

CharacterHub: Open-Domain Character Profiling for LLM Role-play via Deep Search

Anonymous ACL submission

Abstract

Building high-quality character profiles is a foundational prerequisite for developing immersive Role-Playing Language Agents (RPLAs). However, existing profiling methods primarily rely on literature-based extraction or LLM-based generation, which suffer from limited media coverage, high manual costs, and a propensity for factual hallucinations. To address these bottlenecks, we propose CHARACTERHUB, an automated character profiling framework powered by deep search agents. Unlike traditional extractive pipelines, our framework autonomously navigates open web sources to retrieve and aggregate heterogeneous information across multiple dimensions. This agentic approach offers unparalleled scalability, extending high-fidelity profiling beyond literary figures to anime, games, and user-generated characters, without human intervention. To rigorously validate our method, we establish an automatic evaluation protocol using large-scale, human-curated data from Fandom¹ as gold reference. Experimental results demonstrate that our dataset achieves strong alignment with reference sources, notably reaching a 83.13% Support Score in the critical personality dimension, while attaining nearly twice the information density of Fandom references. We will publicly release the dataset and associated resources.

1 Introduction

Recent advances in Large Language Models (LLMs) (Yang et al., 2025; DeepSeek-AI et al., 2025) have facilitated widespread applications across various domains. Among them, Role-Playing Language Agents (RPLAs) (Chen et al., 2024) have gained increasing attention for their ability to construct character-centric interactive AI systems capable of consistently portraying diverse personas or characters. RPLAs enable a wide

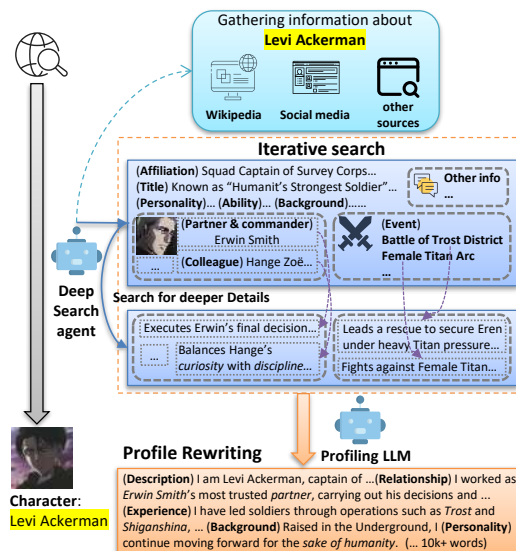


Figure 1: Character Profile Construction Pipeline of CHARACTERHUB. A deep search agent gathers multi-source information about a target character, which is then synthesized by a profiling LLM into a unified character profile.

range of applications, including virtual companions (Sheehy et al., 2024), game NPCs (Park et al., 2023), interactive storytelling (Sun et al., 2023), and educational simulators (Zhao et al., 2025). To build a comprehensive RPLAs system, we require rich types of structured or unstructured character data. Within this data hierarchy, the character profile constitutes the most foundational level, which naturally leads to the research problem of character profiling.

However, constructing high-quality character profiles remains a major bottleneck for RPLA systems. Existing efforts in character profiling can be broadly categorized into two paradigms: **literature-based extraction** (Li et al., 2023; Yuan et al., 2024) and **LLM-based generation** (Zhou et al., 2024, 2025). The former 1) relies on complex pipelines to harvest data from long-form narratives but is often resource-intensive; 2) results in sparse profiles

¹<https://www.fandom.com/>

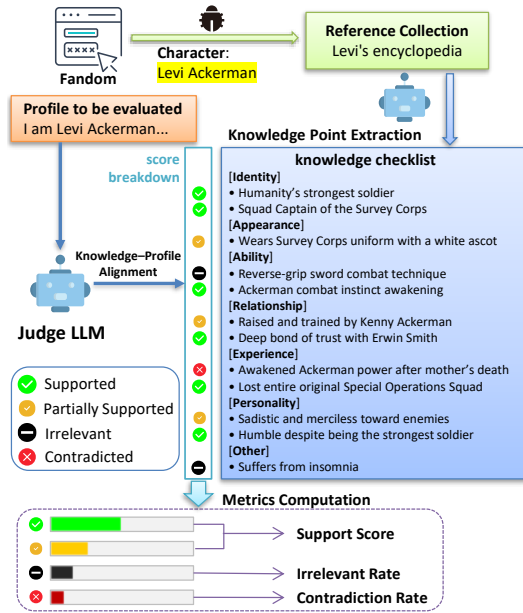


Figure 2: Character profile evaluation pipeline of CHARACTERHUB. Reference character information is extracted from Fandom sources into structured knowledge points, which are then compared against a generated profile by a judge LLM to assign consistency labels and compute evaluation metrics.

due to implicit textual cues like motivations, and 3) fails to generalize beyond characters with extensive written records, like game, animation and virtual characters. The latter offers greater flexibility but frequently suffers from factual hallucinations and a lack of grounding in established lore. Ultimately, these extractive and static approaches struggle to adapt to the evolving universe of characters across diverse media.

To address this challenge, we propose CHARACTERHUB, an automated character profiling framework based on deep search agents that autonomously acquire character information from open web sources. Instead of relying solely on raw narrative texts, our framework actively retrieves structured and semi-structured descriptions from various sources, such as wikis, fan communities, and multimedia-related websites. Through explicit iterative search, the agent constructs comprehensive character profiles covering multiple dimensions such as personality, experiences, relationships, and abilities. Following this pipeline, we further curate a comprehensive character profiling dataset, also named CHARACTERHUB. Overall, our pipeline offers two major advantages: (1) **Broad character applicability**, as it can be applied to a wide spectrum of characters, ranging

from literary figures to game, virtual, and user-generated characters; (2) **Flexible and comprehensive profiling**, as it dynamically integrates heterogeneous sources to produce multi-dimensional and fine-grained character representations.

Furthermore, to assess the quality of our profile, we propose an automatic evaluation protocol to assess character profiles. To rigorously assess the quality of the acquired profiles, we design a set of fine-grained evaluation metrics that measure factual correctness, attribute coverage, and semantic consistency. To build this benchmark, we crawl large-scale character pages from Fandom and treat them as gold references for quantitative comparison. Experimental results demonstrate that our automatically acquired profiles achieve strong alignment with human-curated sources, particularly in the most critical Personality dimension where we achieve a Support Score of 83.13%, while exhibiting substantially higher information density—nearly twice that of Fandom references.

In summary, this work makes three key contributions:

1) We propose a scalable, search-agent-based framework for automatic character profile construction across heterogeneous media domains, along with an LLM-based evaluation protocol that derives knowledge-level metrics from Fandom-grounded references.

2) We release a high-quality and diverse character profile dataset CHARACTERHUB, tailored for character-centric RPLA research, covering a wide range of profile dimensions.

3) Extensive experiments, including dimension-wise analysis, ablation studies, and human evaluation, demonstrate the high quality of the proposed dataset and its clear advantages over directly using raw Fandom content.

2 Related Work

Character Profiling Character profiling is the task of extracting attributes of given character to represent a persona. Current methodologies generally follow two paradigms: **literature-based extraction** and **LLM-based generation**.

The first paradigm focuses on harvesting character knowledge from authoritative sources, including 1) film or animation scripts, such as *ChatHaruhi* (Li et al., 2023) and *RoleLLM* (Wang et al., 2024), 2) literature and fictional works, such as *CoSER* (Wang et al., 2025), *HPD* (Chen et al.,

2023) and *CroSS* (Yuan et al., 2024), 3) Wiki or Baidu, such as *Character-LLM* (Shao et al., 2023) and *CharacterEval* (Tu et al., 2024). These extractive methods are often limited by high manual costs, data sparsity in implicit narratives, and a scope restricted to established literary works.

