Relation Extraction or Pattern Matching? Unravelling the Generalisation Limits of Language Models for Biographical RE

Anonymous ACL submission

Abstract

001 Analysing the generalisation capabilities of relation extraction (RE) models is crucial for assessing whether they learn robust relational patterns or rely on spurious correlations. Our cross-dataset experiments find that RE models struggle with unseen data, even within similar domains. Notably, higher intra-dataset perfor-007 800 mance does not indicate better transferability, instead often signaling overfitting to datasetspecific artefacts. Our results also show that 011 data quality, rather than lexical similarity, is key to robust transfer, and the choice of opti-012 mal adaptation strategy depends on the qual-014 ity of data available: while fine-tuning yields the best cross-dataset performance with highquality data, few-shot in-context learning (ICL) is more effective with noisier data. However, 017 018 even in these cases, zero-shot baselines occasionally outperform all cross-dataset results. 019 Structural issues in RE benchmarks, such as single-relation per sample constraints and nonstandardised negative class definitions, further hinder model transferability.¹

1 Introduction

027

034

035

Relation extraction (RE) is the core information extraction task of identifying the semantic relationship between entities in text. Traditional RE evaluations rely predominantly on in-distribution testing, but this approach often overestimates true model performance by implicitly assuming that individual datasets wholly represent the underlying task (Linzen, 2020; Kovatchev and Lease, 2024). While model generalisation has gained increasing attention in NLP, RE remains relatively unexplored in this context (§ 2).

However, understanding RE generalisation to out-of-distribution (OOD) data is crucial both for the task itself as well as for the robust application of RE systems in downstream tasks like question answering and knowledge-base population (Bassignana and Plank, 2022a). Given the popularity of representing internal language model (LM) knowledge as relational triples (Geva et al., 2023; Hernandez et al., 2024), building robust RE systems beyond the mere memorisation of dataset-specific patterns may also be key to more interpretable and trustworthy models. 040

041

042

045

046

047

048

050

051

052

054

060

061

062

063

064

065

066

067

068

069

070

071

073

074

076

077

This paper systematically analyses how well RE systems generalise across datasets. Due to the limited relation overlap in popular RE datasets, we focus our experiments on biographical relations, which are pervasive in RE settings; this also allows us to include a domain-specific dataset for grounded analysis (§ 3). While prior work explored various robustness tests, including domain shift and adversarial attacks (Chen et al., 2023; Meng et al., 2024), to the best of our knowledge, this work provides the first thorough analysis of cross-dataset generalisation capabilities across sentence-level RE benchmarks.

Through our cross-dataset experiments, this paper makes the following contributions:

- We document key challenges in analysing RE generalisation, including inconsistent relation schemas and highly imbalanced class distributions (§ 3), as well as propose methods for overcoming these issues.
- We find that strong in-distribution RE performance often masks fundamental generalisation failures, with models that excel on intradataset evaluations frequently failing to transfer effectively (§ 5).
- Our cross-dataset analysis shows that data quality dictates the best RE adaptation method: while fine-tuning achieves superior cross-dataset performance on clean data, fewshot in-context learning (ICL) better handles

¹We will release the code upon publication.

161

162

163

164

165

167

168

169

170

171

172

173

174

175

126

127

128

129

130

noisy data. However, zero-shot prompting outperforms all cross-dataset methods in some settings (§ 5.2).

> • We identify structural issues in current RE benchmarks that lead to observed generalisation errors, including single-relation constraints, reliance on external knowledge, and coverage biases (§ 6).

These findings reveal that while current RE systems achieve high in-distribution results, their cross-dataset performance shows critical gaps in genuine relation understanding, limiting their applicability in real-world situations.

2 Related Work

078

090

095

100

101

102

103

104

105

106

109

110

111

112

113

114

115

116

2.1 Approaches for RE

RE is traditionally framed as a classification task, tackled via either a pipeline approach—where subtasks like named entity recognition (NER), coreference resolution, and relation classification (RC) are performed sequentially—or a joint model that processes them simultaneously (Taillé et al., 2020; Bassignana and Plank, 2022b; Saini et al., 2023). It is further categorised into sentence- (Alt et al., 2020; Plum et al., 2022) and document-level RE (Yao et al., 2019; Meng et al., 2024).

Since the introduction of BERT (Devlin et al., 2019), encoder-based models have dominated RE due to their bidirectional attention mechanism, which effectively captures context for classification tasks (Alt et al., 2020; Plum et al., 2022). However, the rise of autoregressive models has led to increasing adoption of decoder-based architectures to RE (Wang et al., 2022; Sun et al., 2023; Xu et al., 2023; Liu et al., 2024; Efeoglu and Paschke, 2024a). While encoder-decoder models have been explored (Huguet Cabot and Navigli, 2021; Li et al., 2023b), our experiments focus on the dominant encoderonly and decoder-only architectures for RE.

2.2 Generalisation Capabilities of RE Models

Recent work advocates for transparent evaluation 117 (Neubig et al., 2019; Liu et al., 2021) and OOD 118 testing (Linzen, 2020; Allen-Zhu and Li, 2024; Qi 119 et al., 2023) to assess model robustness. Com-120 121 mon strategies include cross-dataset (Antypas and Camacho-Collados, 2023; Jang and Frassinelli, 122 2024) and cross-domain (Fu et al., 2017; Liu et al., 123 2020; Bassignana and Plank, 2022a; Calderon et al., 124 2024) experiments, as well as testing on perturbed 125

and adversarially modified sets (Wu et al., 2019; Gardner et al., 2020; Goel et al., 2021; Rusert et al., 2022).

Recent studies have explored various ways to improve RE model robustness. Bassignana and Plank (2022a) introduce a cross-domain RE dataset with broad relation types, while Meng et al. (2024) and Chen et al. (2023) evaluate state-of-the-art (SOTA) document-level RE models on perturbed test sets. Chen et al. (2023) further reveal that even when models predict correctly, they often rely on spurious correlations, highlighting their vulnerability to minor evaluation shifts. To reduce dependence on mere pattern matching, Allen-Zhu and Li (2024) propose augmenting training data with synthetic samples reformulated by an auxiliary model.

3 Methodology

We assess the robustness of RE systems by evaluating their performance on OOD data. Standard evaluations on in-distribution test sets likely overestimate RE performance (Linzen, 2020), as models can exploit spurious cues rather than learning genuine relational patterns (Chen et al., 2023; Meng et al., 2024; Arzt and Hanbury, 2024).

To systematically evaluate generalisation capabilities of RE models, we conduct both intra- and cross-dataset experiments. The intra-dataset experiments act as a control, evaluating RE models on data drawn from the distribution used for model adaptation, while the cross-dataset experiments measure model robustness with OOD test sets derived from a different RE dataset.

For our experiments, we consider three sentencelevel RE datasets: TACRED-RE (Alt et al., 2020), NYT (Riedel et al., 2010), and Biographical (Plum et al., 2022). While TACRED-RE and NYT are general-purpose RE datasets, we focus exclusively on biographical relations, or relations that describe aspects of an individual's life (e.g., *place_of_birth*, studied at, or children). Our focus on biographical RE is motivated by two key factors. First, the relation type overlap between TACRED-RE and NYT is limited to two non-biographical relations (compared to the six overlapping biographical relations). In addition, focusing on biographical relations allows for additional cross-dataset evaluations with the Biographical dataset, which only contains biographical relations. This setup thus allows us to directly compare the generalisation of two popular RE datasets in a third, held-out evaluation setting.

3.1 Data

176

177

179

180

181

183

185

186

188

192

193

194

195

196

197

198

201

202

206

We now briefly describe three RE datasets used.

TACRED-RE (Alt et al., 2020) is a generalpurpose RE dataset with 41 relation types and a 'no_relation' class.² It contains over 106k instances but is highly imbalanced, with $\sim 80\%$ labeled as 'no_relation'. Built from English newswire and web text, it is a revisited version of TACRED (Zhang et al., 2017), where challenging samples were re-annotated by professional annotators to reduce noise from crowdsourced labels. Experimental results demonstrate improved performance on TACRED-RE compared to TACRED (see Tables 11 and 12 in the Appendix), leading to its use in our cross-dataset experiments. Figure 1 shows a TACRED-RE example. We focus on 26 biographical relations from TACRED-RE, including the 'no_relation' class (Table 4, Appendix).

> I I Irene Morgan Kirkaldy, who was born and reared in Baltimore, lived on Long Island and ran a child-care center in Queens with her second husband, Stanley Kirkady.

per:city_of_birth

Figure 1: Example from TACRED-RE dataset. This annotation example is from Zhang et al. (2017).

NYT (Riedel et al., 2010) is a general-purpose RE dataset comprising 24 relations and a 'None' class. The dataset contains over 266k sentences, with 64% of the instances labeled as 'None' and half of positive instances containing a single dominant relation, '/location/location/contains'.³ The NYT dataset was constructed via distant supervision, by applying Freebase (Bollacker et al., 2008) as external supervision on the text of New York Times articles (Sandhaus, 2008). Figure 2 shows an NYT example. We focus on a subset of the NYT dataset with 7 biographical relations, including a 'None' class (Table 5, Appendix).

"I think Amlo truly feels he's the Redeemer of Mexico, but his reign is of
this world," said Enrique Krauze, a historian.
/people/person/nationality

Figure 2: Example from the NYT dataset

Biographical is an RE dataset for the biographical text domain, with 10 relation types (Plum et al.,

2022). Built from Wikipedia articles on prominent individuals and containing 346,257 instances, Biographical was created using a semi-supervised approach.⁴ Named entities were automatically extracted from text using spaCy (Honnibal et al., 2020) and Stanford CoreNLP (Manning et al., 2014). Subsequently, Wikipedia sentences containing the extracted named entities were matched with Pantheon and Wikidata to automatically infer relations between the entities. Figure 3 shows an example⁵ from Biographical. The statistics for Biographical, downsampled to match the size of the TACRED-RE and NYT subsets, are presented in Table 5 in the Appendix. 209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

228

229

230

232

233

234

235

236

237

238

239

240

241

242

243

245

246

247





3.2 Cross-Dataset Comparison: Challenges

Single Relation per Sample: Both TACRED-RE and Biographical restrict each sample to at most two named entities and one relation, even when multiple relations exist within a sentence. For instance, the TACRED-RE example in Figure 1 is labeled with 'per:city_of_birth' but also contains relations like 'per:stateorprovinces_of_residence', 'per:employee_of', and 'per:spouse', all within TACRED-RE's relation set. This constraint may potentially confuse a model trained on such data, as it enforces a single-label assignment per sentence.

