

# MythNER: Multi-Agent NER Extraction and Benchmarking in Chinese Myth Narratives

Anonymous ACL submission

## Abstract

Named entity recognition (NER) performs strongly in well-studied domains, yet mythological narratives pose a long-tail setting where entityhood is defined by referable instances, mentions vary through titles and aliases, and correct extraction requires stable long-span boundaries. We introduce MYTHNER, a culturally grounded NER benchmark for Chinese mythology built from *Ne Zha* subtitle text. MythNER uses four flat labels (PER/LOC/ORG/OBJ) with conservative annotation criteria; in particular, OBJ is restricted to *named* artifacts and fixed technique titles rather than generic concepts. We evaluate a zero-shot Chinese spaCy model, a supervised BERT token-classification baseline, and a multi-agent LLM extraction pipeline with chunk/context tuning on the held-out *Ne Zha Part 2* test set. Results show substantial domain shift for off-the-shelf NER, strong gains from supervised adaptation, and further improvements from agentic extraction under well-chosen constraints. Our analysis characterizes dominant failure modes—boundary drift under exact-span scoring, mythology-specific type ambiguity, and over-extraction of generic nouns as OBJ—highlighting the need for iterative consistency checks in narrative-domain NER.

## 1 Introduction

Named entity recognition (NER) is a foundational component for information extraction, question answering, and knowledge base construction (Li et al., 2022; Jehangir et al., 2023).

While modern NER systems perform strongly in well-studied domains such as newswire, their assumptions often break in narrative text. Narratives routinely exhibit

aliases and titles, implicit references, invented lexicons, and long-range dependencies in which the identity and type of an entity is only resolved after additional context (Bamman et al., 2019; Chu et al., 2020).

In this regime, accurate tagging is less a single-pass string labeling problem and more an exercise in *iterative consistency checks*: predictions must remain coherent across repeated mentions, shifts in viewpoint, and dispersed evidence. This motivates workflows that go beyond local extraction, pairing chunk-level recognition with document-level consolidation and verification to enforce global consistency.

Chinese mythology provides a particularly revealing stress test for narrative NER. Mythological texts, for instance fantasy movie subtitles, contain culturally grounded naming conventions, honorific- and role-heavy mentions, and named artifacts or techniques that behave like proper names, all of which can blur boundaries and types. For Chinese, these challenges are compounded by segmentation ambiguity and compact mentions, increasing the likelihood of span and label confusions (Gui et al., 2019).

At the same time, existing narrative NER resources tend to emphasize either Western myth and fantasy (similar narrative structure but different cultural and linguistic grounding) or Chinese fictional literature (shared language but a different entity ecology), leaving a gap for Chinese mythology. To address this gap, we introduce MYTHNER, a benchmark built from Chinese subtitles of *Ne Zha* films (Parts 1–2), designed to evaluate myth-domain NER under a controlled narrative setting.

We benchmark off-the-shelf zero-shot NER, a supervised encoder baseline, and an agentic extraction pipeline that decomposes long-context NER into chunked extraction, global

consolidation, and verification/correction (Wang et al., 2025b; Zhang et al., 2024). Across this ladder, off-the-shelf tools struggle under myth-domain shift, supervised training improves robustness, and agentic extraction yields further gains by enforcing cross-chunk consistency.

Our contributions are:

- **Dataset:** a culturally grounded Chinese mythology NER benchmark with clear guidelines and a coarse schema (PER/LOC/ORG/OBJ).
- **Benchmark evidence:** a controlled comparison that characterizes the gap between off-the-shelf NER and in-domain modeling for myth narratives.
- **Method:** an agentic extraction pipeline that performs chunked extraction followed by global consolidation and verification/correction to improve document-level consistency.

## 2 Related Work

Named entity recognition has progressed from lexicon- and rule-based systems to statistical sequence labeling, neural encoders, and transformer-based representations (Li et al., 2022; Jehangir et al., 2023; Pakhale, 2023). Despite strong performance on mainstream benchmarks, fictional and mythological narratives remain challenging: entity identity is shaped by evolving roles, aliases, and world-specific ontologies, and errors compound when the same character must be tracked consistently across long contexts. Literary and fiction-oriented resources highlight domain-specific entity definitions and distributions (Bamman et al., 2019), and practical systems for fiction emphasize the need for domain-specific typing and consolidation strategies (Chu et al., 2020). Complementary analyses in other narrative-like domains (e.g., tabletop role-playing corpora) likewise show systematic differences from standard NER settings (Weerasundara and de Silva, 2023), and fine-tuning studies in fantasy settings further support the value of in-domain adaptation (Sivaganeshan and De Silva, 2023).

For Chinese and historically oriented text, persistent issues include segmentation ambigu-

ity and sparse supervision; lexicon-aware modeling remains a common mitigation for Chinese NER (Gui et al., 2019). Adjacent evaluations in ancient Chinese NER further illustrate both the progress and the domain gaps that arise once genre, period, and entity taxonomy shift (Li et al., 2025). Beyond modeling, annotation studies and multi-genre benchmarks underline that entity definitions are not purely technical: label sets and annotator interpretation can vary, affecting both training and evaluation (Tedeschi and Navigli, 2022; Peng et al., 2024).

Recent work explores using large language models (LLMs) directly for NER in zero-/few-shot settings (Wang et al., 2025a). However, narrative-scale extraction stresses long-context tracking and cross-mention consistency, motivating structured agentic alternatives. Multi-agent frameworks decompose extraction into specialized roles (e.g., proposal, verification against an ontology, aggregation), which is especially relevant when outputs must satisfy schema constraints (Tao et al., 2025; Wang et al., 2025b). Long-context collaboration mechanisms provide another path to maintaining state across long documents (Zhang et al., 2024; Zhuang et al., 2025).

