

CARE: Extracting Experimental Findings From Clinical Literature

Anonymous ACL submission

Abstract

Extracting fine-grained experimental findings from literature can provide dramatic utility for scientific applications. Prior work has developed annotation schemas and datasets for limited aspects of this problem, failing to capture the real-world complexity and nuance required. Focusing on biomedicine, this work presents CARE—a new IE dataset for the task of extracting clinical findings. We develop a new annotation schema capturing fine-grained findings as n-ary relations between entities and attributes, which unifies phenomena challenging for current IE systems such as discontinuous entity spans, nested relations, variable arity n-ary relations and numeric results in a single schema. We collect extensive annotations for 700 abstracts from two sources: clinical trials and case reports. We also demonstrate the generalizability of our schema to the computer science and materials science domains. We benchmark state-of-the-art IE systems on CARE, showing that even models such as GPT4 struggle. We release our resources to advance research on extracting and aggregating literature findings.

1 Introduction

It is surely a great criticism of our profession that we have not organised a critical summary, by specialty or sub-specialty, adapted periodically, of all relevant randomised controlled trials. (Archie Cochrane, 1979)

Though this critique focused on clinical trials, the statement arguably applies to much of science today. There is tremendous potential utility in extracting, structuring and aggregating fine-grained information about experimental findings and the conditions under which they were achieved, across scientific studies. Once extracted and aggregated, scientific findings can power many critical applications such as producing literature reviews

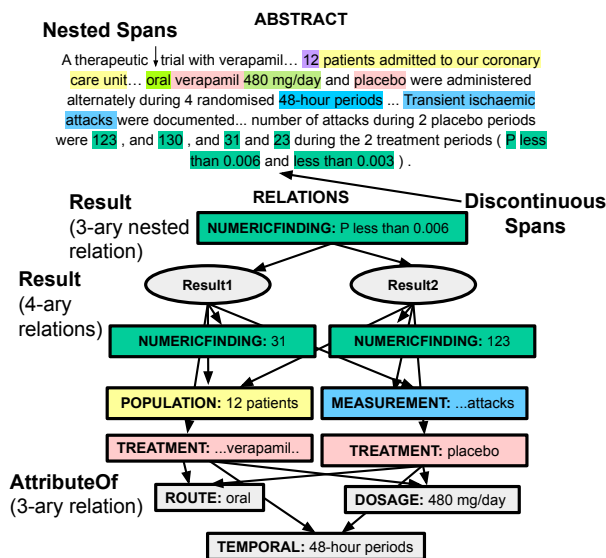


Figure 1: A partial example of entity, attribute and relation annotation using our schema for a clinical trial.

(DeYoung et al., 2021), supporting evidence-based decision-making (Naik et al., 2022), and generating new hypotheses (Wang et al., 2023).

While there have been efforts on building resources and tools to capture findings in various domains such as clinical trials (Lehman et al., 2019), computer science (Jain et al., 2020) and social and behavioral sciences (Magnusson and Friedman, 2021)—a major obstacle has been creating a representation that is expressive enough to capture complex and nuanced information about findings. We propose a new representation schema that makes important progress in capturing the real-world complexity of scientific findings in papers, and use it to build a high-quality annotated dataset focusing on biomedical (clinical) findings. Our schema represents fine-grained information about experimental findings and conditions as n-ary relations between entities and attributes, and includes several structural complexities such as discontinuous span annotation, variable arity in relations and

nestedness in relations. These aspects have been studied individually in previous datasets (Karimi et al., 2015; Tiktinsky et al., 2022), but our schema is the first to unify them. Our dataset also captures *numeric* findings in addition to their interpretation (e.g., significance, utility, etc.); prior datasets typically focus solely on the latter (e.g., Lehman et al. (2019) captures *increases/decreases* in outcomes but not their magnitudes).

To build our dataset, named CARE (Clinical Aggregation-oriented Result Extraction), we collect extensive annotations for 700 abstracts (clinical trials and case reports). We also conduct annotation studies demonstrating that our schema generalizes to computer science and materials science, using minor updates based on analogies between aspects across experimental domains (e.g., *populations/interventions* \rightarrow *tasks/methods* in CS). This reflects the expressive power of our schema to generalize across domains while capturing granular and useful information, making it a strong "backbone schema" for research efforts on result-oriented scientific IE.

We achieve good agreement scores (0.74-0.78 partial F1) comparable to prior work that used simpler schemas that are easier to annotate (Luan et al., 2018; Nye et al., 2018), and at the same time our resulting dataset is larger in size than previous corpora. Our final dataset annotation is extremely rich; at 16.23 relations per abstract, our relation density is nearly 4x that of prior work on annotating findings from clinical trials (Lehman et al., 2019).

We evaluate a wide range of IE models on our dataset, including both extractive systems and generative LLMs. Given the high annotation burden, we test generative LLMs in both fully supervised as well as zero-shot and few-shot settings. Our results demonstrate the difficulty of our dataset, with even SOTA models such as GPT4 struggling to accurately extract clinical findings. As a highly challenging new dataset designed to be reflective of real-world nuance and informational needs, we hope CARE is an important resource for the scientific NLP and IE research community to pursue.

2 Related Work

2.1 Information Extraction from Scientific Literature

Much prior work has focused on information extraction from scientific papers (Luan et al., 2018; Jain et al., 2020), including biomedical literature

(see (Luo et al., 2022a) for a detailed summary). Most relevant to our goal in this work is prior research on extracting findings or results from scientific literature, but it has only explored limited aspects of this problem.

Gábor et al. (2018) and Luan et al. (2018) annotate *associative* relations between entities being compared or producing a result, as part of their broader goal of developing IE resources for computer science, but do not capture any nuance (e.g., directionality, causality, etc. of results). Conversely, Magnusson and Friedman (2021) develop a schema focused solely on capturing associations between experimental variables and evidence. However, their focus on sentence-level annotation from scientific claims limits how much additional nuance about experimental setting can be captured.

