

# HSCODECOMP: A Realistic and Expert-level Agent Benchmark for Hierarchical Rule Application

Tian Lan\*, Yiqian Yang\*, Qianghuai Jia<sup>†</sup>, Li Zhu, Hui Jiang  
Hang Zhu, Weihua Luo, Longyue Wang  
Alibaba Group

## Abstract

Despite recent progress, existing agent benchmarks neglect a fundamental real-world capability: *hierarchical rule application*, a critical requirement in fields like law and medicine where agents must reason from broad categories down to specific exceptions. This introduces significant challenges, particularly in resolving complex logical dependencies and disambiguating fuzzy semantic boundaries. To bridge this gap, we introduce HSCODECOMP, an e-commerce benchmark that requires rigorous interpretation of product attributes and tariff rules to assign unique 10-digit Harmonized System Codes (HS Codes) to products<sup>1</sup>. HSCODECOMP comprises 632 products across 32 categories from e-commerce platforms. Each product contains detailed but noisy attributes mirroring real-world challenges, such as attributes and image. We also compile the official hierarchical tariff rules and knowledge databases. Utilizing these resources, 26 domain experts rigorously annotated the HSCodes. Evaluations of 23 state-of-the-art LLMs, VLMs, and agents reveal a large performance gap: the best-performing agent achieves only 46.8% accuracy versus 95.0% for human experts, and test-time scaling fails to close this gap. Extensive analysis reveals the challenges of hierarchical rule application. For example, excessive reasoning steps often induce "reasoning drift" that degrades accuracy. Furthermore, agents often misapply rules due to insufficient domain knowledge, as well as reasoning hallucinations that lack factual grounding<sup>2</sup>.

## 1 Introduction

Recent advancements have empowered Large Language Model (LLM)-based agents to handle com-

plex tasks involving reasoning and navigating web environments (Wei et al., 2025; Mialon et al., 2023b) or accessing structured databases (Hu et al., 2025a; Yu et al., 2025a). However, a critical gap persists in numerous real-world professional domains (e.g., legal, medical and e-commerce), where decision-making is governed not by retrieved knowledge or internal knowledge of models, but by rigorous expert-written rules (Sadowski and Chudziak, 2025). This capability involves navigating expert-written rules containing vague boundaries and logical dependencies—a challenge largely overlooked by current agent benchmarks (Mialon et al., 2023b; Wei et al., 2025).

To formalize this challenge, we categorize agentic knowledge into three levels of complexity (Figure 1, left): (1) **Level 1: Open-Domain Data**: This level involves understanding and reasoning over massive open-domain knowledge like web pages. Established benchmarks include GAIA (Mialon et al., 2023b) and BrowseComp (Wei et al., 2025); (2) **Level 2: Structured Data**: This level requires precise utilization of structured data like databases and knowledge graphs, as seen in domain-specific agent benchmarks like MedBrowseComp (Chen et al., 2025b) and FinSearchComp (Hu et al., 2025a); and (3) **Level 3: Hierarchical Rule Data**: We position this level as a critical frontier in evaluating current agents, requiring not only open-domain and domain-specific knowledge utilization but also precise navigation through hierarchical expert-written rules. An example of the HSCode hierarchical tariff rules is shown in Figure 10. Applying these rules presents unique challenges: (a) accurate decision-making when rules contain vague natural language descriptions (Sadowski and Chudziak, 2025); (b) reasoning about logical dependencies among rules, such as exception clauses and cross-references (Guha et al., 2023); and (c) precise interpretation and application of hierarchical rules throughout multi-step reasoning.

\* Equal contribution.

<sup>†</sup> Corresponding author: qianghuai.jqh@alibaba-inc.com

<sup>1</sup>HSCode is the global tariff standard that encodes a product's essential character to facilitate cross-border trade.

<sup>2</sup>Resources have been publicly released at <https://github.com/AIDC-AI/Marco-DeepResearch>.

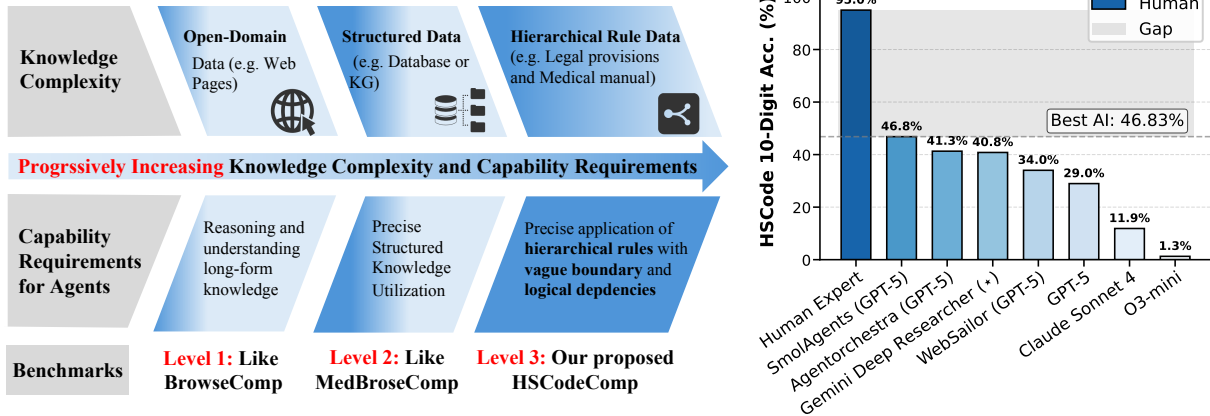


Figure 1. **Left:** Recent benchmarks reveal the increasing knowledge complexity and capability requirements for agents. **Right:** 10-digit HSCode accuracy of state-of-the-art baseline largely lags behind human experts (46.8% < 95.0%), proving the challenges of hierarchical rule application. The closed-source agent (\*) is evaluated on the subset due to API unavailability.

To rigorously evaluate this overlooked Level 3 capability, we introduce HSCODECOMP, the first realistic and expert-level benchmark designed specifically for Hierarchical Rule Application. Grounded in the global standard for international trade (Grainger, 2024) as a representative testbed, the task requires agents to assign 10-digit Harmonized System Codes (HSCodes) to products based on strict hierarchical tariff rules. HSCODECOMP comprises 632 carefully curated product entries from real-world e-commerce platforms across 32 categories. Each product entry contains detailed but noisy features, such as title, price, category, attributes and image. Furthermore, we provide several necessary resources for this task: (1) official tariff classification rules (United States Harmonized Tariff Schedule, HTSUS); (2) expert-written decision rules; and (3) a knowledge base of historical U.S. customs rulings. To guarantee maximum evaluation rigor, we eschew conventional crowdsourcing. Instead, all ground-truth HSCode labels are meticulously annotated and cross-validated by a dedicated panel of 26 domain experts, averaging over five years of professional experience.

We evaluate 23 state-of-the-art LLMs, VLMs and agent systems on HSCODECOMP. The results expose a substantial disparity: As shown in Figure 1 (right), while human experts achieve 95.0% accuracy, the best-performing agent reaches only 46.8% (GPT-5). Crucially, this gap cannot be bridged by test-time scaling strategies (Liu et al., 2025b; Zhu et al., 2025). This performance gap highlights the inadequacy of current agents in handling hierarchical rule applications. Extensive experimental

analysis and failure taxonomy uncover several important findings. For example, excessive reasoning steps frequently induce "reasoning drift", thereby degrading performance. Furthermore, agents often misapply rules due to insufficient domain knowledge, leading to premature decisions on high-level categories, and frequently suffer from reasoning hallucinations that lack factual grounding in official rules. These insights position HSCODECOMP as a challenging testbed for current AI agents.

## 2 Related Works

**HSCode Prediction.** Most previous works treat HSCode prediction as the e-commerce text classification task (Grainger, 2024), using pre-trained BERT models (Liao et al., 2024; Shubham et al., 2022) or Large Language Models (LLMs) prompting (Gholamian et al., 2024). However, these approaches fail to leverage domain-specific knowledge, especially the rules written by human experts (Gholamian et al., 2024; Judy, 2024; Lee et al., 2024; Stassin et al., 2023).

**Benchmarking Level 1 Knowledge.** Numerous benchmarks have been proposed to evaluate agent capabilities in understanding and deep reasoning over long-form open-domain web content (Thomas et al., 2025; Yao et al., 2024; Joshi et al., 2017; Phan et al., 2025). For example, WebArena (Zhou et al., 2023) provides realistic, self-hostable websites with standardized evaluation protocols to assess functional correctness. WebShop (Yao et al., 2022) and ALFWorld (Shridhar et al., 2021) evaluate long-horizon decision-making abilities of

agents in web environments through tool interactions. More recent deep search benchmarks, such as GAIA (Mialon et al., 2023a), BrowseComp (Wei et al., 2025), demand advanced tool-usage capabilities (Zhang et al., 2025; Li et al., 2025b).

**Benchmarking Level 2 Knowledge.** Recent works have focused on how agents utilize structured knowledge in domain-specific applications. Unlike open-domain data, domain-specific knowledge is typically organized into structured formats such as databases and knowledge graphs (Huang et al., 2025; Yu et al., 2025b; Chen et al., 2025a), enabling more precise knowledge retrieval and utilization. To evaluate these capabilities, numerous deep search benchmarks have been proposed, including WebMall (Peeters et al., 2025) and DeepShop (Lyu et al., 2025) for e-commerce, LegalAgentBench for law (Li et al., 2025a), FinSearchComp for finance (Hu et al., 2025a), and MedBrowseComp for medicine (Chen et al., 2025b).

**Neuro-Symbolic Approaches.** Recently, neuro-symbolic architectures have emerged as a promising direction to enforce structured rule adherence by combining the representational flexibility of LLMs with the rigorous logic of symbolic systems (Xu et al., 2025). For example, SymAgent (Liu et al., 2025a) achieves complex reasoning over structured knowledge graphs. MDD-5k (Yin et al., 2024) demonstrates the potential of neuro-symbolic agents in domain-specific applications like medical diagnosis. However, these approaches often struggle with fuzzy, natural-language-based hierarchical rules due to their reliance on well-structured constraints. Our proposed HSCODECOMP directly targets this gap, serving as a realistic testbed to evaluate how effectively neuro-symbolic agent systems can handle ambiguous linguistic boundaries in expert-level tasks.

In summary, while there exists numerous agent benchmarks that target open-domain or domain-specific scenarios, none evaluates hierarchical rule application capability. To address this, we introduce a realistic and expert-level benchmark HSCODECOMP. Our benchmark presents significant challenges even for state-of-the-art closed-source and open-source agent systems.

### 3 Task Formulation

We formalize HSCode prediction as a knowledge-intensive hierarchical reasoning task. Unlike stan-

dard classification, this problem requires an agent not only to retrieve domain knowledge but also to strictly adhere to the expert-written hierarchical rules. Mathematically, the objective is to learn a mapping function  $f : (\mathcal{X}, \mathcal{R}, \mathcal{K}) \rightarrow \mathcal{Y}$ , where the agent predicts a valid HSCode  $y \in \mathcal{Y}$  for product  $x \in \mathcal{X}$  by applying hierarchical rules  $\mathcal{R}$  supported by domain-specific knowledge  $\mathcal{K}$ .

**Input Space ( $\mathcal{X}$ ).** As shown in Figure 2, Each input instance  $x \in \mathcal{X}$  represents a commercial product profile containing real-world and noisy information. We define  $x = (t, A, c, i, p, u)$ , where:  $t$  is the product title;  $A = \{(k_j, v_j)\}_{j=1}^K$  denotes  $K$  structured attributes (e.g., material and package size);  $c$  represents product categories defined by the e-commerce platform;  $i$  is the product image; capturing visual features (e.g., texture, shape) absent in text;  $p$  and  $u$  are the price and currency.

**Hierarchical Rules ( $\mathcal{R}$ ).** The core of HSCODECOMP is the Hierarchical Tariff Rules ( $\mathcal{R}$ ), consisting of two kinds of rules: (1) **Tariff rules** are sourced from official eWTP platform, aggregating authoritative US tariff systems. HSCODECOMP evaluates agents on using US tariff rules. As shown in Figure 2,  $\mathcal{R}$  is characterized by: (a) Hierarchical Structure: organized into Chapters (2-digit), Headings (4-digit), Sub-headings (6-digit) and Country-specific codes (10-digit); (b) Logical Dependencies: frequently contain exclusion clauses and cross-references; and (3) Vague Boundaries: Linguistic descriptions often ambiguous definitions, requiring agents to interpret the boundary between categories; and (2) **Expert-written decision rules** provides the guidelines to resolve conflicts when multiple headings seemingly apply, which is a consolidated HSCode classification protocol distilled from years of accumulated domain expertise. Partial case can be found in Figure 12.

**Domain-Specific Knowledge ( $\mathcal{K}$ ).** To correctly interpret and apply  $\mathcal{R}$ , agents must access external domain knowledge ( $\mathcal{K}$ ), i.e., a repository of historical customs rulings target US rulings, i.e., the CROSS system<sup>3</sup>. These rulings serve as few-shot examples, helping agents navigate edge cases where rules are ambiguous.