While NYT better reflects real-world scenarios by allowing multiple relations per sentence, we filtered it to two-entity, single-relation samples for fair cross-dataset comparison, retaining only those like the example in Figure 2.

Unclear 'Negative' Class: Clear negative samples—instances with entities but no meaningful relation—are crucial for RE systems. While TACRED-RE has an explicit 'no_relation' class, NYT's 'None' class lacks clear definition (Riedel et al., 2010), potentially confusing models about whether it indicates absence of predefined relations or any relation. Similarly, the Biographi-

²TACRED, and therefore TACRED-RE, are licensed by the Linguistic Data Consortium (LDC).

³The dataset is publicly available and accessible at https: //github.com/INK-USC/ReQuest.

⁴It has multiple versions. We use the version without coreference resolution or first Wikipedia sentence skipping (referred to as *m2_normal_final1* in Plum et al. (2022)).

⁵Sample ID 'mS7/1269356' in Biographical

cal dataset lacks an explicit negative class, using
instead an 'Other' class for unspecified relations
(Plum et al., 2022), as shown in Figure 3. This inconsistency between choosing a 'no_relation' versus a 'none_of_the_above' class in RE benchmarks
highlights the general challenge of consistently
defining the boundaries between presence and absence of semantic relations in text (Bassignana and
Plank, 2022b).

Expected Factual Knowledge: The design of RE datasets influences whether models genuinely learn RE or rely on dataset-specific cues. NYT's 259 distant supervision approach incorporates Freebasederived relations not stated in text requiring exter-261 nal world knowledge rather than textual evidence, as shown in Figure 2 where the text lacks explicit 263 information about Enrique Krause's nationality-264 such annotations extend beyond RE's scope and 265 corrupt models trained on such data. Similarly, although manually curated, TACRED-RE encom-267 passes relations like 'per:city_of_birth', which re-268 269 quire factual knowledge from a model, limiting generalisation to instances seen during adaptation. 270

3.3 Cross-Dataset Label Overlap

271

273

275

276

278

279

283

284

290

296

To enable cross-dataset evaluation, a manual label mapping was conducted across TACRED-RE, NYT, and Biographical. Table 7 shows six overlapping biographical relations between NYT and TACRED-RE, with twelve fine-grained TACRED-RE relations mapping to six broader NYT labels (e.g., NYT's 'place_of_birth' encompasses three TACRED-RE location-specific birth relations).

Treating Biographical's 'Other' class as negative—supported by manual analysis of 30 random instances revealing typically negative and not unspecified relations—we find four overlapping relations across all three datasets (Table 9, Appendix). The NYT-Biographical overlap includes these same four relations (Table 6, Appendix), while TACRED-RE and Biographical share nine relations (Table 8, Appendix).

4 Experiments

4.1 Data Format and Standardisation

To facilitate cross-dataset evaluations and focus our experiments on *relation classification*, we standardise our data into a unified format providing the entity spans in each input. Thus, each example's entities are marked with the tags <e1>head entity</e1> and <e2>tail entity</e2>, respectively. To address class imbalance, negative instances were randomly downsampled across all three datasets to balance the number of positive and negative instances, and Biographical (~350,000 instances) was downsampled to match the size of other biographical subsets for computational efficiency (Appendix Tables 3, 4, 5). Given TACRED-RE's fine-grained relations, we mapped them to broader NYT labels during cross-dataset evaluation (see Appendix C for details). 297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

341

342

343

344

345

346

347

4.2 Model Selection, Training, and Evaluation

We consider two types of models: an encoder-only model (DeBERTa-v3-large 304M; He et al. (2021)) and a decoder-only model (an instruction-tuned LLaMA 3.1 8B model; Grattafiori et al. (2024)). We do not include SOTA systems for each dataset; SOTA approaches for the three datasets in question differ–with most SOTA systems for one dataset not evaluated on the others (Orlando et al., 2024; Efeoglu and Paschke, 2024b; Sainz et al., 2024)– and we focus our experiments on the biographical subset of relations from TACRED-RE and NYT. Thus, our results are not directly comparable to prior work on these datasets.

We employ two commonly used model adaptation strategies for RE: fine-tuning and in-context learning (ICL). Specifically, we consider direct finetuning with DeBERTa, fine-tuning LLaMA using low-rank adaptation (LoRA; Hu et al., 2022), and zero-shot and five-shot ICL (Brown et al., 2020) with LLaMA. For few-shot ICL, we perform five runs with different demonstration sets to account for demonstration sensitivity (Zhang et al., 2022; Webson and Pavlick, 2022; Lu et al., 2022). However, due to computational constraints, fine-tuning experiments are limited to a single run. For NYT and TACRED-RE, we conduct experiments in two adaptation settings: adaptation on all biographical relations in each dataset (Appendix Tables 4 and 5) and adaptation on only overlapping relations (Appendix Table 7). This applies to both fine-tuning and ICL, where zero-shot prompts and few-shot demonstrations are selected accordingly.

We then perform two types of evaluations: **intra-dataset**, where models are evaluated on the same dataset they were adapted to; and **cross-dataset**, where the adapted models were tested on OOD data to assess their generalisation capabilities. Further implementation details, including hyperparameter settings, full label sets, and prompting details, are provided in Appendix D.

5 Results

354

367

371

372

373

374

377

391

394

Overview of Reported Results: Table 1 reports intra- and cross-dataset results for NYT and TACRED-RE on six overlapping relations, using models adapted on all biographical relations. For TACRED-RE, which maps its 12 fine-grained labels to NYT's shared label space, we report both dataset-specific and shared label results.

Table 2 shows model generalisation to Biographical across three overlapping relation sets: (1) four relations shared across all datasets, (2) same four relations shared between NYT and Biographical, and (3) nine relations shared between TACRED-RE and Biographical. Models were adapted on each dataset's full biographical relations, with Biographical's intra-dataset results for comparison. We focus on results with models adapted on the full overlap, as they show similar performance to those adapted only on overlap (Table 14, Appendix) while better reflecting real-world scenarios.

While we primarily focus on macro F1 to address class imbalance, additional experimental results, including per-class performance breakdowns, are provided in Appendix F.

5.1 Intra-Dataset Results

We evaluate our RE models on their training data distribution for comparison of cross-dataset generalisation. Unsurprisingly, we find that fine-tuning performs best for intra-dataset evaluations: fine-tuned LLaMA outperforms DeBERTa on TACRED-RE and NYT (Table 1), while the two models perform nearly identically on the full label overlap when fine-tuned on Biographical (Table 2).⁶

Our ICL experiments similarly show expected results, with the five-shot prompting moderately outperforming zero-shot prompting but underperforming full model fine-tuning. In the case of Biographical, this few-shot ICL gain over zero-shot is significant, increasing from 0.24 to 0.53 ± 0.05 with five demonstrations (Table 2).

We also note different performance trends within intra-dataset evaluations of TACRED-RE and NYT. Specifically, we find while fine-tuning yields higher intra-dataset performance on NYT than on TACRED-RE, this performance trend flips for the zero- and few-shot ICL settings, with prompting on NYT performing significantly *worse* than TACRED-RE despite TACRED-RE's finer-grained relation schema. This difference likely stems from differing data quality between datasets: the noisy labeling during NYT creation (Yaghoobzadeh et al., 2017) likely leads to over-fitting during fine-tuning (Tänzer et al., 2022) (rather than learning robust relational patterns), but harms model generalisation to NYT when not fine-tuned for that data distribution. 395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

Comparison with SOTA: As our generalisation analysis focuses on biographical relations across datasets (and prior works rarely report per-class results), we cannot directly compare against SOTA models on the full datasets. However, our best intradataset results on biographical subsets (with finetuned LLaMA) remain close to the overall SOTA performance, trailing the reported scores by 1-2 F1 points on all three datasets (Han et al., 2022; Plum et al., 2022; Orlando et al., 2024).

5.2 Cross-Dataset Results

We now turn to examining the cross-dataset generalisation of our RE systems. Unsurprisingly, we find that performance almost always declines with cross-dataset evaluations. However, models adapted on TACRED-RE exhibit relatively strong generalisation capabilities—the few exceptions of better cross-dataset performance stemming from TACRED-RE models applied to the Biographical dataset—while those adapted on NYT struggle to transfer effectively, likely due to dataset noise.

RE Models Struggle to Generalise across Datasets Cross-dataset evaluations (almost) always perform worse than the comparable intradataset experiment: NYT and TACRED-RE show substantial drops of 25-27 points, while Biographical exhibits smaller decreases of 4-10 points depending on the relation setting (Tables 1 and 2). We also observe somewhat different performance trends across model and adaptation approaches from the intra-dataset experiments; while finetuning LLaMA with TACRED-RE achieves the best cross-dataset performance on NYT, the best TACRED-RE cross-dataset results are obtained using few-shot ICL with NYT demonstrations (rather than fine-tuning). However, these results remain below the zero-shot TACRED-RE baseline.

The cross-dataset experiments on Biographical similarly perform worse than the corresponding intra-dataset experiments in most settings (Table 2); one notable exception is LLaMA prompted

⁶This may indicate a performance ceiling due to data noise limiting further improvement (Alt et al., 2020; Arzt and Hanbury, 2024).