Conversely, the second paradigm leverages the generative capabilities of LLMs to synthesize profiles through 1) self-alignment, like *DITTO* (Lu et al., 2024), 2) instruction-guided expansion, like *CharacterGLM* (Zhou et al., 2024), *Character-Bench* (Zhou et al., 2025), and 3) multimodal synthesis, like *MMRole* (Dai et al., 2025). Although scalable, it frequently struggles with factual hallucinations and inconsistency.

Deep Research The paradigm of information retrieval has evolved from static **RAG** (Lewis et al., 2020) to dynamic **Agentic Deep Research Architectures**, which utilize **Tool Use** (Schick et al., 2023) and **Function Calling** (Li et al., 2024) to interact with external search APIs (e.g., Google Search or Fandom API), effectively bypassing the model’s knowledge cutoff. Central to this evolution is the **ReAct** paradigm (Yao et al., 2023), which interleaves reasoning traces with task-specific actions, enabling autonomous query planning and iterative information synthesis. This capability allows our search-agent framework to target specific character wikis across the open web.

3 Dataset Construction

We construct a large-scale, franchise-aware character corpus through (1) Character Collection with Franchise-Aware Merging, resulting in 1,109 high-precision personas; and (2) Profile Generation, where we employ our CHARACTERHUB pipeline to synthesize comprehensive first-person character profiles for RPLA conditioning.

3.1 Character Collection

Data Source and Coverage We assemble a large-scale, multi-domain character corpus by integrating three popularity rankings: AniList GraphQL API, MyAnimeList Jikan API, and a community-curated Character.AI leaderboard. The combined raw set contains 2,139 distinct fictional personas spanning anime, manga, video-games, movies, and western comics, providing both Japanese and English name variants.

Selection and Deduplication To ensure high precision, we implement a franchise-aware merging

algorithm. Characters are first grouped by their canonical work (e.g., *Jujutsu Kaisen*, *Genshin Impact*). Within each group, we compute name similarity using *SequenceMatcher*². Two entries are merged only if they (i) originate from different sources, (ii) share an exact franchise match, and (iii) exceed a similarity threshold of 0.85. This strategy successfully consolidates alternate spellings (e.g. “Satoru Gojo” vs. “Gojo Satoru”) while preserving cross-franchise homonyms like “Zoro” from *One Piece* and *Demon Slayer*. This process collapsed 32 duplicates, yielding 2,107 unique characters.

Quality Control We filter out characters whose franchise information is missing or cannot be matched to a hand-crafted alias table covering high-impact series, resulting in 1,109 characters that are successfully processed by the CHARACTERHUB pipeline to obtain complete character profiles. The resulting dataset contains character names, normalized franchises, category tags, and cross-platform popularity scores, and serves as the backbone for all character-driven experiments in the remainder of the paper.

3.2 CHARACTERHUB Pipeline

We propose a fully automated search-agent pipeline for open-ended and scalable profile construction for arbitrary characters. As illustrated in Figure 1, the pipeline consists of two stages: character-driven iterative search and profile rewriting.

Character-Driven Iterative Search Given a character name c_i , the system activates a search agent to retrieve character-related information from open web sources. The agent automatically formulates diversified search queries by combining the character name with profiling-oriented keywords (e.g., personality, background, story), and issues them to general search engines and wiki platforms. To improve coverage and reduce ambiguity, the pipeline adopts an iterative search strategy: retrieved information R_i^1 is summarized and fed back into subsequent search rounds for self-refinement, yielding refined information R_i^2 . This multi-round process progressively expands under-represented profile dimensions and enables more comprehensive and fine-grained information acquisition.

²<https://docs.python.org/3/library/difflib.html>

Output Profile Rewriting The refined deep-retrieved information $\{R_i^2\}$ is then provided to a profiling model to generate the final character profile P_i in a first-person narrative form. The profiling prompt instructs the LLM to integrate heterogeneous facts into a coherent persona, resolve potential cross-source conflicts, and rewrite the character description as a unified self-introduction. This design is motivated by RPLA requirements, where first-person self-consistent profiles better support immersive role-playing than third-person factual descriptions.

Finally, for each character c_i , the pipeline outputs $\{c_i, R_i^1, R_i^2, P_i\}$. The entire process is fully automatic, requires no manual annotation, generalizes to arbitrary characters without task-specific retraining, and scales naturally to web-scale character discovery, making it well suited for large-scale open-domain character profile construction.

4 Evaluating Character Profiles

To systematically assess the quality of constructed character profiles, we adopt an **LLM-as-a-Judge** evaluation framework grounded in **external fandom-based reference**. The key idea is to treat large-scale fandom wikis as rich yet noisy factual references, extract fine-grained knowledge points from them, and employ an LLM to judge their consistency with the generated character profiles.

As illustrated in Figure 2, The evaluation pipeline consists of four stages: (1) reference collection, (2) knowledge point extraction, (3) knowledge–profile alignment, and (4) metric computation.

4.1 Reference Collection

To establish a reliable evaluation benchmark, we construct an **external reference** for each character by combining the Fandom Wiki API with the Google Search API.

We first perform **webpage matching** to identify the most authoritative source. Given a character name, the Google Search API is used to retrieve candidate URLs, with priority given to pages within the Fandom Wiki domain. When multiple candidates are available, we select the page with the highest textual relevance and the most comprehensive content structure.

After identifying the target page, we conduct **web crawling and content extraction**. Specifically, we extract the **infobox** to obtain concise

factual attributes (e.g., gender, affiliation, occupation) and parse the **main content** into a structured dictionary that maps section headings (e.g., *Background, Abilities*) to their corresponding narrative paragraphs, preserving the original thematic organization.

Finally, the extracted data is passed through a **normalization** pipeline to remove citation markers, footnotes, navigation artifacts, corrupted Unicode characters, and formatting noise. The cleaned content is standardized into a JSON reference file, serving as a unified factual reference for subsequent knowledge extraction and quantitative evaluation.

4.2 Knowledge Point Extraction

Given the cleaned fandom reference for each character, we employ an LLM to automatically extract **fine-grained knowledge points**, which serve as atomic evaluation units.

Formally, for each character, the LLM outputs a knowledge checklist: $\mathcal{K} = [k_1, k_2, \dots, k_N]$. Each knowledge point k_i is represented as a structured triple: $k_i = \{s_i, e_i, type_i\}$, where s_i denotes the extracted knowledge statement, e_i denotes the evidence from the input text, and $type_i$ denotes the type of knowledge.

Inspired by previous studies (Yuan et al., 2024; Zhou et al., 2025), we facilitate a **multi-dimensional evaluation** by assigning each point a type $type_i \in \{Identity, Appearance, Ability, Relationship, Experience, Personality, Other\}$. Detailed definitions for each dimension are provided in Appendix C.

4.3 Knowledge–Profile Alignment

To quantitatively assess the factual consistency between synthesized profiles and canonical lore, we adopt an **LLM-as-a-Judge** framework to evaluate the alignment between the generated character profile and Fandom knowledge points.

For each character, the evaluator LLM is provided with two inputs: the generated profile P and a set of extracted knowledge checklist $\mathcal{K} = [k_1, k_2, \dots, k_N]$ derived from Fandom. The model evaluates each knowledge point $k_i \in \mathcal{K}$ independently and outputs a judgment pair (y_i, r_i) , where y_i denotes an alignment label and r_i provides a brief textual justification.