**Output Space ( $\mathcal{Y}$ ).** The target  $y \in \mathcal{Y}$  is a unique 10-digit numeric string. A valid prediction must

<sup>3</sup><https://rulings.cbp.gov/home>

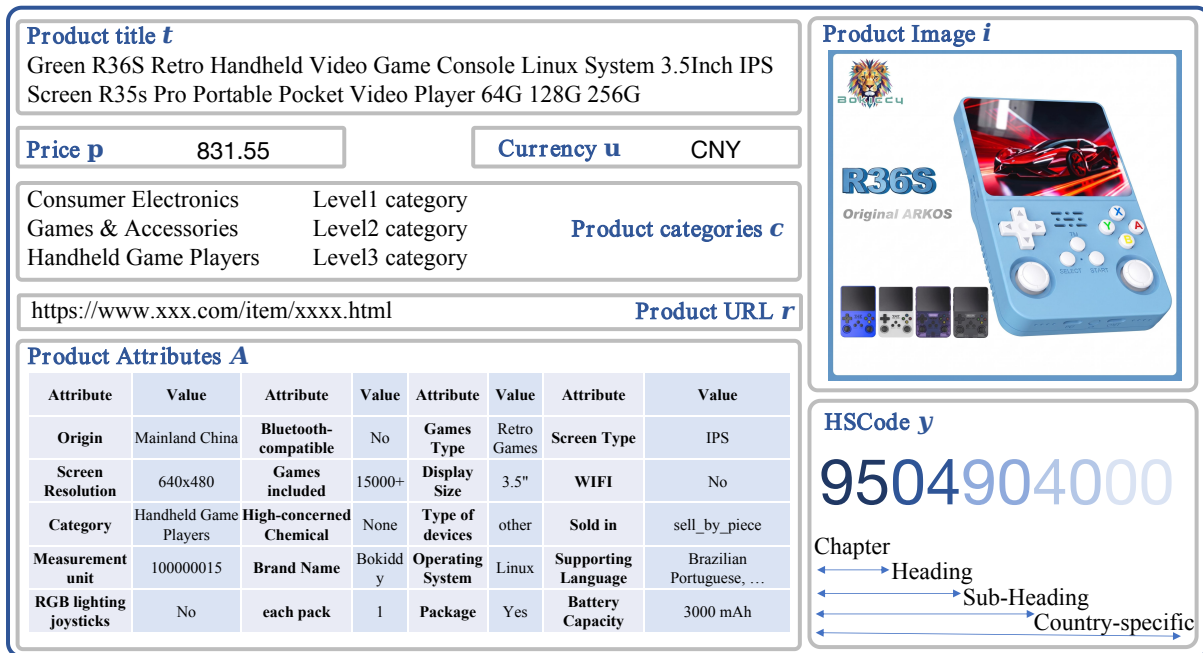


Figure 2. An illustrative example of a game console product profile from HSCODECOMP, which consists of rich metadata, including title, price, categories, attributes and image.

represent a path in the HS taxonomy tree, satisfying all constraints defined in  $\mathcal{R}$ .

## 4 HSCODECOMP Benchmark

### 4.1 Benchmark Construction

We design a rigorous pipeline, ensuring the dataset is diverse, realistic and expert-level: (1) Realistic Sourcing and Filtering; (2) Human Expert Annotation; and (3) Human Expert Validation.

**Realistic Sourcing and Filtering.** We source product entries from a large-scale global e-commerce platform, adopting a sampling strategy anchored to actual sales volume and category statistics to reflect real-world distributions. Distinct from synthetic datasets, we explicitly retain intrinsic noisy metadata (e.g., verbose titles, inconsistent attributes) to preserve the complexity of realistic scenarios. To ensure taxonomic diversity and prevent topical skew, we apply a semantic redundancy filter, discarding instances that share identical categories and HSCodes with existing entries. This ensures the benchmark captures the challenges of valid long-tail distributions rather than being dominated by frequent and trivial items.

**Human Expert Annotation.** We engage human experts specialized in HSCode classification to annotate the HSCode ( $y$ ) for each product profile ( $x$ ). As shown in Figure 3 (left), the annotation

process follows a strict protocol: (1) **Step 1:** two experts gather comprehensive information from the product webpage; (2) **Step 2:** they extract the core structured features of products; (3) **Step 3:** experts query the official ruling databases (CROSS) for related cases. If a related case is found, the corresponding HSCode is then validated on eWTP system, followed by the revision. Otherwise, they refine queries and revisit Step 2 to adjust the extracted features; (4) **Step 4:** for products without any related cases, experts execute expert-written decision rules (Figure 12) to apply tariff rules, and determine the appropriate HSCode; and (5) **Step 5:** experts finally verify the final identified HSCodes on the eWTP website to ensure its validity.

**Human Expert Validation.** To guarantee label reliability, we implement a dual-verification mechanism. As shown in the Figure 3 (Step 6), Two experts independently annotate each sample; if their codes match, the instance is accepted. In cases of disagreement, a senior tariff expert adjudicates the conflict. If no consensus is reached, the sample is discarded. Furthermore, to verify the reliability of our process, we conduct an additional quality review. A fourth senior expert, not involved in the initial annotation, re-annotated a random 10% sample of the dataset. This review shows only a 2% disagreement rate, confirming the effectiveness and consistency of our dataset construction pipeline.

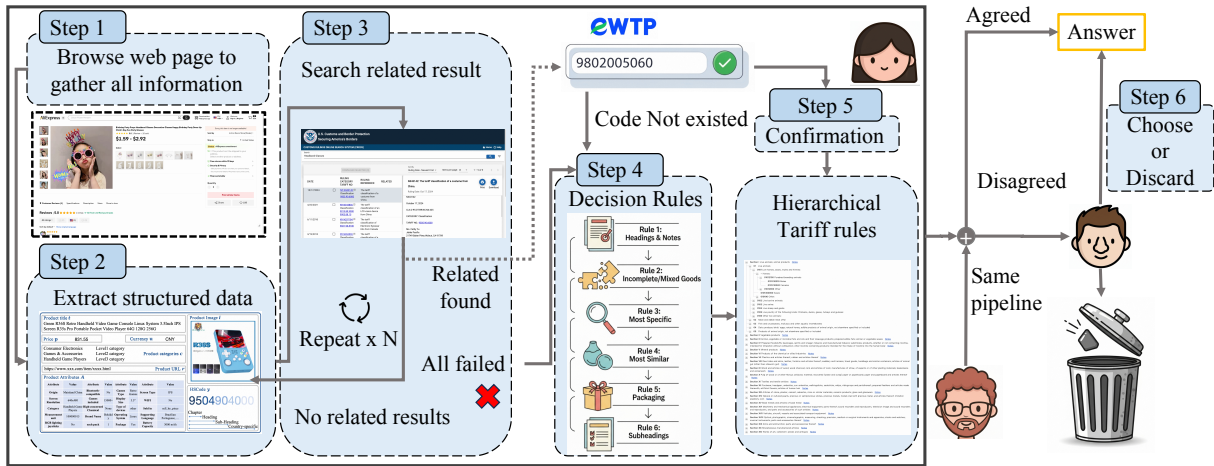


Figure 3. The pipeline for human experts to annotate the HSCodes, including two human experts for HSCode annotation (Step 1 to 5) and one additional expert for quality validation (Step 6).

## 4.2 Evaluation Metric

We adopt Exact Match Accuracy as the primary metric. Given the hierarchical nature of HS Codes, we report accuracy at multiple granularities: 2-digit (Chapter), 4-digit (Heading), 6-digit (Sub-heading), and the full 10-digit (Country-specific) level. The 10-digit accuracy serves as the strictest measure of an agent’s ability to perform deep, end-to-end hierarchical reasoning.

## 5 Experimental Setup

We conduct a comprehensive evaluation across three distinct paradigms.

**Foundation Models.** We test 14 leading LLMs and VLMs (e.g., GPT-5, Gemini-2.5-Pro, Kimi-K2 (Team et al., 2025), Claude-Sonnet-4, DeepSeek series (Guo et al., 2025), Qwen series (Yang et al., 2025), o3-mini, Nemotron-32B (Bercovich et al., 2025) etc. relying solely on internal knowledge. For VLMs, product images are provided to assess visual reasoning capabilities.

**Open-source Agent Frameworks.** We evaluate six advanced frameworks (e.g., SmoAgents (Roucher et al., 2025), Aworld (Yu et al., 2025c), Agentorchestra (Zhang et al., 2025), OWL (Hu et al., 2025b), WebSailor (Li et al., 2025b) and Cognitive Kernel (Fang et al., 2025)) using GPT-5 as the default backbone. All frameworks utilize standardized tools including Google Search, Knowledge and Tariff Rules Retrieval.

**Closed-source Agent Systems.** Following previous works (Li et al., 2025b), we assess the performance of commercial systems (Manus, Gemini

Deep Research, and Grok DeepSearch) on a representative subset of benchmark.

All systems are evaluated using the strict 10-digit Exact Match Accuracy metric. CROSS database is integrated via a web-browsing tool, allowing agents to dynamically query and retrieve relevant rulings during reasoning. More implementation details appear in Appendix G.

## 6 Main Experimental Results

In this section, we present a rigorous evaluation on HSCODECOMP. We analyze performance across three dimensions: (1) the absolute gap between state-of-the-art agents and human experts; (2) the comparison analysis of open-source versus closed-source frameworks; (3) backbone model selection; and (4) test-time scaling analysis.

### 6.1 Human-Agent Performance Gap.

Table 1 reveal a substantial disparity between agents and human experts in hierarchical rule application: (1) **Significant Performance Gap:** The state-of-the-art system, SmoAgents (GPT-5 w/ Vision), achieves only **46.83%** 10-digit accuracy, significantly lagging behind the **95.0%** accuracy of human experts, highlighting that current agents struggle to satisfy the rigorous constraints of expert-level tasks; (2) **Hierarchical Challenges:** Performance decays consistently as taxonomic depth increases. While agents maintain competitive accuracy at the 2-digit Chapter level ( $\sim 82\%$ ), they fail to navigate the fine-grained logical boundaries required for 10-digit country-specific classification. Moreover, to further accurately diagnose where the reasoning breaks, we further analyze the Conditional

Baselines	Type	HSCode Prediction Accuracy				
		2-digit	4-digit	6-digit	8-digit	10-digit
<i>LLM/VLM-Only</i>						
GPT-5	VLM	82.12	70.89	59.97	41.46	29.27
GPT-5	LLM	82.59	69.78	56.33	40.98	28.96
Gemini-2.5-PRO	VLM	<b>82.28</b>	71.04	59.02	40.51	24.21
Gemini-2.5-PRO	LLM	80.54	69.94	58.54	40.35	23.42
GPT-4o	VLM	78.01	64.08	48.10	29.75	18.51
GPT-4o	LLM	75.47	61.55	45.73	30.06	18.35
Claude Sonnet 4	VLM	78.80	64.08	45.25	22.63	11.23
Claude Sonnet 4	LLM	78.80	62.97	44.94	23.58	11.87
Kimi-K2 (Team et al., 2025)	LLM	78.01	62.03	44.15	24.53	12.18
DeepSeek-R1 (Guo et al., 2025)	LLM	77.22	61.71	38.45	16.77	6.65
Qwen-Max (Yang et al., 2025)	LLM	71.52	48.58	24.21	11.23	3.80
O3-mini	LLM	77.22	56.17	24.53	6.65	1.27
<i>Open-source Agent System (GPT-5 Backbone)</i>						
SmolAgents (Roucher et al., 2025)	VLM	82.06	<b>72.06</b>	<b>62.38</b>	<b>52.38</b>	<b>46.83</b>
SmolAgents (Roucher et al., 2025)	LLM	<b>82.28</b>	70.89	59.81	49.05	42.72
Aworld (Yu et al., 2025c)	LLM	<b>82.28</b>	70.41	59.18	48.58	41.30
Agentorchestra (Zhang et al., 2025)	LLM	82.12	70.73	60.44	47.78	41.30
OWL (Hu et al., 2025b)	LLM	72.63	61.87	51.58	41.77	37.34
WebSailor (Li et al., 2025b)	LLM	81.64	70.56	57.27	43.98	35.44
Cognitive Kernel (Fang et al., 2025)	LLM	80.06	69.15	54.59	40.03	26.42

Table 1. Performance comparison of various foundation models and agent systems on the HSCODECOMP benchmark. Accuracy (%) is reported at different HS Code hierarchical levels (2-digit to 10-digit). Standalone LLM/VLM performances are compared against open-source agent frameworks powered by a GPT-5 backbone.

Accuracy. This metric measures the probability of correctly predicting a more fine-grained code (e.g., 4-digit) contingent upon the correctness of its preceding digits (e.g., 2-digit). As shown in Table 10, it can be found that current agents struggle with the complex, multi-step reasoning required for precise rule application at deeper hierarchical levels; and (3) **Agents Outperform LLMs**: Agentic frameworks consistently outperform standalone foundation models, validating the necessity of tool usage like domain-specific knowledge retrieval and tariff rules retrieval.

## 6.2 Open-source vs. Closed-source Agents

We compare open-source frameworks against commercial "Deep Research" systems on a representative subset. Surprisingly, open-source agents demonstrate superior rule adherence. As detailed in Table 2, SmolAgents and AWorld (both ~42.9%) outperform proprietary systems like Gemini Deep Research (40.8%) and Manus (30.6%). Qualitative analysis suggests commercial systems often suffer from premature commitment and information

misprocessing, as detailed in Section 7.3.

System Architecture	10-digit Acc.
<i>Closed-Source Commercial Systems</i>	
Gemini Deep Researcher	40.81
Grok DeepSearch	26.53
Manus	30.61
<i>Open-Source Frameworks (w/ GPT-5)</i>	
SmolAgents	<b>42.86</b>
AWorld	<b>42.86</b>

Table 2. Performance comparison: Commercial Systems vs. Open-Source Agent Frameworks with GPT-5.

## 6.3 Impact of Backbone Models

We evaluate the SmolAgents framework across four state-of-the-art backbone models to analyze the contribution of the underlying backbone models. Results in Table 3 demonstrate that GPT-5 is currently the only viable backbone for this high-complexity task. It achieves 42.72% accuracy, establishing a decisive lead over Gemini-2.5-Pro

Backbone Model	10-digit (%)	$\Delta$
GPT-5 (Default)	<b>42.72</b>	–
Gemini-2.5-Pro	34.49	– <b>8.23</b>
Claude 4 Sonnet	33.70	– <b>9.02</b>
Qwen-Max	17.43	– <b>25.29</b>

Table 3. Ablation on backbone models.  $\Delta$  denotes the performance gap compared to GPT-5.