## 3 Dataset and Annotation

### 3.1 Corpus

MythNER is a culturally-grounded named entity recognition benchmark built from Chinese subtitle text from the animated *Ne Zha* films. The domain is a Chinese mythological narrative world (e.g., *Fengshen*-related traditions and *Journey to the West*-related motifs), which is widely represented across Chinese myth literature and popular media. We use this setting as a representative high-context domain in which entity mentions include characters (PER), places (LOC), organizations/factions (ORG), and named artifacts or technique titles (OBJ).

### 3.2 Task and label schema

We perform coarse-grained NER with four flat labels: PER, LOC, ORG, and OBJ. The task follows standard span-based NER: given a subtitle sequence, the model predicts entity spans and assigns one of the four types.

### 3.3 Data split

We split the corpus by film to evaluate generalization across sequels. Specifically, we use *Ne Zha Part 1* (P1) for training and validation, and *Ne Zha Part 2* (P2) as a held-out test set. Within P1, we randomly split into 80% training and 20% validation. The P1 split is used for training and tuning supervised baselines. The entire P2 test set is held out for evaluation, including fair comparisons for zero-shot spaCy and our LLM-agent methods.

### 3.4 Dataset statistics and entity overlap

Figure 1 summarizes the dataset size and label distribution. Overall, MythNER contains 697 entity mentions (159 unique entities) under the four-label schema. The held-out P2 test set contains 415 mentions (106 unique), with per-label totals of 150 PER, 90 ORG, 75 LOC, and 100 OBJ.

To quantify cross-film generalization difficulty, we measure unique-entity overlap between training and test. Approximately 31.1% of unique entities in the test set also appear in training, implying that most test-set unique entities are unseen during training.

### 3.5 Annotation workflow

MythNER is annotated by three annotators and one adjudicating judge. Annotation is conducted in Doccano, self-hosted on a Google Cloud VM. Annotators label entity spans and coarse types (PER/LOC/ORG/OBJ). Ambiguous cases are resolved via discussion to align on guidelines, after which the judge performs final review and consolidation.

### 3.6 Annotation guidelines

We follow a flat NER schema without nested or overlapping spans. We adopt conservative criteria for entityhood: mentions should refer to an independently identifiable instance in context. We exclude purely event/phenomenon expressions, and treat OBJ as a strict category that covers named artifacts and fixed technique titles rather than generic concepts. When a surface form of a location word can plausibly be interpreted as a place or an institution depending on its syntactic role, we consistently annotate it as LOC. When a full canonical name appears, we prefer the longest

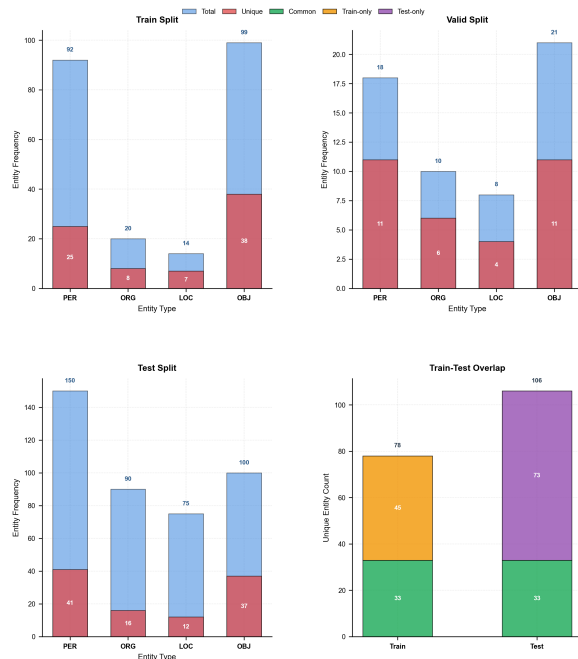


Figure 1: Dataset statistics for MythNER. The figure summarizes entity counts and label distributions across train/validation (P1) and test (P2), and reports the proportion of unique entities in the test set that overlap with training (31.1%).

valid span boundary. Appendix A provides a one-page version of our guidelines with examples.

## 4 Methods

### 4.1 Overview

We compare three families of approaches for PER/LOC/ORG/OBJ schema NER on MythNER: (i) a zero-shot spaCy baseline, (ii) a fine-tuned BERT token-classification baseline, and (iii) our proposed multi-agent LLM extraction system.

Our core hypothesis is that mythological narratives amplify the failure modes of off-the-shelf NER: long-range coreference, culturally specific aliases/titles, dense entity clusters, and heavy use of metaphorical or honorific expressions. The multi-agent design addresses these issues by (a) extracting with local context, (b) enforcing global consistency across chunks, and (c) verifying and selectively correcting uncertain spans.