Some prior efforts have also explored result extraction from biomedical literature. The EBM-NLP (Nye et al., 2018) and Evidence Inference (Lehman et al., 2019) corpora contain annotations for experimental findings from clinical trials, following the well-established PICO (participant, intervention, comparator, outcome) framework (Richardson et al., 1995). Sanchez-Graillet et al. (2022) also develop a PICO-inspired schema-based annotation format for diabetes and glaucoma trials. Chen et al. (2022) focuses on aggregating findings, which are already manually organized in structured format in databases such as AACT (Aggregate Analysis of ClinicalTrials.gov) (Tasneem et al., 2012). However, these efforts are tailored to clinical trials and do not translate easily to other domains. Finally, Luo et al. (2022a) conducted *novelty* annotation for relations, indicating whether they were presented as new observations; however they did not focus on experimental findings.

In contrast, we develop a representation schema expressive enough to capture fine-grained experimental findings, while generalizing across scientific domains. Our schema also contains phenomena challenging for SOTA IE models (§3.2).

2.2 Extracting Numeric Information

Another unique aspect of our schema is our focus on capturing numeric information from experimental findings and setup, which is understudied. Some prior work on open IE has explored extraction and linking of numeric spans (Madaan et al., 2016; Saha et al., 2017), including linking to implied entities (Elazar and Goldberg, 2019) (e.g., "it's worth

| Type | EBM | CTKG | Example |
|-----------------|-----|------|--|
| Population | ✓ | ✓ | This study compared rizatriptan 5 mg and placebo in 1268 outpatients treating a single migraine attack |
| Subpopulation | ✓ | ✓ | We found low-certainty evidence of little or no difference in delirium (RR 1.06, 95% CI 0.55 to 2.06; 2 studies, 800 participants) |
| Treatment | ✓ | ✓ | Dialysate magnesium was 0.375 mM/L for the hemodialysis |
| Measurement | ✓ | ✓ | Headache relief rates after rizatriptan 10 mg were higher |
| Temporal | ✗ | ✓ | After a 48-hour run-in period, oral verapamil 480 mg/day and placebo were administered |
| Numeric Finding | ✗ | ✓ | The number of attacks during treatment periods were 31 and 23 |
| Qualifier | ✗ | ✗ | Pindolol and metoprolol lowered blood pressure to the same extent |

Table 1: Examples of entity types in our schema. EBM and CTKG columns indicate whether these entity types are present in the EBM-NLP and CTKG schemas respectively. EBM-NLP uses IE to extract information according to its schema, while CTKG is a database schema not based on IE.

| Type | EBM | CTKG | Example |
|-------------|-----|------|-------------------------------|
| Age | ✓ | ✗ | for those age 60-67 years |
| Sex | ✓ | ✗ | 210 females |
| Size | ✓ | ✓ | 12 patients |
| Condition | ✓ | ✓ | patients getting hemodialysis |
| Demographic | ✗ | ✗ | A 40's Japanese man |
| Route | ✗ | ✗ | oral verapamil |
| Dosage | ✗ | ✗ | verapamil 480 mg/day |
| Strength | ✗ | ✗ | rizatriptan 5 mg |
| Duration | ✗ | ✗ | for 4 weeks |

Table 2: Examples of attribute types in our schema. EBM and CTKG columns indicate whether these entity types are present in the EBM-NLP and CTKG schemas.

two million” can be linked to currency). However, these models broadly focused on sentence-level extraction and did not evaluate on scientific text.

Within the scientific domain, some studies have focused on numeric information extraction from biomedical/clinical text. Kang and Kayaalp (2013) and Claveau et al. (2017) extract numeric spans from FDA-released decision summaries and clinical trial eligibility criteria respectively. EBM-NLP (Nye et al., 2018) annotates some categories of numeric information associated with cohorts participating in a clinical trial, but ignores trial outcomes and findings. Among non-medical scientific domains, numeric span extraction work has mainly focused on extraction from tables (Hou et al., 2019). None of these studies focus extensively on linking numeric spans with entities that can help in interpreting this information, which is key to our work.

3 Annotation Schema

We develop a new annotation schema to represent fine-grained clinical findings present in biomedical

abstracts, and later demonstrate its broader applicability to domains beyond biomedicine (§6.2). Our schema captures this knowledge via three main elements, commonly used in IE tasks:

1. Entities involved in a study, which are spans of text, either contiguous or non-contiguous, belonging to one of the seven types listed in Table 1.

2. Attributes associated with entities, which are also contiguous or non-contiguous spans of text, belonging to one of the nine types listed in Table 2. The first five attribute types are associated with population and subpopulation entities, while the remaining four types are associated with intervention entities. Other entity types do not have any associated attributes.

3. N-ary Relations linking together various entities and attributes, where N (relation arity) is variable and nesting is allowed. A relation is an n-tuple, where each element can be an entity, attribute or another n-ary relation. Relations are categorized into four types listed in Table 3.

3.1 Comparison to Clinical Schemas

Prior work such as EBM-NLP (Nye et al., 2018) and Evidence Inference (Lehman et al., 2019; DeYoung et al., 2020) has focused on developing IE schemas to represent clinical knowledge appearing in the literature in a structured format. In addition, work such as CTKG (Chen et al., 2022) outside the NLP/IE sphere has built schema for representing clinical information in databases. However, these schemas suffer from a few shortcomings: (i) most are designed for clinical trials; their applicability to other types of biomedical literature is untested, (ii) focus on a small set of broad entity types, which leaves out fine-grained details, (iii) follow strict relation formats, which makes it hard to capture ad-

| Type | Arity | EI | CTKG | Example |
|-----------------|--------|----|------|--|
| AttributeOf | N-ary | ✗ | ✗ | (<i>Subpopulation</i> : 144 had the U-type method, <i>Size</i> : 144) |
| SubpopulationOf | N-ary | ✗ | ✗ | (<i>Population</i> : 285 women, <i>Subpopulation</i> : 144 had the U-type method, <i>Subpopulation</i> : 141 had the H-type method) |
| InterventionOf | Binary | ✗ | ✓ | (<i>Subpopulation</i> : 144 had the U-type method, <i>Intervention</i> : U-type method) |
| Result | N-ary | ✓ | ✓ | (<i>Subpopulation</i> : 144 had the U-type method, <i>Measurement</i> : objective cure rates, <i>NumericFinding</i> : 87.5%) |

