims data set. The results show that tiny models struggle with reliable skill selection, while moderately sized SLMs (approximately 12B–30B parameters) benefit substantially from the Agent Skill approach. Moreover, code-specialized variants at around 80B parameters achieve performance comparable to closed-source baselines while improving GPU efficiency. Collectively, these findings provide a comprehensive and nuanced characterization of the capabilities and constraints of the framework, while providing actionable insights for the effective deployment of Agent Skills in SLM-centered environments[1].

## 1 Introduction

**Agent Skills Concept** Agent skills, initially introduced by a leading code-oriented large language model (LLM) unicorn[2], are widely regarded as an effective framework for agent-centric context engineering. This approach can be viewed as a well-designed "static cheat-sheet" in contrast to a "dynamic cheat-sheet" (Suzgun et al., 2025), enabling LLMs to focus on salient information through progressive context management. This approach substantially optimizes both the length and quality of context windows, enabling LLMs and agents to markedly reduce hallucinations while improving their use of tools, documents, and external knowledge resources. Unlike retrieval-augmented generation (RAG) (Lewis et al., 2020) systems—which rigidly encode textual or visual data into vector-space databases—this design directly leverages the in-context learning (ICL) capabilities of LLMs (Dong et al., 2024a; Brown et al., 2020a) and emergent reasoning to dynamically select the most pertinent information, instructions, and contextual knowledge. Furthermore, its progressive disclosure–oriented design maintains a tightly bounded effective context length(Mei et al., 2025), which, in turn, substantially increases the success rates of skill and tool selection.

**Huge LLMs Dominate** However, these paradigms rely heavily on the smartness of huge-size foundation models, such as ChatGPT or the Claude model series, and are primarily tailored to models specifically optimized for code-related tasks(Zhang et al., 2024a). Most of these models are closed-source and require API calls for interaction, while even their relatively open-source counterparts demand substantial GPU resources to maintain low latency. Moreover, invoking external APIs in industrial applications inevitably raises concerns about data, information, and code security (Yan et al., 2025), which can be particularly problematic in finance or military and other industries that handle sensitive personal information. Accordingly, this paper systematically explores and analyzes the integration of small-sized open-source models within the Agent Skill framework, using multiple open-source

---

[♠]These authors contributed equally to this work.

[1]More results: https://anonymous.4open.science/r/LabAgentSkill-C0E7

[2]https://github.com/anthropics/skills

datasets and real-world industrial applications.

## 2 Related Work

### 2.1 Context Engineering (CE)

Regarding this Agent Skill framework, the concept of CE in its backbone has already garnered substantial research and engineering practice across multiple directions. With the advent of zero-shot or few-shot generalization capabilities in LLMs (Brown et al., 2020b; Dong et al., 2024b), post-training paradigms have gradually been challenged by more convenient, efficient, and cost-effective post-deployment CE approaches for domain adaptation and agent behavior improvement. Moreover, a growing body of research has validated that LLMs exhibit human-like attention limitations and the "Lost in the Middle" / "Context rot" (Anthropic, 2025) phenomenon when confronted with excessively long contexts(Kou et al., 2024; Dai et al., 2024; Du et al., 2025). Contemporary LLMs employ various CE designs, such as hierarchical multi-agent systems to handle complex tasks(Luo et al., 2025); routing steps that direct inputs to designated agents (Yue et al., 2025); and sophisticated management of multi-turn dialogue histories. These are integrated with heterogeneous long and short term agent memory modules (Zhang et al., 2024b; Salama et al., 2025; Hu et al., 2025) leveraging file systems, vector databases, knowledge graphs(Edge et al., 2025), and memory buffers (Xu et al., 2025) to achieve comprehensive enhancements in contextual coherence, personalized learning, and complex task decision-making.

### 2.2 Agent Skills Flourish, Not in SLMs

This framework was first introduced in a blog post by Claude[3], and since then, their design has been widely adopted by major players, including VS-Code, OpenAI, LangChain, and others. Ye et al. conceptualizes CE as an evolvable skill, proposing a two-layer framework termed "Meta CE" to automatically rewrite and optimize skill descriptions. Li demonstrates that, for many reasoning tasks, constructing a single-agent system equipped with a skill library achieves comparable accuracy to multi-agent systems while reducing token consumption and latency by approximately half. Numerous open-source libraries for Agent Skills—such as

DeepAgents[4] and the Agent Skill Collections[5] have grown rapidly (Chen et al., 2026), while also sparking discussions on skill safety and permissions Liu et al.. When using Agent Skills, the default models are typically large-scale proprietary models that require API calls. In the proposed engineering practice, model selection proves crucial for the skill routing of SKILL.md, wherein small-sized models (Li et al., 2025a) often exhibit suboptimal success rates and low performance (Belcak et al., 2025; Li et al., 2025b). Furthermore, research on the feasibility of Agent Skill for small models remains limited, and there is a lack of quantitative evidence on deployment-level efficiency gains—such as VRAM footprint and end-to-end latency.