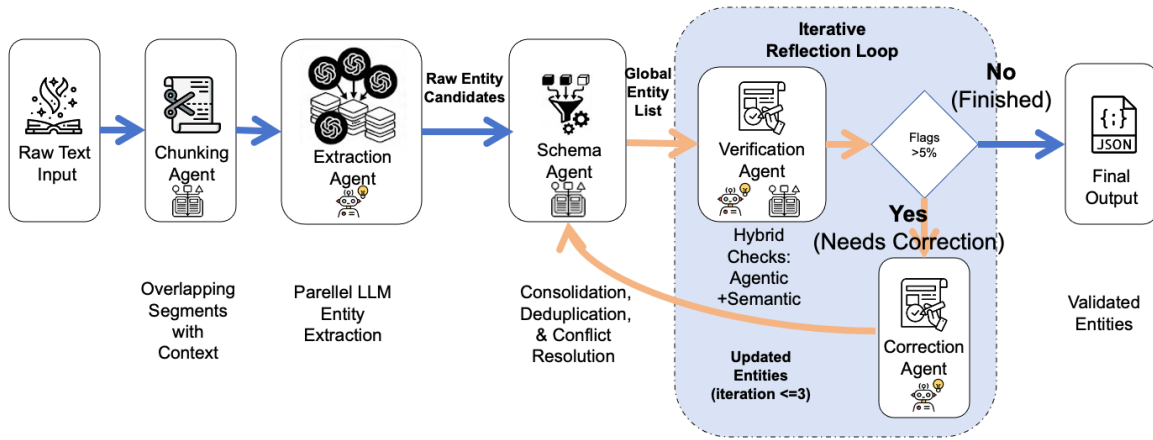


Figure 2: Multi-agent LLM extraction architecture. A text chunking agent creates overlapping segments with context; an extraction agent runs parallel LLM calls to propose entity spans; a consolidation agent deduplicates and resolves conflicts; a verification agent applies hybrid checks (agentic + semantic) and triggers an iterative reflection loop (up to 3 iterations) for targeted correction. The output is a validated entity set.

## 4.2 Proposed method: multi-agent LLM extraction

**Workflow.** The proposed system is a LangGraph-based pipeline<sup>1</sup> with five specialized roles: *chunking* (segmentation), *extraction* (LLM-based span proposal), *consolidation* (global merge across chunks), *verification* (QA over spans and labels), and *correction* (targeted re-extraction for flagged cases). Figure 2 depicts the following workflow:

Input → Chunk → Extract → Consolidate → Verify → [Correct → Re-verify] → Output.

Verification may trigger correction; the loop terminates early if verification passes, and is capped at 3 correction iterations.

**Chunking with context windows.** Given a full subtitle document, we split it into chunks of  $K$  characters and attach left/right context windows of  $C$  characters. This design makes extraction robust to boundary effects (e.g.,

a title introduced in one sentence and referenced in the next) while keeping each LLM call within a bounded context. Chunk boundaries are adjusted to nearby sentence-ending punctuation to reduce mid-phrase truncation.

**Extraction prompting and confidence.** The extraction stage queries an LLM with prompts specialized for (i) output format (document-level spans vs sentence-level spans), (ii) prompt type (a concise zero-shot prompt vs a richer prompt that includes guidelines and examples), and (iii) optional confidence reporting (continuous  $[0, 1]$  or discrete 0–5 ranks).

**Consolidation, verification, and targeted correction.** Chunk-level proposals are merged into a single entity set for the document. A dedicated verifier then checks span validity (e.g., offsets are well-formed and aligned to the text), label plausibility (e.g., preventing obvious confusions such as locations mislabeled as persons), and cross-chunk consistency (e.g., repeated mentions of the same named entity). Only flagged items are sent to the correction stage, which performs targeted fixes rather than re-running the full extraction.

<sup>1</sup>LangGraph v1.0.5, MIT License. Homepage: <https://docs.langchain.com/oss/python/langgraph/overview>

**Reproducibility knobs.** The LLM used is Gemini-2.5-pro We record the chunk/context settings ( $K, C$ ), prompt type, confidence mode, parallelism, and the LLM model identifier used in each run. Prompt templates are provided in Appendix B. These knobs are used in our controlled studies and ablations.

### 4.3 Baselines

#### 4.3.1 Zero-shot spaCy

We use the pretrained zh\_core\_web\_lg spaCy NER model without fine-tuning. Because spaCy’s label set does not match our schema, we apply a deterministic mapping into {PER, LOC, ORG, OBJ}. Unmapped spaCy labels (e.g., DATE/TIME/QUANTITY) are ignored. Predictions are produced at the sentence level.

#### 4.3.2 Fine-tuned BERT

We fine-tune bert-base-chinese for token classification with BIO tagging. The training label set is: {O, B/I-PER, B/I-LOC, B/I-ORG, B/I-OBJ}.

**Training.** HuggingFace Trainer is used for BERT model fine-tuning. We fix 3 random seeds (20261, 12345, 54321) and train each run for up to 10 epochs with early stopping (patience 3), learning rate  $2 \times 10^{-5}$ , weight decay 0.01, and batch size 4 with gradient accumulation 2 (effective batch size 8). The performance reported in Table 1 corresponds to the average of the best-performing checkpoint from each of the three runs.

**Decoding to spans.** At inference time, we decode token-level BIO predictions into entity spans using tokenizer offset mappings, producing sentence-local character spans in the coarse label set.

## 5 Experimental Setup

### 5.1 Preprocessing

**Common format.** Gold annotations are stored as JSONL with sentence-level texts, entity offsets, and entity labels (PER/LOC/ORG/OBJ).

**Label mapping for spaCy.** We map spaCy default labels to coarse labels via:

PERSON→PER;  
{FAC,GPE,LOC}→LOC;

{ORG,NORP}→ORG; {PROD-  
UCT,WORK\_OF\_ART,EVENT}→OBJ.

Unmapped spaCy labels (e.g., DATE/TIME/QUANTITY) are ignored.

**Data splits and standardized format.** We split *Ne Zha* P1 subtitles for train/validation (80/20) and P2 as a held-out test set. To create baseline training splits, we convert JSONL annotations into the spaCy format to standardize model training and evaluation. We use a leakage-aware strategy: sentences are grouped into connected components based on co-occurring entity strings, and entire components are assigned to train vs dev. This reduces leakage from repeated names/aliases across splits.

**BERT data conversion.** We further convert the spaCy format into BERT BIO tags using tokenizer offset mappings. Special tokens [CLS] and [SEP] are assigned loss label  $-100$ . In training/evaluation we truncate to a maximum sequence length of 64 tokens.

### 5.2 Evaluation

We report exact span-match micro-averaged precision/recall/F1 over entity tuples (start, end, label). A predicted entity is counted as correct if and only if its span boundaries and label exactly match a gold entity in the same sentence. We micro-average by summing true positives, false positives, and false negatives across all sentences.