Table 3: Examples of relation types in our schema. EI and CTKG columns indicate whether these relation types are present in the EI and CTKG schemas respectively. While the EI and CTKG datasets contain 4-ary and binary result relations respectively, our n-ary schema allows fine-grained information to be captured more flexibly.

ditional nuance that might be useful for interpreting findings.

Our schema makes several enhancements to tackle these issues. First, it is extensible to other categories of biomedical literature beyond clinical trials, and we demonstrate this by applying our schema to case reports. Second, our schema captures more fine-grained information about various entities than prior work via attributes (see Table 2). Third, allowing for variable arity and nesting in relation annotation provides the flexibility which makes our schema capable of representing both atomic findings (e.g., value of primary outcome observed for a given intervention) as well as composite findings (e.g., outcome improvement observed for intervention vs control groups). Tables 1, 2 and 3 provide a more detailed comparison of our schema with EBM-NLP, EI and CTKG.

3.2 Annotation Complexity

In addition to using an expanded set of entity, attribute and relation types, our annotation schema supports the following phenomena (also illustrated in Figure 1), unifying them all in a single dataset:

Discontinuous spans: Biomedical abstracts often present multiple entities as conjunctive phrases or lists of items, so we allow discontinuous span annotation to capture every entity. For example, given the phrase “maximal diameters and volumes”, our scheme captures two measurement entities: “maximal diameters” and “maximal volumes”, with the latter being a discontinuous span.

Nested/overlapping spans: Attributes, as defined in our annotation scheme, are often present within an entity span or overlap with an entity span. This motivates our decision to allow nested and overlapping spans to be annotated.

Variable arity in relations: Owing to variation in clinical studies, findings are often described in a wide range of formats (e.g., outcome for a single population, outcome for a pair of populations,

outcome for a single population at different time periods, etc.). This diversity motivated our choice of *variable arity* for relation annotation, similar to Tiktinsky et al. (2022).

Nested relations: In addition to outcomes for individual populations/groups, clinical studies often present comparative findings and analyses, such as improvement on an outcome given a pair of interventions. Our scheme allows for annotation of nested relations to link these higher-order observations with their associated atomic findings.

Our complete annotation guidelines are included in the supplementary material. Figure 1 presents partial entity, attribute and relation annotations for an example clinical trial abstract.

4 Dataset Collection

Annotation Tool: We use TeamTat¹ (Islamaj et al., 2020), a web-based tool for team annotation since it allows for n-ary and nested relation annotation, a core component of our schema.

Annotator Background: We recruit two in-house annotators² with backgrounds in data analytics and data science, both having extensive experience in reading and annotating scientific papers. One of our annotators has a background in biology. Both annotators went through several pilot rounds to gain familiarity with our task and schema. Additionally, we used their feedback and insights from pilots to solidify our schema design (see §4.1). We also solicited feedback from two medical students and an MD to validate our final schema.

Data Sources: CARE covers two categories of biomedical literature: (i) clinical trials, and (ii) case reports. Clinical trials are research studies that test a medical, surgical, or behavioral intervention in people to determine whether a new form of treatment or prevention or a new diagnostic device is

¹<https://www.teamtat.org>

²included as co-authors on this paper

| Category | Exact F1 | Partial F1 |
|-----------|----------|------------|
| Entity | 0.5764 | 0.7578 |
| Attribute | 0.6174 | 0.7801 |
| Relation | 0.4209 | 0.7414 |

Table 4: Final inter-annotator agreement scores on a sample of 28 abstracts, measured during full-scale data annotation.

effective. Case reports are detailed reports of the symptoms, signs, diagnosis, treatment, and follow-up of an individual patient, usually motivated by unusual or novel occurrences. We sample clinical trials from the EBM-NLP (Nye et al., 2018) dataset, which consists of 4993 abstracts annotated with PICO spans, only retaining abstracts containing at least one number (4685 in total). To sample case reports, we extract all reports with at least one number in the abstract from PubMed (907,862 in total) and randomly sample from this pool. We sample 350 abstracts from each source, resulting in our final dataset size of 700 abstracts, which is slightly larger than other prior corpora that perform fine-grained annotation (§ 4.3). Further characteristics of our abstract sample are detailed in Appendix C.

4.1 Annotation Pilots

We conducted three pilot rounds with the following goals: (i) training annotators to apply our schema, (ii) evaluating agreement, and (iii) assessing whether our schema captures clinical knowledge of interest. Annotators worked on a fresh set of 5-10 abstracts per round, followed by agreement computation and disagreement discussion. For entity and attribute annotation, agreement is computed as entity-level F1 between annotators, using both strict (entity boundaries match exactly) and partial (entity boundaries overlap on at least one token) matching. For relations, we first align annotations from both annotators by linking pairs of relations which share $\geq 50\%$ of participating entities. Agreement is computed as F1 score between annotators, using both strict (100% of entities match) and partial matching. After achieving reasonable agreement levels by round 3 (partial F1 scores of 0.79, 0.68 and 0.79 for entity, attribute and relation annotation respectively), we started full-scale data annotation (further discussion in Appendix C).