## 3 Problem and Experiments

### 3.1 Agent Skill - Formal Definitions

#### 3.1.1 Mathematical definition

An Agent Skill system can be abstracted as a constrained, information-seeking controller over a partially observed world: at each step, the agent decides whether to (i) **commit** to a skill, (ii) **reveal** more information about that skill (or its references), or (iii) **act** in the environment, trading off task progress against limited attention/context. Considering that the underlying system evolves according to Markov dynamics while the agent cannot directly observe the true (agent–task) state, the agent must maintain a *belief* state that summarizes its interaction history. It aggregates past observations and previously taken actions into a posterior probability distribution over the current latent state. Hence, this setting is naturally modeled as a partially observable Markov decision process (POMDP) (Puterman, 2014) augmented with explicit *information-acquisition* actions (e.g., $\text{reveal}(\rho_k)$) that trade off additional cost for improved observations. Let each skill $k \in \mathcal{K}$ be represented by the abstract triple $k \in \mathcal{K}, k = (d_k, \pi_k, \rho_k)$, where $d_k$ is a textual descriptor (e.g., the skill's name and description), $\pi_k$ is a temporally extended intra-skill policy (i.e., an option-level procedure or workflow that induces a sequence of actions), and $\rho_k$ is a reference mechanism (formalizable as a Markov kernel) that can reveal additional skill-relevant context and tools, potentially including pointers to other skills. Define a POMDP, $\mathcal{M} = \langle S, A, O, T, \Omega \rangle$,

---

[3]https://claude.com/blog/skills

[4]https://github.com/langchain-ai/deepagents
[5]https://github.com/heilcheng/awesome-agent-skills

**1. State** $s_t \in S$: The (possibly hidden) task situation: user intent, task stage/progress, and unretrieved environment facts.

**2. Observation** $o_t \in O$: What the agent can access at time $t$, e.g., the current user message, available history, and any revealed/loaded resources and skills.

**3. Actions** $a_t \in A$: Skill selection $\text{use}(k)$, paid context acquisition $\text{reveal}(\rho_k)$, skill execution $\text{execute}(\pi_k)$, and environment/tool calls (plus internal planning/inference).

**4. Belief** $b_t = P(s_t \mid o_{\leq t}, a_{<t})$: The agent's posterior over $s_t$ given observations so far and past actions (its uncertainty state).

**5. Transition** $T$ (Markov): How the latent state changes after acting (progress, context/skill availability, and inferred inteleant updates): $P(s_{t+1} \mid s_t, a_t) = T(s_{t+1} \mid s_t, a_t)$.

**6. Observation model** $\Omega$: How the next observation is produced from the updated context and action outcomes (e.g., reveal results, tool outputs); the new observation then drives a belief update for the next decision.

### 3.1.2 Intuitive interpretation

When the agent is highly uncertain, as reflected by a diffuse belief state $b_t$; it is often worthwhile to spend additional time and cost to *reveal* relevant skill context before committing to an execution, since the expected value of information can outweigh the reveal cost. When the agent is already nearly certain (i.e., its belief is concentrated after acquiring sufficient evidence), directly executing the appropriate workflow description is typically more cost-effective than further disclosure. This progressive-disclosure behavior in Agent Skill aligns with the classic result that, for finite-horizon POMDPs, the optimal value function over the belief space is piecewise-linear and convex (Kaelbling et al., 1998): different regions of belief correspond to different optimal contingent plans. Consequently, within an agent framework, one can view progressive disclosure as an instance of approximately optimal control with *information actions*: reveal only what is necessary to increase expected utility (by reducing uncertainty) before executing the intra-skill policy, rather than exhaustively loading all context upfront.

### 3.2 Methods

To evaluate Agent Skills on real-world tasks, a temporary Skill repository is constructed for each task by sampling 4–5 distractor Skill entries from a publicly collected skill hub and combining them with the ground-truth Skill. Then, a study is carried out to check how three context-engineering strategies affect agent performance and efficiency in complex decision-making. (1) **Direct Instruction (DI)** uses a minimal prompt to mimic raw user input. (2) **Full-Skill Instruction (FSI)** provides a fixed context containing the entire temporary Skill repository, requiring the model to identify the correct Skill among multiple specifications. (3) **Agent Skill Instruction (ASI)** loads Skill information on demand: the model determines whether additional Skill details are needed, retrieves the relevant Skill, and generates an answer conditioned on that information.

### 3.3 Datasets

Table 1: Dataset overview for the evaluations: average length (words/items), number of labels, domain/topic, and evaluation set size.

| Dataset | Word / Item | # Labels | Topics | # Eval |
|---|---|---|---|---|
| IMDB | 74.05 | 2 | Film reviews | 300 |
| FiNER | 50.43 | 139 | Financial Tags | 403 |
| InsurBench | 710.52 | 2 | Insurance Claims | 200 |

As shown in Table 1, a subset of the **IMDB** dataset derived from the Large Movie Review Dataset v1.0 (Maas et al., 2011) is used. Specifically, reviews are filtered to retain only those with a string length between 300 and 500 characters in order to control for review length. This dataset is used to evaluate the agent's ability to perform binary sentiment classification (positive vs. negative) on movie reviews. In addition, the **FiNER** dataset (Loukas et al., 2022) is used, a comparatively challenging XBRL tagging benchmark. FiNER contains 139 tag types and requires strong domain knowledge in finance as well as robust logical reasoning capabilities. In addition to public benchmarks, a proprietary industry dataset, ***InsurBench***, is employed to examine the applicability and performance of Agent Skills in small language models (SLMs). InsurBench is built from authentic insurance-claim email histories, which are typically long and noisy. The threads include claim details extracted from PDF documents, communication mismatches between claimants and agents, and multilingual interference. The task requires an AI agent to issue a recommendation based on the full thread: whether to continue engaging with the

claim, initiate further actions, or close the case and terminate the process. This task is challenging due to the extended interaction context and the need to infer subtle intents of both parties, while also requiring coherent end-to-end reasoning aligned with the task objective.