The alignment label y_i is selected from four categories reflecting different degrees of consistency. A knowledge point is labeled 1) **Supported(Sup)** if it is explicitly stated or clearly entailed by the profile;

2) **Partially Supported**(Par) if only some aspects are confirmed while others are missing or ambiguous; 3) **Irrelevant**(Irr) if the profile contains no information related to the knowledge point; 4) **Contradicted**(Con) if the profile includes statements that directly conflict with or negate it.

The resulting alignment is represented as a structured mapping:

$$\mathcal{A} = (k_i, y_i, r_i)_{i=1}^N, \quad (1)$$

which serves as the basis for computing quantitative evaluation metrics. This fine-grained formulation enables localized analysis of factual accuracy and coverage across different profile dimensions.

4.4 Evaluation Metrics

Based on the LLM-generated alignment labels, we define three quantitative metrics to reflect different aspects of profile quality.

Support Score. This metric measures overall factual consistency between the profile and knowledge:

$$SS = \frac{|\text{Sup}| + 0.5 \cdot |\text{Par}|}{|\mathcal{K}|}. \quad (2)$$

Irrelevant Rate. This metric reflects **coverage deficiency**, i.e., how much knowledge is missing from the profile:

$$IR = \frac{|\text{Irr}|}{|\mathcal{K}|}. \quad (3)$$

Contradiction Rate. This metric directly captures **factual inconsistency and hallucination risk**:

$$CR = \frac{|\text{Con}|}{|\mathcal{K}|}. \quad (4)$$

5 Experiment

5.1 Settings

Models for Profile Construction For the iterative search stage, we adopt Gemini-2.5-Pro (Comanici et al., 2025) as the core model in our search-agent-based pipeline, where iterative information gathering is explicitly orchestrated by the agent. In parallel, we employ Doubao-seed-1.6 (Seed et al., 2025) as a representative model of the deep research paradigm, which is prompted only at the initial search step and relies on its internal reasoning to implicitly perform iterative information

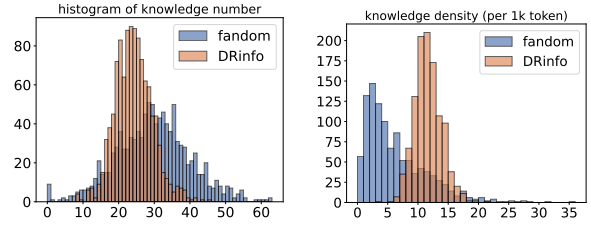


Figure 3: Comparison between Fandom and DRInfo (Gemini) in knowledge quantity and density per character. Left: histogram of the absolute number of extracted knowledge points. Right: histogram of knowledge density, measured as knowledge points number per 1k tokens. DRInfo yields more compact and information-dense character profiles than raw Fandom content.

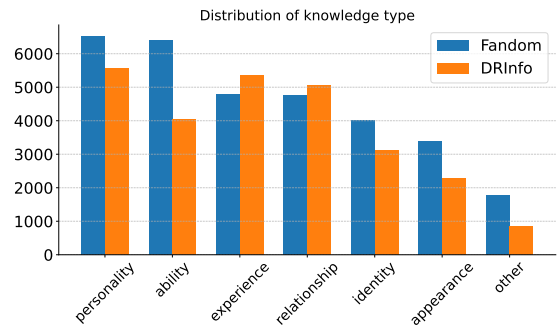


Figure 4: Comparison of knowledge type distributions between Fandom and DRInfo. Fandom exhibits a skewed distribution dominated by personality and ability-related knowledge, whereas DRInfo provides more balanced coverage across experience and relationship-centric dimensions.

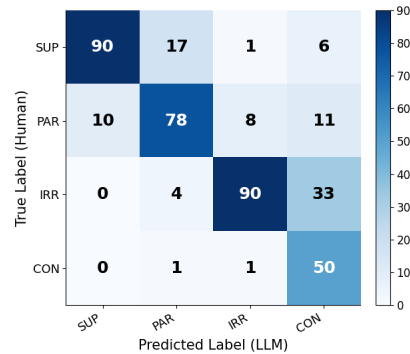


Figure 5: Confusion matrix comparing human annotations and LLM-based judgments for Knowledge-Profile Alignment.

seeking, serving as a comparative setting without explicit agent control. For profile rewriting stage, we use Qwen-plus to transform retrieved evidence into coherent first-person character profiles.

Judge Models For profile evaluation, we adopt a unified LLM-as-a-Judge setup. For knowledge

Type	Source	Gemini			Doubao		
		SS	IR	CR	SS	IR	CR
Personality	DRInfo	81.67	10.00	0.87	78.04	14.16	1.78
	Profile	83.13	9.95	1.58	80.67	10.55	3.51
Ability	DRInfo	61.37	28.69	1.47	53.93	33.52	4.33
	Profile	58.86	31.91	2.05	50.42	39.49	4.14
Experience	DRInfo	60.42	27.02	5.28	50.23	33.02	10.43
	Profile	60.63	27.00	5.90	50.44	32.36	12.18
Relationship	DRInfo	68.32	19.68	3.06	56.05	28.16	7.51
	Profile	68.94	19.71	3.25	56.58	29.00	7.84
Identity	DRInfo	74.66	11.80	2.77	66.21	13.38	8.16
	Profile	73.46	13.22	3.15	65.25	14.90	9.41
Appearance	DRInfo	57.62	23.25	3.06	42.50	26.60	12.76
	Profile	55.83	25.27	4.25	42.12	30.13	13.78
Other	DRInfo	39.72	50.31	2.38	42.36	46.36	3.81
	Profile	40.57	48.58	3.69	42.54	44.67	6.78
Overall	DRInfo	66.54	22.08	2.55	58.36	26.06	6.56
	Profile	66.16	22.65	3.18	58.14	27.04	7.62

Table 1: Evaluation results across seven dimensions and overall performance. Scores are presented as percentages with the % symbol omitted. DRInfo means the raw information retrieved by search agent, Profile means the rewritten profile.

Type	SS (%) \uparrow		IR (%) \downarrow		CR (%) \downarrow	
	Full	w.o. IS	Full	w.o. IS	Full	w.o. IS
Overall	66.54	64.15 (\downarrow 2.39)	22.08	24.22 (\uparrow 2.14)	2.55	2.23 (\downarrow 0.32)
Personality	81.67	80.78 (\downarrow 0.89)	11.78	12.22 (\uparrow 0.44)	0.87	0.83 (\downarrow 0.04)
Ability	61.37	57.48 (\downarrow 3.89)	28.69	32.42 (\uparrow 3.73)	1.47	1.78 (\uparrow 0.31)
Experience	60.42	57.80 (\downarrow 2.62)	27.02	30.23 (\uparrow 3.21)	5.28	4.08 (\downarrow 1.20)
Relationship	68.32	65.63 (\downarrow 2.69)	19.68	21.88 (\uparrow 2.20)	3.06	2.60 (\downarrow 0.46)
Identity	74.66	72.32 (\downarrow 2.34)	11.80	13.31 (\uparrow 1.51)	2.77	2.83 (\uparrow 0.06)
Appearance	57.62	55.94 (\downarrow 1.68)	23.25	24.19 (\uparrow 0.94)	3.06	3.08 (\uparrow 0.02)
Other	39.72	37.29 (\downarrow 2.43)	50.31	53.71 (\uparrow 3.40)	2.40	2.11 (\downarrow 0.29)

Table 2: Ablation study on iterative search. ‘‘Full’’ denotes the complete pipeline with iterative search, while ‘‘w.o. IS’’ removes the iterative search step. Relative changes compared to the full pipeline are shown in parentheses.

point extraction from the Fandom reference, we use Qwen3-235B-A22B-Instruct-2507 (Yang et al., 2025). The same model is further used to assess the consistency between each extracted knowledge point and (1) the generated character profile P_i , and (2) the deep-retrieved information R_i^K (DRInfo). All reported metrics are computed by averaging over all knowledge points across all characters.

