Combining Distantly Supervised Models with In Context Learning for Monolingual and Cross-Lingual Relation Extraction

Vipul Rathore Malik Hammad Faisal Parag Singla Mausam

Indian Institute of Technology New Delhi, India {vipul.rathore, cs1210559, parags, mausam}@cse.iitd.ac.in

Abstract

Distantly Supervised Relation Extraction (DSRE) remains a long-standing challenge in NLP, where models must learn from noisy bag-level annotations while making sentence-level predictions. While existing state-of-the-art (SoTA) DSRE models rely on task-specific training, their integration with in-context learning (ICL) using large language models (LLMs) remains underexplored. A key challenge is that the LLM may not learn relation semantics correctly, due to noisy annotation.

In response, we propose HYDRE – $\underline{\mathbf{H}}\mathbf{Y}$ brid <u>Distantly</u> Supervised <u>Relation</u> Extraction framework. It first uses a trained DSRE model to identify the top-k candidate relations for a given test sentence, then uses a novel dynamic exemplar retrieval strategy that extracts reliable, sentence-level exemplars from training data, which are then provided in LLM prompt for outputting the final relation(s). We further extend HYDRE to cross-lingual settings for RE in low-resource languages. Using available English DSRE training data, we evaluate all methods on English as well as a newly curated benchmark covering four diverse low-resource Indic languages - Oriya, Santali, Manipuri, and Tulu. HYDRE achieves up to 20 F1 point gains in English and, on average, 17 F1 points on Indic languages over prior SoTA DSRE models. Detailed ablations exhibit HYDRE's efficacy compared to other prompting strategies.

1 Introduction

Relation Extraction (RE) is a core task in Information Extraction (IE) that aims to identify semantic relations between entity mentions in text. Given a sentence s containing a head entity e_1 and a tail entity e_2 , the goal is to predict the set of relations ($\hat{r} \subset \mathcal{R}$) expressed between them, where \mathcal{R} is a predefined ontology of relations. RE plays a crucial role in downstream applications such as knowledge base construction, and question answering.

Supervised RE methods depend on sentence-level annotations, which are costly and time-consuming to obtain at scale (Zhang et al., 2017).

Distantly Supervised Relation Extraction (DSRE) (Mintz et al., 2009) alleviates this challenge by aligning text with knowledge base (KB) triples to generate weakly labeled training data. Specifically, DSRE groups all sentences mentioning an entity pair (e_1, e_2) into a bag $B(e_1, e_2)$, which is labeled with all relations $R(e_1, e_2)$ known between (e_1, e_2) in the KB. Although the supervision is weak and bag-level, inference is typically performed at the sentence level (micro-reading (Mitchell et al., 2009)). This mismatch between training and inference granularity – coupled with noisy training labels – often causes state-of-the-art (SoTA) DSRE models to confuse fine-grained relation types, such as Nationality vs. Place_of_Birth or Founder vs. CEO, thereby limiting their overall performance.

Existing DSRE methods (Chen et al., 2021; Rathore et al., 2022; Jian et al., 2024) primarily rely on task-specific fine-tuning of moderate-sized language models such as BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019). However, the potential of Large Language Models (LLMs) for this task remains largely unexplored. Modern LLMs excel at in-context learning (ICL), where the model performs reasoning by conditioning on a few taskspecific exemplars. Yet, directly applying LLMs to DSRE is non-trivial: noisy distant supervision degrades exemplar quality, and the lack of clean, sentence-level exemplars undermines effective ICL. Consequently, prior works in DSRE either ignore LLMs altogether or fail to exploit their reasoning capabilities effectively (Zhao et al., 2024).

In this work, we propose HYDRE: a <u>HY</u>brid <u>DSRE</u> framework that combines the high-recall candidate label selection of fine-tuned DSRE models with the reasoning capabilities of LLMs. Given a query sentence, a fine-tuned DSRE model first

identifies a candidate relation set by filtering out the obvious negatives. These candidates are then passed to an LLM for disambiguation, guided by carefully selected in-context exemplars. To construct these exemplars, we retrieve relevant bags from the DSRE training corpus using a joint scoring function that blends model confidence and semantic similarity to the query. From each selected bag, we extract the most representative sentence to form a dynamic, relation-specific prompt, guiding the LLM to accurately *judge* the correct relation(s).

We further extend HYDRE to the cross-lingual setting, focusing on low-resource languages — a largely underexplored area in DSRE. To facilitate this, we construct a new multilingual benchmark covering four low-resource Indic languages: Oriya, Santhali, Manipuri, and Tulu. Given the limited representation of these languages in LLM pretraining corpora (Singh et al., 2024; Nag et al., 2025), they pose an interesting challenge for evaluating cross-lingual RE in the context of latest LLMs.

We evaluate HYDRE under three transfer settings - (1) *English-only-data*, where no target language data is used; (2) *Translate-train*, where English DSRE training data is translated to target language; and (3) *Translate-test*, where test queries in the target language are translated to English.

Experiments with both open-source and proprietary LLMs show that HYDRE consistently outperforms prior DSRE baselines and naive LLM prompting strategies. Our exemplar retrieval strategy proves robust across both monolingual and cross-lingual setups. Ablation analyses further reveal that removing retrieval components can degrade performance by up to 7 micro-F1 points.

In summary, our key contributions are as follows. (1) We present the first systematic integration of LLMs via In-Context Learning (ICL) into DSRE inference, achieving significant gains over both fine-tuned and prompting-only baselines. (2) We propose a novel retrieval strategy that combines model confidence and semantic similarity to select high-quality ICL exemplars for LLM-based relation disambiguation. (3) We curate and release gold-standard evaluation datasets for relation extraction for four typologically diverse lowresource Indic languages. (4) We propose effective cross-lingual strategies for adapting HYDRE to low-resource languages, demonstrating robustness across multiple transfer scenarios. We commit to releasing our code and data to facilitate further research in multilingual DSRE.

2 Related Work

Distantly Supervised Relation Extraction: DSRE (Mintz et al., 2009) aligns KB triples with text to create bag-level supervision, where labels apply at the bag rather than sentence level. Neural DSRE models typically adopt the multi-instance multi-label (MIML) framework (Surdeanu et al., 2012). Earlier works encoded sentences in a bag using piecewise CNNs (Zeng et al., 2015) or graph CNNs (Vashishth et al., 2018), while recent models employ pre-trained transformers. PARE (Rathore et al., 2022) encodes a bag by treating all bag sentences as a passage, whereas CIL (Chen et al., 2021) uses intra-bag attention and contrastive learning. HiCLRE (Li et al., 2022) introduces hierarchical contrastive learning, and HFMRE (Li et al., 2023b) employs Huffman-tree structures to denoise bags. These advances improve bag-level reasoning but still struggle with fine-grained sentence-level inference.

Multilingual DSRE: Research in multilingual RE has progressed (Ni et al., 2020; Nag et al., 2021), but multilingual DSRE is scarce. DisRex (Bhartiya et al., 2022) introduced a dataset across four European languages, though without manual sentence-level evaluation. No prior work targets typologically diverse low-resource Indic languages, which motivates our benchmark contribution.

LLMs for Relation Extraction: Large language models (LLMs) have recently been explored for RE primarily in supervised settings. Wadhwa et al. (2023) demonstrated that few-shot prompting with GPT-3 can rival fully supervised baselines, particularly when enhanced with chain-of-thought (CoT) reasoning. Other works employ zero-shot prompting with relation label definitions (Zhou et al., 2024) or leverage LLMs to denoise distantly supervised training data (Li et al., 2023a; Jian et al., 2024). However, these approaches either assume access to clean, labeled supervision or utilize LLMs solely for data relabeling purposes. In contrast, HYDRE takes a complementary view by employing LLMs directly at test time, leveraging them as reasoning engines rather than annotation tools.

Diversity-aware exemplar selection: A common in-context exemplar selection strategy is to retrieve the top-K semantically most similar exemplars to the query, but this often results in redundancy amongst the selected exemplars. To address this, Kapuriya et al. (2025) propose maximum and the selection of th

mal marginal relevance (MMR), which balances similarity with diversity, while Wang et al. (2024) employ determinantal point processes (DPPs) to encourage diversity through submodular modeling objectives. These methods highlight the importance of complementing similarity with diversity in exemplar selection, a principle that we extend within our hybrid DSRE–LLM framework.

LLM-as-Judge and Hybrid Models: Recent work explores LLMs as evaluators ("judges") for ranking candidate outputs (Zheng et al., 2023; Bavaresco et al., 2024). Our approach adapts can be seen as an instance, where a DSRE model provides high-recall candidate relations and their exemplars, and the LLM judges among them. Related hybrid paradigms exist in other NLP tasks (Rathore et al., 2024) such as sequence labeling, but not in DSRE.

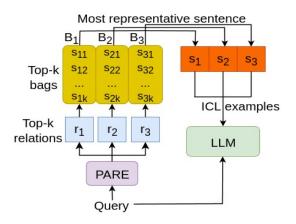


Figure 1: HYDRE overview (shown for k=3): (a) Select top-k candidate labels using PARE confidence, (b) Select best bag for each candidate using combined semantic similarity and label confidence (aggregated over bag sentences), and (c) Select best sentence per bag using aggregate (over bag labels) confidence scoring.

3 HYDRE: <u>HYbrid Distantly supervised</u> Relation Extraction

DSRE models are trained on bags $B(e_1, e_2)$, each containing multiple sentences for an entity pair (e_1, e_2) and annotated with a set of relations $R(e_1, e_2)$. Consequently, such models often generalize poorly to sentence-level queries (Gao et al., 2021; Chen et al., 2021). However, our preliminary analysis (Table 10) reveals that SoTA DSRE models such as PARE exhibit high Recall@k (\approx 85% for English) even for small values of k (E.g. k=5).

This observation motivates a two-stage hybrid approach: use a trained DSRE model to obtain a high-recall candidate relation set, and employ an LLM as a judge to select the final relation(s). We seek to achieve this by providing carefully selected exemplars for these candidate relations to LLM. Effective use of in-context learning (ICL), however, requires clean, sentence-level exemplars—whereas DSRE training examples are noisy and bag-level. To bridge this gap, HYDRE employs a novel three-stage exemplar selection process described below.

Stage 1: Candidate Relation Selection. A trained DSRE model (e.g., PARE) assigns confidence scores $f_{PARE}(q,r)$ to each relation r for a query sentence q. The top-k relations with the highest scores form the candidate relation set \mathcal{R}' .