(34.49%), Claude 3.5 Sonnet (33.70%), and Qwen-Max (17.43%). Therefore, we choose GPT-5 as the default setup for open-source agents. Complete results are provided in Table 12 in Appendix E.3.

#### 6.4 Inefficacy of Test-Time Scaling

While test-time scaling (TTS) has proven effective for reasoning tasks (Guo et al., 2025), our experiments reveal that two standard TTS strategies fail to improve performance on HSCODECOMP (Liu et al., 2025b): (1) **Majority Voting**: We implement majority voting across  $K = \{1, 2, 4, 8, 16\}$  trials (Voting@ $K$ ), using SmolAgent (GPT-5). Figure 4 shows that increasing  $K$  yields negligible improvement; and (2) **Self-Reflection**: We integrate a self-reflection mechanism into SmolAgent (GPT-5), enabling the model to proactively evaluate and revise its reasoning and actions. However, this approach slightly decreases performance from 42.72% to 42.57%. These results highlight the need of more effective test-time scaling strategy for hierarchical rule application.

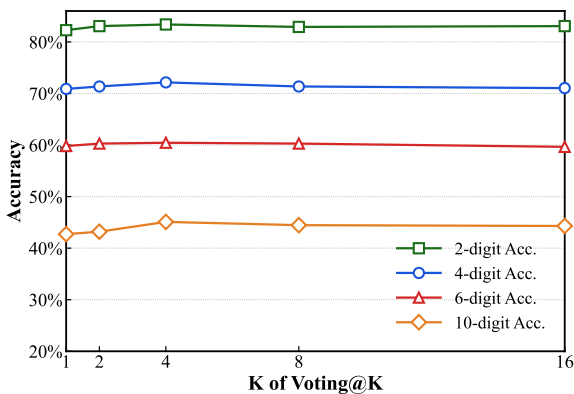


Figure 4. Test-time scaling (majority voting (Zhu et al., 2025)) analysis on HSCODECOMP.

## 7 Analysis on HSCODECOMP

In this section, we systematically diagnose the root causes of the challenges in hierarchical rule application by: (1) investigating the performance degrada-

tion due to "reasoning drift"; (2) conducting ablation studies on the contributions of tariff rules, decision rules, knowledge bases, and visual signals; (3) establishing a qualitative failure taxonomy that categorizes errors into perception, reasoning, and reflection failures; and (4) performing a per-category analysis to validate the benchmark's rigorous difficulty distribution.

### 7.1 Overthinking Decrease Performance

Current agentic paradigms often operate on the assumption that providing richer reasoning improves decision quality. As shown in Table 1, the reasoning-heavy WebSailor framework significantly underperforms the direct-execution SmolAgents (35.44% vs. 42.72%). One case study in Appendix I suggests a phenomenon of "Reasoning Drift": without the grounding of accurate tool feedback (e.g. rules definitions), agents tend to hallucinate constraints or misinterpret vague product details. To validate this, we conducted an ablation study on WebSailor by constraining its reasoning depth: (1) **No-Think** (direct tool invocation); and (2) **Medium-Think** (moderate reasoning controlled by the prompt). As shown in Table 4, restricting reasoning depth yields significant gains, with the No-Think variant approach SmolAgents. This empirical evidence underscores a critical design principle: for domain-specific tasks where valuable information resides in domain knowledge and rules, prioritizing knowledge and rules utilization over internal reasoning is essential to prevent hallucination and reasoning drift.

Agent System	10-digit Accuracy (Depth $\rightarrow$ )			Gain
	Overthink	Medium	No-Think	
SmolAgents	–	–	<b>42.72</b>	–
WebSailor	35.44	37.34	<b>40.82</b>	<b>+5.38</b>

Table 4. Impact of reducing thinking depth. Less thinking leads to higher accuracy for WebSailor.

### 7.2 Ablation Study

To investigate how different knowledge and rules in HSCODECOMP affect performance, we perform ablation studies to isolate their contribution.

**The Naive Rule Injection Gap.** Intuitively, providing expert-written hierarchical decision rules should guide the agent's reasoning. However, Table 5 reveals that injecting these explicit guide-

Agent System	Type	10-digit Accuracy		
		w/o DR	w/ DR	$\Delta$
SmolAgents	VLM	<b>46.83</b>	43.83	<b>-3.00</b>
WebSailor	LLM	<b>35.44</b>	35.34	<b>-0.10</b>

Table 5. The ablation study on expert-written Decision Rules (DR). All systems use GPT-5 as the backbone.

lines degrades performance. For example, SmolAgent with GPT-5 degrades from 46.83%  $\rightarrow$  43.83%. This highlights that current agents struggle to integrate abstract rules with noisy product contexts and primary tariff rules simultaneously, often suffering from information overload or misinterpreting rules as passive text rather than enforceable constraints.

**Hierarchical Tariff Rules are Essential.** Removing the Tariff Rules results in a significant performance drop. SmolAgent with GPT-5 degrades from 46.83%  $\rightarrow$  36.73%. It confirms that HSCode prediction is fundamentally a rule-centric task rather than a retrieval-centric one. Agents cannot rely solely on open-web knowledge or model’s internal knowledge; they must strictly adhere to the specific, hierarchical tariff rules to resolve fine-grained classification paths.

**Supplementary Role of Knowledge Base.** Access to the CROSS knowledge base yields only a marginal performance benefit (SmolAgent w/ GPT-5, 46.53%  $\rightarrow$  46.83%), indicating that while historical precedents provide useful auxiliary context for resolving specific ambiguities, they remain supplementary to the core reasoning process.

**Visual Signal are Helpful.** Table 1 and Table 6 show that incorporating product images yields a significant gain. For example, SmolAgent with GPT-5 improves from 42.72%  $\rightarrow$  46.83% when visual capability is activated. Case studies in Appendix J show that visual attributes that are not present in the textual description (such as material and surface features) are critical for classification.

### 7.3 Qualitative Analysis: Failure Taxonomy

To systematically diagnose agent failures, we contextualize our findings within the taxonomy of model hallucinations (Rawte et al., 2023) and classify the identified error patterns into three categorizations: Perception, Reasoning, and Reflection.

**Perception & Grounding Failures.** Two kinds of errors often originate before reasoning: (1) In-

Backbone Model	10-digit Accuracy		
	w/o Image	w/ Image	$\Delta$
GPT-5	42.72	<b>46.83</b>	<b>+4.11</b>
Claude 4 Sonnet	33.70	<b>34.65</b>	<b>+0.95</b>
GPT-4o	22.03	<b>22.31</b>	<b>+0.28</b>

Table 6. Ablation study on the impact of visual signals. The table compares the 10-digit Exact Match Accuracy of GPT-5, Claude-Sonnet-4 and GPT-4o models evaluated with and without product images.

**formation Misprocessing:** Agents overlook or misinterpret product details. For instance, in Figure 15, the agent becomes fixated on marketing buzzwords ("INS style", "Aesthetic") while overlooking the decisive technical specification ("Material: PVC"), leading to a misclassification based on style rather than substance; (2) **Lack of Domain Knowledge:** As seen in Figure 16, models may misidentify *silicone* products as *rubber* (Chapter 40) instead of *plastics* (Chapter 39), lacking the correct domain-specific ontology to distinguish material properties essential for customs classification.

**Reasoning & Planning Failures.** Even when input is correctly perceived, agents struggle to navigate the logical decision on hierarchical rules: (1) **Premature Decisions:** This is a planning failure where agents lock onto a high-level category too early, as seen in Table 14; (2) **Wrong Rule Application:** As detailed in Figure 17, a prevalent failure is the agent’s inability to dynamically determine the correct rule precedence given the product description; and (3) **Reasoning Hallucination:** Table 13 reveal that agents generate plausible but factually incorrect reasoning (Rawte et al., 2023).

**Reflection.** Finally, agents frequently conduct wrong self-correction: agents predict correct HSCodes initially but revise them incorrectly through excessive critique (Table 13).

### 7.4 Per-category Performance Analysis

We analyze two critical distributions across the 32 first-level product categories to diagnose whether performance correlates with data distribution: (1) **Challenging Product Distribution (CID):** the distribution of hard instances where all baselines failed; (2) **Average Performance Distribution (APD):** the distribution of average 10-digit accuracy across all baselines. Figure 5 reveals two key insights: (1) The CID indicates that the most chal-

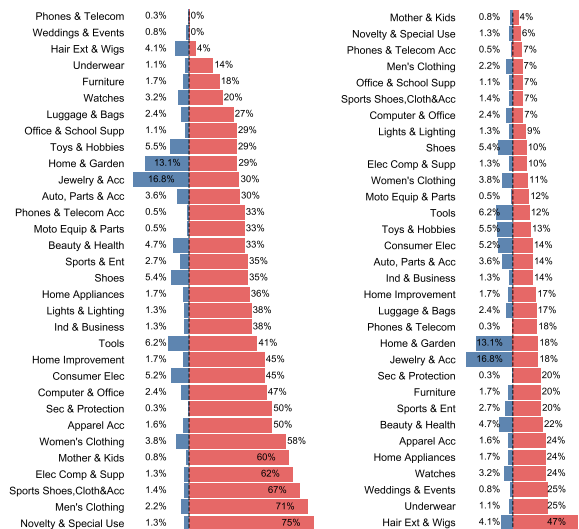


Figure 5. Category distribution on the left blue bars. **Left:** Challenging Product Distribution (CID). **Right:** Average Performance Distribution (APD).

lenging products are concentrated in long-tail categories, such as *Novelty & Special Use* (1.3%) and *Men's Clothing* (2.2%). These categories, characterized by non-standard descriptions and sparse examples, expose the agents' inability to generalize rules to data-sparse domains; (2) The APD shows that average accuracy across most product categories remains below 25%, with only *Hair Extensions & Wigs* achieving a relatively high accuracy of 47%. Importantly, even dominant categories like *Jewelry & Accessories*, *Home & Garden* and *Tools* still stays below 18%. These findings highlight the challenging nature of HSCODECOMP.

## 8 Conclusion

We identified and addressed the critical gap in evaluating agents in hierarchical rule applications. To address this gap, we introduced HSCODECOMP, the first realistic and expert-level benchmark designed to assess agents for multi-hop reasoning with hierarchical tariff rules in e-commerce domain. Our extensive evaluation revealed a substantial performance gap between current state-of-the-art agents (46.8%) and human experts (95.0%), highlighting that hierarchical rule application remains a significant challenge for existing agents.

### Limitation

**Temporal Sensitivity.** As noted in Section 4 and Appendix G, our ground truth is anchored to a specific timestamp (June 30, 2025), which is a standard

implementation in the research community (Wong et al., 2025). While this limits our ability to evaluate an agent's real-time temporal adaptability, it does not compromise the validity of the benchmark. By strictly anchoring the evaluation to a specific tariff version<sup>4</sup>, we effectively isolate the agent's reasoning capability from the volatility of external knowledge, ensuring the benchmark measures the agent's ability to interpret and apply a given set of rules, regardless of their temporal version.

**Trade-off Between Benchmark Quality and Scale.** Diverging from the prevalence of crowd-sourced (Zhou et al., 2025; Phan et al., 2025; Wu et al., 2025), HSCODECOMP prioritizes label quality over scale. The average hourly wage of annotators is  $\approx$  \$34.6/hour. The dataset size reflects a deliberate trade-off necessitated by the resource-intensive nature of verification by domain experts. While our statistical analysis in Appendix C confirms that this sample size provides sufficient power for stable model differentiation, we acknowledge that the benchmark may not exhaustively capture the extreme long-tail of global trade.

**The Scope of Tariff Rules.** Currently, HSCODECOMP is exclusively scoped to the United States Harmonized Tariff Schedule (HTSUS). While the HS system is globally standardized at the 6-digit level, classifications at the 8-to-10 digit level remain country-specific. Future work will need to expand to multi-lingual and multi-jurisdictional scenarios (e.g., EU or China customs) to assess cross-border regulatory generalization.

### Ethical considerations

This research adheres to strict ethical guidelines regarding data privacy and fair labor. The dataset is fully anonymized and contains no personally identifiable information. The hourly wage of our human annotators is over 34.6 USD, which is much higher than average hourly wage 3.13 USD on Amazon Mechanical Turk (Hara et al., 2017). This remuneration structure was designed to provide a fair and competitive wage, acknowledging the expertise and effort required for this task and ensuring that contributors were rewarded appropriately for their work. Moreover, in preparing this manuscript, Qwen and Gemini were used solely as a writing assistant to improve grammar and clarity. The LLMs

<sup>4</sup><https://www.ewtp.com/web>

was not used for generating code, concepts, or any part of the core research methodology.