5.2 Main Results

In the experiments, we wish to answer three research questions: RQ1) Can our method construct high-quality character profiles? RQ2) How does performance vary across different profile dimensions? RQ3) What are the advantages over directly crawling Fandom pages?

5.2.1 Can our method construct high-quality character profiles?

The proposed Search Agent-based method is highly effective in constructing high-quality character profiles with high accuracy and low hallucination rates. As shown in Table 1, the Gemini based pipeline achieves a total Support score of 66.54% for raw info and 66.16% for the final profile. These high scores, coupled with a remarkably low overall Contradict rate of 3.18%, demonstrate that the search-agent approach accurately retrieves knowledge consistent with the Fandom reference.

The Profile Rewriting process introduces only marginal information loss, demonstrating the stability of the persona conversion phase. For Gemini, the Support score shows only a marginal

409 decrease of 0.38% (66.54% to 66.16%) during the
410 rewriting phase, while the Contradict rate remains
411 stable under 4%. This indicates that the rewriting
412 LLM effectively preserves the factual integrity of
413 the retrieved knowledge. This minor trade-off be-
414 tween semantic density and narrative style is ac-
415 ceptable given the significant improvement in the
416 naturalness of the first-person perspective.

417 **Our method significantly outperforms the**
418 **Doubao-based Deep Research architecture.**
419 Gemini maintains an overall Support score that
420 is approximately 8% higher than Doubao, while
421 Doubao’s overall Contradict rate is more than dou-
422 ble that of Gemini. This performance gap can be
423 attributed to the search behavior: Gemini’s tool
424 usage tends to prioritize global information sources
425 like Fandom, whereas Doubao frequently relies
426 on Chinese-centric platforms such as Bilibili and
427 Baidu Baike, which may contain narrative discrep-
428 ancies or different levels of detail compared to the
429 Fandom reference.

430 5.2.2 How does performance vary across 431 different profile dimensions?

432 **Knowledge extraction performance is strongest**
433 **in the Personality dimension, which is consid-**
434 **ered the most critical component for RPLA.**
435 Specifically, the Gemini-based pipeline achieves its
436 highest Support score in the Personality dimension,
437 reaching 81.67% for raw info and 83.13% for the
438 final profile. This high performance ensures that
439 the core behavioral logic and linguistic style of the
440 characters are captured with high precision, lay-
441 ing a solid foundation for immersive role-playing
442 experiences.

443 **The Appearance and Experience dimensions**
444 **present a higher risk of hallucination, likely due**
445 **to the granular nature of the required informa-**
446 **tion.** Doubao’s Contradict rates for Experience
447 and Appearance reach 12.18% and 13.78%, respec-
448 tively, while Gemini also sees its highest contradic-
449 tion rates in these areas. These dimensions often
450 involve minute physical details (e.g., eye color, spe-
451 cific accessories) or complex chronological events,
452 which are more susceptible to model confusion or
453 information overlap.

454 **Stability during the rewriting phase varies by**
455 **dimension, with technical Ability descriptions**
456 **being more prone to conversion loss.** In the Abil-
457 ity dimension, Gemini’s Support score dropped
458 from 61.37% to 58.86% after rewriting. This in-
459 dicates that converting descriptive, often technical

460 ability mechanics into first-person dialogue is more
461 linguistically challenging for the LLM than con-
462 verting straightforward identity or personality facts.

463 5.2.3 What are the advantages over directly 464 crawling Fandom pages?

465 **Directly relying on Fandom introduces popular-**
466 **ity bias, whereas our method alleviates it.** As
467 illustrated in Figure 3, Fandom pages exhibit sub-
468 stantial imbalance in content completeness across
469 characters: while popular protagonists often have
470 rich and detailed entries, characters from niche
471 intellectual properties or peripheral roles are fre-
472 quently under-documented, resulting in sparse and
473 fragmented knowledge distributions. In contrast,
474 our approach provides more consistent coverage
475 across characters, which alleviates this dependency
476 on character popularity and yields more stable and
477 informative profiles even for under-documented
478 characters.

479 **Our dataset achieves substantially higher**
480 **knowledge density.** The right histogram of Figure
481 3 shows that DRInfo consistently yields a higher
482 number of knowledge points per tokens compared
483 to Fandom. This suggests that, despite containing
484 fewer raw tokens, DRInfo encodes character in-
485 formation in a more compact and structured form.
486 In contrast, Fandom pages often include lengthy
487 narrative descriptions, meta information, or loosely
488 related content that dilutes information density.

489 **Our dataset produces a more balanced distri-**
490 **bution of knowledge types and enriches event-**
491 **centric and relational knowledge.** As shown in
492 Figure 4, knowledge points extracted from Fan-
493 dom are heavily skewed toward *personality* and
494 *ability*, reflecting the narrative and fan-driven na-
495 ture of wiki pages, which often emphasize descrip-
496 tive traits and combat abilities over structured life
497 events or relational context. In contrast, our dataset
498 exhibits higher counts in *experience* and *relation-*
499 *ship* categories, which are particularly important
500 for role-playing and long-horizon character behav-
501 ior, as they ground the character in narrative history
502 and social context rather than isolated attributes.

503 These advantages make our dataset more suit-
504 able for scalable and generalizable character-
505 centric applications. The improved compactness,
506 higher information density, and reduced popularity
507 bias demonstrate that our approach offers clear ad-
508 vantages over directly using Fandom pages. This
509 is particularly beneficial for RPLA research, where
510 consistent, concise, and comprehensive character

profiles are critical for controllable and faithful role-playing behavior.

5.3 Ablation Studies

Iterative Search Benefits Profile Quality. As shown in Table 2, removing the iterative search mechanism leads to a consistent drop in Support Score across all knowledge dimensions. The overall Support Score decreases markedly, indicating that a single-round search fails to sufficiently cover diverse and fine-grained character attributes. This confirms that iterative querying and self-refinement are essential for expanding profile coverage and reducing missing information.

The performance degradation is more pronounced on dimensions that require richer contextual evidence. When iterative search is disabled, dimensions such as *Ability*, *Experience*, and *Relationship* exhibit larger declines compared to more factual attributes like *identity* or *appearance*. This suggests that action-oriented, narrative-dependent knowledge benefits more from multi-round evidence aggregation, whereas surface-level attributes can often be captured with fewer search steps. In contrast, more stable attributes such as *Personality* and *Identity* show relatively smaller performance declines.

5.4 Human Evaluation

Due to the inherent label imbalance, we conduct a controlled human evaluation by uniformly sampling 100 instances for each label, resulting in 400 manually annotated cases in total. The details of annotation can be found in Appendix D.

Overall, the LLM-based evaluator shows strong consistency with human judgments. As shown in Figure 5, the confusion matrix indicates high agreement on the supported (SUP) and irrelevant (IRR) categories. Most SUP instances are correctly identified by the LLM, and IRR cases are accurately filtered, demonstrating that the evaluator reliably distinguishes grounded profile content from unrelated information. However, its most prominent failure mode lies in misclassifying *Irrelevant* knowledge as *Contradicted*, indicating that the model sometimes interprets missing evidence as explicit conflict. More detailed analyses can be found in the Appendix E.