Model	Setting	Dataset	Intra-Dataset		Cross-Dataset	
			Shared Labels	Dataset Labels	NYT	TACRED-RE
DeBERTa-v3 large 304M	Fine-tuned on	NYT	0.67	0.67	-	0.27
		TACRED-RE	0.66	0.64	0.50	-
LLaMA 3.1 8B	Fine-tuned on	NYT	0.87	0.87	-	0.45
		TACRED-RE	0.79	0.76	0.62	_
LLaMA 3.1 8B	Zero-Shot	NYT	0.31	0.31	_	_
		TACRED-RE	0.58	0.37	-	_
LLaMA 3.1 8B	5-Shot	NYT	0.45 ± 0.07	0.45 ± 0.07	-	0.52 ± 0.06
		TACRED-RE	0.63 ± 0.06	0.43 ± 0.07	0.39 ± 0.02	

Table 1: Macro F1-scores for intr	a- and cross-dataset	predictions on six	overlapping relations.	Results show both
shared and dataset-specific labels	, with models adapte	d on all biographic	cal relations through fir	ne-tuning or ICL.

Model	Setting	Dataset	Full Overlap	Overlap w. NYT	Overlap w. TACRED-RE
DeBERTa-v3-large 304M	Fine-tuned on	NYT	0.47	0.47	-
		TACRED-RE	0.60	_	0.68
		Biographical	0.79	0.79	0.71
LLaMA 3.1 8B	Fine-tuned on	NYT	0.30	0.30	_
		TACRED-RE	0.69	_	0.70
		Biographical	0.79	0.79	0.74
LLaMA 3.1 8B	Zero-Shot	Biographical	0.24	0.24	0.35
LLaMA 3.1 8B	5-Shot	NYT	0.48 ± 0.04	0.48 ± 0.04	-
		TACRED-RE	0.51 ± 0.04	-	0.58 ± 0.02
		Biographical	0.53 ± 0.05	0.53 ± 0.05	0.54 ± 0.03

Table 2: Evaluation on Biographical Dataset (macro F1-scores). Models adapted on all biographical relations through fine-tuning or ICL.

with five TACRED-RE examples, which outperforms the intra-dataset few-shot experiments on the TACRED-RE/Biographical label overlap. The best Biographical cross-dataset results are achieved with fine-tuning LLaMA on TACRED-RE, though this still underperforms intra-dataset fine-tuning.

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

NYT Models Generalise Worse than TACRED-RE Models NYT-adapted models exhibit significantly poorer generalisation capabilities than those adapted on TACRED-RE (Tables 1 and 2). For example, fine-tuning LLaMA with TACRED-RE surpasses zero- and cross-dataset ICL on NYT (by ~30 and ~20 points, respectively), whereas all cross-dataset experiments transferring from NYT to TACRED-RE underperform zero-shot evaluations with no cross-dataset signal. This is also clear from the evaluations on the held-out Biographical dataset, where transferring from TACRED-RE always performs better than NYT (and occasionally outperforms the intra-dataset performance).

This performance gap is unlikely due to domain differences, as both datasets contain newspaper articles (with TACRED-RE including some NYT newspaper content without instance overlap), while Biographical covers Wikipedia articles. Rather, we attribute it to NYT's distant supervision annotations, which introduce noise and limit model robustness. This is likely why LLaMA fine-tuned on NYT and evaluated on Biographical data (0.30) underperforms DeBERTa (0.47; Table 2)—the overparametrised LLaMA exhibits stronger overfitting to NYT noise and generalises poorly to unseen data, a phenomenon also highlighted by Liu et al. (2022) with corrupted training data. 472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

Effect of Adaptation Strategy on Generalisation While prior work suggests that ICL often generalises more effectively to OOD data than finetuning (Awadalla et al., 2022; Song et al., 2023; Si et al., 2023), our results indicate this advantage depends heavily on data quality. With highquality data like TACRED-RE, fine-tuning consistently achieves the best cross-dataset performance, surpassing few-shot ICL on both NYT and Biographical evaluations. In fact, TACRED-RE adaptations can even perform comparably to intra-dataset ones: DeBERTa (0.68) and LLaMA (0.70) finetuned on TACRED-RE achieve similar results to their Biographical intra-dataset performance (0.71)and 0.74) on the TACRED-RE/Biographical label overlap (Table 2).

However, when adaptation data are noisy, as with NYT, few-shot ICL becomes a more effective strategy: few-shot ICL via NYT consistently performs better than fine-tuning on NYT for both TACRED-



LLAMA3.1 Fine-tuned on: (a) TACRED-RE (b) NYT

Figure 4: Confusion matrices comparing cross-dataset results for LLAMA fine-tuned on TACRED-RE/NYT

RE and Biographical evaluations (Tables 1, 2). This is likely because ICL limits the signal from noisy training data, which in turn reduces the overfitting to dataset-specific artefacts and catastrophic forgetting compared to fine-tuning (Tran et al., 2024; Anonymous, 2024; Kotha et al., 2024).

499

500

501

502

503

504

509

510

511

512

513

514

516

517

518

519

521

522

524

6 Analysing RE Generalisation Failure Cases

Given RE benchmarks' numerous relations and class imbalance, we analyse the strongest performing model, fine-tuned LLaMA, beyond aggregated metrics. We examine per-relation performance (Figure 4) and qualitatively analyse 30 random misclassifications from four evaluation settings to identify their likely underlying causes.⁷ Through these analyses, we find the following causes of RE generalisation mistakes:

Effect of Noisy Supervision on Generalisation Confusion matrices in Figure 4 reveal that NYTadapted models systematically overpredict the 'None' class on both TACRED-RE and Biographical, with manual analysis showing these misclassifications stem primarily from NYT's *distant supervision* (Table 21, Appendix) rather than vocabulary differences (Figure 5) or domain shift. This issue is particularly evident in NYT's location-based relations, where reliance on external knowledge leads to conflicting annotations that hinder pattern learning—for instance, "Henryk Tomaszewski [...] died on Sunday at his home in Warsaw"⁸ is labeled as birthplace despite clear evidence of death location. Similar issues arise with Biographical's *semi-supervised* data, where models adapted on cleaner datasets like TACRED-RE fail to replicate ground truth labels that lack textual evidence. Notably, despite higher lexical overlap between Biographical and NYT (Figure 6, Appendix), TACRED-RE-adapted models perform better on Biographical, indicating that adaptation data quality matters more than lexical similarity. 525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

556

557

558

559

560

561

562

563

564

565

566

567

569

570

571

572

Single Relation Constraint & Negative Class The issue extends beyond noisy supervision to fundamental constraints in RE task design. Models adapted on biographical relations often detect valid but unlabeled relations (marked as 'new_relation' in Figure 4), highlighting limitations of enforcing single-relation per sample. This manifests in cases like "Gross, who is himself Jewish [...] was sent to Cuba"⁹, where only 'place_lived' is labeled while 'religion' is omitted due to TACRED-RE's constraint. Also, unclear or missing negative class affects cross-dataset evaluation: while Biographical's 'Other' class is intended to cover undefined relations, our analysis reveals it contains both instances without meaningful relations and those with valid but undefined relations. This explains the high frequency of 'new_relation' predictions for Biographical when using LLaMA fine-tuned on TACRED-RE with their finer-grained relation schema (Figure 4) and also highlights the fundamental difficulty in defining boundaries between presence and absence of semantic relations. Despite this ambiguity, the class achieves strong cross-dataset performance when mapped to 'None' (0.78 and 0.72)for LLaMA fine-tuned on TACRED-RE and NYT), further underscoring the importance of distinguishing between 'no_relation' and 'none_of_the_above' cases in RE (Bassignana and Plank, 2022b).

Reliance on External Knowledge even in Manually Curated Datasets Even with high-quality manual annotations, RE often requires external knowledge and complex reasoning capabilities. Our analysis reveals this challenge manifests in two key ways: through implicit relations requiring

⁷More misclassification patterns in Appendix Table 21.

⁸NYT instance ID: '/m/vinci8/data1/riedel/projects/relation /kb/nyt1/docstore/nyt-2005-2006.backup/1701917.xml.pb'

⁹TACRED-RE instance ID: '098f6f318be29eddb182'



Figure 5: Vocabulary Overlap (%) per overlapping relation between NYT and TACRED-RE

inference, and through necessary world knowledge for entity interpretation. For example, in "Gross [...] was sent to Cuba as a spy"⁹, the NYT-adapted model predicts 'None' instead of 'place_lived', failing to infer that being sent somewhere as a spy implies residence. While detecting implicit relations is crucial (Geva et al., 2021), ensuring consistent and objective interpretation remains challenging.

573

576

579

581

584

585

586

588

589

590

591

593

594

598

607

Beyond implicit relations, models must also rely on world knowledge for basic entity understandingas in cases like 'Idaho businesswoman'¹⁰, where identifying entity types requires knowing Idaho as a location. TACRED-RE fine-grained relation schema further demonstrates this issue, where even with world knowledge, distinguishing between relations like city of birth and state/province of birth can be ambiguous (e.g., whether New York refers to the city or state). As Chen et al. (2023) note, even human annotators tend to rely on such prior knowledge despite the lack of rationales, motivating the need for finer-grained word evidence annotation.

Dataset Composition and Coverage Biases Analysis of the most shared words across NYT and TACRED-RE demonstrates a strong US-centric coverage bias, likely limiting generalisation to non-US contexts (Appendix Table 22). NYT also exhibits topical skews in specific relations, such as religion being predominantly associated with Islam, potentially leading models to learn narrow, biased representations of relations.

Analysing part-of-speech distributions also reveals distinct patterns across all three datasets (Appendix Table 20). While proper nouns dominate head and tail entities in all datasets (reaching nearly 100% in NYT), TACRED-RE shows more linguis-

¹⁰TACRED-RE instance ID: '098f6bd9fa786293e49d'

tic diversity with 17% of head entities as pronouns and 17% of tail entities as common nouns. Biographical, sourced from Wikipedia, contains a high proportion (26%) of numerical tail entities, primarily dates. These compositional differences, along with TACRED-RE's longer, compound sentences and higher average entity distance (\sim 12 tokens vs NYT's \sim 8 tokens), most likely impact crossdataset performance; NYT-adapted models struggle with these more complex patterns, which are absent from their training data (Appendix Table 21).