Stage 2: Bag Selection. For each candidate relation $r \in \mathcal{R}'$, we identify the most relevant bag B_r from the subset \mathcal{B}_r of bags annotated with relation r as one of it's relations. Each bag $B_j \in \mathcal{B}_r$ is scored using a weighted combination of (i) semantic similarity $\operatorname{sim}(q,B_j)$ and (ii) DSRE model confidence $f(B_j,r)$, computed using max-pooling over its sentence scores. The highest-scoring bag is chosen as B_r for relation r. Essentially, this chosen bag has sentences that are strongly representative of r, while being similar to the test sentence q.

Stage 3: Sentence Selection. From each selected bag B_r , we extract a single sentence s_r that best captures the relation(s) expressed in the bag. For each sentence $s \in B_r$, we compute its coverage $c(s) = \sum_{r_a \in \text{labels}(B_r)} \mathbb{I}[f(s, r_a) > t]$, where t is a confidence threshold. Among sentences with maximum coverage, we select s_r as the one with the highest aggregate confidence $\sum_{r_a} f(s, r_a)$. This encourages selection of sentences that express the largest number of valid bag relations with strong model confidence. In the prompt we place the selected sentences s_r along with their corresponding bag labels $labels(B_r)$.

Our stage 3 offers three key benefits. First, it aids the *multi-label* nature of the problem, since co-occuring relations (which may be present in B_r , but not in \mathcal{R}') are also added to the prompt. This may not have been possible if the method tried to, instead, search for sentences that express only the specific relation. Second, selecting a sentence that has high label confidence over many relations promotes diversity (coverage) of relevant labels in the prompt – this may not occur if only semantic similarity with test sentence were used for selection. Finally, scoring sentences using aggregate confidence over all bag labels helps surface sentences

with stronger overall evidence, not just noisy singlelabel confidence (as justified via ablations shown in subsec. A.15).

Once selected, exemplars are ordered by ascending f(q,r), placing more relevant candidates closer to the query in the prompt. The full selection process is described in Algo. 1 (appendix).

Prompting and Parsing. Each prompt comprises (i) task instruction, (ii) candidate relation ontology along with their definitions, (iii) ICL exemplars, and (iv) the query. The exact prompt template is shown in subsec. A.3. The output is parsed using string matching to obtain predicted labels, including "NA". In case no valid match is found, we label it as "NA".

4 Dataset Curation

To evaluate HYDRE in low-resource settings, we construct gold-standard test sets for four Indic languages – Oriya, Santali, Manipuri, and Tulu. We choose these languages because they are an orthographically diverse set of languages that have limited presence on the Web (3K-25K Wikipedia articles), but have large number of native speakers (over 1 million each). This makes them a good challenge set for testing modern NLP systems.

We begin with a stratified subset of the English NYT-10m test split (Gao et al., 2021), having 538 multi-label queries with a total of 722 labels (incl. 30 "NA"s), ensuring balanced coverage across relation types. Each English sentence is translated into the target language using CODEC (constrained decoding (Le et al., 2024)) with IndicTrans2 (Gala et al., 2023) as the underlying model. CODEC performs joint translation and entity projection to preserve head–tail spans, ensuring that entity mentions remain correctly aligned in the target script. For Tulu, where IndicTrans2 lacks coverage, we use the Google Translate API.

To further enhance entity preservation, we construct synthetic parallel data using entity-spanalignment heuristics (following Chen et al., 2023) and fine-tune IndicTrans2 before translation. All translations are subsequently verified and corrected by native speakers fluent in reading and typing their respective scripts. Annotators confirmed the correctness of both the full sentence translation as well as the $\{\text{head}(e1), \text{tail}(e2)\}$ entity projections. A second annotator independently re-evaluated 100 randomly sampled sentences per language to measure reliability.

Across all languages, over 70% of translations required no correction, while for the corrected ones, their character-level F1 match with original outputs averaged 93%. Inter-annotator agreement exceeded 90% for both sentence translations and entity projections (Table 18), demonstrating strong consistency and translation quality. Detailed statistics on dataset sizes, label distributions, and annotation guidelines are provided in Appendix A.13.

5 Experiments

We do two sets of experiments: monolingual (English) and cross-lingual transfer (in the four low-resource Indic languages). For cross-lingual transfer experiments, we evaluate HYDRE across three standard settings commonly used in multilingual NLP. Let $X_{\rm train}$ denote the translated DSRE training corpus from English to target language X. We define PARE-X and CIL-X as the target-language counterparts of PARE and CIL, fine-tuned respectively on $X_{\rm train}$. The settings are:

- 1. **English-only:** Models are trained and prompted exclusively on English DSRE data. For English test queries, HYDRE uses the PARE model's confidence scores $f_{\text{PARE}}(q,r)$ for candidate ranking and computes semantic similarity using the off-the-shelf encoder e5-large-v2 (Wang et al., 2022). For Indic test queries, since PARE or CIL are not trained on Indic scripts, confidence scores are omitted; HYDRE instead relies solely on the multilingual encoder BGE-M3 (Chen et al., 2024) to compute cross-lingual semantic similarity.
- 2. **Translate-train:** HYDRE_X employs *PARE-X* for confidence estimation and *CIL-X* encoders—trained via contrastive learning to generate high-quality sentence embeddings in X for semantic similarity.
- 3. **Translate-test:** Test queries in the target language X are translated into English (X_{test}), after which the standard English HYDRE pipeline is applied directly.

Setting	Similarity Model	Confidence Model
English-only (En)	e5-large-v2	PARE
English-only (Indic)	BGE-M3	_
Translate-train	CIL-X	PARE-X
Translate-test	e5-large-v2	PARE

Table 1: Models used for semantic similarity and confidence estimation across various evaluation settings.

Table 1 summarizes the model configurations for each setting. In Table 2, we denote a variant as $\mathsf{HYDRE}(M)$ when using LLM M as the judge, and as $\mathsf{HYDRE}_X(M)$ when exemplars are retrieved from X_{train} under the Translate-train setting.

Baselines. We compare HYDRE against the following categories of baselines:

- **Supervised DSRE models:** PARE and CIL for English, and their target-language counterparts PARE-X and CIL-X trained under the *Translate-train* setting.
- 0-shot prompting LLM: We evaluate both open and proprietary large language models: *Qwen3-235B-A22B*, *GPT-4o*, *Llama3.1-8B*, and its fine-tuned variant *Llama3.1-8B-FT*. Here, *Llama3.1-8B-FT* denotes a model fine-tuned using English DSRE data for English experiments, and on *X*_{train} for the *Translate-train* setting. We consider two prompting variants. **Direct:** LLM is provided only with the query sentence and the list of candidate relation labels. **Ontology-based:** LLM additionally receives the definitions of candidate relations to aid disambiguation.
- Few-shot prompting LLM: We further compare against few-shot prompting strategies that differ in exemplar retrieval. Random-K randomly selects K exemplars for each query. TopK-sim retrieves the top-K semantically most similar exemplars to the query sentence (similarity-based selection). LM-MRR employs a diversity-aware selection strategy using a Maximal Marginal Relevance (MMR) objective (Kapuriya et al., 2025) to balance semantic relevance and diversity. Details are provided in Appendix A.2.3.

English Datasets. We use NYT-10m and Wiki-20m datasets (Gao et al., 2021) for English evaluations. Due to space limitations, results and discussion on Wiki-20m are deferred to Appendix A.16.

Implementation details. For obtaining X_{train} and X_{test} in our experiments, we use EasyProject (Chen et al., 2023), a more lightweight joint translation and entity projection tool as compared to CODEC. We use LoRA fine-tuning for LLaMA-3.1, updating only adapter and embedding weights while freezing other parameters (App. A.2.4). For PARE-X training, we continually pretrain mBERT on X_{train} before adapter-based fine-tuning. All

LLMs are run with temperature 0.0, max input tokens 2048, and max generation tokens 256.

Evaluation Metrics. We report micro-F1 and macro-F1. Area-under-curve (AUC) is omitted as LLMs cannot yield calibrated confidence scores. Statistical significance is tested using McNemar's test (McNemar, 1947) for micro-F1.

6 Results & Analysis

Table 2 reports results on English and four Indic languages across the three cross-lingual settings.

English results. Among supervised DSRE models, *CIL* (43 micro-F1) and *PARE* (42 micro-F1) are strongest. Zero-shot prompting with GPT-40 already surpasses them (56 F1), and HYDRE's few-shot prompting further lifts GPT-40 to 63 F1. Smaller open LLMs (Llama 3.1, Qwen 3) gain even more—up to +14 F1—demonstrating that HYDRE's exemplar-guided prompting yields larger benefits for weaker LLMs. These gains highlight the value of HYDRE's label-aware exemplar selection, which leverages DSRE model confidences to guide few-shot reasoning—unlike conventional similarity- or diversity-based retrieval.

Cross-lingual transfer (English-only setting). Since English-trained DSRE models (*PARE*, *CIL*) lack coverage for Indic scripts, their results are not reported in the table. Zero-shot prompting yields low scores (typically 20–30 micro-F1), reflecting the limited representation of these languages in the latest LLMs. Introducing few-shot prompting with English exemplars via HYDRE improves performance by 8 to 11 average F1 points across LLMs. Notably, HYDRE(Llama-FT) reaches from 24/14 to 35/21 — showing that even English exemplars, when semantically aligned, can substantially aid reasoning in low-resource Indic languages.