## References

- Akhiad Bercovich, Itay Levy, Izik Golan, Mohammad Dabbah, Ran El-Yaniv, Omri Puny, Ido Galil, Zach Moshe, Tomer Ronen, Najeeb Nabwani, Ido Shahaf, and Oren Tropp et al. 2025. [Llama-nemotron: Efficient reasoning models](#). *Preprint*, arXiv:2505.00949.
- Kaiyuan Chen, Yixin Ren, Yang Liu, Xiaobo Hu, Haotong Tian, Tianbao Xie, Fangfu Liu, Haoye Zhang, Hongzhang Liu, Yuan Gong, Chen Sun, Han Hou, Hui Yang, James Pan, Jianan Lou, Jiayi Mao, Jizheng Liu, Jinpeng Li, Kangyi Liu, and 14 others. 2025a. [xbench: Tracking agents productivity scaling with profession-aligned real-world evaluations](#). *Preprint*, arXiv:2506.13651.
- Shan Chen, Pedro Moreira, Yuxin Xiao, Sam Schmidgall, Jeremy Warner, Hugo Aerts, Thomas Hartvigsen, Jack Gallifant, and Danielle S. Bitterman. 2025b. [Medbrowsecomp: Benchmarking medical deep research and computer use](#). *Preprint*, arXiv:2505.14963.
- Tianqing Fang, Zhisong Zhang, Xiaoyang Wang, Rui Wang, Can Qin, Yuxuan Wan, Jun-Yu Ma, Ce Zhang, Jiaqi Chen, Xiyun Li, Hongming Zhang, Haitao Mi, and Dong Yu. 2025. [Cognitive kernel-pro: A framework for deep research agents and agent foundation models training](#). *Preprint*, arXiv:2508.00414.
- Sina Gholamian, Gianfranco Romani, Bartosz Rudnikowicz, and Laura Skylaki. 2024. [Llm-based robust product classification in commerce and compliance](#). *arXiv preprint arXiv:2408.05874*.
- Andrew Grainger. 2024. [Customs tariff classification and the use of assistive technologies](#). *World Customs Journal*, 18(1):3–32.
- Neel Guha, Julian Nyarko, Daniel E. Ho, Christopher Re, Adam Chilton, Aditya Narayana, Alex Chohlas-Wood, Austin Peters, Brandon Waldon, and et al. 2023. [Legalbench: A collaboratively built benchmark for measuring legal reasoning in large language models](#). In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shitong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, and Zhibin Gou et al. 2025. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#). *Preprint*, arXiv:2501.12948.
- Kotaro Hara, Abi Adams, Kristy Milland, Saiph Savage, Chris Callison-Burch, and Jeffrey Bigham. 2017. [A data-driven analysis of workers’ earnings on amazon mechanical turk](#). *Preprint*, arXiv:1712.05796.
- Liang Hu, Jianpeng Jiao, Jiashuo Liu, Yanle Ren, Zhoufutu Wen, Kaiyuan Zhang, Xuanliang Zhang, Xiang Gao, Tianci He, Fei Hu, Yali Liao, Zaiyuan Wang, Chenghao Yang, Qianyu Yang, Mingren Yin, Zhiyuan Zeng, Ge Zhang, Xinyi Zhang, Xiyang Zhao, and 4 others. 2025a. [Finsearchcomp: Towards a realistic, expert-level evaluation of financial search and reasoning](#). *Preprint*, arXiv:2509.13160.
- Mengkang Hu, Yuhang Zhou, Wendong Fan, Yuzhou Nie, Bowei Xia, Tao Sun, Ziyu Ye, Zhaoxuan Jin, Yingru Li, Qiguang Chen, Zeyu Zhang, Yifeng Wang, Qianshuo Ye, Bernard Ghanem, Ping Luo, and Guohao Li. 2025b. [Owl: Optimized workforce learning for general multi-agent assistance in real-world task automation](#). *Preprint*, arXiv:2505.23885.
- Kung-Hsiang Huang, Akshara Prabhakar, Sidharth Dhawan, Yixin Mao, Huan Wang, Silvio Savarese, Caiming Xiong, Philippe Laban, and Chien-Sheng Wu. 2025. [Crmarena: Understanding the capacity of llm agents to perform professional crm tasks in realistic environments](#). In *Proceedings of NAACL 2025 (Long Papers)*, Albuquerque, New Mexico. Association for Computational Linguistics.
- Mandar Joshi and 1 others. 2017. [Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension](#). In *ACL*.
- Bryce Judy. 2024. [Benchmarking harmonized tariff schedule classification models](#). *arXiv preprint arXiv:2412.14179*.
- Eunji Lee, Sihyeon Kim, Sundong Kim, Soyeon Jung, Heeja Kim, and Meeyoung Cha. 2024. [Explainable product classification for customs](#). *ACM Transactions on Intelligent Systems and Technology*, 15(2):1–24.
- Haitao Li, Junjie Chen, Jingli Yang, Qingyao Ai, Wei Jia, Youfeng Liu, Kai Lin, Yueyue Wu, Guozhi Yuan, Yiran Hu, Wuyue Wang, Yiqun Liu, and Minlie Huang. 2025a. [Legalagentbench: Evaluating llm agents in legal domain](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Vienna, Austria. Association for Computational Linguistics.
- Kuan Li, Zhongwang Zhang, Huifeng Yin, Liwen Zhang, Litu Ou, Jialong Wu, Wenbiao Yin, Baixuan Li, Zhengwei Tao, Xinyu Wang, Weizhou Shen, Junkai Zhang, Dingchu Zhang, Xixi Wu, Yong Jiang, Ming Yan, Pengjun Xie, Fei Huang, and Jingren Zhou. 2025b. [Websailor: Navigating super-human reasoning for web agent](#). *Preprint*, arXiv:2507.02592.
- Mengjie Liao, Lei Huang, Jian Zhang, Luona Song, and Bo Li. 2024. [Enhanced hs code classification for import and export goods via multiscale attention and ernie-bilstm](#). *Applied Sciences*, 14(22):10267.
- Ben Liu, Jihai Zhang, Fangquan Lin, Cheng Yang, Min Peng, and Wotao Yin. 2025a. [Symagent: A neural-symbolic self-learning agent framework for](#)

- complex reasoning over knowledge graphs. *Preprint*, arXiv:2502.03283.
- Zijun Liu, Peiyi Wang, Runxin Xu, Shirong Ma, Chong Ruan, Peng Li, Yang Liu, and Yu Wu. 2025b. Inference-time scaling for generalist reward modeling. *Preprint*, arXiv:2504.02495.
- Youngang Lyu, Xiaoyu Zhang, Lingyong Yan, Maarten de Rijke, Zhaochun Ren, and Xiuying Chen. 2025. Deepshop: A benchmark for deep research shopping agents. *Preprint*, arXiv:2506.02839.
- Grégoire Mialon, Clémentine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom. 2023a. Gaia: a benchmark for general ai assistants. *arXiv preprint arXiv:2311.12983*.
- Grégoire Mialon, Clémentine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom. 2023b. Gaia: a benchmark for general ai assistants. *Preprint*, arXiv:2311.12983.
- Ralph Peeters, Aaron Steiner, Luca Schwarz, Julian Yuya Caspary, and Christian Bizer. 2025. Web-mall – a multi-shop benchmark for evaluating web agents. *arXiv preprint arXiv:2508.13024*.
- Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Sean Shi, Michael Choi, and 1 others. 2025. Humanity’s last exam. Center for AI Safety and Scale AI (whitepaper / dataset).
- Vipula Rawte, Amit Sheth, and Amitava Das. 2023. A survey of hallucination in large foundation models. *Preprint*, arXiv:2309.05922.
- Aymeric Roucher, Albert Villanova del Moral, Thomas Wolf, Leandro von Werra, and Erik Kaunismäki. 2025. ‘smolagents’: a smol library to build great agentic systems. <https://github.com/huggingface/smolagents>.
- Albert Sadowski and Jarosław A. Chudziak. 2025. Explainable rule application via structured prompting: A neural-symbolic approach. *Preprint*, arXiv:2506.16335.
- Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. 2021. ALFWorld: Aligning Text and Embodied Environments for Interactive Learning. In *ICLR*.
- Shubham, Avinash Arya, Subarna Roy, and Sridhar Jonnala. 2022. An ensemble-based approach for assigning text to correct harmonized system code. *arXiv preprint arXiv:2211.04313*.
- Sédric Stassin, Otmane Amel, Sidi Mahmoudi, and Xavier Siebert. 2023. Similarity versus supervision: Best approaches for hs code prediction.
- Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, and Ruijue Chen et al. 2025. Kimi k2: Open agentic intelligence. *Preprint*, arXiv:2507.20534.
- George Thomas, Alex J. Chan, Jikun Kang, Wenqi Wu, Filippos Christianos, Fraser Greenlee, Andy Toulis, and Marvin Purtorab. 2025. Webgames: Challenging general-purpose web-browsing ai agents. In *arXiv preprint arXiv:2502.18356*.
- Jason Wei, Zhiqing Sun, Spencer Papay, Scott McKinney, Jeffrey Han, Isa Fulford, Hyung Won Chung, Alex Tachard Passos, William Fedus, and Amelia Glaese. 2025. Browsecomp: A simple yet challenging benchmark for browsing agents. *Preprint*, arXiv:2504.12516.
- Ryan Wong, Jiawei Wang, Junjie Zhao, Li Chen, Yan Gao, Long Zhang, Xuan Zhou, Zuo Wang, Kai Xi-ang, Ge Zhang, Wenhao Huang, Yang Wang, and Ke Wang. 2025. Widesearch: Benchmarking agentic broad info-seeking. *Preprint*, arXiv:2508.07999.
- Jialong Wu, Wenbiao Yin, Yong Jiang, Zhenglin Wang, Zekun Xi, Runnan Fang, Linhai Zhang, Yulan He, Deyu Zhou, Pengjun Xie, and Fei Huang. 2025. Webwalker: Benchmarking llms in web traversal. *Preprint*, arXiv:2501.07572.
- Chen Xu, Tian Lan, Yu Ji, Changlong Yu, Wei Wang, Jun Gao, Qunxi Dong, Kun Qian, Piji Li, Wei Bi, and Bin Hu. 2025. Decider: A dual-system rule-controllable decoding framework for language generation. *Preprint*, arXiv:2403.01954.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, and Bowen Yu et al. 2025. Qwen3 technical report. *Preprint*, arXiv:2505.09388.
- Hong Yao, Zipeng Jing, Tianshi Liu, Eric J Wong, Chunyuan Hong, Is Wang, Wen-Loong Wang, Yalda Talebiraad, Rui Wang, Zhaowei Chen, and 1 others. 2024. Webvoyager: Building an end-to-end web agent with large multimodal models. In *ICLR*.
- Shunyu Yao, Howard Chen, John Yang, and Karthik R. Narasimhan. 2022. Webshop: Towards scalable real-world web interaction with grounded language agents. In *NeurIPS*.
- Congchi Yin, Feng Li, Shu Zhang, Zike Wang, Jun Shao, Piji Li, Jianhua Chen, and Xun Jiang. 2024. Mdd-5k: A new diagnostic conversation dataset for mental disorders synthesized via neuro-symbolic llm agents. *Preprint*, arXiv:2408.12142.
- Ailing Yu, Lan Yao, Jingnan Liu, Zhe Chen, Jiajun Yin, Yuan Wang, Xinhao Liao, Zhiling Ye, Ji Li, Yun Yue, Hansong Xiao, Hualei Zhou, Chunxiao Guo, Peng Wei, Junwei Liu, and Jinjie Gu. 2025a. Medresearcher-r1: Expert-level medical deep researcher via a knowledge-informed trajectory synthesis framework. *Preprint*, arXiv:2508.14880.
- Ailing Yu, Lan Yao, Jingnan Liu, Zhe Chen, Jiajun Yin, Yuan Wang, Xinhao Liao, Zhiling Ye, Ji Li, Yun Yue, Hansong Xiao, Hualei Zhou, Chunxiao Guo, Peng Wei, Junwei Liu, and Jinjie Gu. 2025b.

Medresearcher-r1: Expert-level medical deep researcher via a knowledge-informed trajectory synthesis framework. *arXiv preprint arXiv:2508.14880*.

Chengyue Yu, Siyuan Lu, Chenyi Zhuang, Dong Wang, Qintong Wu, Zongyue Li, Runsheng Gan, Chunfeng Wang, Siqi Hou, Gaochi Huang, Wenlong Yan, Lifeng Hong, Aohui Xue, Yanfeng Wang, Jinjie Gu, David Tsai, and Tao Lin. 2025c. *Aworld: Orchestrating the training recipe for agentic ai*. *Preprint*, arXiv:2508.20404.

Wentao Zhang, Liang Zeng, Yuzhen Xiao, Yongcong Li, Ce Cui, Yilei Zhao, Rui Hu, Yang Liu, Yahui Zhou, and Bo An. 2025. *Agentorchestra: A hierarchical multi-agent framework for general-purpose task solving*. *Preprint*, arXiv:2506.12508.

Peilin Zhou, Bruce Leon, Xiang Ying, Can Zhang, Yifan Shao, Qichen Ye, Dading Chong, Zhiling Jin, Chenxuan Xie, Meng Cao, Yuxin Gu, Sixin Hong, Jing Ren, Jian Chen, Chao Liu, and Yining Hua. 2025. *Browsecomp-zh: Benchmarking web browsing ability of large language models in chinese*. *Preprint*, arXiv:2504.19314.

Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. 2023. *Webarena: A realistic web environment for building autonomous agents*. In *arXiv preprint arXiv:2307.13854*.

King Zhu, Hanhao Li, Siwei Wu, Tianshun Xing, Dehua Ma, Xiangru Tang, Minghao Liu, Jian Yang, Jiaheng Liu, Yuchen Eleanor Jiang, Changwang Zhang, Chenghua Lin, Jun Wang, Ge Zhang, and Wangchunshu Zhou. 2025. *Scaling test-time compute for llm agents*. *Preprint*, arXiv:2506.12928.

## A Reproducibility statement

We are committed to the principles of reproducible research. Accordingly, all necessary materials, including code, benchmark dataset and other related resources will be publicly released to promote the development of the deep search agents.

## B Dataset distribution

In this section, we provide following detailed data statistics of our proposed HSCODECOMP: (1) Total unique 10-digit codes is 352; (2) Sample distribution per HSCode is shown in Table 7. Since 72.73% of HSCodes appear only once in our dataset, HSCODECOMP is not intentionally enriched for tricky or dispute-prone items; and (3) HSCode subcategory Hierarchical Coverage is shown in Table 8. Besides, Figure 6 presents the distributions of first-level product categories (left) and HSCode chapter categories (right), which closely mirror real-world product distributions. This alignment confirms that HSCODECOMP accurately reflects practical international trade scenarios, ensuring that model performance evaluations reliably generalize to real-world applications. In summary, our sampling strategy for constructing HSCODECOMP directly follows sales volume and category distributions from a major global E-Commerce platform. This approach ensures the benchmark reflects actual production distributions.

Code	Count	Percentage
1	256	72.73%
2	58	16.48%
3	16	4.55%
4	4	1.14%
5	5	1.42%
6+	13	3.69%

Table 7. Sample Distribution per HSCode.