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

7 Conclusion

This work examines cross-dataset generalisation in language model-based RE systems in biographical settings. We find RE models struggle to generalise even within similar domains, with high intradataset performance often masking spurious overfitting rather than genuine learning of relational patterns. Furthermore, data quality is crucial for robust transfer (with the optimal adaptation strategy depending on data quality); fine-tuning yields the best cross-dataset performance when high-quality data is available, but few-shot ICL performs better in settings with noisy data. However, in some cases, a zero-shot baseline surpasses all cross-dataset results, further underscoring the limitations of current RE systems.

Our analysis also reveals several structural issues in all current RE benchmarks: (1) single-relation constraints that ignore other valid relations between entities in text, (2) the lack of a well-defined negative class with challenging samples (e.g., sentences containing commonly used tokens for relations like 'born' or 'died') to enforce deeper semantic understanding beyond pattern matching, and (3) limited diversity in data sources. These issues, compounded by inconsistent relation definitions and limited overlap across datasets, hinder meaningful evaluations of RE generalisation.

These findings thus highlight the need for more transparent evaluation beyond in-distribution testing and aggregated metrics, as limiting evaluation to these may not reflect genuine improvements in capturing relational patterns or account for class imbalance and the large number of relations in RE benchmarks. We see many promising directions for future work, including testing RE robustness on perturbed evaluation sets and applying interpretability methods to better understand how models infer relational knowledge.

679

694

700

701

702

703

704

707

Limitations

Our cross-dataset analysis is limited to a particular set of biographical relations but reflects a broader 660 challenge in RE evaluation where datasets, even covering the same domain, typically share a small relation overlap. We also constrain our analysis to single-relation examples: while, real-world scenarios often involve multiple relations per instance (and NYT allows multiple relations), we focused 666 on single-relation setting for fair-cross dataset comparison, as TACRED-RE and Biographical are annotated with single relations. Similarly, we exclusively evaluate relation classification (RC) due to dataset constraints: TACRED-RE and Biographical 671 assume a single relation triple per sentence, unlike 672 real-world text where multiple relations can coexist. 673 By focusing on RC with entity tags as guidance, we aim to minimise the prediction of other potential relations present in a sentence, but not between the 676 specified entities. 677

> The adaptation sets we used contain a large class imbalance due to the underlying distributions of the datasets, even after we perform data rebalancing. While this could be viewed as a limitation, it reflects real-world scenarios where models must adapt with limited training data (Bassignana and Plank, 2022a). Finally, we note that our intradataset results cannot be directly compared with reported SOTA performance, as most papers lack detailed relation-based metrics, reporting only aggregated results.

References

- Zeyuan Allen-Zhu and Yuanzhi Li. 2024. Physics of language models: Part 3.1, knowledge storage and extraction. *Preprint*, arXiv:2309.14316.
- Christoph Alt, Aleksandra Gabryszak, and Leonhard Hennig. 2020. TACRED revisited: A thorough evaluation of the TACRED relation extraction task. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1558– 1569, Online. Association for Computational Linguistics.
- Anonymous. 2024. Analyzing and reducing catastrophic forgetting in parameter efficient tuning. In *Submitted to ACL Rolling Review - April 2024*. Under review.
- Dimosthenis Antypas and Jose Camacho-Collados. 2023. Robust hate speech detection in social media: A cross-dataset empirical evaluation. In *The 7th Workshop on Online Abuse and Harms (WOAH)*,

pages 231–242, Toronto, Canada. Association for Computational Linguistics.

- Varvara Arzt and Allan Hanbury. 2024. Beyond the numbers: Transparency in relation extraction benchmark creation and leaderboards. In *Proceedings of the 2nd GenBench Workshop on Generalisation (Benchmarking) in NLP*, pages 120–130, Miami, Florida, USA. Association for Computational Linguistics.
- Anas Awadalla, Mitchell Wortsman, Gabriel Ilharco, Sewon Min, Ian Magnusson, Hannaneh Hajishirzi, and Ludwig Schmidt. 2022. Exploring the landscape of distributional robustness for question answering models. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 5971–5987, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Elisa Bassignana and Barbara Plank. 2022a. CrossRE: A cross-domain dataset for relation extraction. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 3592–3604, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Elisa Bassignana and Barbara Plank. 2022b. What do you mean by relation extraction? a survey on datasets and study on scientific relation classification. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop*, pages 67–83, Dublin, Ireland. Association for Computational Linguistics.
- Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data*, SIGMOD '08, page 1247–1250, New York, NY, USA. Association for Computing Machinery.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA. Curran Associates Inc.
- Nitay Calderon, Naveh Porat, Eyal Ben-David, Alexander Chapanin, Zorik Gekhman, Nadav Oved, Vitaly Shalumov, and Roi Reichart. 2024. Measuring the robustness of NLP models to domain shifts. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 126–154, Miami, Florida, USA. Association for Computational Linguistics.

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

Haotian Chen, Bingsheng Chen, and Xiangdong Zhou.
2023. Did the models understand documents? benchmarking models for language understanding in document-level relation extraction. In *Proceedings* of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 6418–6435, Toronto, Canada. Association for Computational Linguistics.

766

767

771

775

776

777

778

786

790

792

793

794

795

796

797

803

804

807

809

810

811

812 813

814

815

816

817

818

819

821

822

- 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.
 - Sefika Efeoglu and Adrian Paschke. 2024a. Relation extraction with fine-tuned large language models in retrieval augmented generation frameworks. *Preprint*, arXiv:2406.14745.
- Sefika Efeoglu and Adrian Paschke. 2024b. Retrievalaugmented generation-based relation extraction. *Preprint*, arXiv:2404.13397.
- Lisheng Fu, Thien Huu Nguyen, Bonan Min, and Ralph Grishman. 2017. Domain adaptation for relation extraction with domain adversarial neural network. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 425–429, Taipei, Taiwan. Asian Federation of Natural Language Processing.
- Matt Gardner, Yoav Artzi, Victoria Basmov, Jonathan Berant, Ben Bogin, Sihao Chen, Pradeep Dasigi, Dheeru Dua, Yanai Elazar, Ananth Gottumukkala, Nitish Gupta, Hannaneh Hajishirzi, Gabriel Ilharco, Daniel Khashabi, Kevin Lin, Jiangming Liu, Nelson F. Liu, Phoebe Mulcaire, Qiang Ning, Sameer Singh, Noah A. Smith, Sanjay Subramanian, Reut Tsarfaty, Eric Wallace, Ally Zhang, and Ben Zhou. 2020. Evaluating models' local decision boundaries via contrast sets. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1307–1323, Online. Association for Computational Linguistics.
- Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. 2023. Dissecting recall of factual associations in auto-regressive language models. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 12216–12235, Singapore. Association for Computational Linguistics.
- Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of the Association for Computational Linguistics*, 9:346– 361.

Karan Goel, Nazneen Fatema Rajani, Jesse Vig, Zachary Taschdjian, Mohit Bansal, and Christopher Ré. 2021. Robustness gym: Unifying the NLP evaluation landscape. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Demonstrations, pages 42–55, Online. Association for Computational Linguistics. 823

824

825

826

827

829

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, and et al. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.
- Jiale Han, Shuai Zhao, Bo Cheng, Shengkun Ma, and Wei Lu. 2022. Generative prompt tuning for relation classification. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 3170–3185, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. Deberta: Decoding-enhanced bert with disentangled attention. In *International Conference on Learning Representations*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Evan Hernandez, Arnab Sen Sharma, Tal Haklay, Kevin Meng, Martin Wattenberg, Jacob Andreas, Yonatan Belinkov, and David Bau. 2024. Linearity of relation decoding in transformer language models. In *The Twelfth International Conference on Learning Representations*.
- Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. spaCy: Industrialstrength Natural Language Processing in Python.
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.
- Pere-Lluís Huguet Cabot and Roberto Navigli. 2021. REBEL: Relation extraction by end-to-end language generation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2370– 2381, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Hyewon Jang and Diego Frassinelli. 2024. Generalizable sarcasm detection is just around the corner, of course! In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 4238–4249, Mexico City, Mexico. Association for Computational Linguistics.
- Suhas Kotha, Jacob Mitchell Springer, and Aditi Raghunathan. 2024. Understanding catastrophic forgetting

993

994

936

in language models via implicit inference. In *The Twelfth International Conference on Learning Representations*.

879

897

900

901

902

903

905 906

907

908

909

910

911

912

913

914

915

916

917

918

919

923

924

926

927

930

931

932

- Venelin Kovatchev and Matthew Lease. 2024. Benchmark transparency: Measuring the impact of data on evaluation. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 1536–1551, Mexico City, Mexico. Association for Computational Linguistics.
- Alina Leidinger, Robert van Rooij, and Ekaterina Shutova. 2023. The language of prompting: What linguistic properties make a prompt successful? In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9210–9232, Singapore. Association for Computational Linguistics.
 - Guozheng Li, Peng Wang, and Wenjun Ke. 2023a. Revisiting large language models as zero-shot relation extractors. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 6877–6892, Singapore. Association for Computational Linguistics.
 - Shaoshuai Li, Wenbin Yu, Zijian Chen, and Yangyang Luo. 2023b. A joint entity and relation extraction model based on encoder-decoder. In 2023 IEEE 3rd International Conference on Information Technology, Big Data and Artificial Intelligence (ICIBA), volume 3, pages 996–1000.
- Tal Linzen. 2020. How can we accelerate progress towards human-like linguistic generalization? In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5210– 5217, Online. Association for Computational Linguistics.
- Pengfei Liu, Jinlan Fu, Yang Xiao, Weizhe Yuan, Shuaichen Chang, Junqi Dai, Yixin Liu, Zihuiwen Ye, and Graham Neubig. 2021. ExplainaBoard: An explainable leaderboard for NLP. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: System Demonstrations, pages 280–289, Online. Association for Computational Linguistics.
- Sheng Liu, Zhihui Zhu, Qing Qu, and Chong You. 2022. Robust training under label noise by overparameterization. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 14153–14172. PMLR.
- Siyi Liu, Yang Li, Jiang Li, Shan Yang, and Yunshi Lan.
 2024. Unleashing the power of large language models in zero-shot relation extraction via self-prompting.
 In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 13147–13161, Miami, Florida, USA. Association for Computational Linguistics.