4.2 Full-Scale Annotation

The full-scale data annotation process was conducted in six rounds. To continue monitoring agreement, a small agreement set of 5 abstracts (not

| Metric | Train | Dev | Test |
|-------------|---------|--------|--------|
| #Docs | 500 | 100 | 100 |
| #Tokens | 135,363 | 27,120 | 25,219 |
| #Entities | 12022 | 2367 | 2286 |
| #Attributes | 3992 | 804 | 762 |
| #Relations | 8205 | 1594 | 1560 |

Table 5: Statistics for final collected dataset.

| Phenomenon | Train | Dev | Test |
|----------------------|-------|-------|-------|
| #Discontinuous Spans | 8.9% | 10.1% | 9.3% |
| #Nested Spans | 3.4% | 4.3% | 2.5% |
| #Overlapping Spans | 1.6% | 2.0% | 0.7% |
| #Nested Relations | 11.4% | 11.2% | 11.9% |

Table 6: Prevalence of interesting annotation phenomena in final collected dataset.

identified to the annotators) was included in every round. Table 9 in the appendix presents inter-annotator agreement during each annotation round, while Table 4 shows overall agreement scores. Overall and per-round agreement scores continued to remain in the same range as agreement scores from later pilot rounds, demonstrating consistency in annotation quality. Despite the complexity of our schema, our agreement scores are comparable to datasets using simpler schemas like EBM-NLP (entity agreement of 0.62-0.71; Cohen’s kappa) and SciERC (relation agreement of 67.8; kappa score). Appendix C provides additional details about our full-scale annotation setup.

Consensus Annotation: For all abstracts annotated by multiple annotators during pilots or full-scale annotation (55 in total), we construct a “consensus” version post disagreement discussion. The final dataset releases consensus annotations for these abstracts. Since this subset has been annotated by multiple annotators and discussed extensively, we expect annotations to be higher-quality and include all these abstracts in the test set.

4.3 Dataset Statistics

Table 5 gives an overview of statistics for our final collected dataset. Our dataset size is comparable to other prior biomedical corpora which performs exhaustive fine-grained annotation (though not always with a clinical knowledge focus) such as BioRED (Luo et al. (2022a); 600 abstracts) and Sanchez-Graillet et al. (2022) (211 abstracts). Table 6 presents the proportion of various interesting phenomena allowed by our schema in the final dataset. Interestingly, CARE contains 9% discontinuous spans, making it one of the rare datasets

containing a large proportion of discontinuous mentions.³ At 11%, the final data also contains a high proportion of nested relations.

5 Benchmarking IE Models

We benchmark the performance of two categories of models on CARE: (i) extractive models, and (ii) generative LLMs. We also test generative LLMs in two settings: (i) finetuning on the full training set, and (ii) zero-shot and in-context learning.

Experimental Setup: We test each model on the three sub-tasks—entity extraction, attribute extraction and relation extraction—in isolation. Model performance on entity and attribute extraction is evaluated using entity-level F1. Relation extraction performance is evaluated using a relaxed overlap F1 score metric inspired by Tiktinsky et al. (2022), which assigns partial credit to correctly identified subsets of entities in a relation, even if all identified entities do not match. As with agreement score calculation, predicted relations are first aligned with gold relations by choosing the gold relation with highest overlap per predicted relation. Then a partial match score is computed as $\#shared_entities/total_entities$ and used in the F1 computation instead of binary 0/1 score.

5.1 Extractive IE Baselines:

We evaluate the following systems:

- **OneIE** (Lin et al., 2020): A sentence-level joint entity, relation and event extraction system, which extracts an “information network” representation of entities and events (nodes), connected by relations (edges). Beam search is used to find the highest-scoring network.
- **PURE** (Zhong and Chen, 2021): A sentence-level pipelined extraction system, which learns separate contextual representations for entity and relation extraction, using entity representations to further refine relation extraction.
- **LocLabel** (Shen et al., 2021): A sentence-level two-stage named entity recognition (NER) system capable of extracting nested spans. Inspired by object detection work, it produces boundary proposals for candidate entities, then labels them with correct entity types.

³Dai et al. (2020) considers 10% discontinuous spans to be a high proportion, identifying only three biomedical datasets that satisfy this criterion: CADEC (Karimi et al., 2015), ShARe 13 (Pradhan et al., 2013) and ShARe 14 (Mowery et al., 2014).

| Model | Ent F1 | Attr F1 | Rel F1 |
|-----------------------------|--------------|--------------|--------------|
| Extractive Baselines | | | |
| OneIE | 55.07 | 48.84 | – |
| PURE | 55.94 | 61.04 | – |
| LocLabel | 53.69 | 55.25 | – |
| W2NER | 51.84 | 57.98 | – |
| Generative Baselines | | | |
| FLAN-T5 | 45.08 | 23.27 | 33.24 |
| BioGPT | 14.43 | 29.84 | 33.15 |
| BioMedLM | 1.50 | 10.62 | 32.76 |
| GPT-3.5 0-shot | 11.14 | 5.06 | 14.35 |
| GPT-3.5 1-shot | 21.40 | 8.61 | 31.58 |
| GPT-3.5 3-shot | 23.40 | 8.85 | 31.58 |
| GPT-3.5 5-shot | 8.92 | 9.92 | 32.20 |
| GPT-4 0-shot | 26.89 | 9.02 | 32.04 |
| GPT-4 1-shot | 31.07 | 11.82 | 42.81 |
| GPT-4 3-shot | 16.68 | 13.16 | 53.69 |
| GPT-4 5-shot | 5.04 | 13.90 | 55.04 |

Table 7: Performance of all extractive and generative baselines on entity, attribute and relation extraction.

- **W2NER** (Li et al., 2022): A sentence-level unified NER model, capable of extracting nested and discontinuous spans. It recasts NER as word-word relation classification on a 2-D grid of word pairs, then decodes word pair relations into final span extractions.