HSCode Subcategory	Count
2-digit (chapters)	27
4-digit (headings)	151
6-digit (subheadings)	268
8-digit (country-specific)	322
10-digit (country-specific)	352

Table 8. HSCode Hierarchical Coverage.

## C Statistical Stability Analysis

A primary concern in constructing expert-level benchmarks is balancing annotation cost with sta-

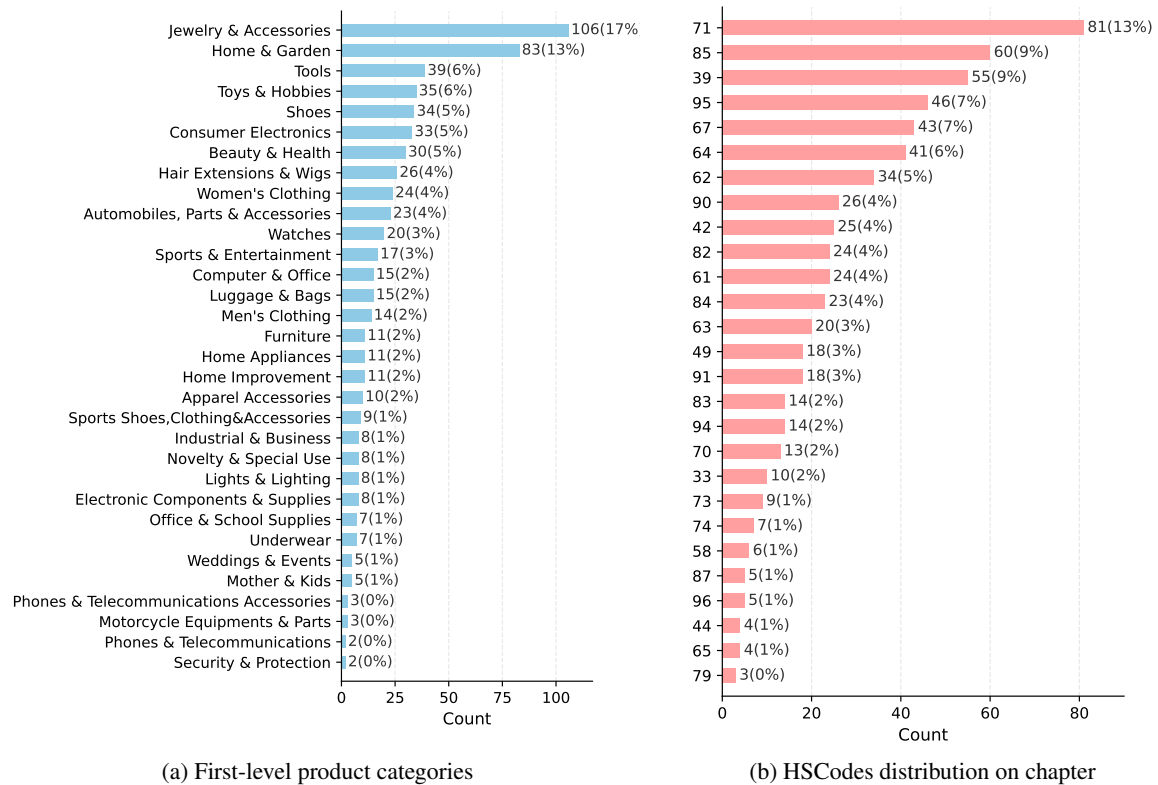


Figure 6. Distributions of the first-level product category and HSCode chapter categories.

tistical significance. To validate the robustness of HSCODECOMP (632 samples), we conducted a bootstrap analysis of the 95% Confidence Intervals (CI) across varying sample sizes for seven representative agent systems.

**Confidence Interval Convergence.** We calculated the 95% CI widths for sample sizes ranging from 200 to 632. As detailed in Table 9, the evaluation exhibits clear diminishing returns in CI width reduction beyond 600 samples. For the best-performing system (SmolAgents with VLM), the CI width decreases by only **0.21%** when expanding the dataset from 600 to 632 samples. This marginal improvement indicates that the evaluation has reached statistical equilibrium at the current dataset size. Further expansion would yield negligible gains in statistical precision.

**Comparison with Community Standards.** Despite the smaller absolute number of samples compared to automated benchmarks, our statistical reliability aligns with or exceeds established agentic benchmarks. For instance, top models on the GAIA benchmark exhibit 95% confidence intervals of approximately  $\pm 8.2\%$  (OpenAI GPT-4 Turbo) to  $\pm 8.7\%$  (Claude 3 Sonnet). In comparison, HSCODECOMP yields narrower or compara-

ble confidence intervals (e.g.,  $\pm 7.78\%$  width for our best agent), demonstrating that 632 expert-validated instances are sufficient to differentiate model performance reliably.

**Resource Allocation Efficiency.** It is also important to note the high cost of expert annotation for hierarchical rule application. As detailed in our Ethics Statement, annotators are domain professionals paid significantly above market rates (\$34.6/hour). Given the statistical stability demonstrated in Table 9, expanding the dataset further provides minimal scientific value relative to the substantial resource cost required for expert validation. We engaged a diverse pool of 26 domain experts (avg. 5+ years of experience). Crucially, these experts are experienced professionals with broad expertise across the Harmonized System (HS), enabling accurate classification across diverse code categories rather than being confined to narrow subsets of HS codes.

**Cost-Benefit Analysis.** We also conduct a cost analysis and found that the average cost per sample for GPT-5 (LLM) is \$0.11 and VLM is \$0.23. In contrast, expert annotation costs approximately \$1.38 per sample (calculated based on an expert rate of \$34/hour). While agents offer a signifi-

cant 6x-12x cost reduction, this comparison further highlights the trade-off between current model affordability and the necessity of high-cost human expertise for accuracy.

## D Semantic Distribution of Tariff Rules

To assess whether the HSCode taxonomy exhibits clear semantic separation, we generate embeddings of the official titles and notes for all HS chapters and sections using a sentence embedding model<sup>5</sup>. We then apply t-SNE to project these embeddings into two dimensions for visualization. As shown in Figure 7, each point represents a chapter, while each star marks a section’s centroid. The visualization reveals significant semantic overlap between adjacent sections: numerous chapters appear closer to neighboring section centroids than to their own section’s centroid, and section centroids themselves form overlapping clusters rather than distinct groupings. This pattern indicates that the semantic structure of hierarchical tariff rules lacks clear boundaries—adjacent sections frequently share similar vocabulary and concepts (*e.g.*, distinctions between raw materials and finished goods, or between component parts and complete articles).

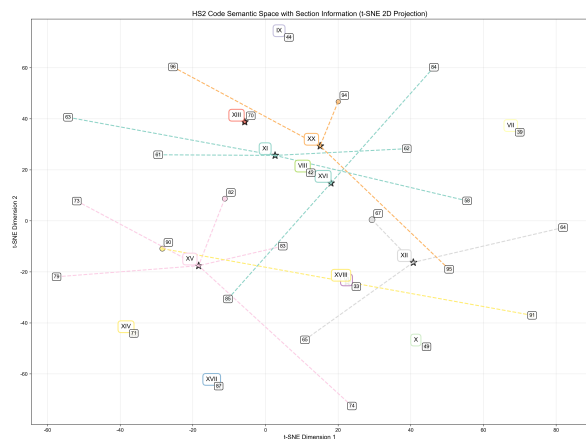


Figure 7. The semantic map of Harmonized System (HS) Codes chapter titles and notes.

## E Supplementary Experimental Results

### E.1 Conditional Performance

As shown in Table 10, we provide the conditional accuracy of partial models on Chapter, Heading, Sub-heading and Country-specific sub-categories code classification. The experimental results reveal

<sup>5</sup><https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>

three key insights: (1) The 2-4 digit heading level shows highest average accuracy (75.16), since its tariff rules are most clearly defined; (2) Intermediate performance drops at the 4-6 digit subheading level due to increased the number of rules and their vague descriptions in tariff rules; and (3) The steepest decline occurs at the country-specific 6-10 digit levels, where rules become highly contextual and country-specific. This hierarchical performance pattern directly supports our main finding: current agents struggle with the complex, multi-step reasoning required for precise rule application at deeper hierarchical levels. The significant accuracy drop from 64.8% (6-digit) to 35.3% (10-digit) demonstrates that correctly navigating the initial hierarchical steps does not guarantee success in the final classification, highlighting the cumulative reasoning challenges in hierarchical rule application.

### E.2 Detailed Ablation Study on Rules

We provide the complete ablation study on used domain-specific knowledge and rules. As shown in Table 11, we remove each used rules or knowledge in WebSailor and Smolagent open-source agent frameworks for solving tasks in HSCODECOMP: (1) **DR**: Human-written Decision Rules; (2) **Tariff**: hierarchical tariff rules; and (3) **CROSS** ruling database. Experimental results reveal three critical insights: (1) Human-written decision rules impair performance: As shown in Table 3 of our paper, adding decision rules (DR) decreases accuracy by 3.00% for Smolagent (GPT-5), demonstrating current agent systems cannot effectively utilize hierarchical rule structures; (2) Tariff rules are essential and contributes most: Removing tariff rules causes a severe 10.10% drop in accuracy. Without tariff rules, agents have to solve the task using its internal knowledge, leading to outdated and hallucination problems; and (3) CROSS provides moderate benefit: The 4.61% accuracy drop when removing CROSS indicates it offers supplementary but non-essential information. These results confirm that tariff rules are the most critical knowledge component, while human-written decision rules present the most significant challenge for current agent architectures to utilize effectively.

### E.3 Experiments on Backbone Models

To investigate how the backbone LLMs affect the performance of the agent system, we conduct more detailed ablation study on four open-source agent systems, replacing the original GPT-5 backbone

Model	200	300	400	500	600	632 (Full)	$\Delta_{600 \rightarrow 632}$
SmolAgents (VLM)	[42.6, 47.6] (11.17%)	[41.6, 51.4] (9.80%)	[46.1, 55.9] (8.75%)	[36.4, 55.4] (13.85%)	[43.3, 51.3] (7.99%)	[42.9, 50.7] (7.78%)	0.21%
AWorld	[37.1, 50.9] (13.76%)	[35.1, 43.7] (10.91%)	[37.9, 42.1] (9.70%)	[31.2, 47.6] (8.57%)	[36.2, 44.1] (7.85%)	[37.5, 45.1] (7.68%)	0.17%
WebSailor	[29.3, 47.8] (13.63%)	[31.1, 40.4] (10.61%)	[34.2, 37.5] (9.39%)	[27.4, 38.0] (8.27%)	[29.6, 37.1] (7.54%)	[31.7, 39.2] (7.46%)	0.08%
GPT-5 (VLM)	[23.6, 35.8] (12.70%)	[26.7, 32.1] (10.23%)	[23.5, 36.4] (9.08%)	[24.3, 33.8] (7.89%)	[25.2, 32.5] (7.25%)	[25.7, 32.8] (7.09%)	0.16%
GPT-5 (LLM)	[24.1, 36.1] (12.58%)	[22.7, 31.9] (9.10%)	[22.9, 33.1] (10.16%)	[26.9, 35.3] (7.87%)	[25.0, 32.3] (7.24%)	[25.4, 32.5] (7.07%)	0.17%
Claude Sonnet 4 (VLM)	[7.5, 14.5] (7.08%)	[7.7, 16.5] (9.01%)	[8.6, 14.9] (6.31%)	[7.5, 13.1] (5.35%)	[8.5, 13.5] (5.01%)	[9.3, 14.4] (5.04%)	-0.03%

Table 9. Analysis of 95% Confidence Interval (CI) widths across varying sample sizes. The " $\Delta$ " column represents the reduction in CI width. The results show that the CI width stabilizes around 600 samples. Note that values in brackets represent the 95% confidence interval range. Values in parentheses represent the CI width.

with Gemini 2.5 Pro. The experimental results in Table 12 indicate that GPT-5 achieves consistently better performance than the advanced Gemini 2.5 pro model on these open-source agent systems.

#### E.4 Ablation Study on Webpage Visit Tool

Augmenting agents with full webpage navigation, for example, visiting real-world e-commerce pages with marketing noise, decrease performance (SmolAgent GPT-5 LLM: 42.72%  $\rightarrow$  42.09%). These webpage content often overwhelms the key information, misleading the agents. Therefore, we remove the webpage visit tool for all open-source agents as the default setup.

#### E.5 Vision Encoder Bottleneck of VLMs

In HS classification, professional-grade visual details such as material textures and assembly states are often decisive, yet it remains unclear whether current Vision Encoders can reliably perceive them. To explicitly quantify this potential visual bottleneck, we perform a consistency analysis between GPT-5 (Text-only) and GPT-5 (Multimodal) to explicitly quantify this bottleneck. While incorporating visual inputs helped correct 6.42% of samples that text-only failed, we also observe a significant "Interference Phenomenon": 6.17% of samples that were correctly classified by text-only became incorrect when images were added. These experimental results confirms that current Vision Encoders in foundation models often misidentify fine-grained details, introducing noise that overrides correct textual reasoning. This indicates that improving the resolution and domain-specificity of Vision En-

coders is an essential research direction.

## F Improvement Gain from Agents

This waterfall in Figure 8 chart reveals that the superior performance of SmolAgent (GPT-5) are primarily from reducing the outdated and hallucination failures, with 56 corrected samples in HSCOD-BENCH. Besides, as shown in Figure 9, it can be found that the rate of outdated and hallucination are significantly reduced in SmolAgent baseline.

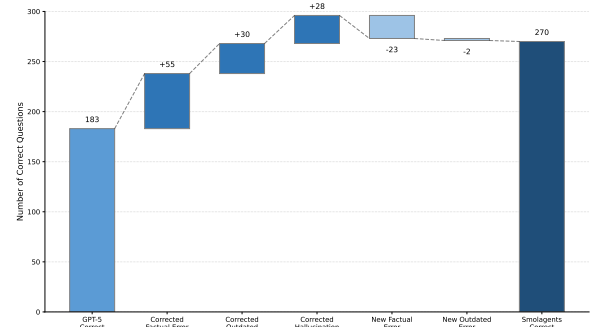


Figure 8. Details of performance gain and loss comparing GPT-5 and Smolagents.