Translate-train setting. Training DSRE models on translated data ($X_{\rm train}$) provides stronger baselines (PARE-X: 30/20 F1). Zero-shot prompting of open LLMs (Qwen3, Llama3.1) still underperforms these supervised models, but five-shot prompting with HYDRE yields significant gains (5 to 8 points) over PARE-X. For Llama-3-FT (fine-tuned Llama), the gains are even more pronounced - improving over PARE-X by +17 micro-F1 and over zero-shot Llama-FT by +9 points. Fine-tuned Llama3-FT achieves the best overall performance (47/29 F1), while GPT-40 follows closely

Model	English		Indic Languages	
· · · · · · · · · · · · · · · · · · ·	n n	English-only	Translate-train	Translate-test
Supervised				
<i>HFMRE</i> (Li et al., 2023b)	33/18	_	22/13	27/16
HiCLRE (Li et al., 2022)	31/18	_	20/13	25/14
CIL (Chen et al., 2021)	43/32	_	26/18	34/24
PARE (Rathore et al., 2022)	42/31	_	30/20	33/23
0-shot (direct)				
Qwen3-235B-A22B	47/39	25/21	25/21	40/34
Llama3.1-8b	24/21	16/12	16/12	22/20
Llama3.1-8B-FT	55/37	26/17	26/17	47/31
GPT-4o	56/55	31/29	31/29	51/49
0-shot (Ontology-based)				
Qwen3-235B-A22B	49/40	25/21	25/21	44/35
Llama3.1-8b	31/17	22/16	22/16	29/25
Llama3.1-8B-FT	60/44	24/14	38/23	49/34
GPT-4o	56/57	33/31	32/30	54/51
few-shot (Random)				
Random-K(Qwen3-235B-A22B)	55/55	21/19	29/25	43/35
Random-K(Llama3.1-8B)	30/24	19/11	11/7	27/21
Random-K(Llama3.1-8B-FT)	55/40	24/14	25/12	46/32
Random-K(GPT-4o)	56/52	31/30	38/31	48/42
few-shot (Similarity-based)				
TopK-sim(Qwen3-235B-A22B)	52/49	22/18	32/26	45/38
TopK-sim(Llama3.1-8B)	33/27	22/14	19/9	30/22
TopK-sim(Llama3.1-8B-FT)	55/40	27/15	29/14	46/30
TopK-sim(GPT-4o)	58/53	29/26	41/32	48/41
few-shot (Diversity-based)				
LM-MRR(Qwen3-235B-A22B)	50/47	24/20	33/28	44/36
LM-MRR(Llama3.1-8B)	34/26	21/12	17/8	30/22
LM-MRR(Llama3.1-8B-FT)	56/40	26/16	28/15	47/32
LM-MRR(GPT-4o)	56/52	28/25	41/33	50/43
few-shot (HYDRE)				
HYDRE(Qwen3-235B-A22B)	63*/62	29/26	38/31	54/46
HYDRE(Llama3.1-8B)	52/47	30/19	35/21	44/36
HYDRE(Llama3.1-8B-FT)	61/45	35/21	47 */29	51/37
HYDRE(GPT-4o)	63* /60	36*/33	38/ 34	56*/54

Table 2: Results for English and Indic languages (across 3 cross-lingual settings). In each entry, we report micro and macro F1 scores. The Indic results are averaged over 4 languages (language-wise results shown in App. Tables 6, 7 and 8). * McNemar's p-value $< 10^{-5}$ (valid for micro-F1 comparison).

(38/34 F1). We also observe interesting language-specific trends (Table 7): GPT-40 performs best on high-resource scripts like Oriya and Tulu, whereas Llama-FT excels on rarer scripts such as Santhali and Manipuri. A plausible explanation is that Tulu shares its script with Kannada, and GPT-40 has likely encountered both Oriya and Kannada during its pretraining or supervised fine-tuning (SFT) phase. In contrast, Llama fine-tuning remains essential for Santhali and Manipuri, where GPT-40 struggles even with few-shot prompting—likely due to their absence from its pretraining data.

Translate-test setting. When the test data are translated to English, all models exhibit a mild degradation compared to their English performance due to translation noise. Nevertheless, HYDRE

continues to offer strong improvements: Llama3.1 and Qwen3 gain +15 and +10 micro-F1 points over their zero-shot variants, while GPT-40 achieves 56/54 F1—the highest among all models. These results demonstrate that HYDRE remains robust across varying LLMs and transfer settings.

HYDRE **vs. other exemplar selection methods:** Across all settings, HYDRE consistently surpasses other few-shot baselines. Compared to diversity-based LM-MMR, HYDRE(GPT-40) achieves average gains of +7 micro-F1 on English and +5–6 micro-F1 across Indic languages. These stem from its hybrid scoring that integrates DSRE confidence with semantic similarity, unlike MMR or TopK-sim, which ignore label information and rely solely on input embedding distance or diversity heuristics.

Translate-train vs. Translate-test: Translate-test consistently outperforms translate-train across all LLMs. This is likely attributed to higher translation quality when sentences are translated from target language (X) to English, as opposed to translating English data into the target language.

6.1 Ablations

To understand the contributions of various components to our proposed three-stage algorithm, we perform systematic ablations as follows¹:

- w/o candidate relation selection (stage 1): Instead of selecting top-k candidate relations in stage 1, we consider all relations in the ontology as candidates.
- w/o bag selection (stage 2): We flatten bags into sentences and assign all bag labels to each of it's sentences. Retrieval is performed directly over sentences.
- w/o sentence selection (stage 3): Entire bags (as passages) are provided to the LLM without representative sentence selection.
- w/o semantic similarity: Retrieval ignores semantic similarity scores; only PARE's confidence scores are used.
- w/o PARE candidate scores: PARE scores are excluded. Bag selection relies solely on semantic similarity, and sentences within bags are chosen randomly.
- *w/o both* (Random): For each candidate label, a random bag containing that label is chosen, followed by a random sentence from it.
- w/o ICL examples: The LLM selects directly from PARE's top-5 candidate relations without being shown their exemplars.

We present results in tables 3 and 4 for the monolingual and cross-lingual transfer experiments, respectively. We make the following observations.

Semantic similarity is most crucial for English, while for low-resource languages *PARE* confidence is critical. We observe that for English, semantic similarity yields with up to 7 F1 point drop in Llama 3.1, while removing the *PARE* confidence only leads a marginal drop (2 to 3 points) for

Ablation	L3.1	Qw3	G4o
Hydre	52/47	63/62	63 /60
w/o sem. sim.	45/39	60/59	60/62
w/o f_{PARE}	49/45	61/60	61/59
w/o both (Random)	45/43	56/56	60/60
w/o candidate selection (stage 1)	39/31	63 /61	62/64
w/o bag selection (stage 2)	48/41	62/59	56/51
w/o sentence selection (stage 3)	46/37	60/59	59/58
w/o ICL	45/37	58/50	59/49

Table 3: English F1 (micro/macro) scores for different ablation variants of HYDRE. Model keys: L3.1 = Llama3.1-8B, Qw3 = Qwen3-235B, G4o = GPT-4o.

Ablation	L3.1	L3.1-FT	Qw3	G4o
HYDRE*	32/ 20	45/28	38/31	38/34
only sem. sim.	26/18	44/27	37/ 31	39 /34
Random	26/16	38/21	37/ 31	36/32
w/o bag selection	33/20	43/25	34/30	38/ 35
(Stage 2)				
w/o ICL	29/19	33/18	33/20	32/21

Table 4: Avg. F1 (micro/macro) over Indic languages for different ablation variants of HYDRE in *translate-train* setting. *HYDRE omits semantic retrieval in *translate-train* setting owing to poor performance of retrieval models for low-resource languages (more details in sec. A.14).

both LLMs. In contrast, *PARE-X* confidence turns to be most critical for cross-lingual setting, with semantic similarity even hurting the performance of HYDRE (more analysis in A.14).

Sentence selection from bags turns out to be crucial step for effective ICL's performance: When we don't subselect representative sentences from top-k bags, and give entire bags-as-exemplars

in prompt to LLM, this hurts the performance by upto 6 micro F1 and 10 macro F1 points as LLMs fail to handle longer context lengths in these cases. In some sense, HYDRE is able to exploit a key property of distant supervision – that all sentences in a bag are not needed for optimal performance.

W/o Canditate relation selection: When we gives exemplars for all 24 relations, it has comparable performance to HYDRE(which uses only top-5 relations) but with almost 5x cost and this method is also not scalable if relation count increases since it requires exemplars for all relations.

ICL exemplars for candidate relations are crucial for both experiments. We observe that just giving PARE top-5 candidate relations to the LLM is not sufficient as the performance with this approach is lesser than the ICL-based approach by up to 7 points in English and 12 F1 points in translate-

¹For low-resource languages, the stage 1 and stage 3 ablations are omitted due to excessive prompt length from high token fertility, often exceeding the LLM's maximum sequence length.

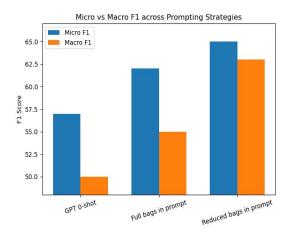


Figure 2: Zero-shot prompting vs full bag ICL vs reduced bag ICL for GPT-40 on English NYT-10m

train setting. This highlights the significance of carefully selected ICL exemplars in guiding LLM for optimal DSRE performance.

6.2 Robustness to bag-level evaluation

In addition to sentence-level evaluation, we seek to evaluate our approach on bag-level test queries as has been done in prior DSRE literature (Gao et al., 2021; Rathore et al., 2024). This finds applications in scenarios where we have a corpus of documents (e.g., news) and we need to populate a KB.

For bag-level queries, we can give bag-level exemplars to LLM by following Algorithm 1 till bag selection stage. Then, instead of finding a single representative sentence from each bag, we find a different representative sentence for each relation in the bag's labelset. Sentences with highest PARE score for each relation in the bag are kept and rest are discarded. As shown in Fig. 2, the reduced bag exemplars improve performance by 3 micro F1 and 7 macro F1 over naively giving the full bag exemplars in the prompt. Further this saves cost with around 2× reduction in prompt tokens.²

6.3 Hydre Sensitivity vs k

We analyze HYDRE's sensitivity to the number of candidate relations (k) selected in Stage 1 before LLM disambiguation. As shown in Figure 17b, micro-F1 initially improves with increasing k, peaking around k=5–10, after which performance plateaus or slightly declines. In contrast, the underlying DSRE model (PARE) continues to show monotonic gains in Recall@k with increasing k (Figure 17a), reaching near-perfect recall beyond

k=20. This discrepancy suggests that presenting too many candidates to LLM may degrade precision likely because LLM receives longer prompts and a complex disambiguation space. So, a moderate k (we set k=5 in all our experiments) offers the best trade-off between recall coverage, F1 performance and inference latency.

6.4 Qualitative Analysis

We conduct confusion analysis of HYDRE against PARE and 0-shot GPT-40 models for English (see appendix Figure 3). Further, we manually identify 16 fine-grained relation pairs amongst which models get confused (details in Fig. 14 and 15 appendix). We observe that HYDRE has stronger ability to resolve between subtle close-by relations such as "Religion" v/s "Ethnicity" (Ex. 4, "Geographic distribution" vs "Nationality" (Ex. 8), etc.

6.5 Error Analysis

We analyze cases where HYDRE fails to output correct relations (see figures 9 to 13 appendix). These include (a) errors due to *position bias* (a prevalent bias observed in LLM-as-Judge (Zheng et al., 2023)) towards PARE's top-1 candidate (nearest exemplar label to query) (fig. 8, 13), (b) low recall for multi-label queries (fig. 12), and (c) low recall for specific labels such as "Ethnicity" (often confused with "Nationality") (fig. 11), "Advisors" (vs "Founders"), etc (detailed analysis in App. A.7).

7 Conclusions and Future Work

We propose HYDRE, a novel hybrid framework that leverages SoTA distantly supervised models for guiding modern day LLMs for the task via efficient In-Context Learning. We show the efficacy of HYDRE across two experiments: (a) Bag-to-sentence transfer in English, (b) Bag-to-sentence transfer in low-resource languages. We propose effective strategies to adapt HYDRE to cross-lingual settings under different data availability settings. Ablations show the efficacy of each component in HYDRE in both monolingual and cross-lingual settings.