Most disagreements occur between supported and partially supported cases, reflecting inherent annotation ambiguity. A noticeable portion

of SUP instances are predicted as PAR and vice versa, suggesting that the boundary between full and partial support is often fuzzy and depends on coverage granularity rather than factual correctness. This observation aligns with the design of our Support Score, where partially supported knowledge is assigned a weight of 0.5, appropriately capturing its intermediate reliability between fully supported and unsupported content.

Contradicted cases are rare and are detected with high recall. The evaluator correctly identifies the majority of contradicted instances, with very few being confused with other categories. This indicates that the LLM is conservative and robust in identifying factual conflicts between knowledge points and character profiles.

These results justify the use of LLM-as-judge for large-scale evaluation. Given its strong alignment with human judgments on clear cases and reasonable behavior on borderline ones, the LLM-based evaluator provides a reliable and scalable alternative to manual annotation for evaluating character profile quality.

6 Conclusion

In this paper, we present a fully automated, search-agent-based framework for large-scale character profile construction across heterogeneous media domains. By leveraging iterative web search and evidence refinement, our approach acquires comprehensive, fine-grained character information without manual annotation or domain-specific retraining. To rigorously assess profile quality, we introduce an LLM-as-a-Judge evaluation framework grounded in external fandom-based knowledge, enabling localized and quantitative analysis of factual alignment and coverage.

Extensive experiments demonstrate that our constructed profiles achieve higher knowledge density, broader attribute coverage, and more balanced representation across profile dimensions compared to fandom-based references, particularly for long-tail and under-represented characters. We further release a high-quality character profile dataset to facilitate future research on role-playing language agents and character-centric generation. We believe this work provides a scalable foundation for open-domain character profiling and reliable evaluation in character-driven AI systems.

Ethical Statement

This work relies exclusively on publicly available web content and fandom wikis, and does not involve any private, sensitive, or personally identifiable information. The constructed dataset focuses on fictional characters from animation and related media, and is released strictly for research purposes.

We acknowledge that fan-curated sources and large language models may contain factual inaccuracies, subjective interpretations, or popularity-driven biases. In addition, while the dataset targets fictional characters, there remains a potential risk of reinforcing stereotypes or misrepresentations embedded in source materials or model generations. We therefore encourage researchers and practitioners to apply critical scrutiny when using the dataset, especially in interactive or user-facing systems.

We also note that large language model-based assistants were used to support language polishing and presentation refinement during manuscript preparation. These tools were not used to generate experimental data, annotations, or evaluation results, and all technical content and conclusions were authored and verified by the researchers.

Overall, we believe the risks associated with this work are limited and manageable. By openly discussing these limitations and releasing the dataset for academic use, we aim to support transparent, responsible research on character-centric language agents and character profiling.

Limitations

Despite its advantages, our approach has several limitations. First, the quality of retrieved evidence depends on the availability and reliability of web sources; for extremely obscure characters, search results may still be sparse or noisy. Second, while the LLM-as-a-Judge framework enables scalable evaluation, it may exhibit systematic biases, such as confusing missing information with contradictions, which we partially mitigate through human evaluation but do not eliminate entirely. Third, our current pipeline focuses on textual character profiling and does not incorporate multimodal evidence such as images or videos, which could further enrich character understanding. Addressing these limitations is left for future work.

References

- Jiangjie Chen, Xintao Wang, Rui Xu, Siyu Yuan, Yikai Zhang, Wei Shi, Jian Xie, Shuang Li, Ruihan Yang, Tinghui Zhu, Aili Chen, Nianqi Li, Lida Chen, Caiyu Hu, Siye Wu, Scott Ren, Ziquan Fu, and Yanghua Xiao. 2024. [From persona to personalization: A survey on role-playing language agents](#). *Transactions on Machine Learning Research*. Survey Certification.
- Nuo Chen, Yan Wang, Haiyun Jiang, Deng Cai, Yuhan Li, Ziyang Chen, Longyue Wang, and Jia Li. 2023. [Large language models meet harry potter: A dataset for aligning dialogue agents with characters](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 8506–8520, Singapore. Association for Computational Linguistics.
- Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke Marris, Sam Petulla, Colin Gaffney, Asaf Aharoni, Nathan Lintz, Tiago Cardal Pais, Henrik Jacobsson, Idan Szpektor, Nan-Jiang Jiang, and 13 others. 2025. [Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities](#). *Preprint*, arXiv:2507.06261.
- Yanqi Dai, Huanran Hu, Lei Wang, Shengjie Jin, Xu Chen, and Zhiwu Lu. 2025. [Mmrole: A comprehensive framework for developing and evaluating multimodal role-playing agents](#). In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net.
- DeepSeek-AI, Aixin Liu, Aoxue Mei, Bangcai Lin, Bing Xue, Bingxuan Wang, Bingzheng Xu, Bochao Wu, Bowei Zhang, Chaofan Lin, Chen Dong, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenhao Xu, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, and 245 others. 2025. [Deepseek-v3.2: Pushing the frontier of open large language models](#). *Preprint*, arXiv:2512.02556.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. [Retrieval-augmented generation for knowledge-intensive NLP tasks](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Cheng Li, Ziang Leng, Chenxi Yan, Junyi Shen, Hao Wang, Weishi MI, Yaying Fei, Xiaoyang Feng, Song Yan, HaoSheng Wang, Linkang Zhan, Yaokai Jia, Pingyu Wu, and Haozhen Sun. 2023. [Chatharuhi: Reviving anime character in reality via large language model](#). *Preprint*, arXiv:2308.09597.
- Zekun Li, Zhiyu Chen, Mike Ross, Patrick Huber, Seungwhan Moon, Zhaojiang Lin, Xin Dong, Adithya Sagar, Xifeng Yan, and Paul A. Crook. 2024. [Large](#)