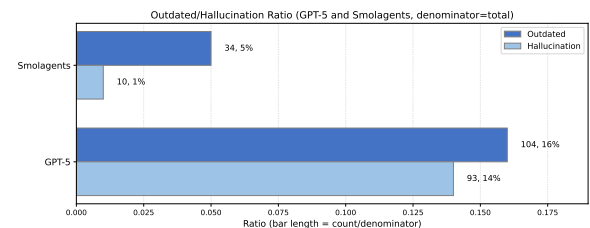


Figure 9. Outdated or hallucination ratio happened on GPT-5 and Smolagent.

Baselines	Model Type	HSCode Prediction Accuracy			
		2-digit	4-digit	6-digit	10-digit
<b>LLM/VLM-Only</b>					
GPT-5	VLM	82.12	86.32	84.60	48.81
Gemini-2.5-PRO	VLM	82.28	86.34	83.08	41.02
Claude Sonnet 4	VLM	78.80	81.32	70.61	24.82
Kimi-K2	LLM	78.01	79.52	71.18	27.59
DeepSeek-V3	LLM	77.06	70.63	59.31	20.11
O3-mini	LLM	77.22	72.74	43.67	5.18
<b>Open-source Agent System (GPT-5 Backbone)</b>					
SmolAgents	VLM	82.06	87.81	86.57	75.07
Aworld	LLM	82.28	85.57	84.05	69.79
Agentorchestra	LLM	82.12	86.13	85.45	68.33
OWL	LLM	72.63	85.19	83.37	72.39
WebSailor	LLM	81.64	86.43	81.16	61.88
Cognitive Kernel	LLM	80.06	86.37	78.94	48.40
Average	-	73.80	75.16	64.78	35.32

Table 10. The conditional performance on Chapter, Heading, Sub-heading and Country-specific codes classification in our proposed HSCODECOMP.

Configuration	10-digit Acc.
<b>Full Framework</b> (w/ <i>Tariff, Rules, CROSS</i> )	43.83%
w/o Decision Rules (DR)	<b>46.83%</b>
w/o Tariff Rules	36.73%
w/o CROSS Knowledge	46.53%

Table 11. **Ablation study on Knowledge Sources.** We evaluate the impact of removing distinct knowledge components from the SmolAgent (GPT-5) backbone.

## G Implementing Details

**Evaluation Details.** All baseline methods are equipped with search tools to access the CROSS database, hierarchical tariff rules, and other related resources, including human-written knowledge bases and hierarchical decision rules. The temperature and context window size of LLMs and agents are set to their default configurations. Moreover, as described in Section 6, the hierarchical decision rules, and webpage visit tool are not used during evaluation, since they do not improve the performance of open-source and closed-source agents. The multi-modal product images are used for open-source agents for SmolAgents. The hierarchical decision rules used in our prompts are illustrated in Figure 12. The hierarchical tariff

rules in the eWTP is shown in Figure 10. It can be found that the red boxes highlight logical dependencies in the tariff rules. These cases demonstrate that rule boundaries are ambiguous, posing significant challenges for accurate rule application by the agent. Moreover, the U.S. Customs Rulings Online Search System (CROSS) interface is shown in Figure 11. As illustrated, the CROSS website contains not only correct precedent results for product HS Codes but also numerous revoked precedents, requiring the agent to carefully evaluate information reliability. Additionally, since the precedent information is presented as plain text emails, the agent must effectively utilize contextual information to perform accurate reasoning.

**Knowledge Alignment.** All these knowledge resources in HSCODECOMP: hierarchical rules come directly from official US HTS legal notes included in eWTP, and the CROSS database contains authentic US Customs rulings. Human-written decision rules are also developed by human experts with US compliance experience.

**Timestamp-based Evaluation.** To ensure the reproducibility of our proposed HSCODECOMP dataset, all our dataset, knowledge and experimental setup are tied to the specific timestamp—June 30, 2024. This is a wide-used community standards

Backbone LLM	HSCode Prediction Accuracy				
	2-digit	4-digit	6-digit	8-digit	10-digit
<i>SmolAgent Framework</i>					
GPT-5	<b>82.28</b>	<b>70.89</b>	<b>59.81</b>	<b>49.05</b>	<b>42.72</b>
Gemini-2.5-Pro	82.19	69.48	57.87	44.04	34.49
Claude 4 Sonnet	80.69	67.09	54.11	42.25	33.70
Qwen-Max	77.34	63.23	42.47	26.62	17.43
<i>Aworld Framework</i>					
GPT-5	<b>82.28</b>	<b>70.41</b>	<b>59.18</b>	<b>48.58</b>	<b>41.30</b>
Gemini 2.5 Pro	79.55	66.97	54.70	38.79	29.24
<i>WebSailor Framework</i>					
GPT-5	<b>81.64</b>	<b>70.56</b>	<b>57.27</b>	<b>43.98</b>	<b>35.44</b>
Gemini 2.5 Pro	78.79	67.58	56.21	42.27	31.21
<i>AgentOrchestra Framework</i>					
GPT-5	82.12	<b>70.73</b>	<b>60.44</b>	<b>47.78</b>	<b>41.30</b>
Gemini 2.5 Pro	<b>82.27</b>	69.39	56.97	41.36	30.61

Table 12. Ablation study of backbone models across different open-source agent systems. Results indicate that GPT-5 consistently outperforms other models across all frameworks.

of established benchmarks like BrowseComp (Wei et al., 2025) and WideSearch (Wong et al., 2025).

**Infeasibility of Flat HSCode Prediction.** A potential alternative to our hierarchical agentic framework is a "flat" classification approach, where a model directly predicts the 10-digit HSCode from the product description. To evaluate the feasibility of this approach, we conducted a preliminary study using a Qwen3-8B model fine-tuned on a dataset of 175,000 historical customs rulings. Despite the substantial training data, the flat prediction model achieved a 10-digit accuracy of only **0.16%**. This catastrophic failure is attributed to the extreme long-tail distribution of the 10-digit classification space (which contains over 19,000 valid codes in the Harmonized Tariff Schedule of the United States (HTSUS)). The model struggled to memorize the specific mappings for rare codes and lacked the reasoning capability to navigate the legal boundaries between similar codes. This result confirms that a flat classification baseline is practically infeasible for this task and underscores the necessity of the hierarchical agentic reasoning adopted in this paper.

## H Case study of Failure Modes

We identify six critical failure modes of open-source and closed-source agent systems in HSCODECOMP: (1) **Premature Decisions:** Agents commit to incorrect classification paths

without collecting sufficient evidence (Figure 14 and Table 14-Grok DeepSearch); (2) **Information Misprocessing:** Agents overlook or misinterpret key product details, indicating challenges with long-context processing (Table 13 and Figure 15); (3) **Unnecessary Self-Correction:** Agents sometimes predict correct HSCodes initially but revise them incorrectly through excessive critique (Table 13, Gemini Deep Research); (4) **Reasoning Hallucination:** Agents generate plausible but factually incorrect reasoning steps (Table 13, Grok DeepSearch); (5) **Wrong Rule Application:** Models frequently miss or misuse relevant tariff rules due to their ambiguous descriptions that confuse the reasoning process, resulting in incorrect classification decisions (Figure 17); and (6) **Lack of Domain Knowledge:** Models exhibit errors due to insufficient domain-specific knowledge, such as misidentify silicone products as rubber instead of plastic (Figure 16). These limitations highlight that HSCODECOMP remains challenging for advanced closed-source and open-source systems.

Moreover, we also provide one case in Figure 13, demonstrating cutting-edge agent’s failure on using human-written hierarchical decision rules. This example demonstrates how explicit rules can encourage mechanical pattern-matching over substantive functional analysis—a critical failure mode in hierarchical rule application. Specifically, the rules-based approach mechanically applied Rule

- Section XVI Machinery and mechanical appliances; electrical equipment; parts thereof; sound recorders and reproducers, television image and sound recorders and reproducers, and parts and accessories of such articles [Notes](#)
- 84 Nuclear reactors, boilers, machinery and mechanical appliances; parts thereof [Notes](#)
  - 85 Electrical machinery and equipment and parts thereof; sound recorders and reproducers, television image and sound recorders and reproducers, and parts and accessories of such articles [Notes](#)
    - 8501 Electric motors and generators (excluding generating sets).
    - 8502 Electric generating sets and rotary converters.
    - 8503 Parts suitable for use solely or principally with the machines of heading 85.01 or 85.02.
    - 8504 Electrical transformers, static converters (for example, rectifiers) and inductors.
    - 8505 Electromagnets; permanent magnets and articles intended to become permanent magnets after magnetisation; electromagnetic or permanent magnet chucks, clamps and similar holding devices; electromagnetic couplings, clutches and brakes; electromagnetic lifting heads.
    - 8506 Primary cells and primary batteries.
    - 8507 Electric accumulators, including separators therefor, whether or not rectangular (including square).
    - 8508 Vacuum cleaners.
    - 8509 Electromechanical domestic appliances, with selfcontained electric motor other than vacuum cleaners of heading 85.08.
    - 8510 Shavers, hair clippers and hairremoving appliances, with selfcontained electric motor.
    - 8511 Electrical ignition or starting equipment of a kind used for sparkignition or compressionignition internal combustion engines (for example, ignition magnetos, magnetodynamos, ignition coils, sparking plugs and glow plugs, starter motors); generators (for example, dynamos, alternators) and cutouts of a kind used in conjunction with such engines.
    - 8512 Electrical lighting or signalling equipment (excluding articles of heading 85.39), windscreen wipers, defrosters and demisters, of a kind used for cycles or motor vehicles.
    - 8513 Portable electric lamps designed to function by their own source of energy (for example, dry batteries, accumulators, magnetos), other than lighting equipment of heading 85.12.
    - 8514 Industrial or laboratory electric furnaces and ovens (including those functioning by induction or dielectric loss); other industrial or laboratory equipment for the heat treatment of materials by induction or dielectric loss.
    - 8515 Electric (including electrically heated gas), laser or other light or photon beam, ultrasonic, electron beam, magnetic pulse or plasma arc soldering, brazing or welding machines and apparatus, whether or not capable of cutting; electric machines and apparatus for hot spraying of metals or cermets.
    - 8516 Electric instantaneous or storage water heaters and immersion heaters; electric space heating apparatus and soil heating apparatus; electrothermic hairdressing apparatus (for example, hair dryers, hair curlers, curling tong heaters) and hand dryers; electric smoothing irons; other electrothermic appliances of a kind used for domestic purposes; electric heating resistors, other than those of heading 85.45.
    - 8517 Telephone sets, including smartphones and other telephones for cellular networks or for other wireless networks; other apparatus for the transmission or reception of voice, images or other data, including apparatus for communication in a wired or wireless network (such as a local or wide area network), other than transmission or reception apparatus of heading 84.43, 85.25, 85.27 or 85.28.
    - 8518 Microphones and stands therefor; loudspeakers, whether or not mounted in their enclosures; headphones and earphones, whether or not combined with a microphone, and sets consisting of a microphone and one or more loudspeakers; audiofrequency electric amplifiers; electric sound amplifier sets.
    - 8519 Sound recording or reproducing apparatus.
    - 8520 [deleted]
    - 8521 Video recording or reproducing apparatus, whether or not incorporating a video tuner.
    - 8522 Parts and accessories suitable for use solely or principally with the apparatus of heading 85.19 or 85.21.
    - 8523 Discs, tapes, solidstate nonvolatile storage devices, «smart cards» and other media for the recording of sound or of other phenomena, whether or not recorded, including matrices and masters for the production of discs but excluding products of Chapter 37.

Figure 10. The case of the US HSCode hierarchical tariff rules (Section XVI). The red boxes highlight implicit logical relationships, such as exclusion clauses and cross-references. The blue boxes indicate vague decision boundaries characterized by open-ended illustrative examples, requiring agents to perform semantic interpretation beyond literal matching.

DATE	RULING CATEGORY TARIFF NO	RULING REFERENCE	RELATED
8/22/1995	<input type="checkbox"/> <a href="#">NY 813517</a> Classification <a href="#">8461.40.5070</a>	The tariff classification of a spur gear honing machine from Germany	
4/25/1995	<input type="checkbox"/> <a href="#">NY 809261</a> Classification <a href="#">8461.40.1020</a>	The tariff classification of gear hobbing machine from Japan or Belgium	
5/29/2003	<input type="checkbox"/> <a href="#">NY J84453</a> Classification <a href="#">8483.40.5010</a>	The tariff classification of gear boxes for turbines from ItalyDear Mr. Schwartz:	
8/22/1995	<input type="checkbox"/> <a href="#">NY 813513</a> Classification <a href="#">8461.40.5070</a>	The tariff classification of gear finishing machines from Germany.	
4/30/2012	<input type="checkbox"/> <a href="#">NY N212762</a> Classification <a href="#">8708.40.7580</a>	The tariff classification of transmission parts from an unspecified country.	
3/2/1993	<input type="checkbox"/> <a href="#">HQ 953365</a> Classification <a href="#">8501.10.40</a>	Protest No. 2704-93-100147; Gear Motor, Explanatory Note 85.01(1)(A); HQ 952500; 8432.80.00	References: <a href="#">952500</a>
7/3/2007	<input type="checkbox"/> <a href="#">NY N012496</a> Classification <a href="#">9817.00.6000</a> <a href="#">8483.40.5010</a>	The tariff classification of gear boxes from Germany and Italy.	
4/25/1995	<input type="checkbox"/> <a href="#">NY 809259</a> Classification <a href="#">8461.40.1020</a>	The tariff classification of bevel gear cutting machine from Japan or Belgium	
4/25/1995	<input type="checkbox"/> <a href="#">NY 809260</a> Classification <a href="#">8461.40.5070</a>	The tariff classification of gear grinder from Japan or Belgium	
8/27/1992	<input type="checkbox"/>	Protests 5301-9-000340, 000359 - 000366; Gear Boxes; Parts of Gear Boxes; Agricultural Use	

**J84453: The tariff classification of gear boxes for turbines from ItalyDear Mr. Schwartz:**

Ruling Date: May 29, 2003

NY J84453

May 29, 2003

CLA-2-84:RR:NC:1:102 J84453

CATEGORY: Classification

TARIFF NO.: [8483.40.5010](#)

Mr. David M. Schwartz  
Thompson Coburn LLP  
1909 K Street, N.W. (Suite 600)  
Washington, D.C. 20006-1167

RE: The tariff classification of gear boxes for turbines from Italy

Dear Mr. Schwartz:

In your letter dated April 30, 2003 you requested a tariff classification ruling on behalf of your client FiatAvio.