HYDRE lays the foundation for further research in DSRE in context of the latest LLMs. Possible future works involve a more complex agentic workflow wherein the agents interact iteratively until convergence to arrive at the correct answer. Moreover, applying HYDRE to newly added entries in Wikipedia or to local news articles for regional languages would be an interesting and useful future direction.

²Avg. tokens/prompt: full bag - 2451, reduced bag - 1242

8 Limitations

One potential limitation of our work is the high cost or latency of retrieval from large bags in the training corpus. This problem is escalated further for low-resource languages due to their high token fertility even w.r.t. to latest LLMs.

Further our framework is not tested for low-resource domains such as biomedical or finance domains involving more complicated semantic relationships between the evolving entities. Similarly, we have only performed experiments on four Indic languages, and have not been able to perform experiments on other low-resource language families due to unavailability of relation extraction data.

References

- Anna Bavaresco, Raffaella Bernardi, Leonardo Bertolazzi, Desmond Elliott, Raquel Fernández, Albert Gatt, E. Ghaleb, Mario Giulianelli, Michael Hanna, Alexander Koller, André F. T. Martins, Philipp Mondorf, Vera Neplenbroek, Sandro Pezzelle, Barbara Plank, David Schlangen, Alessandro Suglia, Aditya K Surikuchi, Ece Takmaz, and Alberto Testoni. 2024. Llms instead of human judges? a large scale empirical study across 20 nlp evaluation tasks.
- Abhyuday Bhartiya, Kartikeya Badola, and Mausam. 2022. DiS-ReX: A multilingual dataset for distantly supervised relation extraction. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, Dublin, Ireland. Association for Computational Linguistics.
- Jianly Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2024. Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation.
- Tao Chen, Haizhou Shi, Siliang Tang, Zhigang Chen, Fei Wu, and Yueting Zhuang. 2021. CIL: Contrastive instance learning framework for distantly supervised relation extraction. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6191–6200, Online. Association for Computational Linguistics.
- Yang Chen, Chao Jiang, Alan Ritter, and Wei Xu. 2023. Frustratingly easy label projection for cross-lingual transfer. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 5775–5796.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of* the North American Chapter of the Association for

- Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Jay Gala, Pranjal A Chitale, A K Raghavan, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar M, Janki Atul Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, Pratyush Kumar, Mitesh M Khapra, Raj Dabre, and Anoop Kunchukuttan. 2023. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages. Transactions on Machine Learning Research.
- Tianyu Gao, Xu Han, Yuzhuo Bai, Keyue Qiu, Zhiyu Xie, Yankai Lin, Zhiyuan Liu, Peng Li, Maosong Sun, and Jie Zhou. 2021. Manual evaluation matters: Reviewing test protocols of distantly supervised relation extraction. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1306–1318, Online. Association for Computational Linguistics.
- Zhaorui Jian, Shengquan Liu, Wei Gao, and Jianming Cheng. 2024. Distantly supervised relation extraction based on non-taxonomic relation and self-optimization. In 2024 International Joint Conference on Neural Networks (IJCNN), pages 1–9. IEEE.
- Janak Kapuriya, Manit Kaushik, Debasis Ganguly, and Sumit Bhatia. 2025. Exploring the role of diversity in example selection for in-context learning. In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '25, page 2962–2966, New York, NY, USA. Association for Computing Machinery.
- Duong Minh Le, Yang Chen, Alan Ritter, and Wei Xu. 2024. Constrained decoding for cross-lingual label projection. In *The Twelfth International Conference on Learning Representations*.
- Dongyang Li, Taolin Zhang, Nan Hu, Chengyu Wang, and Xiaofeng He. 2022. Hiclre: A hierarchical contrastive learning framework for distantly supervised relation extraction. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2567–2578.
- Junpeng Li, Zixia Jia, and Zilong Zheng. 2023a. Semiautomatic data enhancement for document-level relation extraction with distant supervision from large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5495–5505, Singapore. Association for Computational Linguistics.
- Min Li, Cong Shao, Gang Li, and Mingle Zhou. 2023b. Hfmre: Constructing huffman tree in bags to find excellent instances for distantly supervised relation extraction. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 12820–12832.

- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Quinn McNemar. 1947. Note on the sampling error of the difference between correlated proportions or percentages. *Psychometrika*, 12(2):153–157.
- Mike Mintz, Steven Bills, Rion Snow, and Dan Jurafsky. 2009. Distant supervision for relation extraction without labeled data. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pages 1003–1011
- Tom M. Mitchell, Justin Betteridge, Andrew Carlson, Estevam R. Hruschka Jr., and Richard C. Wang. 2009. Populating the semantic web by macro-reading internet text. In *The Semantic Web ISWC 2009, 8th International Semantic Web Conference, ISWC 2009, Chantilly, VA, USA, October 25-29, 2009. Proceedings*, volume 5823 of *Lecture Notes in Computer Science*, pages 998–1002. Springer.
- Arijit Nag, Soumen Chakrabarti, Animesh Mukherjee, and Niloy Ganguly. 2025. Efficient continual pretraining of LLMs for low-resource languages. In Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 3: Industry Track), pages 304–317, Albuquerque, New Mexico. Association for Computational Linguistics.
- Arijit Nag, Bidisha Samanta, Animesh Mukherjee, Niloy Ganguly, and Soumen Chakrabarti. 2021. A data bootstrapping recipe for low-resource multilingual relation classification. In *Proceedings of the 25th Conference on Computational Natural Language Learning*, pages 575–587.
- Jian Ni, Taesun Moon, Parul Awasthy, and Radu Florian. 2020. Cross-lingual relation extraction with transformers. *arXiv preprint arXiv:2010.08652*.
- Vipul Rathore, Kartikeya Badola, Parag Singla, et al. 2022. Pare: A simple and strong baseline for monolingual and multilingual distantly supervised relation extraction. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics* (Volume 2: Short Papers), pages 340–354.
- Vipul Kumar Rathore, Aniruddha Deb, Ankish Kumar Chandresh, Parag Singla, and Mausam . 2024. SSP: Self-supervised prompting for cross-lingual transfer to low-resource languages using large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 15081–15102, Miami, Florida, USA. Association for Computational Linguistics.

- Harman Singh, Nitish Gupta, Shikhar Bharadwaj, Dinesh Tewari, and Partha Talukdar. 2024. IndicGen-Bench: A multilingual benchmark to evaluate generation capabilities of LLMs on Indic languages. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11047–11073, Bangkok, Thailand. Association for Computational Linguistics.
- Mihai Surdeanu, Julie Tibshirani, Ramesh Nallapati, and Christopher D. Manning. 2012. Multi-instance multi-label learning for relation extraction. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 455–465, Jeju Island, Korea. Association for Computational Linguistics.
- Shikhar Vashishth, Rishabh Joshi, Sai Suman Prayaga, Chiranjib Bhattacharyya, and Partha Talukdar. 2018. RESIDE: Improving distantly-supervised neural relation extraction using side information. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1257–1266, Brussels, Belgium. Association for Computational Linguistics.
- Somin Wadhwa, Silvio Amir, and Byron C Wallace. 2023. Revisiting relation extraction in the era of large language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15566–15589.
- Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022. Text embeddings by weakly-supervised contrastive pre-training. *arXiv* preprint *arXiv*:2212.03533.
- Peng Wang, Xiaobin Wang, Chao Lou, Shengyu Mao, Pengjun Xie, and Yong Jiang. 2024. Effective demonstration annotation for in-context learning via language model-based determinantal point process. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 1266–1280, Miami, Florida, USA. Association for Computational Linguistics.
- Daojian Zeng, Kang Liu, Yubo Chen, and Jun Zhao. 2015. Distant supervision for relation extraction via piecewise convolutional neural networks. In *Proceedings of the 2015 conference on empirical methods in natural language processing*, pages 1753–1762.
- Yuhao Zhang, Victor Zhong, Danqi Chen, Gabor Angeli, and Christopher D. Manning. 2017. Position-aware attention and supervised data improve slot filling. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 35–45, Copenhagen, Denmark. Association for Computational Linguistics.
- Xiaoyan Zhao, Yang Deng, Min Yang, Lingzhi Wang, Rui Zhang, Hong Cheng, Wai Lam, Ying Shen, and

- Ruifeng Xu. 2024. A comprehensive survey on relation extraction: Recent advances and new frontiers. *ACM Computing Surveys*, 56(11):1–39.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623.
- Sizhe Zhou, Yu Meng, Bowen Jin, and Jiawei Han. 2024. Grasping the essentials: Tailoring large language models for zero-shot relation extraction. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 13462–13486, Miami, Florida, USA. Association for Computational Linguistics.

A Appendix

A.1 Language details

Please refer to table 5.

Language	Script	Language code
English	Latin	eng_Latn
Oriya	Oriya	ory_Orya
Santhali	Ol Chiki	sat_Olck
Manipuri	Meitei	mni_Mtei
Tulu	Kannada	tcy_Tulu

Table 5: Scripts and language codes for our test languages

A.2 Implementation Details

We make our code and data public at this link - https://anonymous.4open.science/r/AC2f-pool/.

A.2.1 LLMs

We use (1) unsloth/Meta-Llama-3.1-8B-Instruct-unsloth-bnb-4bit version of Llama 3.1 for local inference as well as for fine-tuning, (2) TogetherAI's unsloth/gemma-3-4b-it-unsloth-bnb-4bit for Gemma3, (3) TogetherAI's Qwen3 235B A22B Instruct 2507 FP8 Throughput for Qwen3 and (4) OpenAI's gpt-4o-2024-05-13 for GPT-4o.

For local inference and fine-tuning, we use a single NVIDIA A100 40GB GPU node.

A.2.2 Hyperparameters

For HYDRE, we set k (no. of exemplars) = 5 following the sensitivity analysis in section A.10, model threshold t = 0.5 following *PARE*'s original implementation.

For LLM inference, we set temperature (τ) = 0.0, max. input length as 2048 tokens and max. output tokens as 256.

A.2.3 Baselines

For few-shot prompting baselines, we flatten bags to sentences (assign all bag labels to each of it's sentences) and then perform exemplar selection using (a) Random, (b) Top-k similarity-based, and (c) MMR-based retrieval.

Following (Kapuriya et al., 2025), we implement MMR using iterative selection as follows:

$$MMR(q, s, S_{t-1}) = \alpha \cdot Sim(q, s) - (1 - \alpha) \cdot max_{s' \in S_{t-1}} Sim(s, s'),$$

where S_{t-1} denotes the candidate set selected so far at time step t. Here, α trades-off the relevance with diversity and is tuned on English dev set. The optimal value of α is found to be 0.3 and is used across all experiments.