713	language models as zero-shot dialogue state tracker through function calling. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 8688–8704. Association for Computational Linguistics.	770
714		771
715		772
716		773
717		774
718		775
719	Keming Lu, Bowen Yu, Chang Zhou, and Jingren Zhou. 2024. Large language models are superpositions of all characters: Attaining arbitrary role-play via self-alignment. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 7828–7840. Association for Computational Linguistics.	776
720		777
721		778
722		779
723		780
724		781
725		782
726		783
727	Joon Sung Park, Joseph C. O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. Generative agents: Interactive simula- cra of human behavior. In <i>Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology, UIST 2023, San Francisco, CA, USA, 29 October 2023- 1 November 2023</i> , pages 2:1–2:22. ACM.	784
728		785
729		786
730		787
731		788
732		789
733		790
734		791
735	Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	792
736		793
737		794
738		795
739		796
740		797
741		798
742		799
743	ByteDance Seed, :, Jiaze Chen, Tiantian Fan, Xin Liu, Lingjun Liu, Zhiqi Lin, Mingxuan Wang, Chengyi Wang, Xiangpeng Wei, Wenyuan Xu, Yufeng Yuan, Yu Yue, Lin Yan, Qiyang Yu, Xiaochen Zuo, Chi Zhang, Ruofei Zhu, Zhecheng An, and 255 others. 2025. Seed1.5-thinking: Advancing superb reasoning models with reinforcement learning. <i>Preprint</i> , arXiv:2504.13914.	800
744		801
745		802
746		803
747		804
748		805
749		806
750		807
751	Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. Character-LLM: A trainable agent for role-playing. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 13153–13187, Singapore. Association for Computational Linguistics.	808
752		809
753		810
754		811
755		812
756		813
757	Lisa Sheehy, Stéphane Bouchard, Anupriya Kakkar, Rama El Hakim, Justine Lhoest, and Andrew Frank. 2024. Development and initial testing of an artificial intelligence-based virtual reality companion for people living with dementia in long-term care. <i>Journal of Clinical Medicine</i> , 13(18):5574.	814
758		815
759		816
760		817
761		818
762		819
763		820
764		821
765		822
766		823
767		824
768		825
769		826
		827
	Quan Tu, Shilong Fan, Zihang Tian, Tianhao Shen, Shuo Shang, Xin Gao, and Rui Yan. 2024. CharacterEval: A Chinese benchmark for role-playing conversational agent evaluation. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 11836–11850, Bangkok, Thailand. Association for Computational Linguistics.	
	Noah Wang, Z.y. Peng, Haoran Que, Jiaheng Liu, Wangchunshu Zhou, Yuhan Wu, Hongcheng Guo, Ruitong Gan, Zehao Ni, Jian Yang, Man Zhang, Zhaoxiang Zhang, Wanli Ouyang, Ke Xu, Wenhao Huang, Jie Fu, and Junran Peng. 2024. RoleLLM: Benchmarking, eliciting, and enhancing role-playing abilities of large language models. In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 14743–14777, Bangkok, Thailand. Association for Computational Linguistics.	
	Xintao Wang, Heng Wang, Yifei Zhang, Xinfeng Yuan, Rui Xu, Jen tse Huang, Siyu Yuan, Haoran Guo, Jiangjie Chen, Shuchang Zhou, Wei Wang, and Yanghua Xiao. 2025. CoSER: Coordinating LLM-based persona simulation of established roles. In <i>Forty-second International Conference on Machine Learning</i> .	
	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025. Qwen3 technical report. <i>Preprint</i> , arXiv:2505.09388.	
	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan Cao. 2023. React: Synergizing reasoning and acting in language models. In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	
	Xinfeng Yuan, Siyu Yuan, Yuhan Cui, Tianhe Lin, Xintao Wang, Rui Xu, Jiangjie Chen, and Deqing Yang. 2024. Evaluating character understanding of large language models via character profiling from fictional works. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 8015–8036, Miami, Florida, USA. Association for Computational Linguistics.	
	Jiajia Zhao, Jingru Zhang, and Yuhe Lu. 2025. Enhancing design historical education through AI virtual characters role-playing narratives in serious games. <i>Int. J. Gaming Comput. Mediat. Simulations</i> , 17(1):1–20.	
	Jinfeng Zhou, Zhuang Chen, Dazhen Wan, Bosi Wen, Yi Song, Jifan Yu, Yongkang Huang, Pei Ke, Guanqun Bi, Libiao Peng, JiaMing Yang, Xiyao Xiao, Sahand Sabour, Xiaohan Zhang, Wenjing Hou, Yijia Zhang, Yuxiao Dong, Hongning Wang, Jie Tang, and Minlie Huang. 2024. CharacterGLM: Customizing social characters with large language models. In	

828 *Proceedings of the 2024 Conference on Empirical*
829 *Methods in Natural Language Processing: Industry*
830 *Track*, pages 1457–1476, Miami, Florida, US. Asso-
831 ciation for Computational Linguistics.

832 Jinfeng Zhou, Yongkang Huang, Bosi Wen, Guanqun
833 Bi, Yuxuan Chen, Pei Ke, Zhuang Chen, Xiyao
834 Xiao, Libiao Peng, Kuntian Tang, Rongsheng Zhang,
835 Le Zhang, Tangjie Lv, Zhipeng Hu, Hongning Wang,
836 and Minlie Huang. 2025. **Characterbench: Bench-**
837 **marking character customization of large language**
838 **models**. In *AAAI-25, Sponsored by the Association*
839 *for the Advancement of Artificial Intelligence, Febru-*
840 *ary 25 - March 4, 2025, Philadelphia, PA, USA*, pages
841 26101–26110. AAAI Press.

842	A Experimental Setup		
843	The hyperparameters of profile construction and		
844	evaluation can be see in Table 3. For evaluation,		
845	we deploy Qwen3-235B-A22B-Instruct-2507 us-		
846	ing data-parallel inference on eight NVIDIA A800		
847	GPUs with 80GB memory each.		
848	B Dataset Statistics		
849	This section provides a brief statistical overview		
850	of the CHARACTERHUB character dataset used		
851	throughout the paper. The dataset contains 1,109		
852	high-precision characters spanning 399 franchises		
853	across multiple media domains. Most characters		
854	in the dataset are drawn from anime franchises,		
855	with the remainder covering games, movies, and		
856	novels <i>etc.</i> Each character is associated with a		
857	comprehensive first-person profile covering seven		
858	core dimensions. On average, profiles in CHAR-		
859	ACTERHUB exhibit substantially higher informa-		
860	tion density than their Fandom references, while		
861	maintaining strong alignment with human-curated		
862	references. Table 4 summarizes key statistics of the		
863	dataset.		
864	C Knowledge Type Dimensions		
865	To provide a more granular understanding of the		
866	evaluation framework, this section details the seven		
867	semantic dimensions used for character knowledge		
868	categorization. These dimensions are designed to		
869	cover the full spectrum of a character’s persona,		
870	ranging from static factual identifiers to latent psy-		
871	chological traits. Each category serves as a distinct		
872	axis for evaluating the comprehensiveness and fac-		
873	tual integrity of the generated profiles.		
874	Identity Identity refers to factual and relatively		
875	stable attributes that define who the character is.		
876	This includes the character’s name, origin, birth-		
877	place, nationality, species or race, as well as formal		
878	titles, ranks, aliases, nicknames, professions, and		
879	social or institutional affiliations such as organiza-		
880	tions, teams, clans, or factions. These attributes		
881	typically serve as the core identifiers of a charac-		
882	ter and are often explicitly stated in encyclopedic		
883	sources.		
884	Appearance Appearance captures the character’s		
885	physical and visual traits, including body features,		
886	facial characteristics, clothing style, colors, acces-		
887	sories, and other distinctive visual elements. This		
888	category focuses on how the character looks, and is		
	commonly used in visual media such as animation,	889	
	films, and games.	890	
	Ability Ability describes the character’s compe-	891	
	tencies and capabilities, including skills, powers,	892	
	techniques, magical abilities, combat styles, and	893	
	special talents. This category encompasses both	894	
	innate abilities and learned skills, and reflects what	895	
	the character is capable of doing in practice, espe-	896	
	cially in action-driven narratives.	897	
	Relationship Relationship covers the character’s	898	
	social connections and interpersonal ties, such as	899	
	family members, friends, allies, rivals, enemies,	900	
	romantic partners, or other significant associations.	901	
	These relationships define the character’s position	902	
	within a social network and are central to narra-	903	
	tive development and character interaction, a dis-	904	
	tinction commonly emphasized in character-centric	905	
	narrative analysis.	906	
	Experience Experience represents the charac-	907	
	ter’s past actions and life events, including impor-	908	
	tant battles, missions, training experiences, achieve-	909	
	ments, failures, and major story arcs. Unlike iden-	910	
	tity, which is relatively static, experience is event-	911	
	driven and reflects what the character has done	912	
	over time, contributing to narrative progression and	913	
	character growth.	914	
	Personality Personality characterizes the char-	915	
	acter’s internal traits and behavioral tendencies,	916	
	such as temperament, motivations, desires, goals,	917	
	attitudes, emotional patterns, moral stance, and ha-	918	
	bitual ways of thinking or acting. These attributes	919	
	are often implicit and inferred from behavior or di-	920	
	alogue rather than explicitly stated, and have been	921	
	widely studied in computational character model-	922	
	ing.	923	
	Other Other is a residual category for miscella-	924	
	neous information that does not fit the above di-	925	
	mensions, such as personal preferences, hobbies,	926	
	or trivial habits.	927	
	D Human Annotation Instructions	928	
	To examine the alignment between LLM-based	929	
	judgments and human judgments, we conduct a	930	
	small-scale manual annotation study on a subset	931	
	of knowledge–profile pairs. All annotations are	932	
	performed internally by the authors for validation	933	
	purposes.	934	
	For each instance, annotators are provided with a	935	
	generated character profile and a single knowledge	936	

Parameter	Gemini	Doubao	Qwen-plus	Qwen3-235B
Temperature	1.0	1.0	0.85	0.7
Top- p	0.95	0.8	0.8	0.8
Top- k	64	50	50	20

Table 3: Hyperparameter settings used for different models.