All methods are evaluated in a unified representation as sentence-local character spans; if a method produces document-level offsets internally, they are mapped back to sentence-local offsets before scoring.

## 6 Results

### 6.1 Main results

Table 1 compares a general-purpose Chinese spaCy model (zero-shot), a supervised BERT token-classification baseline, and our multi-agent extraction pipeline.

The zero-shot spaCy baseline substantially underperforms (F1=0.185), reflecting the domain shift in mythological narratives. Fine-tuning a strong encoder improves performance (BERT-base-chinese: F1=0.655), while the multi-agent pipeline achieves the best overall

Model	P	R	F1	Setting	Variant	P	R	F1	
Zero-shot spaCy	0.271	0.141	0.185	<b>Prompting</b>	Zero-shot (con- tinuous)	0.459	<b>0.822</b>	0.589	
Fine-tuned BERT	0.613	<b>0.704</b>	0.655			Few-shot + guide (continuous)	<b>0.770</b>	0.691	<b>0.728</b>
Multi-agent LLM	<b>0.770</b>	0.691	<b>0.728</b>						
				<b>Confidence Scoring</b>		Few-shot + guide (continuous)	<b>0.770</b>	0.691	<b>0.728</b>
						Few-shot + guide (discrete)	0.758	<b>0.697</b>	0.726
						<b>Example Richness</b>			
					Few-shot + guide (discrete)	<b>0.758</b>	0.697	<b>0.726</b>	
					Fully annotated (discrete)	0.562	<b>0.757</b>	0.645	

Table 1: Main NER results on *NE ZHA*. We report exact span match micro-P/R/F1 on the held-out *NE ZHA* P2 test set.

Chunk	Context	P	R	F1
200	250	0.779	0.667	0.718
250	200	0.770	<b>0.691</b>	<b>0.728</b>
1000	500	0.768	0.496	0.603
2000	1000	0.773	0.407	0.533
4000	1000	<b>0.788</b>	0.310	0.445

Table 2: Chunk/context tuning for the multi-agent pipeline (few-shot+guide with extracted-entity examples). We report exact span match micro-P/R/F1 on the held-out *NE ZHA* P2 test set.

Table 3: Ablation slices for the multi-agent pipeline at the best chunk/context setting. We separate prompting style, confidence scoring, and example richness for readability.

result (F1=0.728), improving over BERT by +0.073 F1 and over the zero-shot baseline by +0.543 F1.

## 6.2 Configuration study

We next analyze which design choices drive performance (Table 2; Table 3).

**Chunk/context tuning.** With the prompting strategy held fixed (few-shot + guide using extracted-entity examples), moderate chunk and context sizes yield the strongest results; the best configuration is a chunk size of 250 with a context window of 200.

**Prompting, confidence, and example richness.** Using only label definitions (agent zero-shot) yields high recall but much lower F1 than few-shot+guide prompting. Confidence scoring (continuous vs discrete) has a small effect at the best chunk/context. For example richness, we observe that fully annotated examples underperform simple few-shot examples in our current prompting setup; we treat this as a prompt-formatting effect rather than a definitive statement about annotation richness.

## 7 Analysis and Discussion

### 7.1 The challenge of myth-domain NER

Myth-domain NER is challenging because our label space is defined around *referable in-*

*stances* rather than surface words. Following our annotation criteria, a span should be annotated only if it can answer the core questions “*who / what / where*”. Consequently, many expressions that look “important” in myth narratives are deliberately excluded, including events or phenomena such as “*天劫*” (*heavenly tribulation*) and “*封神大战*” (*Fengshen War*).

The most ambiguous label is OBJ. In our setting, OBJ is reserved for *named* artifacts or fixed technique titles (e.g., a spell name invoked as a unit), while generic concepts such as “*命运/因果/力量*” (*fate/karma/power*) are not annotated unless they behave as a proper name in context. This “named-object” constraint is crucial for dataset consistency, but it also makes the boundary between “object” and “non-entity” particularly sharp and error-prone.

Subtitles further amplify ambiguity through pervasive titles and references. Generic titles like “*师尊/师弟*” (*master/junior disciple*) are excluded, while unique, referential variants such as “*李大人*” (*Lord Li*), “*吒儿*” (*Zha’er, a nickname*), and “*三公子*” (*third young master*) may be annotated as PER when they uniquely identify a character.

Finally, we adopt longer span preference : when a full mention is available, we annotate the complete name (e.g., “*乾元山金光洞太乙真人*” / *Taiyi Zhenren of Qianyuan Mountain*

*Jinguang Cave*). While this reduces fragmentation in the gold labels, it increases the probability of exact-span mismatches for models.

This design choice interacts directly with our evaluation protocol. Because we use exact span match, boundary near-misses are penalized maximally: a prediction that captures the right referent but trims or extends a title-heavy mention is still counted as an error. In practice, this makes span normalization a central challenge in MythNER, and it partly explains why boundary-related errors dominate across both conventional taggers and LLM-based systems.

## 7.2 Traditional baseline errors

We observe several recurring failure modes in traditional baselines. First, zero-shot models frequently over-predict “entity-like” common words, producing false positives from adjectives, verbs, and generic nouns that do not satisfy the *referable-instance* criterion. Second, both zero-shot and supervised baselines are sensitive to long, title-heavy mentions, often fragmenting a gold span into smaller pieces. Under exact-span evaluation, these near-misses are counted as errors even if a partial substring is correct.

Type confusion is also prominent in mythology-specific contexts where the same surface form can plausibly be a person-like being, a collective, or a named object. For example, role-title composites and faction-like phrases can trigger drift between PER/ORG/LOC. Finally, OBJ remains a residual bottleneck: models may miss named objects/techniques entirely (recall errors) or incorrectly promote abstract substances or generic items to OBJ (precision errors), reflecting the intrinsic ambiguity of “named artifact vs. generic concept” in myth discourse.