A.2.4 Fine-tuning details

Translate-Train: PARE Adaptation: In the translate-train setting, we first adapt mBERT to each target language X by pretraining it on monolingual corpora derived from X_{train} . This results in a language-specific variant denoted $mBERT_X$. For pretraining, we extend the vocabulary with up to 10,000 randomly initialized new tokens for language X. Subsequently, we fine-tune the PARE model on X_{train} using $mBERT_X$ as the encoder using a task-specific adapter to obtain $PARE_{total}$.

LLaMA 3.1 Fine-Tuning: We use Low-rank adaptation (LoRA) with $lora_alpha = 64$, $lora_r = 16$, $lora_dropout = 0.0$, Learning rate scheduler as Cosine and warmup ratio as 10%. For training, we use $per_device_train_batch_size = 8$, $gradient_accumulation_steps = 4$ and maximum training steps of 5000 on single NVIDIA A100 40GB GPU.

The NYT-10m training data consists of 41624 bags and so the number of effective training epochs is = $8*4*5000/41624 \approx 4$ epochs.

Semantic Retriever Training: As we do not have an off-the-shelf retriever supporting our target languages, we seek to use a task-specific retriever. Specifically, we leverage the sentence-encoder of *CIL-X* for embedding queries and examples with cosine similarity as their similarity scores.

A.3 Prompt details

Task Description: Choose all applicable relations between head and tail entities from the set below. Print each relation in a new line. If none of the relations are applicable, output 'NA'.

/people/person/nationality : head entity is a person and tail entity is a country

/time/event/locations : head entity is an event and tail entity is a location

/people/person/children: head entity is a person and tail entity is another person (child)

/business/company/advisors : head entity is a company and tail entity is a person (advisor)

/business/location : head entity is a business and tail entity is a location

/business/company/majorshareholders: head entity is a company and tail entity is a person or organization (major shareholder)

/people/person/place_lived : head entity is a person and tail entity is a location

/business/company/place_founded : head entity is

a company and tail entity is a location

/location/neighborhood/neighborhood_of: head entity is a neighborhood and tail entity is a larger location (city, town)

/people/deceasedperson/place_of_death: head entity is a deceased person and tail entity is a location /film/film/featured_film_locations: head entity is a film and tail entity is a location

/location/region/capital: head entity is a region and tail entity is a city (capital)

/business/company/founders: head entity is a company and tail entity is a person (founder)

/people/ethnicity/geographic_distribution: head entity is an ethnicity and tail entity is a location where the ethnicity is commonly found

/location/country/administrative_divisions: head entity is a country and tail entity is a subdivision (state, province)

/people/deceasedperson/place_of_burial : head entity is a deceased person and tail entity is a burial site

/location/country/capital: head entity is a country and tail entity is a city (capital)

/business/person/company: head entity is a person and tail entity is a company they are associated with

/location/location/contains: head entity is a larger location and tail entity is a smaller location within it

/location/administrative_division/country: head entity is an administrative division (state, province) and tail entity is a country

/location/us_county/county_seat: head entity is a U.S. county and tail entity is the county seat /people/person/religion: head entity is a person and tail entity is a religion

/people/person/place_of_birth : head entity is a person and tail entity is a location (birthplace)

/people/person/ethnicity: head entity is a person and tail entity is an ethnicity

NA: no relation from the set exists between the given entity pair

Input format:

Input: {sentence}
Output format:
Output: {relation}

Verbalizer:

Extract the relation based on exact string match A sample format of input and output is shown in Figure 4.

A.4 HYDRE Algorithm

Please refer to Algorithm 1.

A.5 Detailed Language-wise results

Please refer to tables 6, 7 and 8.

A.6 Confusion Analysis

We present confusion matrix (aggregated over all labels) in Figure 3. We manually identify 16 categories of relation pairs amongst which confusion is observed and present their confusion matrices in Figures 14 and 15. We depict some qualitative examples, in which HYDRE correctly identifies the relation while 0-shot and PARE make a mistake, in the following subsection.

Confusion Matrix Comparison HYDRA Correct 223 HYDRA Wrong -25 272 PARE wrong PARE correct **HYDRA** Correct 312 83 HYDRA Wrong 30 267 0-shot correct 0-shot wrong

Figure 3: Confusion analysis of HYDRE vs the baselines (aggregated over all relations)

A.6.1 Qualitative examples

We select examples which are correctly predicted by HYDRE but misclassified by 0-shot prompting in order to study the impact of ICL examples. An example is shown in figure 4 where HYDRE correctly predicts "Religion" as opposed to "Ethnicity" by 0-shot prompting. ICL examples depict that "Ethnicity" is PARE's top-ranked candidate (closest to the query) while "Religion" is ranked second and HYDRE is able to subtly distinguish between the 2 based on their exemplars. Example 5 shows another case where HYDRE correctly predicts "Advisors" v/s "Company" predicted by 0-shot. This is a case where the predicted label

Algorithm 1 Exemplar Selection for HYDRE

```
Input: Query sentence q; bags of sentences \mathcal{B} = \{B_j\}; trained DSRE model scores f(s,r) for sentence s and relation r; semantic similarity function \sin(q, B_j); confidence threshold t; number of exemplars k.
```

```
Output: Ordered list of selected exemplar sentences.
 1: Compute f(q, r) for all relations r in the ontology.
 2: \mathcal{R}' \leftarrow \text{top-}k relations ranked by f(q,r)
                                                                                                                                         3: Initialize exemplar list \mathcal{E} \leftarrow [].
 4: for each r \in \mathcal{R}' do
           \mathcal{B}_r \leftarrow \{B_j \in \mathcal{B} \mid r \in labelset(B_j)\}.
 5:
          \begin{array}{l} B_r \leftarrow \arg\max_{B_j \in \mathcal{B}_r} [\sin(q,B_j) + f_{PARE}(B_j,r)]. \\ \text{For each } s \in B_r, \text{ compute } c(s) = \sum_{r_a \in \mathsf{labels}(B_r)} \mathbb{I}[f(s,r_a) > t]. \end{array}
 6:
                                                                                                                                      \triangleright Select most relevant bag for r
 7:
 8:
           Let v_{\max} \leftarrow \max_{s \in B_r} c(s).
                                                                                                          ⊳ max. no. of relations with confidence > threshold
           \mathcal{S}' \leftarrow \{s \in B_r \mid c(s) = v_{\max}\}.
                                                                                                               ⊳ candidate sentences with max. label coverage
           s^* \leftarrow \arg\max_{s \in \mathcal{S}'} \sum_{r_a \in labels(B_r)} f(s, r_a).
10:
                                                                                                            ▶ Pick sentence with highest aggregate confidence
           Add (s^*, r_a) to \mathcal{E}.
12: end for
13: Sort \mathcal{E} in ascending order of f(q, r).

    b keep most informative examples at the last (closer to query)

14: return Ordered list of exemplar sentences \{s^*\} from \mathcal{E}.
```

Model	eng_Latn	ory_Orya	sat_Olck	mni_Mtei	tcy_Tulu	Avg. (Micro) [†]	Avg. (Macro)†
Supervised							
PARE (Rathore et al., 2022)	42/31	_	_	_	_	_	_
CIL (Chen et al., 2021)	43/32	_	_	_	_	_	_
zero-shot							
Qwen3-235B-A22B	49/40	46/40	6/5	2/1	45/39	25	21
Llama3.1-8b	31/17	26/21	21/14	10/7	29/21	22	16
Llama 3.1 -8B-FT $_{En}$	60/44	33/23	15/8	11/4	36/21	24	14
GPT-4o	56/57	56/56	10/5	11/7	55/ 55	33	31
5-shot							
HYDRE(Qwen3-235B-A22B)	63*/62	52/46	9/7	2/2	53/48	29	26
HYDRE(Llama3.1-8B)	52/47	37/24	24*/16	21*/11	38/25	30	19
$HYDRE(Llama 3.1-8B-FT_{En})$	61/45	47/32	24* /13	21* /8	49/31	35	21
HYDRE(GPT-4o)	63* /60	57/57	18/10	11/8	56/55	36*	33

Table 6: Results for English-only data setting. In each entry, we report micro and macro F1 scores. †The reported average scores are over non-English languages. * McNemar's p-value $< 10^{-5}$ (valid for micro-F1 comparison).

"Advisors" is not listed in PARE's top-5 candidates and despite this the HYDRE is able to predict it. This is attributed to ICL's ability to predict labels beyond what are included in the in-context examples. Other scenarios where HYDRE dominates 0-shot prompting include multi-label prediction as depicted in Ex. 7, where HYDRE predicts "/place_of_burial" as well as "/place_lived" while 0-shot misses the obvious label "/place_lived" while only predicting "/place_of_burial". this is not always the case as shown in Ex. 11, where PARE top-1 is correct label (" /administrative_division/country") but HYDRE predicts a wrong one (" /location/location/contains") which is not in PARE top-5 (ICL Examples). Ex. 12 depicts multi-label query where HYDRE predicts 2 out of 3 labels correctly while 0-shot predicts only one of those. Ex. 13 is another case where HYDRE correctly predicts a more nuanced label "/place_founded" but misses an obvious label

"/business/location", which is not covered in ICL exemplars.

A.7 Error Analysis

We next analyze cases where HYDRE fails to recall correct relations. One case is where HYDRE gets biased to simply predict the PARE's top-1 candidate and misses others in multi-label queries. An example is shown in Figure 9 where HYDRE predicts "/geographic_dsitribution" (PARE's top-1) while missing the label "/nationality" (PARE's top-3), which is predicted by 0-shot prompting.

We further divide the testset facts into 2 categories - (a) for which gold label is in PARE's top-5, and (b) gold label is out of PARE top-5. The individual micro F1 scores on both categories is presented in Table 9. We observe that missing the gold label in PARE top-5 significantly hurts the downstream LLM F1 score compared to hitting it. This shows the positive impact of candidate selection step on

Model	ory_Orya	sat_Olck	mni_Mtei	tcy_Tulu	Avg. (Micro)	Avg. (Macro)
Supervised						
PARE-X	31/20	29/20	28/19	30/19	30	20
CIL-X	29/18	22/15	25/19	29/20	26	18
zero-shot						
Qwen3-235B-A22B	46/40	6/5	2/1	45/39	25	21
Llama3.1-8b	26/21	21/14	10/7	29/21	22	16
GPT-4o	56/ 56	8/6	8/6	56/53	32	30
Llama 3.1 -8B-FT $_{En}$	33/23	10/7	2/1	36/21	20	13
Llama3.1-8B-FT _X	51/37	33/18	24/11	44/26	38	23
5-shot						
$HYDRE_X(Qwen3-235B-A22B)$	56/53	23/12	14/6	57/52	38	31
$HYDRE_X(Llama 3.1-8B)$	38/22	34/19	25/13	42/28	35	21
$HYDRE_X(GPT-40)$	58*/56	17/11	20/15	57/55	38	34
$HYDRE_X(Llama 3.1-8B-FT_{En})$	48/33	27/15	21/8	47/31	36	22
$HYDRE_X(Llama3.1-8B-FT_X)$	56/39	45*/24	35*/16	51/36	47*	29

Table 7: Results for Translate-train setting. * McNemar's p-value $< 10^{-5}$.