For comparability and better adaptation to our dataset, we replace BERT-based encoders in all systems with PubmedBERT (Gu et al., 2021), and follow best-reported hyperparameters per system (see Appendix E). Table 7 presents their performance on entity and attribute extraction. Unfortunately, applying these systems to our relation extraction task is infeasible, since none of them are designed for document-level relation extraction or n-ary relations. Tiktinsky et al. (2022) modify PURE for n-ary relation extraction with variable arity. However, given a set of candidate entities, they consider all possible n-ary combinations and predict relationships per cluster. This is tractable for their work on sentence-level extraction of single-type (drug interaction) relations, but not tractable for document-level multi-type n-ary relation extraction.⁴ Therefore, we do not test extractive models on relation extraction.

Another caveat with extractive models is that they do not identify discontinuous spans (except W2NER). To assess how this impacts model performance, we compute an additional entity-level F1 score which merges span predictions linked in gold

⁴On limiting combination size to 10, every abstract produces 500,000 candidate combinations

447 annotation (i.e., we assume oracle span merging),
448 and observe that this does not significantly improve
449 performance (avg. increase of ~ 1.5 F1). Therefore,
450 Table 7 reports F1 scores without merging.

451 5.2 Generative IE Baselines:

452 Motivated by recent work demonstrating LLM ca-
453 pabilities on information extraction (Wadhwa et al.,
454 2023), we assess the ability of LLMs on our tasks,
455 in both finetuning and zero-shot/in-context learning
456 settings.

457 We evaluate the following finetuned LLMs:

- 458 • **FLAN-T5** (Chung et al., 2022): Enhanced ver-
459 sion of T5 (Raffel et al., 2020) finetuned on a
460 large mixture of tasks, but not specifically pre-
461 trained for biomedicine. We use FLAN-T5-XL,
462 which has 3B parameters.
- 463 • **BioGPT** (Luo et al., 2022b): A 1.6B autoregres-
464 sive model, pretrained from scratch on 15M ab-
465 stracts and titles from PubMed with a custom
466 Pubmed-trained tokenizer.
- 467 • **BioMedLM**⁵: A 2.7B autoregressive model, pre-
468 trained from scratch on all PubMed abstracts and
469 full-texts from the Pile (Gao et al., 2020) with a
470 custom PubMed-trained tokenizer.

471 When training and testing on attribute and re-
472 lation extraction, these models are provided gold
473 entities and attributes by surrounding them with
474 entity markers (`< ent >< /ent >`) in the input.

475 We evaluate GPT3.5 and GPT4 in zero-shot and
476 in-context learning settings. We provide our IE
477 schema and example outputs and prompt the model
478 to produce extractions in a clean JSON format that
479 adheres to the schema. Additionally, for our in-
480 context learning experiments, we follow (Liu et al.,
481 2021) and select the k most similar examples from
482 the training set for every test instance according to
483 similarity computed by the SPECTER v2.0 (Singh
484 et al., 2022) PRX model trained on scientific titles
485 and abstracts. Selected examples are appended to
486 the prompt in decreasing order of similarity, with
487 later examples dropped if they don't fit. We run
488 experiments for the $k = 1, 3, 5$ most similar exam-
489 ples. Further hyperparameter details for all models
490 are provided in Appendix E.

491 Table 7 shows the performance of all generative
492 models. One caveat with GPT3.5/4 is that model
493 outputs sometimes contain correct entity/attribute
494 spans assigned to the wrong type (e.g., a subpop-
495 ulation misclassified as a population entity in a

496 result relation). Since we are evaluating the perfor-
497 mance of relation extraction in isolation, we do not
498 consider such mistyping as errors.

499 5.3 End-to-End Evaluation:

500 In addition to evaluating SOTA systems on each
501 sub-task in isolation, we assess the feasibility of
502 end-to-end extraction. Table 7 shows that PURE
503 is the best-performing system on entity and at-
504 tribute extraction. On the other hand, GPT4 5-shot
505 and FLAN-T5 perform best on relation extraction
506 (GPT3.5 5-shot and BioGPT are close). We test
507 out a hybrid end-to-end extraction system in which
508 entities and attributes are detected using PURE,
509 then input text marked up with these extractions is
510 provided to FLAN-T5 for relation extraction. This
511 hybrid system achieves an F1 score of 33.58, very
512 similar to RE performance with gold markup. Hy-
513 pothesizing that this might be an indication that
514 finetuned LLMs ignore entity/attribute markup dur-
515 ing RE, we run an additional experiment in which
516 we train FLAN-T5 to extract relations from raw
517 text (no markup). This setup achieves an F1 score
518 of 33.07, showing that entity/attribute markup does
519 not provide significant benefit.

520 6 Discussion

521 6.1 How much does strict evaluation 522 underestimate LLM performance?

523 Table 7 shows that even fully-supervised generative
524 models severely lag behind much smaller extractive
525 models on entity and attribute extraction. However,
526 prior work (Wadhwa et al., 2023) has observed
527 that strict IE evaluation metrics underestimate the
528 performance of LLMs since their outputs often
529 contain minor variations from gold annotations,
530 which could still be correct. Therefore, we conduct
531 a human evaluation of a subset of FLAN-T5 and
532 GPT4 5-shot predictions on entity and attribute
533 extraction for a more accurate assessment.