- Zihan Liu, Yan Xu, Tiezheng Yu, Wenliang Dai, Ziwei Ji, Samuel Cahyawijaya, Andrea Madotto, and Pascale Fung. 2020. Crossner: Evaluating cross-domain named entity recognition. In *AAAI Conference on Artificial Intelligence*.
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. Fantastically ordered prompts and where to find them: Overcoming fewshot prompt order sensitivity. In *Proceedings of the* 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 8086–8098, Dublin, Ireland. Association for Computational Linguistics.
- Christopher Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven Bethard, and David McClosky. 2014. The Stanford CoreNLP natural language processing toolkit. In *Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 55–60, Baltimore, Maryland. Association for Computational Linguistics.
- Shiao Meng, Xuming Hu, Aiwei Liu, Fukun Ma, Yawen Yang, Shuang Li, and Lijie Wen. 2024. On the robustness of document-level relation extraction models to entity name variations. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 16362–16374, Bangkok, Thailand. Association for Computational Linguistics.
- Graham Neubig, Zi-Yi Dou, Junjie Hu, Paul Michel, Danish Pruthi, and Xinyi Wang. 2019. compare-mt: A tool for holistic comparison of language generation systems. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations), pages 35–41, Minneapolis, Minnesota. Association for Computational Linguistics.
- Riccardo Orlando, Pere-Lluís Huguet Cabot, Edoardo Barba, and Roberto Navigli. 2024. ReLiK: Retrieve and LinK, fast and accurate entity linking and relation extraction on an academic budget. In *Findings of the Association for Computational Linguistics: ACL* 2024, pages 14114–14132, Bangkok, Thailand. Association for Computational Linguistics.
- Alistair Plum, Tharindu Ranasinghe, Spencer Jones, Constantin Orasan, and Ruslan Mitkov. 2022. Biographical semi-supervised relation extraction dataset. In Proceedings of the 45th International ACM SI-GIR Conference on Research and Development in Information Retrieval, SIGIR '22, page 3121–3130, New York, NY, USA. Association for Computing Machinery.
- Ji Qi, Chuchun Zhang, Xiaozhi Wang, Kaisheng Zeng, Jifan Yu, Jinxin Liu, Lei Hou, Juanzi Li, and Xu Bin. 2023. Preserving knowledge invariance: Rethinking robustness evaluation of open information extraction. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5876–5890, Singapore. Association for Computational Linguistics.

Sebastian Riedel, Limin Yao, and Andrew McCallum. 2010. Modeling relations and their mentions without labeled text. In Machine Learning and Knowledge Discovery in Databases, pages 148-163, Berlin, Heidelberg. Springer Berlin Heidelberg.

995

997 998

999

1000

1001

1002

1005

1006

1007

1008

1010

1012

1013

1016

1018

1019

1020

1021

1022

1023

1024 1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

- Jonathan Rusert, Zubair Shafiq, and Padmini Srinivasan. 2022. On the robustness of offensive language classifiers. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 7424-7438, Dublin, Ireland. Association for Computational Linguistics.
- Pratik Saini, Samiran Pal, Tapas Nayak, and Indrajit Bhattacharya. 2023. 90% f1 score in relation triple extraction: Is it real? In Proceedings of the 1st GenBench Workshop on (Benchmarking) Generalisation in NLP, pages 1-11, Singapore. Association for Computational Linguistics.
- Oscar Sainz, Iker García-Ferrero, Rodrigo Agerri, Oier Lopez de Lacalle, German Rigau, and Eneko Agirre. 2024. GoLLIE: Annotation guidelines improve zero-shot information-extraction. In The Twelfth International Conference on Learning Representations.
- Evan Sandhaus. 2008. The new york times annotated corpus. Linguistic Data Consortium, Philadelphia, 6(12):e26752.
- Chenglei Si, Zhe Gan, Zhengyuan Yang, Shuohang Wang, Jianfeng Wang, Jordan Lee Boyd-Graber, and Lijuan Wang. 2023. Prompting GPT-3 to be reliable. In The Eleventh International Conference on Learning Representations.
- Yisheng Song, Ting Wang, Puyu Cai, Subrota K. Mondal, and Jyoti Prakash Sahoo. 2023. A comprehensive survey of few-shot learning: Evolution, applications, challenges, and opportunities. ACM Comput. Surv., 55(13s).
- Xiaofei Sun, Xiaoya Li, Jiwei Li, Fei Wu, Shangwei Guo, Tianwei Zhang, and Guoyin Wang. 2023. Text classification via large language models. In Findings of the Association for Computational Linguistics: EMNLP 2023, pages 8990–9005, Singapore. Association for Computational Linguistics.
- Bruno Taillé, Vincent Guigue, Geoffrey Scoutheeten, and Patrick Gallinari. 2020. Let's Stop Incorrect Comparisons in End-to-end Relation Extraction! In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 3689-3701, Online. Association for Computational Linguistics.
- Michael Tänzer, Sebastian Ruder, and Marek Rei. 2022. Memorisation versus generalisation in pre-trained language models. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 7564-7578, Dublin, Ireland. Association for Computational Linguistics.

Komal Teru. 2023. Semi-supervised relation extrac-1051 tion via data augmentation and consistency-training. 1052 In Proceedings of the 17th Conference of the European Chapter of the Association for Computational *Linguistics*, pages 1112–1124, Dubrovnik, Croatia. Association for Computational Linguistics.

1054

1055

1058

1059

1060

1061

1062

1064

1065

1066

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

1080

1081

1083

1084

1085

1087

1088

1089

1090

1091

1092

1094

1095

1096

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

- Quyen Tran, Nguyen Xuan Thanh, Nguyen Hoang Anh, Nam Le Hai, Trung Le, Linh Van Ngo, and Thien Huu Nguyen. 2024. Preserving generalization of language models in few-shot continual relation extraction. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 13771-13784, Miami, Florida, USA. Association for Computational Linguistics.
- Shubham Vatsal and Harsh Dubey. 2024. A survey of prompt engineering methods in large language models for different nlp tasks. Preprint, arXiv:2407.12994.
- Chenguang Wang, Xiao Liu, Zui Chen, Haoyun Hong, Jie Tang, and Dawn Song. 2022. DeepStruct: Pretraining of language models for structure prediction. In Findings of the Association for Computational Linguistics: ACL 2022, pages 803-823, Dublin, Ireland. Association for Computational Linguistics.
- Qing Wang, Kang Zhou, Qiao Qiao, Yuepei Li, and Qi Li. 2023. Improving unsupervised relation extraction by augmenting diverse sentence pairs. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 12136-12147, Singapore. Association for Computational Linguistics.
- Albert Webson and Ellie Pavlick. 2022. Do promptbased models really understand the meaning of their prompts? In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 2300-2344, Seattle, United States. Association for Computational Linguistics.
- Tongshuang Wu, Marco Tulio Ribeiro, Jeffrey Heer, and Daniel Weld. 2019. Errudite: Scalable, reproducible, and testable error analysis. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 747–763, Florence, Italy. Association for Computational Linguistics.
- Xin Xu, Yuqi Zhu, Xiaohan Wang, and Ningyu Zhang. 2023. How to unleash the power of large language models for few-shot relation extraction? In Proceedings of The Fourth Workshop on Simple and Efficient Natural Language Processing (SustaiNLP), pages 190-200, Toronto, Canada (Hybrid). Association for Computational Linguistics.
- Yadollah Yaghoobzadeh, Heike Adel, and Hinrich Schütze. 2017. Noise mitigation for neural entity typing and relation extraction. In Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers, pages 1183–1194, Valencia, Spain. Association for Computational Linguistics.

Yuan Yao, Deming Ye, Peng Li, Xu Han, Yankai Lin, 1109 1110 Zhenghao Liu, Zhiyuan Liu, Lixin Huang, Jie Zhou, and Maosong Sun. 2019. DocRED: A large-scale 1111 1112 document-level relation extraction dataset. In Pro-1113 ceedings of the 57th Annual Meeting of the Associa-1114 tion for Computational Linguistics, pages 764–777, Florence, Italy. Association for Computational Lin-1115 1116 guistics.

1117

1118

1119

1120

1121

1122

1123

1124

1125 1126

1127

1128

- Yiming Zhang, Shi Feng, and Chenhao Tan. 2022. Active example selection for in-context learning. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 9134–9148, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- 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.

A Dataset Statistics: Class Distribution

Table 3:	Balanced	Biographical	Dataset
		Ør	

Relation	# of Samples
Other	10,000
birthdate	2914
bplace_name	2845
dplace_name	1138
occupation	1105
deathdate	1011
parent	394
educatedAt	339
child	136
sibling	118
Positive Samples	10,000
Negative Samples	10,000
All	20,000

Table 4: Balanced TACRED-RE Subset with Biographi-cal Relations (26 relations)

Relation	# of Samples
no_relation	14,192
per:title	3805
per:employee_of	2104
per:age	818
per:countries_of_residence	695
per:cities_of_residence	596
per:origin	652
per:stateorprovinces_of_residence	444
per:spouse	463
per:date_of_death	343
per:children	347
per:cause_of_death	318
per:parents	282
per:charges	270
per:other_family	241
per:siblings	238
per:schools_attended	219
per:city_of_death	204
per:religion	145
per:alternate_names	132
per:city_of_birth	107
per:stateorprovince_of_death	100
per:date_of_birth	99
per:stateorprovince_of_birth	77
per:country_of_death	57
per:country_of_birth	45
Positive Samples	12,801
Negative Samples	14,192
All	26,993

Table 5: Balanced NYT Subset with Biographical Relations after removal of instances with multiple labels (7 relations)

Relation	# of Samples
None	5068
/people/person/nationality	2160
/people/person/place_lived	2016
/people/person/place_of_birth	437
/people/deceased_person/place_of_death	284
/people/person/children	147
/people/person/religion	24
Positive Samples	5068
Negative Samples	5068
All	10,136