## 3.4 Small Language Models

What qualifies as a "small" model remains debatable. However, in industrial settings, we often cannot rely on proprietary models via API calls due to information-security constraints; accordingly, the evaluations are restricted to open-source models that span a wide range of scales, from 270M to 80B parameters, as shown in Table 2. Multiple models are also benchmarked at comparable scales to capture model size and training-objective differences. In particular, *code* variants of LLMs are specialized or enhanced for programming-centric behaviors such as code generation, completion, and repair, whereas *reasoning* models are optimized for multi-step deliberation, planning, and verification, and tend to excel at tasks requiring decomposition and derivation. Additionally, evaluations are also done on the closed-source model gpt-4o-mini as a baseline.

Table 2: Model inventory for evaluation across scales (270M–80B) and model specializations, including size, estimated VRAM, and release information.

| Model | Size | VRAM (GB) | Release |
|---|---|---|---|
| Gpt-4o-mini | - | - | 07/2024 (OpenAI) |
| Gemma-3-270m-it | 0.27B | 1 | 07/2025 (Google, 2025b) |
| Gemma-3-4b-it | 4B | 10 | 03/2025 (Google, 2025c) |
| Gemma-3-12b-it | 12B | 29 | 03/2025 (Google, 2025a) |
| Qwen3-30B-Instruct | 30B | 72 | 07/2025 (Qwen Team, 2025a) |
| Qwen3-80B-Instruct | 80B | 192 | 09/2025 (Qwen Team, 2025b) |
| Qwen3-80B-Thinking | 80B | 192 | 09/2025 (Qwen Team, 2025c) |
| Qwen3-80B-Coder | 80B | 192 | 01/2026 (Qwen Team) |

## 3.5 Experimental Settings

The experiments use an agent–skill framework to assess whether small models can approach the performance of much larger ones. In practice, most open-source models are primarily optimized for chat (with only a subset explicitly trained for stronger reasoning), and robust native support for tool invocation remains limited. Furthermore, introducing external components (e.g., tool invocation or MCP-style connectors) can confound evaluation, since end-to-end performance may depend strongly on tool availability and execution correctness rather than on the agent–skill framework itself. Accordingly, the experiments are structured

to emphasize two core aspects: (i) *skill selection*: selecting (routing to) the appropriate skill for a given task and (ii) the agent's execution correctness after the selected skill is obtained, thereby assessing the validity of the agent–skill framework. For evaluation, **Cls ACC (classification accuracy)** and **Cls F1 (F1 score)** are used to quantify classification performance, together with **Skill ACC (skill-selection accuracy)** which measures routing quality. To account for practical efficiency in industrial settings, **Avg GT (min)** is added which is defined as the average processing time per task in minutes. Finally, **Avg VRAM Time (GB·min)**, defined as the average GPU memory, time cost per task is also proposed. This metric design is motivated by common production billing practices based on GPU-hours, under which both wall-clock latency and memory residency translate directly into operational cost. Furthermore, given fixed VRAM budgets, memory occupancy can constitute a primary throughput bottleneck: once GPU memory is saturated by a workload, other jobs may be prevented from running concurrently, an effect not adequately reflected by conventional compute-centric measures such as FLOPS (floating-point operations per second).

# 4 Results

## 4.1 Main Performance

**Skill Returns in SLMs** As shown in Table 3, most SLMs exhibit clear performance improvements while maintaining a high skill-selection accuracy. The gains are particularly pronounced for relatively larger small models; for example, on FiNER, Qwen3-80B-Instruct improves from 0.198 to 0.654 compared with Direction Instruction. In contrast, smaller models such as Gemma-3-4B-IT and Gemma-3-270M-IT show more limited improvements using Agent Skill. It is also observed that for simpler tasks (e.g., IMDB), the benefits of Agent Skills are modest. For more challenging benchmarks such as FiNER and InsurBench, however, the results highlight the necessity of Agent Skills—underscoring the importance of CE. InsurBench further strengthens this conclusion, as its closed-source nature reduces the likelihood of dataset contamination (i.e., that the benchmark has already been seen during training).

**Tiny Models Fail at Skill Routing** Across these three datasets, each evaluation includes 4–6 distracting (irrelevant) skills, which would ostensibly

Table 3: Main performance on IMDB, FiNER, and InsurBench. gpt-4o-mini was not evaluated on InsurBench due to data privacy and security constraints; the corresponding entries are left blank.