534 For every setting, we collect all abstracts with
535 one or more wrong predictions and randomly sam-
536 ple ten to evaluate. We go over all false positives
537 per abstract marking ones that could be considered
538 correct. Our evaluation shows that for FLAN-T5,
539 35 out of 73 entity and 12 out of 32 attribute er-
540 rors are marked correct. For GPT4, these numbers
541 are worse; 38 out of 126 entity and 20 out of 79
542 attribute errors are marked correct. This indicates
543 that LLMs indeed struggle with our span extraction
544 tasks, and their poor performance is not simply a

⁵<https://crfm.stanford.edu/2022/12/15/biomedlm.html>

| Original Type | Generalized Type | Description |
|-----------------|--------------------------|---|
| Population | Research Problem Context | Setting/scenario in which the authors are testing their hypothesis (e.g., task or dataset being studied in ML/NLP). |
| Subpopulation | Problem Stages/Sub-parts | Subgroups or subsamples of overall setting (e.g., dataset splits in ML/NLP). |
| Treatment | Technique/Method | Key technique being proposed or investigated and other techniques being compared (e.g., model or metric in ML/NLP). |
| SubpopulationOf | Sub-PartOf | Links together problem context entities to stage/sub-part entities (e.g., for ML/NLP, this relation would link the overall task to low-data and fully supervised settings). |
| TreatmentOf | AppliedTo | Links together a technique to all the problem contexts/sub-parts it is being tested in. |

Table 8: Changes required to construct a generalized version of our original schema developed for clinical finding extraction, which we use to test whether it applies to other domains such as computer science and materials science

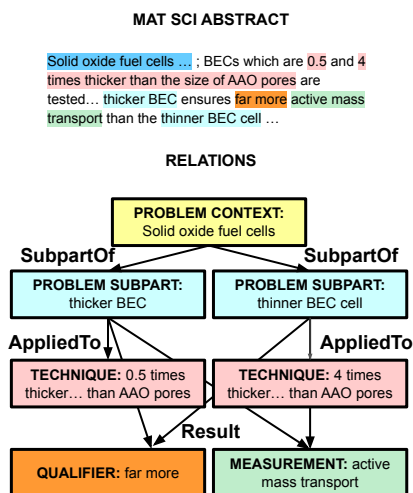


Figure 2: A partial example of entity, attribute and relation annotation using our generalized schema for a materials science abstract.

consequence of strict evaluation.

6.2 How easily can we extend our schema to other domains?

Though we focus on extracting clinical findings from biomedical literature during schema design, we try to incorporate enough flexibility to allow our schema to be easily adapted to other scientific domains. To demonstrate this flexibility, we conduct small-scale pilots in two additional domains: (i) Computer Science, and (ii) Materials Science.

We first develop a *generalized* version of our proposed schema for these studies. Of the three elements in our schema, entities and relations are largely transferable and only require minor renaming. Table 8 provides an overview of changes made to entity/relation nomenclature. Attributes on the other hand, were tailored more closely to our goal of extracting clinical findings. Therefore, we drop all attributes and ask our annotators to propose candidate attributes as they go through the annotation

process. We use the same annotators who participated in dataset create, to leverage their existing familiarity with our schema, assigning one annotator to each domain. Their task is to annotate ten abstracts each while documenting: (i) potential attributes that can be added to the schema, and (ii) important experimental information missed by the generalized schema.

After completing the task, annotators reported that it was feasible to apply our proposed schemas to these scientific domains. Computer science posed some difficulty due to the presence of lots of relative results and references in the abstract, which made entity annotation ambiguous. However, there were no important aspects of experimental information, aside from potential attribute proposals, that our current schema could not account for.

7 Conclusion

In this work, we presented CARE, a new IE dataset for the task of extracting clinical findings from biomedical literature. To collect this dataset, we first developed a new annotation schema capable of capturing fine-grained information about experimental findings, which unified several challenging IE phenomena such as discontinuous spans, nested relations and variable arity n-ary relations. Using this annotation scheme, we collected an extensively annotated dataset of 700 abstracts from clinical trials and case reports. Our benchmarking experiments showed that state-of-the-art extractive and generative LLMs including GPT4 still struggle on this task, particularly on relation extraction. We release both our annotation schema and CARE as a challenging new resource for the IE community and to encourage further research on extraction and representation of findings from scientific literature.

8 Limitations

Despite being a cornerstone of our work, the richness and complexity of our newly proposed annotation schema also poses some limitations. Annotators needed some prior experience with reading and understanding complex scientific text, and had to undergo multiple rounds of additional training before they were able to accurately apply our schema and start full-scale annotation. Though these stringent expertise and training requirements and heavy reliance on human annotators helped us collect a high-quality resource in CARE, they simultaneously limit the scalability of our collection protocol and make it difficult to construct large-scale benchmarks for this task, spanning multiple domains/fields of science.

Our annotated corpus, CARE, is based on RCTs and case reports. While our schema is broad and expressive enough to generalize to other experimental domains with minor adaptations, our generalization annotation studies were comparatively small and preliminary, limited to testing the schema on computer science and material science papers. In addition, while our schema covers many types of experimental finding information, the richness and huge variety of scientific experiments necessarily means that more types of findings could be added. In the future, more studies should be performed on using our schema in other domains, and on extending our schema with more types of informations (entities, attributes, relations). CARE also focuses on English-language papers only, and in the future it would be interesting and important to extend our schema and dataset to cover biomedical/clinical studies in other languages, to capture important scientific findings that are potentially missed when only looking at papers in English.

Finally, a limitation of our current benchmarking effort is the lack of more flexible evaluation metrics, particularly when assessing the performance of generative LLMs. We try to provide supplementary human evaluation for some models to overcome this issue, but this is not scalable and would require ongoing/continuous evaluation efforts. This is not a major focus for our current work, but developing more flexible automated evaluation is an important future direction for IE research.