## 7.3 LLM-agent errors: what changes and what remains

The LLM-agent pipeline changes the error profile rather than eliminating errors altogether. In an agent *zero-shot* configuration (high-recall, lightly constrained), the dominant issue is over-extraction: the agent tends to label abstract substances and generic nouns as OBJ, e.g., “天地灵气” (*heaven-and-earth aura*), “仙气” (*immortal aura*), and “魔气”

(*demonic aura*), which are not annotated in the gold standard. It also exhibits precision collapse for PER by eagerly labeling context-dependent addresses (e.g., “猪兄” / *Brother Pig*, “仙长” / *immortal elder*) that our dataset often treats as non-entities.

Meanwhile, exact-span sensitivity remains: long mentions are frequently split into overlapping subspans (a typical false-positive/false-negative pair under span-level matching). Even in the extracted, tuned chunk/context settings, we still observe boundary drift where the agent extends a named object with nearby descriptive material (e.g., predicting “宝莲仙气” / *lotus aura* instead of the gold “宝莲” / *lotus*).

Compared to agent zero-shot, tuned extraction configurations drastically reduce false positives but increase false negatives, illustrating a clear precision–recall tradeoff. In practice, many of the remaining misses concentrate on short, frequently mentioned myth objects/techniques such as “魔丸” (*Demon Pill*) and “天雷” (*heavenly thunder*), which require stable span decisions across fragmented subtitle context.

The configuration study in our results further supports a simple principle: context helps until it hurts. Too little context deprives the model of cross-mention cues for disambiguation, while too much context (or overly large chunks) increases boundary drift and encourages “entity-like” generic nouns, especially for OBJ. This reinforces the role of constraint design and iterative verification in agentic extraction.

## 7.4 Practical takeaway

Across baselines, the dominant challenges are (i) **boundary alignment** under exact-span scoring, (ii) **type ambiguity** in mythology-specific language, and (iii) the particularly sharp **scope constraint** of OBJ (named artifact/technique vs. generic concept). The LLM-agent approach reduces some forms of type drift via stronger contextual reasoning, but introduces *over-extraction* when constraints are weak and remains highly sensitive to span normalization, motivating careful confidence scoring and post-processing for deployment.

Overall, MythNER is less a test of mem-

555 orizing surface patterns than a test of de-  
556 ciding entityhood and stabilizing boundaries  
557 under culturally grounded naming conven-  
558 tions. The dataset therefore complements  
559 existing benchmarks by emphasizing (i) con-  
560 servative referable-instance criteria, (ii) long-  
561 span boundary decisions, and (iii) ontology-  
562 like constraints implicit in OBJ.

## 563 8 Limitations

564 MythNER targets a specific long-tail setting:  
565 Chinese mythological narratives in subtitle  
566 form. As a result, models trained and eval-  
567 uated on MythNER may not transfer directly  
568 to other genres (e.g., classical prose, novels)  
569 or other media with different discourse conven-  
570 tions.

571 Our annotation design also imposes delib-  
572 erate scope constraints. We adopt a flat,  
573 non-overlapping schema with four coarse types  
574 (PER/LOC/ORG/OBJ), and we use conser-  
575 vative entityhood criteria that exclude many  
576 events and phenomena. In particular, OBJ is  
577 restricted to *named* artifacts or fixed technique  
578 titles; this improves consistency but creates a  
579 sharp boundary between named objects and  
580 generic concepts.

581 Evaluation uses exact span match, which  
582 is appropriate when downstream use requires  
583 boundary-faithful extraction, but it penalizes  
584 near-miss boundary predictions maximally.  
585 This interacts with our longer-span prefer-  
586 ence and makes span normalization a primary  
587 source of error.

588 Finally, our LLM-agent pipeline introduces  
589 practical constraints: performance is sensitive  
590 to prompting and chunk/context choices, and  
591 deployment may incur higher cost and latency  
592 than conventional taggers.

## 593 References

594 David Bamman, Sejal Papat, and Sheng Shen.  
595 2019. [An annotated dataset of literary entities](#).  
596 In *Proceedings of the 2019 Conference of the*  
597 *North American Chapter of the Association for*  
598 *Computational Linguistics: Human Language*  
599 *Technologies, Volume 1 (Long and Short Papers)*,  
600 pages 2138–2144, Minneapolis, Minnesota. Asso-  
601 ciation for Computational Linguistics.

602 Cuong Xuan Chu, Simon Razniewski, and Ger-  
603 hard Weikum. 2020. [ENTYFI: A System for](#)  
604 [Fine-grained Entity Typing in Fictional Texts](#).

In *Proceedings of the 2020 Conference on Em-  
605 pirical Methods in Natural Language Processing:  
606 System Demonstrations*, pages 100–106, Online.  
607 Association for Computational Linguistics. 608

Tao Gui, Yicheng Zou, Qi Zhang, Minlong Peng,  
609 Jinlan Fu, Zhongyu Wei, and Xuanjing Huang. 610  
2019. [A Lexicon-Based Graph Neural Network  
611 for Chinese NER](#). In *Proceedings of the 2019  
612 Conference on Empirical Methods in Natural  
613 Language Processing and the 9th International  
614 Joint Conference on Natural Language Process-  
615 ing (EMNLP-IJCNLP)*, pages 1040–1050, Hong  
616 Kong, China. Association for Computational  
617 Linguistics. 618

Basra Jehangir, Saravanan Radhakrishnan, and  
619 Rahul Agarwal. 2023. [A survey on Named En-  
620 tity Recognition — datasets, tools, and method-  
621 ologies](#). *Natural Language Processing Journal*,  
622 3:100017. 623