| Model Name | Method | IMDB | | | | | FiNER | | | | | InsurBench | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Cls ACC.(↑) | Cls F1(↑) | Skill ACC.(↑) | AVG GT (min)(↓) | AVG VRAM Time (GB·min)(↓) | Cls ACC.(↑) | Cls F1(↑) | Skill ACC.(↑) | AVG GT (min)(↓) | AVG VRAM Time (GB·min)(↓) | Cls ACC.(↑) | Cls F1(↑) | Skill ACC.(↑) | AVG GT (min)(↓) | AVG VRAM Time (GB·min)(↓) |
| Gpt-4o-mini | DI | 0.970 | 0.970 | - | - | - | 0.484 | 0.444 | - | - | - | - | - | - | - | - |
| Qwen3-80B-Instruct | DI | **0.968** | **0.968** | - | 0.027 | 5.259 | 0.198 | 0.172 | - | 0.054 | 10.425 | 0.498 | 0.331 | - | 0.049 | 9.379 |
| | FSI | 0.989 | 0.989 | - | 0.026 | 5.022 | 0.196 | 0.163 | - | 0.055 | 10.547 | 0.530 | 0.530 | - | 0.043 | 8.201 |
| | ASI | 0.953 | 0.953 | 0.997 | **0.022** | **4.195** | **0.654** | **0.647** | 0.978 | **0.027** | **5.242** | **0.620** | **0.601** | 0.915 | **0.028** | **5.321** |
| Qwen3-80B-Thinking | DI | **1.000** | **1.000** | - | **0.088** | **16.822** | 0.478 | 0.506 | - | **0.109** | 20.948 | 0.498 | 0.331 | - | **0.203** | 33.893 |
| | FSI | 0.990 | 0.990 | - | 0.212 | 40.811 | 0.398 | 0.377 | - | 0.188 | 36.076 | 0.435 | 0.426 | - | 0.208 | 40.017 |
| | ASI | 0.973 | 0.973 | 1.000 | 0.494 | 94.802 | **0.717** | **0.699** | 0.960 | 0.492 | 94.503 | **0.545** | **0.495** | 0.945 | 0.943 | 181.003 |
| Qwen3-80B-Coder | DI | 0.938 | 0.938 | - | **0.020** | **3.833** | 0.309 | 0.307 | - | 0.047 | 8.928 | 0.498 | 0.331 | - | 0.058 | 11.176 |
| | FSI | 0.965 | 0.965 | - | 0.173 | 33.289 | 0.350 | 0.340 | - | 0.049 | 9.381 | 0.541 | 0.521 | - | 0.055 | **10.555** |
| | ASI | **0.969** | **0.969** | 1.000 | 0.025 | 4.729 | **0.657** | **0.646** | 1.000 | **0.033** | **6.359** | **0.660** | **0.658** | 0.990 | **0.057** | 10.975 |
| Qwen3-30B-Instruct | DI | 0.948 | 0.948 | - | 0.023 | 1.646 | 0.184 | 0.182 | - | 0.023 | 1.683 | 0.500 | 0.333 | - | 0.051 | 3.680 |
| | FSI | **0.976** | **0.976** | - | 0.033 | 2.386 | 0.082 | 0.095 | - | 0.157 | 11.313 | **0.510** | **0.427** | - | 0.041 | 2.923 |
| | ASI | 0.950 | 0.950 | 1.000 | **0.015** | **1.083** | **0.564** | **0.528** | 0.995 | **0.023** | **1.678** | 0.450 | 0.332 | 0.990 | **0.030** | **2.153** |
| Gemma-3-12b-it | DI | 0.965 | 0.965 | - | 0.441 | 12.693 | 0.433 | 0.442 | - | 0.457 | 13.147 | 0.500 | 0.333 | - | 0.034 | 0.964 |
| | FSI | **0.969** | **0.969** | - | 0.874 | 25.178 | 0.314 | 0.331 | - | 0.676 | 19.455 | 0.460 | 0.446 | - | 0.040 | 1.149 |
| | ASI | 0.968 | 0.968 | 1.000 | **0.245** | **7.069** | **0.503** | **0.486** | 1.000 | **0.149** | **4.298** | **0.575** | **0.540** | 0.990 | **0.025** | **0.713** |
| Gemma-3-4b-it | DI | 0.950 | 0.950 | - | 0.037 | 0.354 | 0.216 | **0.200** | - | 0.013 | 0.123 | 0.500 | 0.333 | - | 0.034 | 0.333 |
| | FSI | 0.941 | 0.941 | - | 0.053 | 0.506 | 0.041 | 0.032 | - | 0.058 | 0.553 | 0.515 | **0.507** | - | 0.041 | 0.389 |
| | ASI | **1.000** | **1.000** | 1.000 | **0.015** | **0.145** | 0.219 | 0.175 | 1.000 | **0.008** | **0.081** | **0.525** | 0.389 | 0.780 | **0.016** | **0.150** |
| Gemma-3-270m-it | DI | **0.637** | **0.609** | - | **0.002** | **0.006** | 0.034 | 0.007 | - | **0.028** | **0.067** | **0.500** | **0.333** | - | **0.007** | **0.016** |
| | FSI | 0.500 | 0.418 | - | 0.006 | 0.015 | 0.000 | 0.000 | - | 0.061 | 0.147 | **0.500** | **0.333** | - | 0.017 | 0.041 |
| | ASI | 0.500 | 0.418 | 0.000 | 0.060 | 0.144 | **0.250** | **0.192** | 0.050 | 0.051 | 0.123 | 0.415 | 0.293 | 0.190 | 0.021 | 0.051 |

make skill identification relatively straightforward for the model. However, we find that extremely small models such as Gemma-3-4B-it and Gemma-3-270M-it still struggle to retrieve the appropriate skill. In particular, Gemma-3-270M-it appears to largely miss the objective of skill retrieval, and Gemma-3-4B-it achieves only a 0.78 success rate on InsurBench. These results suggest that, within an agent-skill framework, models below 4B parameters often lack even the basic capability to identify the correct skill—let alone reliably carry out the subsequent, more complex execution steps.
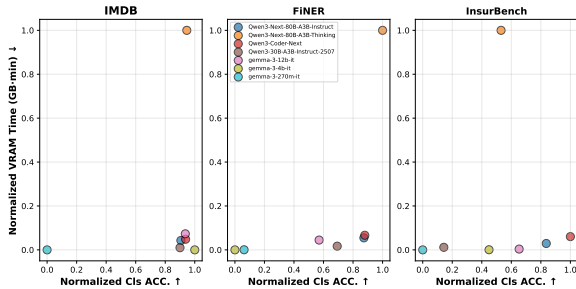


Figure 1: Normalized average VRAM-time vs. task performance across different model variants on three datasets. Lower-left region indicates superior performance; both axes are normalized.