References

- Ziqi Chen, Bo Peng, Vassilis N Ioannidis, Mufei Li, George Karypis, and Xia Ning. 2022. A knowledge graph of clinical trials (ctkg). *Scientific reports*, 12(1):4724.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.
- Vincent Claveau, Lucas Emanuel Silva Oliveira, Guillaume Bouzillé, Marc Cuggia, Claudia Maria Cabral Moro, and Natalia Grabar. 2017. Numerical eligibility criteria in clinical protocols: annotation, automatic detection and interpretation. In *Artificial Intelligence in Medicine: 16th Conference on Artificial Intelligence in Medicine, AIME 2017, Vienna, Austria, June 21-24, 2017, Proceedings 16*, pages 203–208. Springer.
- Xiang Dai, Sarvnaz Karimi, Ben Hachey, and Cecile Paris. 2020. An effective transition-based model for discontinuous NER. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5860–5870, Online. Association for Computational Linguistics.
- Jay DeYoung, Iz Beltagy, Madeleine van Zuylen, Bailey Kuehl, and Lucy Lu Wang. 2021. MS²: Multi-document summarization of medical studies. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7494–7513, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Jay DeYoung, Eric Lehman, Benjamin Nye, Iain Marshall, and Byron C Wallace. 2020. Evidence inference 2.0: More data, better models. In *Proceedings of the 19th SIGBioMed Workshop on Biomedical Language Processing*, pages 123–132.
- Yanai Elazar and Yoav Goldberg. 2019. Where’s my head? definition, data set, and models for numeric fused-head identification and resolution. *Transactions of the Association for Computational Linguistics*, 7:519–535.
- Kata Gábor, Davide Buscaldi, Anne-Kathrin Schumann, Behrang QasemiZadeh, Haïfa Zargayouna, and Thierry Charnois. 2018. SemEval-2018 task 7: Semantic relation extraction and classification in scientific papers. In *Proceedings of the 12th International Workshop on Semantic Evaluation*, pages 679–688, New Orleans, Louisiana. Association for Computational Linguistics.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. 2020. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*.

| | | | | |
|-----|--|-----|---|-----|
| 703 | Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2021. Domain-specific language model pretraining for biomedical natural language processing. <i>ACM Transactions on Computing for Healthcare (HEALTH)</i> , 3(1):1–23. | 757 | Yi Luan, Luheng He, Mari Ostendorf, and Hannaneh Hajishirzi. 2018. Multi-task identification of entities, relations, and coreference for scientific knowledge graph construction . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3219–3232, Brussels, Belgium. Association for Computational Linguistics. | 758 |
| 704 | | 759 | | 760 |
| 705 | | 761 | | 762 |
| 706 | | 763 | | |
| 707 | | | | |
| 708 | | | | |
| 709 | Yufang Hou, Charles Jochim, Martin Gleize, Francesca Bonin, and Debasis Ganguly. 2019. Identification of tasks, datasets, evaluation metrics, and numeric scores for scientific leaderboards construction . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 5203–5213, Florence, Italy. Association for Computational Linguistics. | | Ling Luo, Po-Ting Lai, Chih-Hsuan Wei, Cecilia N Arighi, and Zhiyong Lu. 2022a. Biored: a rich biomedical relation extraction dataset. <i>Briefings in Bioinformatics</i> , 23(5):bbac282. | 764 |
| 710 | | 765 | | 766 |
| 711 | | 767 | | |
| 712 | | | | |
| 713 | | | | |
| 714 | | | | |
| 715 | | | | |
| 716 | | | | |
| 717 | Rezarta Islamaj, Dongseop Kwon, Sun Kim, and Zhiyong Lu. 2020. Teamtat: a collaborative text annotation tool. <i>Nucleic acids research</i> , 48(W1):W5–W11. | | Renqian Luo, Liai Sun, Yingce Xia, Tao Qin, Sheng Zhang, Hoifung Poon, and Tie-Yan Liu. 2022b. Biogpt: generative pre-trained transformer for biomedical text generation and mining. <i>Briefings in Bioinformatics</i> , 23(6):bbac409. | 768 |
| 718 | | 769 | | 770 |
| 719 | | 771 | | 772 |
| 720 | Sarthak Jain, Madeleine van Zuylen, Hannaneh Hajishirzi, and Iz Beltagy. 2020. SciREX: A challenge dataset for document-level information extraction . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7506–7516, Online. Association for Computational Linguistics. | | Aman Madaan, Ashish Mittal, Ganesh Ramakrishnan, Sunita Sarawagi, et al. 2016. Numerical relation extraction with minimal supervision. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 30. | 773 |
| 721 | | 774 | | 775 |
| 722 | | 776 | | 777 |
| 723 | | | | |
| 724 | | | | |
| 725 | | | | |
| 726 | | | | |
| 727 | Yanna Shen Kang and Mehmet Kayaalp. 2013. Extracting laboratory test information from biomedical text. <i>Journal of pathology informatics</i> , 4(1):23. | | Ian H Magnusson and Scott E Friedman. 2021. Extracting fine-grained knowledge graphs of scientific claims: Dataset and transformer-based results. <i>arXiv preprint arXiv:2109.10453</i> . | 778 |
| 728 | | 779 | | 780 |
| 729 | | 781 | | |
| 730 | Sarvnaz Karimi, Alejandro Metke-Jimenez, Madonna Kemp, and Chen Wang. 2015. Cadec: A corpus of adverse drug event annotations. <i>Journal of biomedical informatics</i> , 55:73–81. | | Danielle L Mowery, Sumithra Velupillai, Brett R South, Lee Christensen, David Martinez, Liadh Kelly, Lorraine Goeuriot, Noemie Elhadad, Sameer Pradhan, Guergana Savova, et al. 2014. Task 2: Share/clef ehealth evaluation lab 2014. In <i>Proceedings of CLEF 2014</i> . | 782 |
| 731 | | 783 | | 784 |
| 732 | | 785 | | 786 |
| 733 | | 787 | | |
| 734 | Eric Lehman, Jay DeYoung, Regina Barzilay, and Byron C Wallace. 2019. Inferring which medical treatments work from reports of clinical trials. In <i>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)</i> , pages 3705–3717. | | Aakanksha Naik, Sravanthi Parasa, Sergey Feldman, Lucy Lu Wang, and Tom Hope. 2022. Literature-augmented clinical outcome prediction . In <i>Findings of the Association for Computational Linguistics: NAACL 2022</i> , pages 438–453, Seattle, United States. Association for Computational Linguistics. | 788 |
| 735 | | 789 | | 790 |
| 736 | | 791 | | 792 |
| 737 | | 793 | | |
| 738 | | | | |
| 739 | | | | |
| 740 | | | | |
| 741 | Jingye Li, Hao Fei, Jiang Liu, Shengqiong Wu, Meishan Zhang, Chong Teng, Donghong Ji, and Fei Li. 2022. Unified named entity recognition as word-word relation classification. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 36, pages 10965–10973. | | Benjamin Nye, Junyi Jessy Li, Roma Patel, Yinfei Yang, Iain Marshall, Ani Nenkova, and Byron C Wallace. 2018. A corpus with multi-level annotations of patients, interventions and outcomes to support language processing for medical literature. In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 197–207. | 794 |
| 742 | | 795 | | 796 |
| 743 | | 797 | | 798 |
| 744 | | 799 | | 800 |
| 745 | | 801 | | |
| 746 | | | | |
| 747 | Ying Lin, Heng Ji, Fei Huang, and Lingfei Wu. 2020. A joint neural model for information extraction with global features . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7999–8009, Online. Association for Computational Linguistics. | | Sameer Pradhan, Noemie Elhadad, Brett R South, David Martinez, Lee M Christensen, Amy Vogel, Hanna Suominen, Wendy W Chapman, and Guergana K Savova. 2013. Task 1: Share/clef ehealth evaluation lab 2013. <i>CLEF (working notes)</i> , 1179. | 802 |
| 748 | | 803 | | 804 |
| 749 | | 805 | | 806 |
| 750 | | | | |
| 751 | | | | |
| 752 | | | | |
| 753 | Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2021. What makes good in-context examples for gpt-3? <i>arXiv preprint arXiv:2101.06804</i> . | | Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>The Journal of Machine Learning Research</i> , 21(1):5485–5551. | 807 |
| 754 | | 808 | | 809 |
| 755 | | 810 | | 811 |
| 756 | | 812 | | |