The articles in question are described as an inlet gear box, part number 9185M71G28, and a transfer/accessory gear box, part number L24644G06. You indicate that the gear boxes are imported exclusively for use in the manufacture and operation of an LM6000 aero-derivative turbine. You also offer that the gearboxes are properly classified as parts of such turbines in heading 8411, Harmonized Tariff Schedule of the United States (HTSUS). Although descriptive literature on the turbine was submitted, literature specific to the gear boxes was not made available. We assume from your explanation of what the gear boxes do within the turbine assembly that they are essentially fixed ratio speed changers.

We agree the gear boxes may be, prima facie, parts of the turbine provided for in HTSUS heading 8411. However, heading 8411 falls within chapter 84 of the HTSUS. Chapter 84 is in section XVI, HTSUS. Note 2(a) to section XVI provides in relevant part that parts which are goods included in any of the headings of chapter 84 are in all cases to be classified in their respective headings. If the gear boxes are articles of any other heading of chapter 84, then they cannot be classified as parts in heading 8411, HTSUS.

We note that gear boxes are specifically provided for in HTSUS heading 8483. The Explanatory Notes (EN) to the Harmonized Tariff System relevant to HTSUS heading 8483 explain that the heading includes all types of gears and assemblies of such gears used to change speed and/or the direction in which power is transmitted. Accordingly, we find that the gear boxes in question are articles of HTSUS heading 8483. Accordingly, as required by note 2(a) to HTSUS section XVI, the gear boxes fall to be classified in heading 8483, rather than heading 8411, HTSUS.

The applicable subheading for the gear boxes will be [8483.40.5010](#), HTSUS, which provides for other fixed ratio speed changers. The rate of duty will be 2.5 percent ad valorem.

Figure 11. The case of CROSS knowledge database that contains the products rulings.

1 literal description, classifying the collar under heading 8526 (radio navigational equipment) because it contains a GPS receiver. It cited ruling NY N006896 for “GPS pet locator devices”, treating GPS presence as determinative of classification. However, this misapplies Rule 3 Essential Character principle. The product’s essential function is not navigation but electronic pet containment and training. GPS merely serves as the technical means to enforce boundaries. The non-rules version correctly recognized this functional distinction, classifying it under heading 8543 (electrical machines with individual functions).

### **I Case Study of Overthinking**

Figure 18 and Figure 19 exhibit the overthinking case study of WebSailor framework.

### **J Case Study of Multi-modal Information**

Table 15 presents experimental case studies comparing SmolAgent (GPT-5) performance with and without product image processing.

### **K Case Study of Step-by-Step Failure Trace**

Figure 20 visualizes the step-by-step failure trace of advanced GPT-5 model.

## Decision Rules

### **The six decision rules for hierarchical tariff rules application**

The following six rules must be applied progressively from Rule 1 to Rule 6, without skipping.

#### **Rule 1: Priority of Headings and Notes**

The classification of goods shall be determined primarily according to the terms of the headings (4/6-digit HS codes) and any related Notes. Subsequent rules shall only be applied if the terms of the headings and the Notes do not suffice for classification.

#### **Rule 2: Incomplete/Unfinished Articles and Extension to Materials/Substances**

Rule 2(a): An incomplete or unfinished article (e.g., a bicycle missing wheels), if it has the essential character of the complete article, is to be classified as the complete article.

Rule 2(b): An article consisting of a certain material or substance, which retains its original character after the addition of other materials/substances (e.g., a plastic cup with a metal base), is to be classified according to the original material.

Example: An unassembled computer motherboard (which already has the function of a motherboard) is classified under heading 8473 (parts of computers).

#### **Rule 3: Decision Logic for Goods Classifiable Under Multiple Headings**

When goods are classifiable under two or more headings, classification shall be effected as follows, in order of priority:

Specificity (The more specific description shall be preferred to a more general description);  
Essential Character (Determined by the main material, function, or use of the goods);  
Last in Numerical Order (If classification cannot be determined otherwise, classify under the heading which occurs last in numerical order).

Example: An electric toothbrush (which has the attributes of both a "household appliance" and an "oral hygiene tool"): Specificity: Classified as a "domestic electro-mechanical appliance" (heading 8509) rather than a "toothbrush" (heading 9603).

#### **Rule 4: Principle of Closest Analogy**

When goods cannot be classified by applying the preceding three rules, they shall be classified under the heading appropriate to the goods to which they are most similar.

Example: Imitation leather made from a new material (not listed in the HS) is classified as "artificial leather" (heading 3921).

#### **Rule 5: Packing Materials and Containers**

Rule 5(a): Packing materials/containers presented with the goods (e.g., a jewelry box), if normally sold with the goods, are classified with the goods; otherwise, they are classified separately.

Rule 5(b): Reusable packing containers (e.g., metal gas cylinders) are classified separately.

Example: A glass bottle presented with perfume is classified under the heading for perfume (3303); however, a glass bottle sold separately is classified under 7010.

#### **Rule 6: Hierarchical Classification at the Subheading Level**

The classification of goods in the subheadings (6-digit and subsequent HS codes) of a heading shall be determined level by level, first determining the 1-dash subheading (5-6 digits), and then successively the lower-level subheadings. At each level, classification must take into account any Subheading Notes and the relationship between subheadings at the same level.

Example: After classifying goods under heading 6205 (men's shirts), the subheading is chosen based on material (cotton, man-made fibers, etc.): 620520 (of cotton) or 620530 (of man-made fibres).

Figure 12. The expert-written decision rules that used during expert annotation.

### Failure of Using Human-written Decision Rules (Task#31)

#### Task

**Product Title:** GPS-enabled dog fence collar that triggers alerts when pets cross virtual boundaries.

**Ground-truth HSCode:** 8526910040

**SmolAgent + GPT-5 (without rules):** Correctly predicted 8526910040 (Rule 1)

**SmolAgent + GPT-5 (with rules):** Incorrectly classified as 8543709860 (Rule 3)

**Used Rule 1 and 3 can be found in Figure G.1**

#### **Failure Analysis:**

The rules-based approach mechanically applied Rule 1's literal description, classifying the collar under heading 8526 (radio navigational equipment) because it contains a GPS receiver. It cited ruling NY N006896 for "GPS pet locator devices", treating GPS presence as determinative of classification. However, this misapplies Rule 3's Essential Character principle. The product's essential function is not navigation but electronic pet containment and training. GPS merely serves as the technical means to enforce boundaries. The non-rules version correctly recognized this functional distinction, classifying it under heading 8543 (electrical machines with individual functions).

Figure 13. The case of failure of agent on applying hierarchical human-written rules.

**Error: Premature Decisions**

**Task**

**Product Title:** PEN-F PU Leather Half Case for Olympus PEN-F Digital PENF Camera Brown/Black/Coffee

**Product Attributes:**

Origin: Mainland China	Measurement unit: 100000015
Use: Mirrorless System Camera	Package size - length (cm): 15
<b>Material: PU</b>	Model Number: PEN-F
Brand Name: NiYi	Type: Camera Bags, Hard Bag
each pack: 1	Package size - width (cm): 8
Package weight: 0.200	Style: handbags
Package: No	Package size - height (cm): 5

**Product Price:** 9.8 USD

**Category:** Consumer Electronics → Accessories & Parts → Camera Bags & Cases

**Search Query (Premature Focus)**

**Search query:** HTSUS camera case **4202.92** United States 10-digit code camera cases outer surface of plastic sheeting.

*Note: The agent immediately narrowed the search to Heading 4202 (Trunks, Suitcases, etc.) based on the keyword "Case", ignoring the specific material implications of "PU Leather" (Plastics).*

**HSCode Description**

**Wrong Direction (Heading 4202):**

42 Articles of leather; saddlery and harness...  
4202 Trunks, suit-cases, vanity-cases... camera cases...:  
**4202.92 ... (Incorrect for PU/Plastic composition in this context)**

**True Label (Chapter 39):**

39 Plastics and articles thereof  
3926 Other articles of plastics...:  
392690 Other:  
39269099 Other  
**3926909989 Other**

**Analysis**

The accurate code is **3926.90.9989**. This first search query leads to a wrong direction since it decides the item is under **4202** and ignored it is **PU leather** which is chemically **plastic**. The agent committed a "Premature Decision" error by locking onto the "Camera Case" function (Heading 4202) before verifying if the material composition (PU) excluded it from that chapter, and later turns did not realize this ignorance.

Figure 14. The case of early wrong search query leading to incorrect classification.

**Error: Real-world noise**

Task

**Product Title:** 10/30/60PCS INS Blue Color PVC Sticker **Aesthetic** Hand Accounting DIY Decoration **Scrapbooking** **Korean Stationery** Supplies

**Product Attributes:**

Origin : Mainland China	Shape : malformed
Size : M	Package size - length (cm) : 20
Material : <b>Plastic</b>	Model Number : sticky PVC stickers
LOT_NUM : 1	High-concerned Chemical : None
PRODUCT_UNIT : 100000015	Brand Name : NoEnName_Null
PACKAGE_TYPE : sell_by_piece	Package size - width (cm) : 15
Package weight : 0.030	Package size - height (cm) : 3
Age : >3 years old	

**Product Price:** 0.77 USD

**Category:** Office & School Supplies → Stationery Sticker → Stationery Sticker

Search query and Prediction

**Search query:** site:rulings.cbp.gov PVC stickers classification 4911.91

**Prediction:**

4911.91.4040

**Classification path:** Chapter 49 → 4911 → 4911.91 → 4911.91.4040

**Decision rationale (Excerpt):** The product consists of decorative, self-adhesive PVC stickers... CBP consistently classifies decorative stickers ... under 4911.91 when the printing is the primary purpose...

HSCode Description

**Predicted (Wrong):**

49 Printed books...  
4911 Other printed matter...  
491191 Pictures, designs...  
**4911914040 Other**

**True Label:**

39 Plastics and articles thereof  
3919 Selfadhesive plates...  
391910 In rolls...  
**3919102055 Other**

Analysis

The accurate code is 3919.10.2055. The material is PVC plastic... The misclassification error stems from **real-world noise** in the product title. This noise takes the form of subjective marketing terms (e.g., **Aesthetic**, **INS**, **Korean Stationery**) and application-focused words... This noise creates a misleading narrative that the printed design is the product's essential character.

Figure 15. The case of real-world noise leading to misclassification.

**Error: Lack domain knowledge**

Task

**Product Title:** 20/50/100PCS Pink Nitrile Gloves Disposable Latex Powder Free Gloves for Household Cleaning Beauty Salon Working Kitchen Gloves

**Product Attributes:**

Origin : Mainland China	Brand Name : INMAICH
Type : Disposable	each pack : 1
Measurement unit : 100000015	Use : Cleaning
Classification : Oversleeve	Style : Brief
<b>Material : PVC</b>	Model Number : Nitrile Disposable Gloves
Sold in : sellbypiece	

**Product Price:** 11.56 USD

**Category:** Home & Garden → Household Merchandises → Disposable Household Cleaning Products

Search query and Prediction

**Search query:** HTSUS nitrile disposable gloves HS code 10-digit 4015.19 CBP ruling

**Prediction:**

4015.19.1010

**Decision rationale (Excerpt):** Material and nature: The product is described and marketed as “nitrile gloves.”... This places the goods within Chapter 40 (rubber), not plastics...

**Notes:** If the goods were actually made of PVC/vinyl (plastic) rather than nitrile rubber, they would instead fall in Chapter 39... However, based on the stated product title... the correct US HTSUS is 4015.19.1010.

HSCode Description

**Predicted (Wrong & Obsolete):**

40 Rubber and articles thereof  
4015 Articles of apparel... of vulcanized rubber:  
- Gloves, mittens and mitts:  
401519 Other:  
40151911 Seamless  
**4015191110 Disposable (Obsolete)**

**True Label:**

39 Plastics and articles thereof  
3926 Other articles of plastics...:  
392620 Articles of apparel...:  
39262010 Seamless  
**3926201020 Disposable**

Analysis

This failure is a classic example of **lacking critical domain knowledge**, demonstrated by an inability to resolve **contradictory product data** and the use of an **outdated HTS code**. The analysis was fundamentally flawed because it prioritized the marketing-oriented *Product Title (Nitrile Gloves)* over the explicit technical attribute, *Material: PVC*. This **critical error** led the analyst to misclassify the item as a **rubber** product under **Chapter 40**. A domain expert would recognize that a **specific material attribute is a more reliable basis** for classification... Compounding this mistake, the predicted code, 4015.19.1010, is an **obsolete tariff number**...

Figure 16. The case of domain knowledge deficiency.