**Code Model is Better** Qwen3-NEXT provides three 80B variants. Holding model size constant, it is seen that after incorporating Agent Skills, the code-oriented models consistently outperform their instruction-tuned counterparts, suggesting that under an agent-skill framework, code models can achieve higher accuracy than instruct models. On InsurBench, the code model even surpasses the thinking variant, although on FiNER its accuracy (0.657) remains below that of the thinking model (0.717). However, as shown in Figure 1, when VRAM consumption and end-to-end inference time are taken into account, code models substantially reduce GPU VRAM-time. Therefore, the conclusion drawn out is that for industrial deployment, combining Agent Skills with a code model constitutes the most VRAM-efficient choice. This may also help explain, at least in part, why the Claude family has seen particularly rapid adoption in agent-skill applications.

## 4.2 Hitting Paradigm? Small Vs Large

Prior experiments indicate that tiny SLMs exhibit noticeable performance degradation in model selection tasks for Gemma-3-4b-it, even when interference is limited to only 4–6 competing skills. To further investigate this phenomenon, the robustness of SLMs is evaluated under larger skill hubs, reflecting more realistic autonomous agent development scenarios that require extensive skill repertoires (e.g., exceeding 50 skills to achieve full project autonomy). As shown in Figure 2, tiny models exhibit rapid accuracy decline beyond $N = 10$–$20$, whereas models exceeding 12B parameters demonstrate exceptional robustness, maintaining high precision even at $N = 100$. In particular, the code-specialized variant outperforms its counterparts in skill selection tasks. SLMs struggle to capture hierarchical skill-revealing structures, whereas only large-scale models reliably handle nested depen-

dencies within a single SKILL.md. Even proprietary models like GPT-4o-mini occasionally falter in interpreting these relationships accurately, as evidenced in several LangChain DeepAgent CLI experiments, where solely Claude-Opus models consistently achieved near-100% success rate in identifying referenced skills within SKILL.md descriptions.
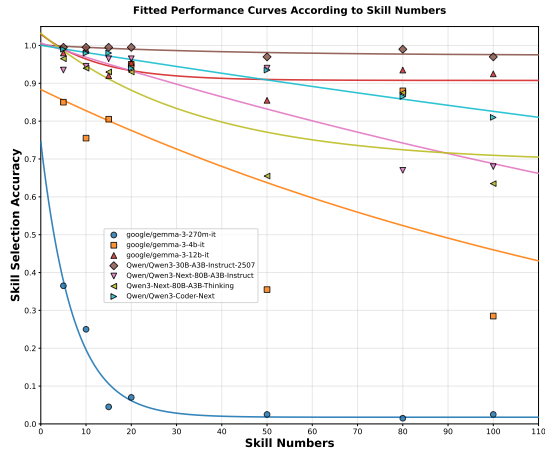


Figure 2: Fitted decay curves (solid lines) and empirical data points (markers) for skill selection accuracy across skill number $N = 5$ to 100.

### 4.3 Post-hoc Exploration

**Chat History Matter?** As shown in Table 4, the impact of incorporating conversational history on model performance in InsurBench is evaluated. To avoid context-window overflow, the dialogue before each call is truncated, retaining the system prompt and the most recent 3–4 turns. Results indicate that history benefits are largest for very small models (e.g., Gemma-3-4b-it and Gemma-3-270m-it), while larger SLMs improve only slightly. However, chat history increases VRAM-time usage—for Qwen3-80B-Instruct, the cost rises from 5.321 to 10.035 GB·min per item. Therefore, enabling chat-history processing is mainly recommended for lightweight SLMs in Agent Skills deployments.

**Replacing "Skill" with "others"?** An exploratory experiment is conducted to investigate whether replacing the keyword "Skill" with its synonyms affects the efficiency and accuracy of agent tasks. As shown in Table 5, four synonyms were tested. Their impact is shown to be minimal on the performance. Notably, "Expertise" consistently outperformed "Skill" across metrics, suggesting it as a potentially superior alternative. Additionally, "Know-how" demonstrated substantial improvements in

| Model Name | Methods | Cls ACC. ↑ | Avg VRAM Time (GB·min) ↓ |
|---|---|---|---|
| Qwen3-80B-Instruct | ASI | 0.620 | 5.321 |
| | ASIH | 0.535 | 10.035 |
| Qwen3-30B-Instruct | ASI | 0.450 | 2.153 |
| | ASIH | 0.500 | 2.243 |
| Gemma-3-12b-it | ASI | 0.575 | 0.713 |
| | ASIH | 0.585 | 0.7521 |
| Gemma-3-4b-it | ASI | 0.525 | 0.150 |
| | ASIH | 0.660 | 0.152 |
| Gemma-3-270m-it | ASI | 0.415 | 0.051 |
| | ASIH | 0.525 | 0.058 |

Table 4: Performance of Various SLMs on InsurBench: ASI vs. ASIH (ASI with Chat History) on Qwen3-80B-Instruct.

| Method | Keyword Name | Cls ACC. ↑ | Cls F1 ↑ | Skill ACC. ↑ | AVG GT (min) ↓ | Avg VRAM Time (GB·min) ↓ |
|---|---|---|---|---|---|---|
| ASI | Skill | **0.620** | 0.601 | 0.915 | 0.028 | 5.321 |
| | Capability | 0.595 | 0.594 | 0.915 | 0.032 | 6.147 |
| | Expertise | 0.610 | **0.608** | **0.930** | 0.032 | 6.184 |
| | Proficiency | 0.580 | 0.578 | 0.925 | 0.031 | 5.946 |
| | Know-how | 0.580 | 0.579 | 0.910 | **0.022** | **4.302** |
| FSI | Skill | 0.530 | 0.506 | - | 0.043 | 8.201 |
| | Capability | 0.485 | 0.436 | - | 0.055 | 10.458 |
| | Expertise | **0.570** | **0.539** | - | 0.056 | 10.739 |
| | Proficiency | 0.510 | 0.462 | - | 0.055 | 10.486 |
| | Know-how | 0.535 | 0.497 | - | **0.040** | **7.459** |

Table 5: Performance comparison of different skill synonyms under ASI and FSI frameworks on Qwen3-80B-Instruct tested on InsurBench. Bold values indicate the best performance within each method group.