Bin Li, Bolin Chang, Ruilin Liu, Xue Zhao,  
624 Si Shen, Lihong Liu, Yan Zhu, Zhixing  
625 Xu, Weiguang Qu, and Dongbo Wang. 2025.  
626 [Overview of EvaHan2025: The First Interna-  
627 tional Evaluation on Ancient Chinese Named  
628 Entity Recognition](#). In *Proceedings of the Sec-  
629 ond Workshop on Ancient Language Processing*,  
630 pages 156–164, The Albuquerque Convention  
631 Center, Laguna. Association for Computational  
632 Linguistics. 633

Jing Li, Aixin Sun, Jianglei Han, and Chenliang Li.  
634 2022. [A Survey on Deep Learning for Named En-  
635 tity Recognition](#). *IEEE Transactions on Knowl-  
636 edge and Data Engineering*, 34(1):50–70. 637

Kalyani Pakhale. 2023. [Comprehensive Overview  
638 of Named Entity Recognition: Models, Domain-  
639 Specific Applications and Challenges](#). *Preprint*,  
640 arXiv:2309.14084. 641

Siyao Peng, Zihang Sun, Sebastian Loftus, and  
642 Barbara Plank. 2024. [Different Tastes of En-  
643 tities: Investigating Human Label Variation in  
644 Named Entity Annotations](#). In *Proceedings  
645 of the Third Workshop on Understanding Im-  
646 plicit and Underspecified Language*, pages 73–81,  
647 Malta. Association for Computational Linguistics.  
648 649

Aravinth Sivaganeshan and Nisansa De Silva. 2023.  
650 [Fine Tuning Named Entity Extraction Models  
651 for the Fantasy Domain](#). In *2023 Moratuwa  
652 Engineering Research Conference (MERCon)*,  
653 pages 346–351. 654

Xinli Tao, Xin Dong, and Xuezhong Zhou.  
655 2025. [OEMA: Ontology-Enhanced Multi-Agent  
656 Collaboration Framework for Zero-Shot Clin-  
657 ical Named Entity Recognition](#). *Preprint*,  
658 arXiv:2511.15211. 659

660	Simone Tedeschi and Roberto Navigli. 2022. <a href="#">Multi-NERD: A Multilingual, Multi-Genre and Fine-Grained Dataset for Named Entity Recognition (and Disambiguation)</a> . In <i>Findings of the Association for Computational Linguistics: NAACL 2022</i> , pages 801–812, Seattle, United States. Association for Computational Linguistics.	<b>Labels.</b> PER: a specific character or person-like being (humans, immortals, demons, named creatures). LOC: a named place/location. ORG: a named faction/organization acting as a collective. OBJ: a <i>named</i> artifact/item, or a fixed technique title invoked as a unit.	713 714 715 716 717 718 719
667	Shuhe Wang, Xiaofei Sun, Xiaoya Li, Rongbin Ouyang, Fei Wu, Tianwei Zhang, Jiwei Li, Guoyin Wang, and Chen Guo. 2025a. <a href="#">GPT-NER: Named Entity Recognition via Large Language Models</a> . In <i>Findings of the Association for Computational Linguistics: NAACL 2025</i> , pages 4257–4275, Albuquerque, New Mexico. Association for Computational Linguistics.	<b>Span rules.</b> (1) <b>No overlaps / no nesting</b> (flat NER). (2) <b>Prefer the longest valid boundary</b> when a full canonical mention is available. (3) Exclude pure pronouns and generic role nouns.	720 721 722 723 724
675	Zihan Wang, Ziqi Zhao, Yougang Lyu, Zhumin Chen, Maarten de Rijke, and Zhaochun Ren. 2025b. <a href="#">A Cooperative Multi-Agent Framework for Zero-Shot Named Entity Recognition</a> . In <i>Proceedings of the ACM on Web Conference 2025, WWW '25</i> , pages 4183–4195, New York, NY, USA. Association for Computing Machinery.	<b>Titles, aliases, and references.</b> Generic, context-dependent addresses (e.g., 师尊/师兄/大人) are not annotated. However, if a variant uniquely identifies a character (contains a surname, name fragment, rank/epithet that is used as a unique handle), annotate it as PER.	725 726 727 728 729 730
683	Gayashan Weerasundara and Nisansa de Silva. 2023. Comparative Analysis of Named Entity Recognition in the Dungeons and Dragons Domain. In <i>Proceedings of the 14th International Conference on Recent Advances in Natural Language Processing</i> , pages 1225–1233, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.	<b>OBJ is conservative by design.</b> Annotate OBJ only when the mention behaves like a proper name or a stable technique title. Do not annotate generic substances/abstract concepts (e.g., 神力/力量/因果/命运) unless they are explicitly treated as a named unit in context.	731 732 733 734 735 736 737
690	Yusen Zhang, Ruoxi Sun, Yanfei Chen, Tomas Pfister, Rui Zhang, and Sercan Ö Arik. 2024. <a href="#">Chain of Agents: Large Language Models Collaborating on Long-Context Tasks</a> . <i>Advances in Neural Information Processing Systems</i> , 37:132208–132237.	<b>Exclude events and phenomena.</b> Named or salient events/rituals/astronomical phenomena are not entities (e.g., 天劫, 封神大战, 生辰宴).	738 739 740 741
696	Yufan Zhuang, Xiaodong Yu, Jialian Wu, Ximeng Sun, Ze Wang, Jiang Liu, Yusheng Su, Jingbo Shang, Zicheng Liu, and Emad Barsoum. 2025. <a href="#">Self-Taught Agentic Long Context Understanding</a> . <i>Preprint</i> , arXiv:2502.15920.	<b>LOC vs. ORG (surface-syntax heuristic).</b> Some names can denote both a place and an institution. We default to LOC for consistency.	742 743 744 745
701	<b>A Annotation Guide</b>	<b>Examples.</b>	746
702	We annotate <i>flat</i> , non-overlapping named entities in Chinese mythological subtitles with four coarse types: PER, LOC, ORG, OBJ. The annotation target is not “words”, but <i>referable instances in context</i> .	<b>PER (unique name).</b> 我乃 [太乙真人] PER。 [李大人] PER 来了。 但是 [师尊] <i>not annotated</i> .	747 748 749
703		<b>LOC (named place).</b> 我们去 [陈塘关] LOC。 他在 [昆仑山] LOC 修行。	750 751
704		<b>ORG (faction acting collectively).</b> [天庭] ORG 下令捉拿。 [妖族] ORG 来犯。	752 753
705		<b>OBJ (named artifact / technique title).</b> 哪吒祭出 [乾坤圈] OBJ。 他使出 [龙族秘术] OBJ。 但 [神力/力量] <i>not annotated</i> .	754 755 756 757
706		<b>Long-span preference.</b> [乾元山金光洞太乙真人] PER (not 乾元山 + 金光洞 + 太乙真人).	758 759 760
707	<b>Entityhood test (fast rule).</b> Annotate a span only if, in its local context, it can reasonably answer <i>who / what (thing) / where</i> . If an expression only answers “what event/state/phenomenon is happening”, do <i>not</i> annotate it.	<b>Exclude events/phenomena.</b> [天劫] <i>not annotated</i> 。 [封神大战] <i>not annotated</i> .	761 762
710			
711			
712			