**Error: Wrong rule application**

Task

**Product Title:** Fashion TYPE 2 **DIABETIC** Bracelet Soft Silicone Bracelet With Type 2 Diabetes Medical Warning Bracelet Hot Sale

**Product Attributes:**

Origin : Mainland China	Shape pattern : Round
Measurement unit : 100000015	Package size - length (cm) : 15
Fine or Fashion : Fashion	<b>Material : Silicone</b>
Gender : Unisex	CN : Fujian
Item Type : BRACELETS	Model Number : SL-268
Metals Type : None	High-concerned Chemical : None
Brand Name : NPKDS	Function : None
each pack : 1	Package size - width (cm) : 15
Package weight : 0.013	Style : TRENDY
Setting Type : None	Clasp Type : None
Compatibility : All Compatible	Package size - height (cm) : 2
Sold in : sell_by_piece	Bracelets Type : Hologram Bracelets
Chain Type : Rope Chain	

**Product Price:** 1.71 USD

**Category:** Jewelry & Accessories → Fashion Jewelry → Bracelets

Search query and Prediction

**Search query:** US HTS 7117.90 silicone bracelet imitation jewelry HS code medical alert bracelet

**Prediction:**

7117.90.7500

**Decision rationale (Excerpt):** Product is a bracelet made of silicone (a plastic), marketed as a fashion/medical warning bracelet... It is not of precious metal... functions as personal adornment... fitting the terms of "imitation jewelry" of heading 7117...

HSCode Description

**Predicted (Wrong):**

71 Natural or cultured pearls... imitation jewellery...

7117 Imitation jewelry:

711790 Other:

**7117907500 Of plastics**

**True Label:**

39 Plastics and articles thereof

3926 Other articles of plastics...:

392690 Other:

39269099 Other

**3926909989 Other**

Analysis

This failure is a classic case of **Wrong Rule Application**, where a basic rule was incorrectly prioritized over a decisive, higher-level one. The agent mistakenly applied the simple rule of **classification by apparent form**, seeing a *bracelet* and defaulting to the seemingly specific **Chapter 71 (Imitation Jewelry)**. However, the correct, overriding rule that should have been applied is the Rule of **Essential Character**, which embodies the important principle of **Function Over Form** in customs classification. In this case, the text "**TYPE 2 DIABETIC**" was the key that defined its core **function**, fundamentally changing the product's essence from an ornament to a **utilitarian medical alert device**. An article with such an overriding utilitarian function is legally excluded from the scope of **Chapter 71**. Therefore, the correct rule application mandates that the item be classified by its **function** and material under **Chapter 39**, not by its misleading physical form.

Figure 17. The case of wrong rule application due to confusion between imitation jewelry and functional articles.

Analysis Dimension	SmolAgent	Gemini DeepResearch	Grok DeepSearch	Manus
<b>Final HTSUS Code</b>	8431.20.0000 (Correct)	7326.90.8688 (Incorrect)	8428.90.0290 (Incorrect)	8427.90.0020 (Incorrect)
<b>Core Logic Explained</b>	Based on the core principle of HTSUS Section XVI, Note 2, the cage is a <b>part</b> as it is 'solely or principally for use with' a forklift (heading 8427). Its design, function, and identity are entirely dependent on the forklift.	The argument is based on the <b>'part vs. accessory' distinction</b> . It posits the cage is not an <b>'indispensable' part</b> , but an optional 'accessory'. Since accessories are precluded from 8431, classification defaults to its constituent material (steel).	Characterizes the cage as a <b>functional piece of machinery</b> . The rationale is that it enables a new function and <b>incorrectly compares it to complex attachments</b> with their own mechanics (e.g., rotators, clamps).	Characterizes the cage as a <b>complete 'aerial work platform'</b> . The core argument is that its <b>'4 universal wheels' constitute a 'mobile base'</b> per the legal notes, thus assembling it into a complete vehicle.
<b>Key Flaw Analysis</b>	This approach correctly identifies the product's primary use and applies the controlling legal note directly, which is the standard and most reliable method for classification.	This is overthinking because the model became <b>fixated on a complex, secondary legal nuance</b> (part vs. accessory) while <b>ignoring the more direct, primary rule</b> ('solely or principally for use with'), leading to an unnecessarily complicated and incorrect conclusion.	This is an analysis hallucination because the model <b>invents characteristics the product lacks</b> , effectively <b>treating a passive structure as an active machine</b> . The entire analysis is built on this fabricated, non-existent product feature.	The model misses key product information by <b>misunderstanding the function of a key feature</b> (the wheels). It correctly identifies the wheels but misses their trivial context (ground convenience), instead <b>mistaking them for a vehicle's chassis</b> , which invalidates the entire classification.
<b>Failure Modes</b>	<b>Correct</b>	<b>Unnecessary Self-Correction</b>	<b>Reasoning Hallucination</b>	<b>Information Misprocessing</b>

Table 13. Comparative Analysis of HSCode Prediction for a Forklift Safety Cage mentioned above.

Analysis Dimension	SmolAgent	Gemini DeepResearch	Grok DeepSearch	Manus
<b>Final HTSUS Code</b>	8487.90.0080 (Correct)	8412.21.0075 (Incorrect)	8302.49.6085 (Incorrect)	8412.31.0080 (Incorrect)
<b>Core Logic Explained</b>	Based on the hierarchical structure of HTSUS Chapter 84, the shock absorber is a <b>generic machinery part</b> . After systematically eliminating more specific headings, it correctly classifies the item in the <b>residual heading 84.87</b> for parts "not elsewhere specified."	Characterizes the product as an <b>active hydraulic motor</b> . The logic is that because it is a linear-acting hydraulic device, it must be a "motor" under heading 84.12, which is an apparatus that <b>generates force or motion</b> .	It correctly identifies the passive function but then classifies it as a <b>simple base metal fitting</b> . The argument is that since it's not a motor, its classification defaults to a general heading for <b>common hardware and accessories</b> .	Characterizes the product as an <b>active pneumatic motor</b> . The rationale is based on the keyword "Pneumatic Cylinder" in the title, concluding it must be an <b>actuator that performs work</b> under heading 84.12.
<b>Key Flaw Analysis</b>	This approach correctly identifies the product's non-specific nature and applies the HTSUS's hierarchical structure and residual headings, which is the standard and most reliable method for such goods.	This is a fundamental misunderstanding of the product's function, as the model <b>mistakes a passive energy-dissipating device (a damper) for an active power-generating device (a motor)</b> . It confuses braking with accelerating.	The model makes a decision with insufficient information about HTSUS structure because it <b>fails to consider the critical distinction between Chapter 83 (simple fittings) and Chapter 84 (machinery)</b> and thus <b>underestimates the product's nature</b> as a piece of machinery.	This is a fundamental misunderstanding of the product's function, as the model is <b>misled by an inaccurate keyword</b> and <b>mistakes a passive damper for an active motor</b> , ignoring contradictory product attributes.
<b>Failure Modes</b>	<b>Correct</b>	<b>Information Misprocessing</b>	<b>Premature Decisions</b>	<b>Information Misprocessing</b>

Table 14. Comparative Analysis of AI Model Classifications for a Cylinder Shock Absorber mentioned above.

**Case 1: WebSailor Overthink.**

**Task**

**Product Title:** Men's and Women's Fashion Printed T-shirts, Men's and Women's Fashion Shirts, WWE John Scene, Winning This Day, Summer, New

**Product Attributes:**

Origin : Mainland China	Package size - length (cm) : 10
Gender : Men	Fabric Type : Broadcloth
High-concerned Chemical : None	Special Features : Breathable,Lightweight
Applicable Scene : Daily	Place Of Origin : China (mainland)
each pack : 1	Applicable Season : Spring And Summer
Style : Casual	Sleeve Length(cm) : short (4-16inch)
Package size - height (cm) : 10	Collar : O-Neck
Measurement unit : 100000015	Tops Type : TEES
Material : POLYESTER	Item Type : tops
Craft of Weaving : Knit	CN : Jiangxi
Brand Name : NoEnNameNull	Package size - width (cm) : 10
Product Care Instructions : Machine Wash	Hooded : No
Pattern Type : 3D	Package weight : 0.200
Closure Type : Pull On	Sold in : sellbypiece

**Product Price:** 74.96 CNY

**Category:** Men's Clothing → Tops & Tees → T-Shirts

**First step**

**Websailor:**

*We need the precise 10-digit HTSUS classification for a men's knit T-shirt of man-made fiber (polyester)... The common 10-digit code used by importers ... is **6109.90.1050 (men's or boys')**. I will confirm using CBP rulings...*

name : search , arguments : { query : [ HTSUS **6109.90.1050 men t-shirt** man-made fibers ... ]}

**Smolagents:** search query: 'HTSUS 6109.90 T-shirts of man-made fibers 10-digit statistical suffix men

**HSCode Description**

**Websailor wrong prediction:**

61 Articles of apparel...  
 6109 Tshirts, singlets...  
 610990 Of other textile materials:  
 61099010 Of manmade fibers  
**6109901050 Women's (639)**

**Smolagents correct prediction:**

Section XI Textiles...  
 61 Articles of apparel...  
 6109 Tshirts...  
 610990 Of other textile materials:  
 61099010 Of manmade fibers  
**6109901007 Men's (638)**

Figure 18. The first case of HS Code classification comparing Websailor and Smolagents reasoning.

**Case 2: WebSailor Overthink and Directly Answer.**

**Task**

**Product Title:** SPRAY GUN with adapter 400cc Mix tank 1.3/1.8mm nozzle air spray gun with quick-connect wash-free tank paint Mixing Cup sprayer

**Product Attributes:**

Origin : Mainland China	Power Type : Pneumatic
Measurement unit : 100000015	Battery Included : No
Use : Commercial Manufacture	Cup Volume : 400cc
Package size - length (cm) : 30	Nozzle Diameter : 1.3mm
Model Number : GTI	Hign-concerned Chemical : None
Feed Type : Gravity	Certification : CE
Brand Name : NTOOLS	Type : HVLP
each pack : 1	Package size - width (cm) : 22
Package weight : 1.300	No-Load Speed : 1200rpm
Application : Paint Spray Gun	Package size - height (cm) : 20
Sold in : sellbypiece	

**Product Price:** 98.0 USD

**Category:** Tools → Power Tools → Spray Gun

**First step**

**Websailor:**

*We need the precise 10-digit US HS (HTSUS) code for a pneumatic paint spray gun... Specifically, subheading 8424.20 covers "Spray guns and similar appliances."... The 10-digit breakouts... fall in 8424.20.10 at the 8-digit level... Therefore, the full 10-digit code is **8424.20.1000** or **8424.20.1090**... Now present classification path...* 8424.20.1000

...

**SmolAgents:**

search query: HTS 842420 paint spray guns 10 digit statistical suffix

**HSCode Description**

**Websailor wrong prediction:**

84 Nuclear reactors...  
 8424 Mechanical appliances...  
 842420 Spray guns and similar appliances:  
**8424201000 Simple piston pump sprays and powder bellows**

**Smolagents correct prediction:**

84 Nuclear reactors...  
 8424 Mechanical appliances...  
 842420 Spray guns and similar appliances:  
**8424209000 Other**

Figure 19. The second case of HS Code classification comparing Websailor and Smolagents reasoning.





Image	Product Information	Non-vision HS (10-digit) + brief rationale	With-vision HS (10-digit) + brief rationale	Vision advantage (reason category)
	<p><b>Title:</b> Royal Wedding Bouquet Rhinestone Bride and Bridesmaid Hand Flowers Handmade Bridal Bouquet.  <b>Category:</b> Artificial Decorations.  <b>Attributes:</b> Material: <b>Silk (natural textile)</b>, not plastics; bouquet of artificial flowers with rhinestones.</p>	<p><b>6702.90.3500</b> (Artificial flowers of <i>man-made fibers</i>). Assumes “silk flowers” = polyester; treats bouquet as MMF artificial flowers.</p>	<p><b>6702.90.6500</b> (Artificial flowers, <i>other than man-made fibers</i>). Images/text indicate <b>silk fabric (natural fiber)</b>; rhinestones do not change essential character (GRI 3(b)).</p>	<p><b>Material/fiber identification</b> (natural silk vs assumed polyester).</p>
	<p><b>Title:</b> OTF knife parts, aviation aluminum tactical handle kit.  <b>Category:</b> Hand Tools / Knife accessories.  <b>Attributes:</b> <b>Handle-parts only</b>, no blade present; aluminum body with clip, actuator, screws.</p>	<p><b>8211.93.0035</b> (Folding/pocket knives). Interprets listing as a <i>complete folding knife</i> rather than components.</p>	<p><b>8211.95.9000</b> (Handles of base metal, other). Visuals confirm <b>no blade</b>; essential character is the base-metal handle assembly (parts provision applies).</p>	<p><b>Completeness/parts identification</b> (parts-only vs whole article).</p>
	<p><b>Title:</b> Acrylic buckle, beaded/openwork shoulder bag with inner pouch.  <b>Category:</b> Women’s Handbags.  <b>Attributes:</b> Outer surface is <b>beads/openwork lattice (ABS/acrylic)</b>, not plastic sheeting; linen pouch is interior.</p>	<p><b>4202.22.1500</b> (Handbags with outer surface of <i>sheeting of plastics</i>). Assumes exterior is plastic sheeting based on “ABS.”</p>	<p><b>4202.29.9000</b> (Other). Images show <b>beads/rings openwork</b>, not “sheeting of plastics”; classification by outer surface (Additional U.S. Note 2 to Ch. 42).</p>	<p><b>Structural outer-surface feature</b> (beads/openwork vs plastic sheeting).</p>
	<p><b>Title:</b> The Simpsons character collectible paper cards.  <b>Category:</b> Printed matter / collectibles.  <b>Attributes:</b> Cards bear <b>pictures only</b>; no rules-based deck specified; <b>no lithographic process evidence</b>.</p>	<p><b>4911.99.6000</b> (Other printed matter, often linked to lithographic printing). Assumes lithographic process without explicit evidence.</p>	<p><b>4911.91.4040</b> (Pictures, designs and photographs; other). Visuals confirm <b>picture-only cards</b> and lack of lithographic evidence/criteria; not a playing-card game.</p>	<p><b>Process/evidence-based exclusion</b> (no lithography evidence; pictures-only).</p>

Table 15. Vision vs Non-vision: Four Representative Cases with Key Evidence and Reasons



Figure 20. The case of GPT-5’s step-by-step reasoning on HSCODECOMP.