Statistic	Value
Number of Characters	1109
Number of Franchises	399
Media Domains	Mainly Anime
Avg. Profile Length (words)	10372
Avg. Profile tokens	2187
Avg. Knowledge Points per Character	24
Information Density (per 1k token)	11.68

Table 4: Statistical overview of the CHARACTERHUB dataset.

point extracted from the Fandom-based ground truth. Annotators are asked to assess the factual relationship between the two using the same four labels as in the LLM evaluation: *Supported*, *Partially Supported*, *Irrelevant*, and *Contradicted*.

Judgments must be based only on the information explicitly contained in the given profile. Annotators should not rely on external knowledge, future events, or information from other versions or spin-offs unless such content is explicitly mentioned in the profile.

All annotated content concerns fictional characters from publicly available sources and involves no personal or sensitive information. As annotations are conducted by the authors without external participants, no recruitment, compensation, or risk disclosure is required.

E Case Study

To better understand the failure modes of automated character profile construction and evaluation, we conduct a case study focusing on knowledge points that are judged as *Contradicted* by the LLM-based evaluator. We choose this subset for analysis because contradicted cases represent the most explicit and verifiable form of errors, where the extracted or generated knowledge directly conflicts with canonical evidence, making them particularly suitable for qualitative inspection. By examining these samples in detail and comparing LLM judgments with human annotations, we aim to both find common reasons of factual inconsistency in character profiling and assess the reliability and limitations of LLM-as-a-Judge in identifying such

errors.

E.1 Contradicted Types

From those cases in which both the LLM judge and human annotators consistently identify knowledge points as *Contradicted*, we observe five contradiction types: (1) *Cross-version Confusion* refers to contradictions caused by mixing information across different canons, adaptations, spin-offs, or alternate versions of the same character, where events or attributes valid in one version are incorrectly projected onto another. (2) *Ability Misassignment* occurs when a character is attributed with abilities or powers belonging to a different character or continuity. (3) *Objective Attribute Error* captures fine-grained but verifiable factual mistakes, such as incorrect species, titles, or fixed descriptors. (4) *Temporal Inconsistency* captures cases where a knowledge point is factually correct in a later narrative stage but conflicts with the profile due to an implicit mismatch in temporal scope. (5) *Fabricated Narrative Event* denotes hallucinated plot developments that directly contradict canonical story outcomes.

As illustrated in Table 5, these contradiction types manifest as clear, localized conflicts with fandom reference. Taking *Cross-version Confusion* as an example, *Haruhi Suzumiya* is described as having created an alternate world due to her dissatisfaction with reality; however, this event occurs in other derivative or alternate-universe narratives and does not take place in *Nyoron! Churuya-san*. In the case of *Yuu Ishigami*, he does eventually become the Student Council Vice President later in the narrative; however, the generated profile only includes his present role as Treasurer.

E.2 Bias of LLM-Judge

From those cases where the LLM judge labels as *Contradicted* while human annotators assign a different label, we conclude 3 major bias types exhibited by the LLM judge. (1) *Absence-as-Contradiction Bias*, where missing or unspecified information in the profile is incorrectly treated as a factual contradiction. (2) *Temporal Misinterpre-*

Knowledge Point	LLM Evidence	Contradiction Type
<i>Haruhi Suzumiya (Nyoron! Churuya-san):</i> Haruhi created an alternate world due to her dissatisfaction with reality.	She does not create alternate worlds in this series; the world remains static and mundane.	Cross-version Confusion
<i>Illyasviel von Einzbern (Carnival Phantasm):</i> She was attacked by Zouken, leading to Berserker’s corruption.	In this version, Berserker remains under Illya’s control.	
<i>Jousuke Higashikata (JoJolion):</i> He exploits Crazy Diamond’s restoration ability for tactical advantage.	His Stand is Soft & Wet, not Crazy Diamond.	Ability Misassignment
<i>Ein (Cowboy Bebop):</i> Ein is a Cardigan Welsh Corgi.	The text explicitly identifies Ein as a Pembroke Welsh Corgi.	Objective Attribute Error
<i>Yuu Ishigami (Kaguya-sama):</i> He is the Student Council Vice President and former Treasurer.	The text states he is the treasurer, not Vice President.	Temporal Inconsistency
<i>Lelouch Lamperouge (Code Geass):</i> He resurrects after death and returns to protect the world.	The text presents his death as final; no resurrection is mentioned.	Fabricated Narrative Event

Table 5: Representative **contradicted** knowledge points where LLM-judge and human annotations agree. Each row is separated by horizontal rules for clarity, and all entries are vertically centered.

1013 *tation Bias*, where the LLM fails to distinguish between current, future, or alternative timeline facts.
1014
1015 (3) *External Lore Intrusion Bias*, where the LLM
1016 relies on external or non-local canon knowledge
1017 not present in the evaluated profile.

1018 As shown in Table 6, a dominant source of disagreement is *Absence-as-Contradiction Bias*, arising from the LLM judge’s tendency to treat missing information as explicit contradiction. For instance,
1019 in the cases of *Hei* and *Rei Kiriyama*, the profiles
1020 simply omit the described events, which human annotators correctly label as Irrelevant, whereas the
1021 LLM judge interprets absence as negation.
1022
1023
1024
1025

1026 F Prompt Templates

1027 Specifically, Table 7 present the prompts for the
1028 iterative search stage, where the first prompt
1029 initializes diversified character-oriented queries and
1030 the second prompt refines subsequent searches by
1031 incorporating summarized evidence from the
1032 previous round. Table 8 shows the prompt used for
1033 profile rewriting, which integrates multi-source
1034 evidence into a coherent first-person character
1035 self-introduction tailored for role-playing agents. For
1036 evaluation, Table 9 details the prompt for extracting
1037 fine-grained knowledge points from fandom
1038 reference, while Table 10 presents the LLM-as-a-Judge
1039 prompt used to assess the alignment between
1040 individual knowledge points and the generated profiles.
1041 Together, these prompts operationalize the full au-

1042 tomatic pipeline from information acquisition to
1043 quantitative profile evaluation.

Knowledge Point	LLM Evidence (Excerpt)	Human Label	Bias Type
<i>Hei (Darker than Black):</i> Hei struggled with the moral conflict of potentially having to kill his sister Bai.	No indication he ever considered or was required to kill her.	Irrelevant	
<i>Rei Kiriyama (March Comes in Like a Lion):</i> He lacks common sense, shown by attempting to register his engagement.	The text does not mention any attempt to register an engagement or lack of common sense related to legal procedures.	Irrelevant	Absence-as-Contradiction
<i>Nagisa Furukawa (Clannad):</i> She has two M-shaped ahoge inherited from both parents.	Ahoge is described, but no M-shape or inheritance is mentioned.	Partially Supported	
<i>Hiroto Suwa (Orange):</i> In an alternate timeline, Hiroto marries Naho and has a child.	Marriage occurs in the original timeline, not the alternate one.	Supported	Temporal Misinterpretation
<i>Asuka Langley Souryuu (Evangelion 2.0):</i> Her mother committed suicide after an EVA experiment.	Rebuild version presents a different backstory.	Irrelevant	External Lore Intrusion
<i>Princess Leia (Star Wars):</i> She survives in space using the Force.	Such an ability is neither mentioned nor canonically supported.	Irrelevant	

Table 6: Representative cases where the LLM judge marks knowledge points as **Contradicted** while human annotators assign other labels. The disagreements reveal systematic LLM judge biases.