Model	ory_Orya	sat_Olck	mni_Mtei	tcy_Tulu	Avg. (Micro)	Avg. (Macro)
Supervised						
PARE	35/25	31/21	33/25	32/22	33	23
CIL	34/23	33/21	34/25	36/25	34	24
zero-shot						
Qwen3-235B-A22B	44/37	43/34	43/35	44/34	44	35
Llama3.1	29/25	26/22	29/25	30/26	29	25
Llama 3.1 -8B-FT $_{En}$	49/35	45/29	49/35	52/39	49	34
GPT-40	53/51	50/45	55/ 57	56/51	54	51
5-shot						
HYDRE(Qwen3-235B-A22B)	52/46	51/42	55/47	56/48	54	46
HYDRE(Llama3.1-8B)	46/38	41/31	46/38	44/38	44	36
HYDRE(Llama $3.1-8B-FT_{En}$)	51/37	50/33	50/37	54/40	51	37
HYDRE(GPT-4o)	55*/54	53*/49	56/57	59*/55	56*	54

Table 8: Results for Translate-test setting. * McNemar's p-value $< 10^{-5}$.

the downstream performance of HYDRE..

Category	No. of Facts	micro F1 score
gold label in top-5	594	67
gold label out of top-5	98	34

Table 9: micro F1 scores for test-set categories

A.8 Recall@5 for Supervised DSRE models

We present Recall@5 for both *PARE-X* and *CIL-X* models in table 10 for English and other languages.

Language	PARE	CIL
English	84	82
Oriya	75	_ 72
Santhali	72	64
Manipuri	69	69
Tulu	69	71
Avo *	71	69

Table 10: Recall@5 scores for *PARE* and *CIL* models on all languages. *Averaged over non-English languages

A.9 Additional supervised DSRE baselines

We include a comparison of *PARE* and *CIL* with HiCLRE and HFMRE baselines in Table 11.

Baseline	Micro F1	Macro F1
HiCLRE	31	18
HFMRE	32	18
PARE		31
CIL	43	32

Table 11: Comparison of additional English baselines with PARE and CIL.

A.10 Scalability w.r.t. No. of candidates k

We analyze how PARE's Recall@k and the HYDRE performance vary as a function of k (no. of candidate relations) on the NYT-10m English dev set. Results are shown in Figure 17. While the recall@k keeps on increasing even beyond k=5, the downstream F1 performance saturates at k=5 showing that increasing k may not only increase LLM's prompt length but also confuse the LLM

since the number of candidates is too large to filter the correct relation(s). Therefore, we fix k=5 consistently for all our experiments.

A.11 Detailed Language-wise Ablations

Results are presented in tables 12, 13, 14 and 15 for each of the 4 LLMs.

A.12 Dataset statistics

A.12.1 Test data split

We construct the test split from NYT-10m (Gao et al. 2021) using stratified sampling, ensuring a minimum of 30 instances per relation, except when a relation contains fewer than 30 instances in total. This results in 538 sentences, with the distribution of relation counts shown in Table 16. Since NYT-10m is a multi-label dataset, the total number of relation instances in the test split exceeds the number of sentences, amounting to 722.

A.12.2 Training data

We take the original training data from NYT-10m (Gao et al. 2021) and ensure that "NA" bags do not exceed 10% of total bags to avoid model overfit on "NA" label. This leads to a total number of 41624 training bags.

A.13 Data annotation for Indic languages

A.13.1 Annotator details

We conduct human verification for four Indic languages using native speakers—either students or IT professionals—who were proficient in reading and typing in their respective scripts. Each annotator was compensated approximately \$60 for verifying translations of 538 sentences. Prior to annotation, the speakers were informed that the task was intended solely for research purposes and posed no risk to them.

Each annotator was presented with the following questionnaire, with binary (YES/NO) responses and rectifications requested in case of a NO:

- 1. **Q1.** Is the translation of the given English sentence correct?
- 2. **Q2.** Is the *head entity* correctly translated into your native language?
- 3. **Q3.** Is the *head entity* correctly projected in your native language?
- 4. **Q4.** Is the *tail entity* correctly translated into your native language?

5. **Q5.** Is the *tail entity* correctly projected in your native language?

A.13.2 Quality Assessment and Interannotator Agreement

We first assess the quality of the system-generated translations presented to annotators. Native speakers across all languages found the translations to be generally decent—likely due to high-quality output from IndicTrans2 (and Google Translate in the case of Tulu).

To quantify this, Table 17 reports:

- 1. The percentage of translations that required no human correction
- 2. The character-level F1 score (Char-F1) between the original system-generated translation and the human-corrected version (for cases requiring rectification).

To further evaluate annotation reliability, we conducted an inter-annotator agreement study. A second native speaker independently judged the quality of translations for 100 randomly sampled sentences in each language. Agreement is reported as the percentage of samples where the second speaker's judgment matched that of the first. Results are shown in Table 18.

In summary, on average, 74% of system-generated translations were accepted as correct by native speakers, and for the remaining, the human-corrected outputs had a 93% Char-F1 match with the original translations. Inter-annotator agreement for human-corrected outputs averaged 91%, indicating strong consistency and translation quality across languages.

A.14 Semantic retrieval hurts performance for low-resource languages

We try to use both off-the-shelf retriever BGE-m3 and a fine-tuned retriever (*CIL-X*'s sentence encoder) for HYDRE in translate-train setting. We observe (Table 19) that though the fine-tuned retriever's performance is significantly better than BGE-m3, they both are worse compare to the variant of HYDRE that only uses *PARE-X*'s confidence for candidate scoring (results in table). This analysis suggest that semantic similarity is not useful in translate-train settings for our target languages.

Ablation variant	Ory	Sat	Mni	Tcy	mi.	mu.
Hydre	38/22	24/16	25/13	42/28	32	20
w. sem. sim.	34/24	21/16	14/10	35/22	26	18
w/o PARE (Random)	33/22	24/14	17/10	30/18	26	16
w/o stage 1	39/26	33/20	23/10	38/25	33	20
w/o ICL	32/21	28/19	25/16	30/19	29	19

Table 12: Average F1 scores for **Llama3.1-8b** over our target languages for different ablation variants of our few-shot approach in *translate-train* setting

Ablation variant	Ory	Sat	Mni	Tcy	mi.	mu.
Hydre	56 /39	39/22	35/16	51/36	45	28
w. sem. sim.	54/ 40	38/ 22	33/14	49/32	44	27
w/o PARE (Random)	53/34	31/15	21/8	47/28	38	21
w/o stage 1	54/35	35/21	33/13	48/30	43	25
w/o ICL	50/30	27/15	10/6	44/22	33	18

Table 13: Average F1 scores for **Llama3.1-8b-FT** $_X$ over our target languages for different ablation variants of our few-shot approach in *translate-train* setting

A.15 Additional ablation to justify aggregate confidence for sentence selection

HYDRE selects the representative sentence based on aggregated confidence over all labels of the bag. This helps surface sentences with stronger overall evidence, and not just noisy single-label confidence. We justify our design choice with an additional ablation *Candidate-only scoring* in which we use only candidate label's confidence score for stage 2 retrieval instead of aggregate score over all bag labels. Results are shown in Table 15. While candidate-only scoring performs comparably for GPT-40, aggregate scoring yields substantial gains for smaller LLMs like Llama and Qwen3, validating its broader effectiveness.

A.16 Additional results and ablations on Wiki-20m

We expand our evaluation on another English dataset Wiki-20m (Gao et al., 2021), containing 2386 manually annotated test sentences and 81 relations (including "NA"). Fig. 21 presents the results and ablations for GPT-40 and Llama-3.1-8b-Instruct models. We observe that HYDRE achieves huge gains compared to 0-shot and ablations, beating it's closest competitor by 10 and 20 micro F1 points respectively for GPT-40 and Llama-3.1, showing that our approach also scales well when the number of relations in the ontology is large (80 in this case).

Ablation variant	Ory	Sat	Mni	Tcy	mi.	mu.
Hydre	56 /53	23/12	14/6	57/52	38	31
w. sem. sim.	55/53	24/15	15/5	55/52	37	31
w/o PARE (Random)	55/ 55	19/11	15/5	58/53	37	31
w/o stage 1	51/52	17/9	16/7	52/50	34	30
w/o ICL	54/36	15/8	8/4	53/33	33	20

Table 14: Average F1 scores for **Qwen3-235B-A22B** over our target languages for different ablation variants of our few-shot approach in *translate-train* setting

Ablation variant	Ory	Sat	Mni	Tcy	mi.	mu.
Hydre	58 /56	17/11	20/15	57 /55	38	34
w. sem. sim.	57/ 57	18/ 12	22/15	56/52	38	34
w/o PARE (Random)	55/55	15/8	16/10	57 /55	36	32
w/o stage 1	57/56	20 /11	20/15	57/57	39	35
w/o ICL	54/38	12/6	12/7	48/32	32	21

Table 15: Average F1 scores for **GPT-40** over our target languages for different ablation variants of our few-shot approach in *translate-train* setting

Relation	Count
/people/person/place_lived	73
/people/person/nationality	34
/business/person/company	34
/people/person/place_of_birth	31
/location/location/contains	102
/location/country/administrative_divisions	58
/business/location	33
/location/administrative_division/country	31
/business/company/advisors	37
/business/company/founders	31
/business/company/majorshareholders	4
/location/neighborhood/neighborhood_of	30
/location/country/capital	30
/film/film/featured_film_locations	1
/location/us_county/county_seat	6
/people/person/children	30
/people/deceasedperson/place_of_death	30
/people/deceasedperson/place_of_burial	4
/people/ethnicity/geographic_distribution	30
/location/region/capital	12
/business/company/place_founded	4
/people/person/religion	21
/time/event/locations	3
/people/person/ethnicity	23
NA	30
Total	538

Table 16: Label-wise statistics of our evaluation data. Total number of sentences (538) in our test split is different from total number of labels (722) due to multi-label characteristics of NYT-10m dataset.