GPU memory efficiency with negligible performance degradation.

## 5 Conclusion & Limitations

This paper evaluates SLMs in regulated settings with data-security and GPU-VRAM constraints using the Agent Skill framework, measuring efficiency through GPU VRAM time. Very small models struggle with large skill hubs (50–100 skills), while moderately larger SLMs significantly improve selection accuracy and robustness. Code-specialized variants demonstrate the highest VRAM efficiency. Also, The impact of model scale varies across subtasks: mid-sized models ($\approx 13B - 30B$) are better suited to skill selection, whereas larger models ($\approx 80B$) attain execution quality comparable to GPT-4o-mini. In addition, this paper reports exploratory experiments that provide guidance for future Agent Skill deployment. Evaluations are limited to a narrow set of tasks, primarily classification and tagging. The underlying causes of SLM difficulties in sustained or recursive reasoning under Progressive Disclosure remain unclear, as does the observed accuracy and

VRAM efficiency of code-oriented LLMs. The optimal structure and representation of Skill.md also remain open questions.

# References

Anthropic. 2025. Effective context engineering for ai agents. Accessed: February 6, 2026.

Peter Belcak, Greg Heinrich, Shize Diao, Yonggan Fu, Xin Dong, Saurav Muralidharan, Yingyan Celine Lin, and Pavlo Molchanov. 2025. Small language models are the future of agentic ai. *Preprint*, arXiv:2506.02153.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020a. Language models are few-shot learners. *CoRR*, abs/2005.14165.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020b. Language models are few-shot learners. *Preprint*, arXiv:2005.14165.

Tianyi Chen, Yinheng Li, Michael Solodko, Sen Wang, Nan Jiang, Tingyuan Cui, Junheng Hao, Jongwoo Ko, Sara Abdali, Leon Xu, Suzhen Zheng, Hao Fan, Pashmina Cameron, Justin Wagle, and Kazuhito Koishida. 2026. Cua-skill: Develop skills for computer using agent. *Preprint*, arXiv:2601.21123.

Hui Dai, Dan Pechi, Xinyi Yang, Garvit Banga, and Raghav Mantri. 2024. Deniahl: In-context features influence llm needle-in-a-haystack abilities. *Preprint*, arXiv:2411.19360.

Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. 2024a. A survey on in-context learning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 1107–1128, Miami, Florida, USA. Association for Computational Linguistics.

Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Tianyu Liu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. 2024b. A survey on in-context learning. *Preprint*, arXiv:2301.00234.

Yufeng Du, Minyang Tian, Srikanth Ronanki, Subendhu Rongali, Sravan Bodapati, Aram Galstyan, Azton Wells, Roy Schwartz, Eliu A Huerta, and Hao Peng. 2025. Context length alone hurts llm performance despite perfect retrieval. *Preprint*, arXiv:2510.05381.

Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, Dasha Metropolitansky, Robert Osazuwa Ness, and Jonathan Larson. 2025. From local to global: A graph rag approach to query-focused summarization. *Preprint*, arXiv:2404.16130.

Google. 2025a. google/gemma-3-12b-it. https://huggingface.co/google/gemma-3-12b-it. Hugging Face model card. Accessed: 2026-02-15.

Google. 2025b. google/gemma-3-270m-it. https://huggingface.co/google/gemma-3-270m-it. Hugging Face model card. Accessed: 2026-02-15.

Google. 2025c. google/gemma-3-4b-it. https://huggingface.co/google/gemma-3-4b-it. Hugging Face model card. Accessed: 2026-02-15.

Mengkang Hu, Tianxing Chen, Qiguang Chen, Yao Mu, Wenqi Shao, and Ping Luo. 2025. HiAgent: Hierarchical working memory management for solving long-horizon agent tasks with large language model. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 32779–32798, Vienna, Austria. Association for Computational Linguistics.

Leslie Pack Kaelbling, Michael L. Littman, and Anthony R. Cassandra. 1998. Planning and acting in partially observable stochastic domains. *Artificial Intelligence*, 101(1):99–134.

Bonan Kou, Shengmai Chen, Zhijie Wang, Lei Ma, and Tianyi Zhang. 2024. Do large language models pay similar attention like human programmers when generating code? *Proceedings of the ACM on Software Engineering*, 1(FSE):2261–2284.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474. Curran Associates, Inc.

Lujun Li, Lama Sleem, Niccolo' Gentile, Geoffrey Nichil, and Radu State. 2025a. Small language models in the real world: Insights from industrial text classification. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 6: Industry Track)*, pages 971–982, Vienna, Austria. Association for Computational Linguistics.