## B LLM Prompt Templates

This appendix documents the prompt templates used by the multi-agent pipeline. We report (i) the prompts by agent role and (ii) the prompt variants used in our ablation settings (Table 3). Templates contain placeholders such as {chunk\_text} and {start\_char}.

### B.1 Agent prompts (by role)

**Chunking agent (no LLM call).** Chunking is programmatic: a document is split into chunks of length  $K$  and each chunk is provided with left/right context windows of length  $C$ . These parameters change the *inputs* to the extraction prompt (the chunk text and its context), but do not change the prompt text itself.

**Extraction agent: zero-shot, label definitions only. Original (ZH).**

你是一个命名实体识别 (NER) 专家。  
任务：从中文神话文本中提取命名实体，遵循以下 schema：

**【实体类型定义】** - PER: 人物，包括人类、神仙、妖魔、龙族等能说话、行动的角色。- LOC: 地点，具体的地理位置或场所名称。- ORG: 组织，有组织性的群体、种族、派系、领域。- OBJ: 物品，法器、武器、法术名、魔法物品、坐骑等名。

输出格式：必须返回有效的 JSON：  
{ "entities": [ { "text<sub>span</sub>" :  
" ", "start<sub>char</sub>" : 0, "end<sub>char</sub>" :  
2, "type" : "PER", "confidence" :  
1.0} ] }

重要：start\_char 和 end\_char 是相对于完整原文档的绝对位置！

请从以下 **【文本片段】** 中提取所有命名实体。

**【文本片段】** (第 {start\_char} 到 {end\_char} 字符): {chunk\_text}

**【上文】** (仅供参考): {context\_before}

**【下文】** (仅供参考): {context\_after}

提取要点：1. 只从 **【文本片段】** 提取，位置必须是文档绝对位置。2. 只要输出 entities 列表。

#### English rendering.

You are an expert in Named Entity Recognition (NER).

Task: Extract named entities from Chinese mythology text, using this schema: - PER: characters/person-like beings. - LOC: named locations. - ORG: named organizations/factions. - OBJ: named artifacts/weapons/spell or technique titles/mounts.

Output must be valid JSON:  
{ "entities": [ { "text<sub>span</sub>" :  
"ENTITY", "start<sub>char</sub>" :  
0, "end<sub>char</sub>" : 2, "type" :  
"PER", "confidence" : 1.0} ] }

IMPORTANT: start\_char/end\_char are absolute offsets in the full document.

Extract all named entities from the following [Text Fragment] only. [Text Fragment] (characters {start\_char} to {end\_char}): {chunk\_text}

[Previous Context] (reference only): {context\_before}

[Following Context] (reference only): {context\_after}

**Extraction agent: few-shot + guide; continuous confidence.** This variant adds explicit boundary/type rules and in-prompt examples, and requests a continuous confidence score  $\in [0, 1]$ .

**Original (ZH).**

你是一个中国神话文本的命名实体识别 (NER) 专家。

任务：从中文神话文本中提取命名实体，遵循 Stage 1 schema：

**【实体类型定义】**- PER (人物)：能说话、行动的角色，包括人类、神仙、妖魔、龙族 包括：太乙真人、哪吒、敖丙、李靖、申公豹、龙王、海夜叉、咤儿、太乙、天尊 包括带地点的完整称谓：“乾元山金光洞太乙真人”整体是 PER 不包括：通用称呼 (仙长、师父)、代词、“X 夫妇”中只标注人名

- LOC (地点)：具体地理位置或场所名称 包括：陈塘关、龙宫、昆仑山、东海、乾元山、金光洞、海底炼狱 不包括：天庭、仙界 (组织/领域)、方位词 (南面、北面)