Language	No Rectification Needed (%)	Char-F1 Match (if rectified)
Oriya	69	92
Santhali	72	88
Manipuri	83	96
Tulu	73	95
Average	74	93

Table 17: Percentage of system translations requiring no rectification, and character-level F1 match for rectified translations.

Language	Translation (%)	Head entity projection	Tail entity projection
Oriya	92	92	92
Santhali	89	84	87
Manipuri	85	98	100
Tulu	96	92	88
Average	91	92	92

Table 18: Inter-annotator agreement for 100 samples in each language

Ablation variant	Llama3.1	Llama3.1-ft	GPT-40
HYDRE (using only	32	45	38
PARE-X confidence)			
sem. sim. (BGE-M3)	24	42	36
sem. sim. (CIL-X)	26	44	39

Table 19: Average F1 scores over our target languages for HYDRE (using only *PARE-X* confidence) and sem. sim. variants (both off-the-shelf and fine-tuned) in *translate-train* setting.

Ablation variant	GPT-40	Qwen3	Llama3.1
HYDRE (aggregate scoring)	63 /60	63/62	52/47
Candidate-only scoring	63/62	62/60	47/35

Table 20: Ablations of confidence aggregation step for sentence selection in stage 2.

Ablation	Llama3.1	GPT-40
0-shot	11/16	55/48
HYDRE	54/53	67/63
Ablations		
w/o semantic similarity	25/29	34/36
w/o PARE confidence	21/16	40/36
w/o both (Random)	34/34	57/54
w/o stage 2	39/39	38/37
w/o ICL	3/4	12/4

Table 21: English F1 (micro/macro) scores on Wiki-20m evaluation set.

Input: But Mr. Wallace , 45 , ... in partnership for a time with <Tail> Mark Thatcher </Tail> , son of the former British Prime Minister <Head> Margaret Thatcher </Head> , has detractors ,

Output: /children

Input: In the wake of the torture and killing in February of <Head> Ilan Halimi </Head> , a 23-year-old Jew , attention has focused on an undeniable problem : anti-Semitism among <Tail> France </Tail> 's second-generation immigrant youth ,

Output: /nationality

Input: Another concern is that <Tail> Ethiopia </Tail> and Eritrea , bitter enemies that recently fought ... , with Eritrea suspected of ... to support the Islamists and <Head> Ethiopian </Head> officials now admitting ,

Output: /geographic_distribution

Input: The case has fueled feelings here of an assault against <Tail> Islam </Tail> , coming after ... and , more recently , cartoons ... that mocked the Prophet <Head> Muhammad </Head> .

Output: /religion

Input: ... , Mr. Delli Colli worked with generations of <Tail> Italian </Tail> directors , including <Head> Pier Paolo Pasolini
</Head> ,

Output: /ethnicity

Query:

Input: A jury convicted a 27-year-old British <Tail> Muslim </Tail> , <Head> Umran Javed </Head> , of soliciting murder and inciting racial hatred

Outputs:

HYDRE: /religion

GPT-40 0-shot: /ethnicity

Figure 4: **Example 4**: HYDRE correctly predicts the "*Religion*" label while 0-shot confuses it with "*Ethnic-ity*".

ICL Exemplars (HYDRE)

Input: Mr. White , 54 , of <Tail> Centerport </Tail> , has held top environmental posts with <Head> Suffolk County </Head> and the Town of Huntington .

Output: /location/location/contains

Input: Being able to combine Loudeye services with <Head> Nokia </Head> terminals provides ... , " said Ilkka Raiskinen , senior vice president for multimedia experiences at Nokia , based in <Tail> Espoo </Tail> , Finland .

Output: /place_founded

Input: When people say '...', Jerry was just as responsible for that as my dad, "' said <Tail> Brian Henson </Tail> , <Head> Jim Henson </Head> 's son , who serves with his sister Lisa as chairman and chief executive of the Jim Henson Company.

Output: /children

Output: /business/company/founders, /business/company/majorshareholders

Input: The service has more ways to shield
users 'identities , said <Tail> Antony
Brydon </Tail> , <Head> Visible Path
</Head> 's chief executive , above ."

Output: /business/company/founders

Query:

Input: "<Head> MySpace </Head> is dedicated to ensuring that ..., "<Tail> Chris DeWolfe </Tail> , the chief executive of MySpace , said in a statement .

Outputs:

HYDRE: /business/company/advisors 0-shot: /business/person/company

Figure 5: **Example 5**: HYDRE correctly predicts the "*Advisors*" label while 0-shot misses it.

Input: Mr. Marek started the <Head> Fairfield </Head> Theater Company in 2001 ... at <Tail> Fairfield University </Tail>.

Output: /location/location/contains

Input: But Mr. Wallace , 45 , a businessman and an investor in partnership for a time with <Tail> Mark Thatcher </Tail> , son of the former British Prime Minister <Head> Margaret Thatcher </Head> , has detractors , including"

Output: /people/person/children

Input: Behind the News – Founded eight years ago in a Silicon Valley garage by two <Tail> Stanford University </Tail> graduate students , <Head> Google </Head> went public two years ago at \$ 85 a share .

Output: /business/company/place_founded

Input: A notch or two down-market , the 777 's , 767 's and 757 's are often coveted by corporate titans , among them Larry Page and <Tail> Sergey Brin </Tail> , the co-founders of <Head> Google </Head> , who bought

Output: /business/company/founders /business/company/majorshareholders

Input: When he was 32, <Tail> Bill Hambrecht </Tail> was a co-founder of <Head> Hambrecht & Quist </Head>, a West Coast investment bank that

Output: /business/company/founders

Query:

Input: Mr. Bechtolsheim, one of the first investors in Head> Google , co-founded Kealia in 2001 with Tail> David
Cheriton Total>, a Stanford professor who was another early Google investor.

Outputs:

HYDRE: /busi-

ness/company/majorshareholders 0-shot: /people/person/company

Figure 6: **Example 6:** HYDRE correctly predicts the "*Majorshareholders*" relation while 0-shot fails to recognize it.

ICL Exemplars (HYDRE)

Input: ... Mr. Koppel to take another look at a once-unknown man , <Head> Morrie Schwartz </Head> , a <Tail> Brandeis University </Tail> professor who

Output: /business/person/company

Input: ... as you race away from the pleasant corporate maw of <Tail> Seattle </Tail> , from Starbucks and <Head> Boeing </Head> , Amazon and Microsoft . **Output**: /business/company/place founded

Input: ... <Head> Ernest Hemingway </Head> was born in <Tail> Oak Park </Tail> in 1899 and lived here through high school.

Output: /place_of_birth /people/person/place_lived

Input: ... Mr. Narayanan 's body will be cremated ... in <Tail> New Delhi </Tail> , near the funeral ground of <Head> Jawaharlal Nehru </Head> , India 's first prime minister

Output: /place_of_death /place_lived

Input: The easiest ... way to see <Tail> Philadelphia </Tail> is to stick with the older , central parts of town , emulate <Head> Benjamin Franklin </Head>

Output: /place_of_death /place_lived

Query:

Input: Poe, Evermore A mystery man arrived at <Head> Edgar Allan Poe </Head> 's grave at the Westminster Burial Grounds in <Tail> Baltimore </Tail> on Friday morning, as he has on Poe 's birthday (Jan. 19) every year since 1949,

Outputs:

HYDRE: /place_of_burial, /place_lived 0-shot: /place_of_burial

Figure 7: **Example 7:** HYDRE correctly predicts both "place_of_burial" and "place_lived" while 0-shot partially predicts only one of them.

Input: ... Mendelssohn , who was born Jewish and converted to Christianity , and <Head> Otto Klemperer </Head> , who converted to Christianity and then back to <Tail> Judaism </Tail> .

Output: /religion

Input: As the <Head> <Tail> Baltimore </Tail> Orioles </Head> return home ... , Baltimore embraces its rich sports and maritime history .

Output: /business/location

Input: <Head> Lorenzo Da Ponte </Head> , a Bridge From <Tail> Italy </Tail> to New York " includes three vocal recitals , beginning tonight with

Output: /nationality

Input: The Museum of Modern Art 's exhibition of four films starring the <Tail> Italian </Tail> actress <Head> Laura Morante </Head> concludes this weekend with four films, including

Output: /ethnicity

Input: ... <Head> Madhesi </Head> ethnic group, which by some estimates represents as much as a third of <Tail> Nepal </Tail> 's population of 29 million, has been granted citizenship rights

Output: /geographic_distribution

Query:

Input: Among the performances of note : ... the <head> Italian </head> dancer ALESSANDRA FERRI gives her final performance with the company on Saturday night , with ROBERTO BOLLE , a guest artist also from <Tail> Italy </Tail> .

Outputs:

Correct: /geographic_distribution, /nation-

ality

HYDRE: /geographic_distribution

0-shot: /nationality

Figure 8: **Example 8:** HYDRE partially predicts "/geo-graphic_distribution" while 0-shot partially predicts the other label "/nationality".

ICL Exemplars (HYDRE)

Input: ... owned by a <Tail> Cincinnati </Tail> company , American Financial Group , whose chairman and chief executive officer is Carl H. Lindner III , who is also an owner of the <Head> Cincinnati Reds </Head> .

Output: /business/location

Input: A biomedical research institute in <Tail> Chengdu </Tail> , <Head> China </Head> , is planning to show true commitment to scientific principles

Output: /location/location/contains /location/country/administrative_divisions

Input: <Tail> Robert Bigelow </Tail>
, the founder of <Head> Budget Suites
of America </Head> , is likely to push
forward

Output: /founders

Input: A notch or two down-market, the 777's, 767's and 757's are often coveted by corporate titans, among them Larry Page and <Tail> Sergey Brin </Tail>, the co-founders of <Head> Google </Head>, who bought

Output: /founders /majorshareholders

Input: The N.F.L. ... is very popular, ..., " said Tony Ponturo, vice president for global media and sports marketing at <Head> Anheuser-Busch </Head> in <Tail> St. Louis </Tail>,"

Output: /place_founded

Query:

Input: <Head> Nestlé </Head> , based in
<Tail> Vevey </Tail> , Switzerland ,

Outputs:

Correct: /business/location HYDRE: /place_founded 0-shot: /place_founded

Figure 9: **Example 9:** Both HYDRE and 0-shot misclassify the "/business/location" label as "/place_founded"

Input: But Mr. Wallace ... in partnership for a time with <Tail> Mark Thatcher </Tail> , son of the former British Prime Minister <Head> Margaret Thatcher </Head> , has detractors

Output: /people/person/children

Input: Because the Newbery is open only to American citizens or residents , one enormously popular writer who is n't in the running is <Head> Cornelia Funke </Head> , who lives in <Tail> Germany </Tail> and whose books appear here in translation

Output: /people/person/nationality

Input: <Head> Keith Jarrett </Head> lives on the New Jersey side of the Pennsylvania border , within an hour 's drive of his childhood home of <Tail> Allentown </Tail> , Pa. .