Lujun Li, Lama Sleem, Niccolo' Gentile, Geoffrey Nichil, and Radu State. 2025b. Exploring the impact of temperature on large language models: Hot or cold? *Procedia Computer Science*, 264:242–251. International Neural Network Society Workshop on Deep Learning Innovations and Applications 2025.

Xiaoxiao Li. 2026. When single-agent with skills replace multi-agent systems and when they fail. *Preprint*, arXiv:2601.04748.

Yi Liu, Weizhe Wang, Ruitao Feng, Yao Zhang, Guangquan Xu, Gelei Deng, Yuekang Li, and Leo Zhang. 2026. Agent skills in the wild: An empirical study of security vulnerabilities at scale. *Preprint*, arXiv:2601.10338.

Lefteris Loukas, Manos Fergadiotis, Ilias Chalkidis, Eirini Spyropoulou, Prodromos Malakasiotis, Ion Androutsopoulos, and Paliouras George. 2022. FiNER: Financial Numeric Entity Recognition for XBRL Tagging. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics.

Junyu Luo, Weizhi Zhang, Ye Yuan, Yusheng Zhao, Junwei Yang, Yiyang Gu, Bohan Wu, Binqi Chen, Ziyue Qiao, Qingqing Long, Rongcheng Tu, Xiao Luo, Wei Ju, Zhiping Xiao, Yifan Wang, Meng Xiao, Chenwu Liu, Jingyang Yuan, Shichang Zhang, and 7 others. 2025. Large language model agent: A survey on methodology, applications and challenges. *Preprint*, arXiv:2503.21460.

Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA. Association for Computational Linguistics.

Lingrui Mei, Jiayu Yao, Yuyao Ge, Yiwei Wang, Baolong Bi, Yujun Cai, Jiazhi Liu, Mingyu Li, Zhong-Zhi Li, Duzhen Zhang, Chenlin Zhou, Jiayi Mao, Tianze Xia, Jiafeng Guo, and Shenghua Liu. 2025. A survey of context engineering for large language models. *Preprint*, arXiv:2507.13334.

OpenAI. Gpt-4o mini model | openai api. https://developers.openai.com/api/docs/models/gpt-4o-mini. Accessed: 2026-02-15.

Martin L Puterman. 2014. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons.

Qwen Team. Qwen3-coder-next technical report. Technical report. Accessed: 2026-02-03.

Qwen Team. 2025a. Qwen/qwen3-30b-a3b-instruct-2507. https://huggingface.co/Qwen/Qwen3-30B-A3B-Instruct-2507. Hugging Face model card. Accessed: 2026-02-15.

Qwen Team. 2025b. Qwen/qwen3-next-80b-a3b-instruct. https://huggingface.co/Qwen/Qwen3-Next-80B-A3B-Instruct. Hugging Face model card. Accessed: 2026-02-15.

Qwen Team. 2025c. Qwen/qwen3-next-80b-a3b-thinking. https://huggingface.co/Qwen/Qwen3-Next-80B-A3B-Thinking. Hugging Face model card. Accessed: 2026-02-15.

Rana Salama, Jason Cai, Michelle Yuan, Anna Currey, Monica Sunkara, Yi Zhang, and Yassine Benajiba. 2025. MemInsight: Autonomous memory augmentation for LLM agents. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 33136–33152, Suzhou, China. Association for Computational Linguistics.

Mirac Suzgun, Mert Yuksekgonul, Federico Bianchi, Dan Jurafsky, and James Zou. 2025. Dynamic cheatsheet: Test-time learning with adaptive memory.

Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao, Juntao Tan, and Yongfeng Zhang. 2025. A-mem: Agentic memory for llm agents. *Preprint*, arXiv:2502.12110.

Biwei Yan, Kun Li, Minghui Xu, Yueyan Dong, Yue Zhang, Zhaochun Ren, and Xiuzhen Cheng. 2025. On protecting the data privacy of large language models (llms) and llm agents: A literature review. *High-Confidence Computing*, 5(2):100300.

Haoran Ye, Xuning He, Vincent Arak, Haonan Dong, and Guojie Song. 2026. Meta context engineering via agentic skill evolution. *Preprint*, arXiv:2601.21557.

Yanwei Yue, Guibin Zhang, Boyang Liu, Guancheng Wan, Kun Wang, Dawei Cheng, and Yiyan Qi. 2025. MasRouter: Learning to route LLMs for multi-agent systems. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15549–15572, Vienna, Austria. Association for Computational Linguistics.

Kechi Zhang, Jia Li, Ge Li, Xianjie Shi, and Zhi Jin. 2024a. Codeagent: Enhancing code generation with tool-integrated agent systems for real-world repo-level coding challenges. *Preprint*, arXiv:2401.07339.

Zeyu Zhang, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Quanyu Dai, Jieming Zhu, Zhenhua Dong, and Ji-Rong Wen. 2024b. A survey on the memory mechanism of large language model based agents. *Preprint*, arXiv:2404.13501.

## Acknowledgments

## A Progressive Disclosure Difficulty

The feasibility of progressive disclosure for small-scale models was also evaluated in several trial, but poor results were obtained. Cross-Skill references within Skill descriptions were often not detected unde SLMs senarios, which prevented the system from triggering intra-Skill calls (e.g., when a Skill description in `SKILL.md` referenced another Skill that needed to be disclosed). Even with GPT-4o-mini, a low hit rate for such references was observed. The official LangChain CLI was also tested with its default system-prompt configuration, and higher hit rates were observed only with Claude Sonnet 4.5 or Claude Opus 4.5. Therefore, intra-Skill invocation (i.e., calling one Skill from within another) was excluded from our experiments, since skill selection rates yielded by open-source models were too low for meaningful comparison; this exclusion is also recommended by Anthropic.