B Relation Type Overlap

1131

Table 6: NYT/Biographical Relation Overlap

NYT	Biographical
/people/person/children	child
/people/person/place_of_birth	bplace_name
/people/deceased_person/place_of_death	dplace_name
None	Other

Table 7: NYT/TACRED-RE Relation Overlap

NYT	TACRED-RE
None	no_relation
/people/person/children	per:children
/people/person/religion	per:religion
/people/person/place_lived	per:stateorprovinces_of_residence per:countries_of_residence per:cities_of_residence
/people/person/place_of_birth	per:stateorprovince_of_birth per:country_of_birth per:city_of_birth
/people/deceased_person/place_of_death	per:stateorprovince_of_death per:country_of_death per:city_of_death

Table 8: Biographical/TACRED-RE Relation Overlap

Biographical	TACRED-RE
bplace_name	per:stateorprovince_of_birth per:country_of_birth per:city_of_birth
birthdate	per:date_of_birth
deathdate	per:date_of_death
parent	per:parents
educatedAt	per:schools_attended
dplace_name	per:stateorprovince_of_death per:country_of_death per:city_of_death
sibling	per:siblings
child	per:children
Other	no_relation

Table 9: Biographical/TACRED-RE/NYT Overlap

Biographical	TACRED-RE	NYT
child	per:children	/people/person/children
bplace_name	<pre>per:stateorprovince_of_birth per:country_of_birth per:city_of_birth</pre>	/people/person/place_of_birth
dplace_name	<pre>per:stateorprovince_of_death per:country_of_death per:city_of_death</pre>	/people/deceased_person/ place_of_death
Other	no_relation	None

1133

1134

1135

1136

1137

1138

1139

1140

1144

1145

1146

1147

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1168 1169

1179

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

C Data Formatting Details

TACRED-RE's fine-grained relations (e.g., "per:city_of_birth" and "per:country_of_birth") were mapped to broader categories (e.g., "place_of_birth") used in NYT and Biographical datasets, as shown in Tables 7 and 8. Cross-dataset results are reported using NYT label names (Table 15).

D Implementation Details

Listing 1: System prompt for LLaMA 3.1 8B zero-shot, few-shot, and fine-tuning experiments, shown here with NYT relation inventory

```
system_message = {
     'role":
               'system"
    "content":
                 (
              "You are an intelligent assistant
                    specializing in identifying
                    relations between entities in a
                    sentence
               Ouestion: What is the relation between
                    two tagged entities <e1>entity1</e1
                    > and <e2>entity2</e2> in the
                    following sentence? "
              "Choose one relation from the list: "
"['/people/person/children', '/people/
               'Choose one relation from
['/people/person/children', '/people/
(retionality', '/people/
                    person/nationality',
person/place_lived',
                '/people/person/place_of_birth', '/
                     people/deceased_person/
                     place_of_death',
                                          '/people/person/
                     religion',
               "'None'].
               'Rules: Select exactly one relation from
                     the list. If none of the listed
                    relations apply, select 'None
                'Output must strictly follow this
                     format: <relation_type>. Provide
                     no additional text or explanation
    )
}
```

We fine-tuned DeBERTa-v3-large¹¹ for 10 epochs employing early stopping. Following prior work in RE (Teru, 2023; Wang et al., 2023), we extended the deberta-v3-large tokeniser with entity marker tokens, namely, <e1> and </e1> for the head entity, and <e2> and </e2> for the tail entity. For LLAMA-3.1¹², we used LoRA fine-tuning (r = 8) over three epochs, applying it to attention and feedforward modules. Both models were fine-tuned using HuggingFace's Trainer class. For evaluation of LLaMA 3.1, predictions were considered correct only if they matched ground-truth labels exactly (Hendrycks et al., 2021).

For prompting, we used vanilla prompting vanilla (Li et al., 2023a; Vatsal and Dubey,



Figure 6: Vocabulary Overlap (%) Between Biographical and TACRED-RE/NYT Relations following lemmatisation and stop word removal using spaCy's en_core_web_trf

2024)and tested several RE-specific prompt designs (Leidinger et al., 2023; Li et al., 2023a; Efeoglu and Paschke, 2024a), given LLaMA's sensitivity to prompt formulation (Leidinger et al., 2023). The prompt 1 performed best and was used across all datasets with adapted label sets. Further prompt optimisation techniques were not considered, as they were beyond the scope of this paper. Hyperparameter settings for all experiments are detailed in Table 19.

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1206

Due to Biographical's ambiguous 'Other' class (Section 3), we use it only for evaluation, excluding these relations during adaptation.

All experiments with Deberta-v3-large were run on a single NVIDIA[®] TITAN RTX 24GB GPU; all experiments with Llama-3.1-8B-Instruct were run on a single NVIDIA[®] A100 80GB GPU. All experiments were performed with a fixed random seed for reproducibility.

E Vocabulary Overlap with Biographical

Figure 6 depicts vocabulary overlap between Bio-
graphical and TACRED-RE as well as Biographical1207and NYT per overlapping relation.1208

¹¹https://huggingface.co/microsoft/

deberta-v3-large

¹²https://huggingface.co/meta-llama/Llama-3.

¹⁻⁸B-Instruct

F Results

F.1 Intra-Dataset Results

Model	Deber	ta-v3-large	e 304M	Llama	-3.1 8B ze	ro-shot	Llan	na-3.1 8B 5	5-shot	Llama-3.1 8B fine-tuned		
	P	R	F1	P	R	F1	Р	R	F1	P	R	F1
/people/deceased person/place of death	0.81	0.76	0.78	0.93	0.82	0.87	0.75	0.09	0.16	0.93	0.82	0.87
/people/person/children	0.80	0.86	0.83	0.87	0.93	0.90	0.77	0.71	0.74	0.87	0.93	0.90
/people/person/nationality	0.95	1.00	0.97	0.99	0.99	0.99	0.96	0.10	0.18	0.99	0.99	0.99
/people/person/place_lived	0.88	0.91	0.90	0.84	0.95	0.89	0.36	0.58	0.44	0.84	0.95	0.89
/people/person/place_of_birth	0.54	0.56	0.55	0.88	0.46	0.60	0.07	0.06	0.07	0.88	0.46	0.60
/people/person/religion	0.00	0.00	0.00	1.00	1.00	1.00	0.83	1.00	0.91	1.00	1.00	1.00
None	0.98	0.95	0.96	0.97	0.97	0.97	0.62	0.76	0.68	0.97	0.97	0.97
macro avg	0.71	0.72	0.71	0.92	0.87	0.89	0.62	0.47	0.45	0.92	0.87	0.89
micro avg	0.92	0.92	0.92	0.30	0.22	0.25	0.52	0.52	0.52	0.94	0.94	0.94
weighted avg	0.92	0.92	0.92	0.94	0.94	0.94	0.63	0.52	0.47	0.94	0.94	0.94

Table 10: NYT Results

Model	Deber	ta-v3-large	304M	Llama	a-3.1 8B zei	ro-shot	Llan	na-3.1 8B 5	-shot	Llama	-3.1 8B fine	e-tuned
	Р	R	F1									
macro avg micro avg weighted avg	0.67 0.83 0.83	0.64 0.83 0.83	0.64 0.83 0.83	0.49 0.31 0.68	0.38 0.29 0.29	0.34 0.30 0.32	0.37 0.51 0.59	0.30 0.51 0.51	0.29 0.51 0.49	0.76 0.87 0.87	0.71 0.87 0.87	0.73 0.87 0.87

Model	DeBE	RTa-v3-la	rge 304M	Llama	-3.1 8B z	ero-shot	Llam	a-3.1 8B	5-shot	Llama	-3.1 8B fi	ne-tuned
	P	R	F1	P	R	F1	P	R	F1	P	R	F1
no_relation	0.96	0.79	0.87	0.81	0.01	0.02	0.80	0.61	0.69	0.97	0.87	0.92
per:age	0.91	0.99	0.95	0.97	0.32	0.48	0.93	0.65	0.77	0.95	1.00	0.97
per:cause_of_death	0.74	0.74	0.74	0.46	0.26	0.33	0.53	0.24	0.33	0.65	0.95	0.77
per:charges	0.79	0.95	0.86	0.76	0.30	0.43	0.67	0.43	0.52	0.87	0.99	0.93
per:children	0.54	0.73	0.62	0.23	0.14	0.17	0.19	0.38	0.25	0.96	0.73	0.83
per:cities_of_residence	0.51	0.95	0.66	0.38	0.55	0.45	0.51	0.34	0.41	0.61	0.92	0.73
per:city_of_birth	0.75	0.50	0.60	0.50	0.67	0.57	0.43	0.50	0.46	1.00	0.50	0.67
per:city_of_death	0.78	0.44	0.56	0.33	0.31	0.32	0.28	0.56	0.38	0.43	0.56	0.49
per:countries_of_residence	0.44	0.81	0.57	0.33	0.45	0.38	0.34	0.31	0.32	0.59	0.91	0.71
per:country_of_death	0.00	0.00	0.00	0.00	0.00	0.00	0.50	0.56	0.53	0.50	0.44	0.47
per:date_of_birth	0.70	1.00	0.82	0.26	0.86	0.40	0.71	0.71	0.71	0.86	0.86	0.86
per:date_of_death	0.76	0.80	0.78	0.53	0.42	0.47	0.62	0.12	0.21	0.74	0.93	0.82
per:employee_of	0.76	0.90	0.82	0.10	0.97	0.19	0.17	0.73	0.27	0.86	0.89	0.88
per:origin	0.70	0.83	0.76	0.65	0.12	0.21	0.45	0.13	0.20	0.82	0.79	0.80
per:other_family	0.39	0.49	0.43	0.00	0.00	0.00	0.11	0.59	0.19	0.61	0.97	0.75
per:parents	0.79	0.94	0.86	0.33	0.05	0.09	0.41	0.39	0.40	0.87	0.87	0.87
per:religion	0.65	0.75	0.70	0.67	0.35	0.46	0.65	0.75	0.70	0.71	0.93	0.80
per:schools_attended	0.95	0.70	0.81	1.00	0.04	0.07	0.86	0.22	0.35	1.00	0.81	0.90
per:siblings	0.69	0.85	0.77	0.40	0.60	0.48	0.45	0.69	0.54	0.98	0.94	0.96
per:spouse	0.77	0.97	0.86	0.26	0.69	0.38	0.20	0.58	0.30	0.89	0.81	0.85
per:stateorprovince_of_birth	1.00	0.67	0.80	0.00	0.00	0.00	0.50	0.33	0.40	0.50	0.67	0.57
per:stateorprovince of death	0.50	0.20	0.29	0.00	0.00	0.00	1.00	0.10	0.18	0.44	0.70	0.54
per:stateorprovinces_of_residence	0.53	0.84	0.65	0.40	0.03	0.06	0.43	0.52	0.47	0.63	0.84	0.72
per:title	0.90	0.97	0.93	0.88	0.01	0.03	0.97	0.08	0.14	0.94	0.97	0.96
macro avg	0.69	0.74	0.70	0.43	0.30	0.25	0.53	0.44	0.41	0.77	0.83	0.78
micro avg	0.83	0.83	0.83	0.20	0.14	0.17	0.53	0.51	0.52	0.89	0.89	0.89
weighted avg	0.87	0.83	0.84	0.70	0.14	0.11	0.72	0.51	0.54	0.91	0.89	0.89