Output: /people/person/place_of_birth /people/person/place_lived

Input: ..., Pollan finds his hero in <Head> Joel Salatin </Head> , an "alternative "farmer in <Tail> Virginia </Tail>

Output: /people/person/place_lived

Input: Charles G. Taylor, the former president of Liberia ... arrived in handcuffs in the <Tail> Netherlands </Tail> on Tuesday, and was immediately taken to the jail near <Head> The Hague </Head>

Output: /administrative_division/country

Query:

Input: Mr. Meshal lives in exile in <head> Damascus </head> , <Tail> Syria </Tail> , where the government has rebuffed all Western requests to close his office .

Outputs:

Correct: /administrative_division/country
HYDRE: /location/location/contains
0-shot: /location/location/contains

Figure 10: **Example 10:** Both HYDRE and 0-shot misclassify the "/administrative_division/country" label as "/location/location/contains"

ICL Exemplars (HYDRE)

Input: ... the epic poem "Paterson" by
<Head> William Carlos Williams </Head> ,
a native of <Tail> Rutherford </Tail> .

Output: /people/person/place_of_birth

Input: It goes on to list notable <Tail>
Mississippi </Tail> writers including
William Faulkner , <Head> Richard Wright
</Head> ,

Output: /people/person/place_lived

Input: ... , including Giocangga , the founder of the <Head> Manchu </Head> dynasty in <Tail> China </Tail> , and Niall of the Nine Hostages ,"

Output: /peo-ple/ethnicity/geographic_distribution

Input: <Head> Anthony Trollope </Head>
, the brilliant depicter of the 19th-century
social strata in <Tail> England </Tail> ,

Output: /people/person/nationality

Input: ... , Mr. Delli Colli worked with generations of <Tail> Italian </Tail> directors , including <Head> Pier Paolo Pasolini </Head> ,

Output: /people/person/ethnicity

Query:

Input: Lukacs, a distinguished historian of 20th-century Europe, makes very large claims for his subject in "George Kennan": He was "a better writer and a better thinker "than <Head> Henry Adams </Head>; he was "the best and finest <Tail> American </Tail> writer about Europe "in the interwar years, better than Hemingway.

Outputs:

Correct: /people/person/ethnicity HYDRE: /people/person/nationality 0-shot: /people/person/nationality

Figure 11: **Example 11:** Both HYDRE and 0-shot misclassify the "/people/person/ethnicity" label as "/people/person/nationality"

Input: Halliburton operates in Iran through a unit based in nearby <Head> Dubai </Head> , <Tail> United Arab Emirates </Tail> ,

Output: /administrative_division/country

Input: ... when the State Department
sponsored its tour to Damascus , Homs and
<Tail> Lattakia </Tail> , <Head> Syria
</Head> .

Output: /location/location/contains /country/administrative_divisions

Input: That is one reason that <Head> Hunan </Head> 's fast-growing provincial capital , <Tail> Changsha </Tail> , is beginning to

Output: /location/location/contains/region/capital

Input: It has outfitted the World Financial Center and the new Bloomberg L.P. head-quarters in New York as well as the <Tail> Cheung Kong Center </Tail> in <Head> Hong Kong </Head>

Output: /location/location/contains

Input: There is no shortage of cruises that
stop in either <Tail> Tallinn </Tail> ,
<Head> Estonia </Head> , or Riga , Latvia ,
or both .

Output: /location/location/contains /location/country/capital

Query:

Input: Last year , it opened offices in Warsaw and <Tail> Bucharest </Tail> , the capital of <Head> Romania </Head> .

Outputs:

Correct: /country/capital, /country/administrative_divisions, /location/location/contains

HYDRE: /location/location/contains,

/country/capital

0-shot: /country/capital

Figure 12: **Example 12:** HYDRE misses the "/country/administrative_divisions" label while 0-shot misses both "/country/administrative_divisions" and /location/location/contains.

ICL Exemplars (HYDRE)

Input: Mr. Schwartz, 32, teaches English ... as a <Tail> New York City </Tail> teaching fellow at Intermediate School 349 in <Head> Bushwick </Head>, Brooklyn.

Output: /neighborhood_of

Input: <Tail> Nyack </Tail> 's mayor, John Shields, said ... to locate a theater company in the area if efforts to save the <Head> Helen Hayes </Head> theater fail.

Output: /place_of_death

Input: "..." said Mr. Jordan , who just signed on with <Head> Mark Warner </Head> , a Democrat and the former governor of <Tail> Virginia </Tail> who is considering a run for president .

Output: /place_lived

Input: <Head> John David Booty </Head>
's decision to skip his senior year at ... in
<Tail> Shreveport </Tail> did not

Output: /place_of_birth, /place_lived

Input: Jack Emmert had already earned his master 's degree ... to become creative director at <Head> Cryptic Studios </Head> , a game company based in <Tail> Los Gatos </Tail> , Calif. , where he

Output: /place_founded

Query:

Input: But before Mr. Joffe, ..., accepted a job, a friend suggested he check out a <Tail> San Francisco </Tail> start-up, <Head> Powerset </Head>, which is trying to build a rival search engine.

Outputs:

Correct: /business/location, /place_founded

HYDRE: /place_founded 0-shot: /place_founded

Figure 13: **Example 13:** Both HYDRE and 0-shot GPT-40 miss the obvious label "/business/location" while correctly predict a more nuanced label "/place_founded".

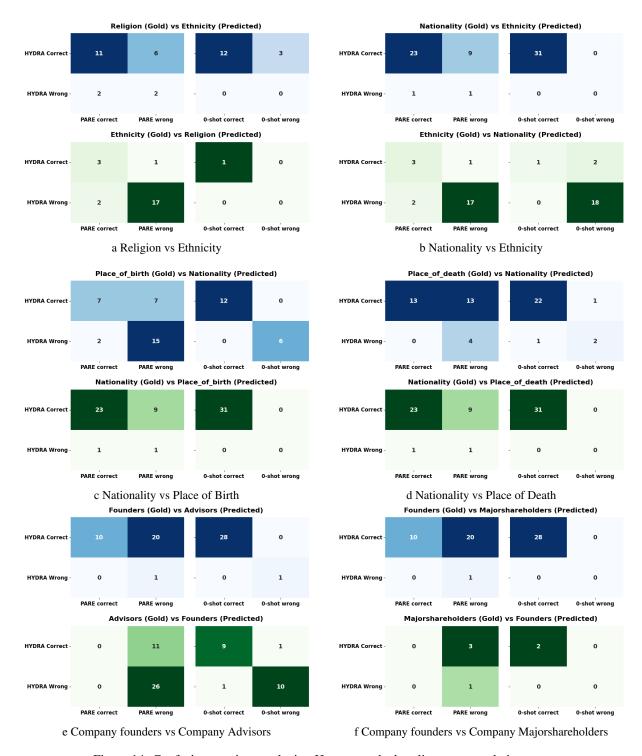


Figure 14: Confusion matrices analyzing HYDRE vs the baselines across relation types

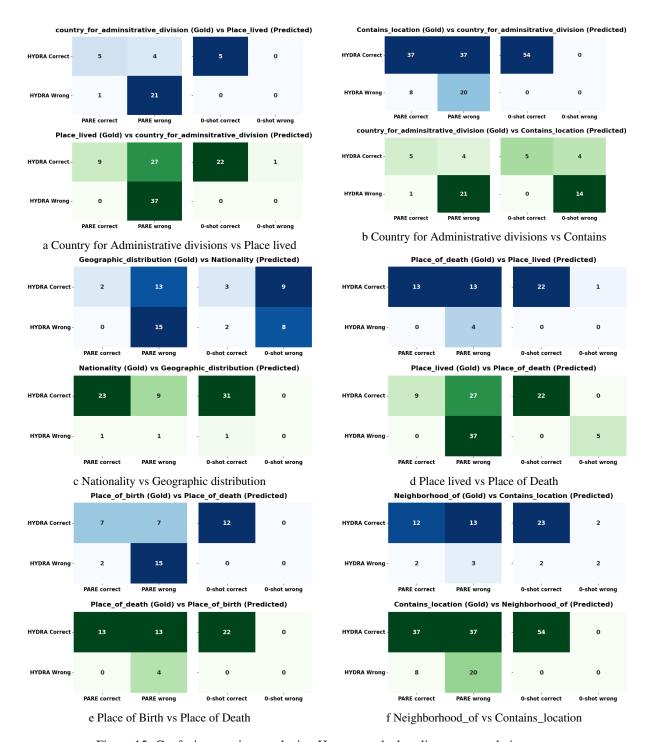


Figure 15: Confusion matrices analyzing HYDRE vs the baselines across relation types

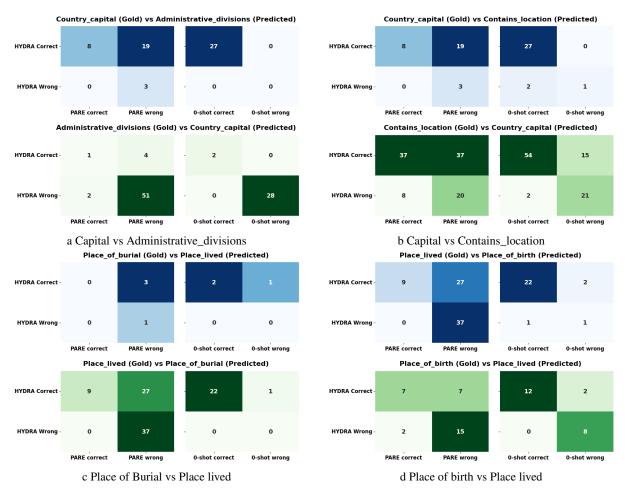


Figure 16: Confusion matrices analyzing HYDRE vs the baselines across relation types

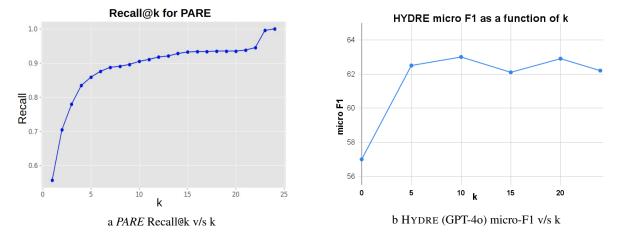


Figure 17: Analysis of PARE's Recall@k and HYDRE (GPT-40) downstream performance on NYT-10m English dev set