846	- ORG (组织): 有组织性的群体、种				
847	族、派系、领域 包括: 龙族、天庭、				
848	十二金仙、妖族、阐教、仙界、伏魔				
849	帮、人族				
850	- OBJ (物品): 法器、武器、法术名、				
851	魔法物品、坐骑 包括: 混元珠、灵				
852	珠、魔丸、乾坤圈、火尖枪、混天绫、				
853	天劫咒、结界、虚空之门、风火轮 不				
854	包括: 咒语念词 (如"急急如律令")				
855	<b>【关键提取规则】</b> 1. 边界规则: 提取				
856	核心实体名, 不要过度扩展 2. 只提				
857	取文本片段中的实体 3. 同一实体的				
858	不同形式都要提取 4. 不标注: 通用				
859	称呼、方位词、代词、咒语念词				
	输出格式: 必须返回有效的 JSON:				
	{ "entities": [ { "text <sub>s</sub> pan" :				
	" ", "start <sub>c</sub> har" : 45, "end <sub>c</sub> har" :				
	47, "type" : "PER", "confidence" :				
	0.95} ] }				
860	重要: start_char 和 end_char 是文档绝对位				
861	置!				
862	—				
863	<b>【文本片段】</b> (第 {start_char} 到 {end_char}				
864	字符): {chunk_text}				
865	<b>【上文】</b> (参考): {context_before}				
866	<b>【下文】</b> (参考): {context_after}				
867	提取要点: 1. 只从文本片段提取 2. 区分:				
868	天庭/仙界 =ORG; 陈塘关/龙宫 =LOC; 结				
869	界/虚空之门 =OBJ 3. 提取核心实体名				
870	<b>English rendering.</b>				
871	You are an expert NER annotator for				
872	Chinese mythology text.				
873	Task: Extract named enti-				
874	ties under a Stage-1 schema				
875	(PER/LOC/ORG/OBJ).				
876	Boundary rule: extract the core en-				
877	tity name; do not over-extend spans.				
	Output must be valid JSON with				
	absolute offsets: { "entities": [				
	{ "text <sub>s</sub> pan" : "...", "start <sub>c</sub> har" :				
	0, "end <sub>c</sub> har" : 2, "type" :				
	"PER", "confidence" : 0.95} ] }				
878	Extract entities from the Text Fragment				
879	(characters {start_char} to {end_char}):				
880	{chunk_text}				
	<b>Extraction agent: discrete confidence</b>				
	<b>variant (0–5 ranks). Original (ZH).</b>				
	置信度等级 (0–5): - 5: 绝对确定 -				
	4: 高置信度 - 3: 中等置信度 - 2: 低				
	置信度 - 1: 极低置信度 - 0: 不确定				
	输出格式: { "entities": [ {				
	"text <sub>s</sub> pan" : "...", "start <sub>c</sub> har" :				
	45, "end <sub>c</sub> har" : 47, "type" :				
	"PER", "confidence <sub>r</sub> ank" :				
	5, "confidence <sub>r</sub> eason" : " " ] }				
	追加要求: 每个实体提供 confidence_rank 与				
	confidence_reason.				
	<b>English rendering.</b>				
	Discrete confidence (0–5). Output				
	JSON must include: { "entities": [ {				
	"text <sub>s</sub> pan" : "...", "start <sub>c</sub> har" : 0, "end <sub>c</sub> har" :				
	2, "type" : "PER", "confidence <sub>r</sub> ank" :				
	5, "confidence <sub>r</sub> eason" :				
	"shortjustification" ] }				
	<b>Extraction agent: fully annotated exam-</b>				
	<b>ples. Original (ZH).</b>				
	输出格式: JSON 数组, 每句一个对				
	象: [ { "id": 1, "text": "这就是我万				
	人敬仰的太乙真人", "label": [[9, 13,				
	"PER"]], "Comments": [] } ]				
	重要: label 偏移量相对于句子 text!				
	<b>English rendering.</b>				
	Output format: JSON array, one ob-				
	ject per sentence. Offsets are relative				
	to the sentence text.				
	<b>Verification agent. Original (ZH).</b>				
	你是一个 NER 标注质量检查专家。				
	任务: 验证实体提取是否符合 Stage				
	1 规则: 1. 非重叠 2. 最大跨度 3. 类				
	型正确 4. OBJ 必须是专有名称 5.				
	边界准确				
	输出格式: { "summary": { "to-				
	tal <sub>e</sub> ntities" : 0, "approved <sub>c</sub> ount" :				
	0, "flagged <sub>c</sub> ount" : 0, "issues" :				
	[] }, "detailed <sub>f</sub> indings" :				
	[ { "entity <sub>i</sub> d" : "...", "issue <sub>t</sub> ype" :				
	"boundary <sub>e</sub> rror", "description" :				
	"...", "suggested <sub>f</sub> ix" : "..." } ] }				
	原文 (前 5000 字符): {raw_text} 提取的实				
	体: {entities_json}				

914 **English rendering.**

915 You are a quality-control agent for  
916 NER. Return JSON with summary  
917 and detailed findings.

918 **Correction agent. Original (ZH).**

919 你是一个 NER 标注修正专家。  
920 任务：根据验证报告修正实体。  
921 输出格式：{ "corrections": [...], "re-  
922 buttals": [...] }  
923 被 标 记 的 实 体：  
924 {flagged\_entities\_json} 验 证 报  
925 告：{verification\_report\_json} 扩 展  
926 上 下 文：{expanded\_contexts} 原 文  
927 (前 10000 字符)：{raw\_text}

928 **English rendering.**

929 You are a correction agent. Apply  
930 fixes or rebuttals and return JSON.

931 **B.2 Prompt variants used in ablations**

932 **Prompt type.** *Zero-shot* uses concise label  
933 definitions only. *Few-shot + guide* adds ex-  
934 plicit rules and examples.

935 **Confidence scoring.** *Continuous*: real-  
936 valued confidence in  $[0, 1]$ . *Discrete*: integer  
937 confidence\_rank plus confidence\_reason.

938 **Example richness.** *Extracted-entity exam-*  
939 *ples*: entity-level prompts. *Fully annotated ex-*  
940 *amples*: sentence-level labeled format.

941 **Chunk/context.** Chunk/context settings  
942  $(K, C)$  determine the injected chunk\_text and  
943 context\_before/after.