Table 11: TACRED Results

Table 12: TACRED-RE Results

Model	DeBE	RTa-v3-la	rge 304M	Llama	-3.1 8B z	ero-shot	Llam	a-3.1 8B	5-shot	Llama	-3.1 8B fi	ne-tuned
	Р	R	F1	P	R	F1	P	R	F1	P	R	F1
Other	0.96	0.89	0.92	0.78	0.29	0.42	0.94	0.54	0.69	0.92	0.93	0.93
birthdate	0.99	1.00	0.99	0.53	0.91	0.67	0.83	0.91	0.87	1.00	1.00	1.00
bplace_name	0.82	0.92	0.87	0.82	0.08	0.14	0.69	0.82	0.75	0.87	0.85	0.86
child	0.44	0.44	0.44	0.02	0.56	0.05	0.00	0.00	0.00	0.57	0.44	0.50
deathdate	0.95	0.96	0.96	0.73	0.82	0.77	0.65	0.85	0.74	0.96	0.99	0.98
dplace name	0.64	0.71	0.67	0.39	0.25	0.31	0.50	0.55	0.52	0.56	0.81	0.66
educatedAt	0.72	0.90	0.80	0.06	1.00	0.12	0.19	1.00	0.32	0.00	0.00	0.00
occupation	1.00	0.99	1.00	0.83	0.41	0.55	0.77	0.94	0.85	1.00	0.98	0.99
parent	0.58	0.91	0.71	0.38	0.59	0.46	0.26	0.84	0.39	1.00	0.68	0.81
sibling	0.00	0.00	0.00	0.10	0.13	0.11	0.42	0.33	0.37	1.00	0.67	0.80
macro avg	0.71	0.77	0.74	0.46	0.50	0.36	0.53	0.68	0.55	0.79	0.74	0.75
micro avg	0.90	0.90	0.90	0.41	0.40	0.41	0.69	0.69	0.69	0.91	0.91	0.91
weighted avg	0.91	0.90	0.90	0.70	0.40	0.43	0.80	0.69	0.71	0.90	0.91	0.90

Table 13: Biographical Results

F.2 Cross-Dataset Results

Model	Setting	Dataset	Intra-	Dataset	Cross	-Dataset
			Shared Labels	Dataset Labels	NYT	TACRED-RE
DeBERTa-v3-large 304M	Fine-tuned on	NYT	0.62	0.62	-	0.27
		TACRED-RE	0.71	0.67	0.49	-
LLaMA 3.1 8B	Fine-tuned on	NYT	0.83	0.83	-	0.45
		TACRED-RE	0.79	0.76	0.58	-
	Shot Setting					
LLaMA 3.1 8B	Zero-Shot	NYT	0.35	0.35	-	-
		TACRED-RE	0.64	0.54	-	-
LLaMA 3.1 8B	5-Shot	NYT	0.46 ± 0.04	0.46 ± 0.04	-	0.50 ± 0.06
		TACRED-RE	0.59 ± 0.03	0.45 ± 0.05	0.44 ± 0.04	-

Table 14: Macro F1-scores for intra- and cross-dataset predictions for the overlap of six relations. Results are reported for shared and dataset-specific labels in both fine-tuned and shot settings. For fine-tuning, models are fine-tuned the overlap of six relations; for shot setting, also only the overlap of six relations is used for prompt or random shot samples.

F.3 Cross-Dataset Per-Class Results

Model	DeBE	DeBERTa-v3-large 304M LLaMA-3.1 8B 5-shot				LLaMA-3.1 8B fine-tuned			
	P	R	F1	Р	R	F1	P	R	F1
/people/deceased_person/place_of_death	0.75	0.36	0.49	0.54 ± 0.07	0.39 ± 0.11	0.45 ± 0.08	0.71	0.45	0.56
/people/person/children	0.50	0.43	0.46	0.47 ± 0.05	0.29 ± 0.10	0.34 ± 0.08	0.69	0.64	0.67
/people/person/place_lived	0.62	0.93	0.75	0.27 ± 0.06	0.22 ± 0.15	0.23 ± 0.11	0.63	0.86	0.73
/people/person/place_of_birth	0.33	0.06	0.11	0.12 ± 0.05	0.07 ± 0.04	0.09 ± 0.04	0.17	0.02	0.04
/people/person/religion	1.00	0.20	0.33	0.64 ± 0.03	0.72 ± 0.23	0.66 ± 0.11	1.00	0.80	0.89
None	0.94	0.78	0.85	0.84 ± 0.03	0.46 ± 0.07	0.59 ± 0.06	0.90	0.82	0.86
macro avg	0.69	0.46	0.50	0.48 ± 0.03	0.36 ± 0.04	0.39 ± 0.02	0.68	0.60	0.62
weighted avg	0.80	0.75	0.75	0.62 ± 0.04	0.37 ± 0.03	0.45 ± 0.03	0.78	0.76	0.76

Table 15: Adapted on TACRED-RE with all biographical relations present in TACRED-RE. Evaluations on NYT. The labels, although borrowed from NYT dataset, reflect the shared labels between NYT and TACRED-RE. More fine-grained TACRED-RE were mapped to broader shared labels to enable cross-dataset evaluation comparison.

Model	DeBE	RTa-v3-la	rge 304M	LI	aMA-3.1 8B 5-s	hot	LLaMA-3.1 8B fine-tuned			
	P	R	F1	Р	R	F1	P	R	F1	
/people/deceased_person/place_of_death /people/person/children /people/person/place_lived /people/person/place_of_birth /people/person/religion None	0.47 0.50 0.69 0.01 0.00 0.81	0.20 0.14 0.17 0.08 0.00 0.90	0.28 0.21 0.27 0.02 0.00 0.85	$\begin{array}{c} 0.66 \pm 0.11 \\ 0.49 \pm 0.15 \\ 0.65 \pm 0.09 \\ 0.18 \pm 0.10 \\ 0.86 \pm 0.04 \\ 0.85 \pm 0.04 \end{array}$	$\begin{array}{c} 0.33 \pm 0.20 \\ 0.48 \pm 0.20 \\ 0.38 \pm 0.06 \\ 0.88 \pm 0.26 \\ 0.64 \pm 0.17 \\ 0.78 \pm 0.12 \end{array}$	$\begin{array}{c} 0.40 \pm 0.17 \\ 0.47 \pm 0.17 \\ 0.48 \pm 0.05 \\ 0.28 \pm 0.11 \\ 0.72 \pm 0.11 \\ 0.81 \pm 0.05 \end{array}$	0.62 1.00 0.71 0.38 0.94 0.79	0.23 0.14 0.31 0.25 0.40 0.95	0.33 0.24 0.43 0.30 0.56 0.86	
macro avg weighted avg	0.41 0.73	0.25 0.66	0.27 0.66	$\begin{array}{c} 0.62 \pm 0.05 \\ 0.79 \pm 0.01 \end{array}$	$\begin{array}{c} 0.58 \pm 0.06 \\ 0.66 \pm 0.09 \end{array}$	$\begin{array}{c} 0.52 \pm 0.06 \\ 0.71 \pm 0.05 \end{array}$	0.74	0.38 0.74	0.45 0.72	

Table 16: Adapted on NYT with all biographical relations present in NYT. Evaluations on TACRED-RE. The labels, although borrowed from NYT dataset, reflect the shared labels between NYT and TACRED-RE. More fine-grained TACRED-RE were mapped to broader shared labels to enable cross-dataset evaluation comparison.

Model	DeBE	RTa-v3-la	rge 304M	LL	aMA-3.1 8B 5-s	hot	LLaMA-3.1 8B fine-tuned			
	P	R	F1	Р	R	F1	P	R	F1	
None /people/person/place_of_birth /people/person/children /people/deceased_person/place_of_death	0.92 0.90 0.57 0.92	0.63 0.73 0.44 0.23	0.75 0.80 0.50 0.36	$ \begin{vmatrix} 0.87 \pm 0.10 \\ 0.84 \pm 0.02 \\ 0.14 \pm 0.14 \\ 0.87 \pm 0.05 \end{vmatrix} $	$\begin{array}{c} 0.49 \pm 0.09 \\ 0.70 \pm 0.08 \\ 0.13 \pm 0.14 \\ 0.36 \pm 0.05 \end{array}$	$\begin{array}{c} 0.62 \pm 0.04 \\ 0.76 \pm 0.05 \\ 0.14 \pm 0.14 \\ 0.51 \pm 0.05 \end{array}$	0.90 0.90 0.71 0.95	0.68 0.75 0.56 0.38	0.78 0.82 0.63 0.54	
macro avg weighted avg	0.83 0.91	0.51 0.61	0.60 0.72	$\begin{array}{c} 0.68 \pm 0.06 \\ 0.85 \pm 0.06 \end{array}$	$\begin{array}{c} 0.42 \pm 0.04 \\ 0.55 \pm 0.01 \end{array}$	$\begin{array}{c} 0.51 \pm 0.04 \\ 0.65 \pm 0.01 \end{array}$	0.87 0.90	0.59 0.67	0.69 0.76	

Table 17: Adapted on TACRED-RE with all biographical relations present in TACRED-RE. Evaluations on Biographical with four biographical relations (full overlap between three datasets). The labels, although borrowed from NYT dataset, reflect the shared labels between NYT, TACRED-RE, and Biographical.