## B Additional Experimental Settings

In all experiments, the models' default decoding configurations (e.g., temperature, top-$p$, and top-$k$) were used to avoid confounding effects from additional hyperparameter tuning. All methods were implemented using the LangChain agent framework to ensure a consistent agentic workflow across settings. For inference, vLLM was used as the serving engine, and the context length was fixed to 10240 tokens.

To manage conversational history under the context-window constraint, a deterministic message-trimming policy was applied. Specifically, the first message (typically the system prompt) was always retained, and only the most recent 3 or 4 messages were kept depending on the parity of the total message count. Core instructions were preserved while the prompt length was bounded, thereby reducing truncation risk and stabilizing inference cost in longer interactions. For formatting and metadata extraction from `SKILL.md`, the publicly available codebase is primarily used.[6]

## C Prompts

### C.1 Direct Instruction Prompts

This is the Direct Instruction prompt setting on the IMDB dataset, where the task is to classify a given review as either positive or negative. Since the task is relatively simple, introducing the skill-based mechanism does not yield a significant performance improvement.

---

> **Direct Instruction Of IMDB:**
>
> Task:
> Classify the following movie review as positive or negative sentiment.
>
> Review:
>
> <<<Review Content>>>

---

This is the Direct Instruction prompt setting for the FiNER dataset. The objective is to identify and classify XBRL tags. As shown, we provide not only the sentence but also the target numerical value appearing in that sentence. Because the 139 candidate tags correspond to specialized financial terminology, this task demands **strong** logical reasoning as well as substantial domain knowledge in finance.

---

> **Direct Instruction Of FiNER:**
>
> Given a sentence from financial documents, a target numeric entity from that sentence, and a list of candidate XBRL tags, choose the single best-matching XBRL tag.
>
> Tag List (candidates): [InterestExpense, ....]
>
> Inputs:
>
> Sentence: <<<Sentence Content>>>
> Target entity: <<<Numeric Entity>>>

---

This is the Direct Instruction prompt setting for InsurBench. The task requires the agent to identify the key email(s) within a very long email thread and to make decisions on behalf of an insurance company. This setting evaluates the agent model's ability to localize salient information and to perform complex decision-making under long-context conditions.

---

## Direct Instruction Of InsurBench:

Task:

Given the full email thread in Email History, decide whether the insurance company must take action to reply.

Email History:

<<<Email History>>>

## Skill Selection System Prompt (Part 2):

1. **Think step-by-step** about the user's request:
- What is the main task?
- Which skill domain does it match?

2. **Skill matching rules:**
- **Multiple skills**: List names separated by commas (e.g., "langgraph-docs, sales-analytics")
- **No skill matches**: Use empty list '[]'

3. **Generate response based on skills found:**

| Skills Found | Message Content |
|————|————|
| 1+ skills | "Yes I need to read the skill information first because ..." |
| No skills | "I didn't find the right skill." |

## C.2 Skill Selection Prompts

In this setting, we configure the model at initialization via the system prompt to prioritize skill selection rather than answering the user query directly. This system prompt is adapted from the Agent Skill system-agent implementation in the LangChain DeepAgent CLI, and we further refine it to better suit our use case.

## Skill Selection System Prompt (Part 1):

In order to complete the objective that the user asks of you, you have access to a number of skills.

## Skills System

**Available Skills:**
Skill Context

## Step-by-Step Process

**ALWAYS follow these exact steps:**

## Skill Selection System Prompt (Part 3):

**Final Output Format (Strict JSON)**

Your final response must **always** follow this JSON structure:

{
"Message": "Your complete response to the user query goes here.",
"Skills": ["List of skills selected to complete the request, e.g., ['sales-analytics','sentiment-analytics']"]
}

## Examples

## Skill Selection System Prompt (Part 4):

## Examples


**User:** "How do I use LangGraph StateGraph?"
**Skills:** ["langgraph-docs"]
**JSON:** '"Message": "es I need to read the skill information first because I need details on StateGraph from LangGraph docs.", "Skills": ["langgraph-docs"]'

**User:** "Write a poem about cats"
**Skills:** []
**JSON:** '"Message": "I didn't find the right skill.", "Skills": []'


**User:** "Write a report about LangGraph usage in sales."
**Skills:** ["langgraph-docs", "sales-analytics"]
**JSON:** '"Message": "Yes I need to read the skill information first because it involves LangGraph documentation and sales analytics.", "Skills": ["langgraph-docs", "sales-analytics"]'

## Skill Execution System Prompt (Part 1):

In order to complete the objective that the user asks of you, you have access to a number of skills description.


## Skills System
You have access to a skills library that provides specialized capabilities and domain knowledge.


**How to Use Skills:**


1. **Read the skill's full instructions**:
2. **Follow the skill's instructions**: contains step-by-step workflows, best practices, and examples


**When to Use Skills:**
- User's request matches a skill's domain (e.g., "research X" -> web-research skill)
- A skill provides proven patterns for complex tasks

## Skill Execution System Prompt (Part 2):

**Example Workflow:**
User: "Can you research the latest developments in quantum computing?"
1. Check available skills description
2. Read the skill
3. Follow the skill's research workflow (search -> organize -> synthesize)
4. Make the final decision.

**Skill Information Collected**

{{Skill Context}}

Remember: Skills make you more capable and consistent. When in doubt, check if a skill exists for the task!

## Output instructions

**Final Output Format (Strict JSON)**

{ "Message": Your message here.
}