First Search Prompt: Please help me collect all relevant knowledge about `{{entity}}` and integrate it into a detailed encyclopedic document. Use markdown format. Be truthful, comprehensive, and thorough, without omitting important information. You must verify the authenticity of information and the credibility of sources. Do not include or fabricate false information.

You should include the character’s personality (very important), background, physical description, core motivations, notable attributes, relationships, key experiences, major plot involvement and key decisions or actions, character arc or development throughout the story, if there is any information about these aspects.

Iterative Search Prompt: Please continue searching to verify and expand on the knowledge you have already collected, such as searching for closely related people, events, objects, and other entities.

Table 7: First-round and iterative search prompt used in the CHARACTERHUB pipeline.

Prompt: Please completely rewrite all the above information from `{{entity}}`’s first-person perspective. Ensure that the information is comprehensive and accurate. You need to focus on the character’s personality, which you could also analyze based on the characters’ experiences. Besides, include include the character’s, background, physical description, core motivations, notable attributes, relationships, key experiences, major plot involvement and key decisions or actions, character arc or development throughout the story, and other important details, if they appear in the information that you obtained.

Table 8: Prompt for rewriting refined evidence into a first-person character profile.

Prompt: You will receive a piece of character data (in JSON format) crawled from the Fandom website. The data includes the character’s basic information (infobox) and textual descriptions (content).

Please parse the input JSON and extract all **knowledge items** related to the character, and focus on *meaningful, distinctive, and role-relevant* information that contributes to understanding or role-playing this character, labeling each item with its type.

Each knowledge item should represent a semantically complete factual statement or a clearly defined attribute. For example: “He is a passionate young man determined to XXX” or “He is humanity’s strongest soldier.”

Output the results in a unified structured format.

Please IGNORE any information that is irrelevant to character role-playing, including:

- Field echo information that only restates a label or name (e.g. “His name is XXX.”).
- Voice actors, performers, production companies, publishers, episode or chapter numbers (e.g. “voiced by ...”).
- Language, romanization, or spelling variants (e.g. Kanji, Romaji, kana notation).
- Generic or obvious facts that do not distinguish this character from ordinary people (e.g. “He is human” or “He is male”), unless they have clear narrative or psychological significance (e.g. “He is human but can transform into a giant”—at this point, “human” information has semantic value).
- Statistical or external information (e.g. number of appearances, popularity rankings, fan titles).
- External commentary or meta descriptions about the character (e.g. “first appeared in episode X”).
- Any meta-level information not belonging to the character’s in-universe context.

Output Format

Please output your result in JSON format as follows:

```
{
  "entity": "Character Name",
  "knowledge_points": [
    {
      "type": "identity | appearance | ability | relationship | experience | personality | other",
      "evidence": "Original supporting text or sentence from the input",
      "knowledge": "Extracted knowledge statement"
    },
    ...
  ]
}
```

Rules

1. Each knowledge item should be semantically clear, self-contained, and useful for **understanding or role-playing the character**.
2. Merge and deduplicate overlapping or equivalent facts. Keep only one representative version of each unique fact. Do not output two items that share the same meaning even if worded differently. If two statements express similar meaning, combine them into one richer description.
3. **Do NOT invent or infer** beyond the text. If something is only hinted at or ambiguous, exclude it.
4. Ensure the output is valid JSON with complete and logically consistent fields.
5. If certain categories are not present, you may omit them — do not fill with “unknown”.
6. You must only use the following categories for type:
 - identity: factual identity, origin, birthplace, nationality, species, titles, nicknames, ranks, aliases, professions, affiliations with organizations, teams, clans, groups, and roles or statuses associated with them.
 - appearance: physical traits, looks, body, clothing, colors, and distinctive visual features.
 - ability: skills, powers, techniques, magic, fighting capability, special talents, or combat style.
 - relationship: family members, friends, rivals, lovers, allies, enemies, or any significant social or character connections.
 - experience: key life events, battles, training, missions, achievements, historical background, or things the character has done.
 - personality: temperament, behavior patterns, motivations, desires, goals, attitudes, emotional tendencies, moral stance, or thinking style.
 - other: anything that does not reasonably fit into the above categories. Do **not** create or use any other category.

Input Data:

```
{{input text}}
```

Table 9: Prompt for extracting structured knowledge points from reference character descriptions.

Prompt: You are an expert evaluator for role-playing character consistency.

Task

You will receive two items:

1. `knowledge_list`: a JSON array of knowledge objects about a character, each describing a factual aspect of the character. Each object has:

- `"id"`: a unique integer identifier for the knowledge item,
 - `"knowledge"`: a concise factual statement about the character.
2. `text`: a piece of character-related text. This may be:
- a first-person profile written from the character's perspective,
 - a third-person descriptive text about the character,
 - or structured JSON data crawled from a source such as Fandom.
- Treat all three forms uniformly as sources of factual claims.

Your goal is to determine, for each knowledge item in `knowledge_list`, determine whether the provided text:

- **supported**: The text explicitly expresses or clearly implies this knowledge point.
- **partially supported**: The text supports some but not all aspects of the knowledge point (e.g., mentions part of a list or partial causation).
- **irrelevant**: The text does not mention or relate to this knowledge point.
- **contradicted**: The text states or implies something that conflicts with this knowledge point.

Do NOT use any labels other than the four above.

Important Rules

1. Focus on **semantic meaning**, not wording. Example: if the knowledge says "He is calm and collected" and the text says "I rarely lose my temper," that is **supported**.
2. A **contradicted** label should be used only if the text explicitly or implicitly rejects or reverses something of the fact.
3. Missing details or partial lack of coverage should be labeled **irrelevant**, not contradicted.
4. **Do NOT treat narrative perspective or tense as evidence**. Do NOT interpret first-person narration ("I") or present tense ("I am", "I do") as evidence that the character is alive or active in the current timeline. Only mark a contradiction if the text explicitly denies or reverses the fact.
5. When uncertain, choose the most reasonable interpretation based on the content of the given text.
6. For each evaluation, provide a brief quote or reasoning from the text that justifies your judgment.

Example

knowledge_list:

```
[
  {"id": 1, "knowledge": "Eren was born in the Shiganshina District."},
  {"id": 2, "knowledge": "Eren is calm and patient."}
]
```

text:

"I grew up behind the walls of Shiganshina, dreaming of the world beyond. I've always been reckless and driven by anger – patience was never my strength. But since that day, I've carried the Titan power within me."

Expected Output:

```
[
  {
    "id": 1,
    "knowledge": "Eren was born in the Shiganshina District.",
    "evaluation": "supported",
    "evidence": "He says he grew up in Shiganshina."
  },
  {
    "id": 2,
    "knowledge": "Eren is calm and patient.",
    "evaluation": "contradicted",
    "evidence": "He describes himself as reckless and impatient."
  }
]
```

Input

`knowledge_list`: {{knowledge_list}}

`text`: {{character_text}}

Table 10: Prompt for knowledge–profile alignment using an LLM-as-a-Judge framework.