Model	DeBE	RTa-v3-la	arge 304M	LL	.aMA-3.1 8B 5-s	hot	LLaMA-3.1 8B fine-tuned			
	P	R	F1	Р	R	F1	P	R	F1	
/people/person/place_of_birth /people/person/children /people/deceased_person/place_of_death None	0.76 0.27 0.62 0.65	0.39 0.44 0.22 0.83	0.51 0.33 0.32 0.73	$ \begin{array}{c} 0.82 \pm 0.02 \\ 0.19 \pm 0.06 \\ 0.76 \pm 0.09 \\ 0.85 \pm 0.09 \end{array} $	$\begin{array}{c} 0.71 \pm 0.14 \\ 0.51 \pm 0.06 \\ 0.11 \pm 0.08 \\ 0.62 \pm 0.09 \end{array}$	$\begin{array}{c} 0.75 \pm 0.09 \\ 0.27 \pm 0.07 \\ 0.19 \pm 0.12 \\ 0.71 \pm 0.06 \end{array}$	0.92 0.00 0.74 0.57	0.14 0.00 0.13 0.96	0.25 0.00 0.22 0.72	
macro avg weighted avg	0.57	0.47 0.59	0.47 0.59	$\begin{array}{c} 0.65 \pm 0.02 \\ 0.82 \pm 0.04 \end{array}$	$\begin{array}{c} 0.49 \pm 0.04 \\ 0.59 \pm 0.05 \end{array}$	$\begin{array}{c}\textbf{0.48}\pm\textbf{0.04}\\\textbf{0.66}\pm\textbf{0.04}\end{array}$	0.56 0.71	0.31 0.55	0.30 0.48	

Table 18: Adapted on NYT with all biographical relations present in NYT. Evaluations on Biographical with four biographical relations (full overlap between three datasets). The labels, although borrowed from NYT dataset, reflect the shared labels between NYT, TACRED-RE, and Biographical.

	1		LL MA 219D	LL MA 219D	LL MA 219D
Setting	Parameter	large Fine-tuned	Zero-Shot	Five-Shot	Fine-tuned
	# of Epochs	10	-	-	3
	seed	42	42	42	42
	Loss	Cross-Entropy Loss	-	-	Cross-Entropy Loss
	Optimiser	AdamW	-	-	AdamW
Common	Batch Size	8	-	-	4
Common	Gradient Accumulation	4	-	-	-
	Early Stopping Patience	2	-	-	2
	Temperature	-	0.1	0.1	-
	Nucleus Sampling	-	0.9	0.9	-
	Lora Settings [†]	-	-	-	8/32/0.1
TACRED DE	Learning Rate	$5 \times 10^{-6} / 5 \times 10^{-5}$	-	-	5×10^{-5}
IACKED-KE	Max Length	-	_	_	800/384
	Max New Tokens	-	40	256	_
	Cross-Validation	–/5-fold	-	-	-
	Learning Rate	5×10^{-6}	_	_	1×10^{-4}
NYT	Max Length	-	-	-	384
	Max New Tokens	-	256	256	-
	Cross-Validation	–/5-fold	-	-	-
	Learning Rate	5×10^{-6}	_	_	1×10^{-4}
Biographical	Max Length	_	-	-	384
0.1	Max New Tokens	-	40	256	-

Table 19: Hyperparameter settings across datasets. Two values (x/y) indicate *All/Overlap* relation experiment settings respectively (if these differ), where *All* indicates experiments with the whole set of biographical relations in each dataset and *Overlap* uses only the intersection. Biographical experiments are performed only with the whole set of biographical relations. [†]Lora Settings: Rank/Alpha/Dropout.

	TACRE	ED-RE	NY	T	Biogra	phical
POS	Head Entity	Tail Entity	Head Entity	Tail Entity	Head Entity	Tail Entity
PROPN	77.4	55.6	98.4	98.8	87.7	60.5
PRON	16.8	6.2	-	-	0.1	0.2
NOUN	2.4	16.8	0.2	0.4	1.7	5.4
ADJ	1.2	4.5	0.2	-	0.7	0.9
ADP	0.7	1.6	0.2	0.3	0.3	1.1
NUM	0.0	8.5	-	-	6.4	25.5
DET	0.3	1.0	0.2	0.1	0.8	2.3
VERB	0.4	0.7	-	-	0.1	0.1

Table 20: (Top 8) POS Distribution Across TACRED-RE, NYT, and Biographical Test Sets with all Biographical Relations (%). POS tags are obtained with spaCy's transformer-based en_core_web_trf model.

G Misclassification Analysis

Issue	Description	Representative Example	Misclassifications on	Models Affected
Overpredicting 'None'	Overpredicting 'None' and strug- gling with even clear relations with cues like 'born' or'died'	"My name is <e1>Ruben</e1> and I am from <e2>Holland</e2> " (GT: <i>place_lived</i> , Pred: <i>None</i> ; TACRED-RE sample ID: '098f6f318bc468878bbb')	TACRED-RE and Biographical	NYT-adapted models
Failure to Cap- ture Implicit Re- lations	Models struggling to detect implicit relations requiring reasoning	" <e1>Gross</e1> [] was sent to <e2>Cuba</e2> as a spy" (GT: <i>place_lived</i> , Pred: <i>None</i> ; TACRED-RE sample ID: '098f6f318be29eddb182')	TACRED-RE	NYT-adapted models
Expected world knowledge	For NYT and Bio- graphical this issue is also frequently paired with detat- able ground truth la- bels	" <e1>Augustus</e1> also amassed an impressive art collection and built lavish baroque palaces in Dresden and <e2>Warsaw</e2> " (GT: <i>dplace_name</i> , Pred: <i>None</i> ; Biographical sample ID: 'mS2/247724')	NYT, TACRED- RE, Biographical	models adapted on all 3 datasets affected
Relation Present in Sentence but Not Between Specified Entities	This issue raises concerns about the framing of the RE task itself	"Jan Malte, [] resident of <e1>Bridgehampton</e1> , died [] in <e2>San Fran- cisco</e2> " (GT: None, Pred: place_of_death; NYT article ID: '/m/vinci8/data1/riedel/projects /relation/kb/nyt1/docstore/nyt- 2005-2006.backup /1777142.xml.pb')	NYT, TACRED- RE, Biographical	models adapted on all 3 datasets affected
Debatable ground truth (GT) labels	Caused by distantly or semi-supervised manner in which NYT and Biograph- ical were created	" <e1>Ida Freund</e1> was born in <e2>Austria</e2> " (GT: <i>Other</i> , Pred: <i>place_of_birth</i> ; Biographical sample ID: 'mS10/37387826')	TACRED-RE, Bi- ographical	NYT-adapted models
Single-Label An- notation Limita- tion	Sentences labeled with a single re- lation may contain additional relations that remain unla- beled	" <e1>Gross</e1> , who is himself Jewish [] was sent to <e2>Cuba</e2> " (GT: <i>place_lived</i> , Pred: <i>None</i> ; TACRED-RE sample ID: '098f6f318b69f98c850c')	NYT, TACRED- RE, Biographical	models adapted on all 3 datasets affected
Relation missing in annotation schema	Lack of granularity needed to fully cap- ture an individual's biography	"Wen was detained in August and accused of protecting the gang operations masterminded by his sister-in-law, <e1>Xie Caiping</e1> , 46, known as the "godmother" of the <e2>Chinese</e2> under- world (GT: <i>place_lived</i> , Pred: <i>nationality</i> ; TACRED-RE sample ID: '098f637935e6e6d1d093')	NYT, TACRED- RE, Biographical	
Failure to Cap- ture Relations in Long, Compound Sentences	Models struggling with long-term rela- tional dependencies	"Ecoffey told jurors that he and another federal agent met with <e1>Graham</e1> in April 1994 in Yellowknife, the city in northwest <e2>Canada</e2> where Graham lived at the time" (GT: <i>place_lived</i> , Pred: <i>None</i> ; TACRED-RE sample ID: '098f6f318b3ea9531448')	TACRED-RE	NYT-adapted models

Table 21: Common Misclassification Patterns Across TACRED-RE, NYT, and Biographical

Relation	NYT	TACRED-RE	Biographical
None	year, york, united, mr,	year, national, president,	release, contract, an-
	like, states, president,	group, include, state,	nounce, song, star,
	company, work, city	percent, million, amer-	award, series, role, sign,
		ican, china	championship
children	father, son, higgins,	son, daughter, grand-	daughter, son, child, li,
	clark, favre, richard,	child, survive, wife,	father, mother, wife,
	mary, daughter, carol,	year, child, include,	give, marry, actor
	daley	gude, jr	
religion	islam, muhammad,	jewish, al, islam, shiite,	-
	prophet, religion, con-	christian, group, mus-	
	vert, leader, school, al,	lim, sunni, mohammed,	
	church, close	tantawi	
place_lived	senator, republican,	year, state, die, home,	-
	state, year, representa-	york, city, president,	
	tive, gov, democrat, city,	live, iran, old	
	john, mr		
place_of_birth	city, year, orleans,	bear, grow, family, child,	bear, raise, née, grow,
	chicago, bear, bill,	york, year, native, july,	family, youth, york, cal-
	attorney, general, mr,	old, son	ifornia, city, mother
	california		
place_of_death	die, year, home, city,	die, home, hospital, can-	die, home, paris, age,
	london, los, angeles, mr,	cer, paris, wednesday,	california, near, october,
	yesterday, paris	sunday, find, early, dead	london, live, move

Table 22: Top 10 tokens per overlapping relation in NYT, TACRED-RE, and Biographical datasets, following lemmatisation and stop word removal using spaCy's transformer-based en_core_web_trf model.