| | | | |
|-----|--|---|-----|
| 813 | W Scott Richardson, Mark C Wilson, Jim Nishikawa, | Zexuan Zhong and Danqi Chen. 2021. A frustratingly | 870 |
| 814 | Robert S Hayward, et al. 1995. The well-built clinical | easy approach for entity and relation extraction . In | 871 |
| 815 | question: a key to evidence-based decisions. <i>Acp</i> | <i>Proceedings of the 2021 Conference of the North</i> | 872 |
| 816 | <i>j club</i> , 123(3):A12–A13. | <i>American Chapter of the Association for Computa-</i> | 873 |
| | | <i>tional Linguistics: Human Language Technologies</i> , | 874 |
| 817 | Swarnadeep Saha, Harinder Pal, and Mausam. 2017. | pages 50–61, Online. Association for Computational | 875 |
| 818 | Bootstrapping for numerical open IE . In <i>Proceedings</i> | <i>Linguistics</i> . | 876 |
| 819 | <i>of the 55th Annual Meeting of the Association for</i> | | |
| 820 | <i>Computational Linguistics (Volume 2: Short Papers)</i> , | | |
| 821 | pages 317–323, Vancouver, Canada. Association for | | |
| 822 | Computational Linguistics. | | |
| | | A Schema Definitions | 877 |
| 823 | Olivia Sanchez-Graillet, Christian Witte, Frank Grimm, | A.1 Entity Types | 878 |
| 824 | and Philipp Cimiano. 2022. An annotated corpus of | Entities can belong to one of the following seven | 879 |
| 825 | clinical trial publications supporting schema-based | types: | 880 |
| 826 | relational information extraction. <i>Journal of Biomed-</i> | 1. Population: Patient groups/cohorts studied in | 881 |
| 827 | <i>ical Semantics</i> , 13(1):1–18. | an article. | 882 |
| | | 2. Subpopulation: Slices/sub-groups of a popula- | 883 |
| 828 | Yongliang Shen, Xinyin Ma, Zeqi Tan, Shuai Zhang, | tion entity sharing some underlying characteris- | 884 |
| 829 | Wen Wang, and Weiming Lu. 2021. Locate and label: | tic. | 885 |
| 830 | A two-stage identifier for nested named entity | 3. Treatment: Treatment regimens, procedures, | 886 |
| 831 | recognition . In <i>Proceedings of the 59th Annual Meet-</i> | therapies etc. prescribed and/or tested to allevi- | 887 |
| 832 | <i>ing of the Association for Computational Linguistics</i> | ate a population’s conditions/symptoms. | 888 |
| 833 | <i>and the 11th International Joint Conference on Natu-</i> | 4. Measurement: Tests used to assess population | 889 |
| 834 | <i>ral Language Processing (Volume 1: Long Papers)</i> , | status and outcomes of the tested intervention. | 890 |
| 835 | pages 2782–2794, Online. Association for Computa- | 5. Temporal: Temporal information such as time | 891 |
| 836 | tional Linguistics. | points at which outcomes are measured. | 892 |
| 837 | Amanpreet Singh, Mike D’Arcy, Arman Cohan, Doug | 6. Numeric Finding: All numeric information as- | 893 |
| 838 | Downey, and Sergey Feldman. 2022. Scirepeval: | sociated with study findings (e.g., p-values, haz- | 894 |
| 839 | A multi-format benchmark for scientific document | ard ratios, etc.). | 895 |
| 840 | representations. <i>arXiv preprint arXiv:2211.13308</i> . | 7. Qualifier: Non-numeric information associated | 896 |
| | | with study findings that provides important per- | 897 |
| 841 | Asba Tasneem, Laura Aberle, Hari Ananth, Swati | spective for interpreting them (e.g., phrases in- | 898 |
| 842 | Chakraborty, Karen Chiswell, Brian J McCourt, and | dicating evidence directionality). | 899 |
| 843 | Ricardo Pietrobon. 2012. The database for aggre- | | |
| 844 | gate analysis of clinicaltrials. gov (aact) and subse- | A.2 Attribute Types | 900 |
| 845 | quent regrouping by clinical specialty. <i>PLoS one</i> , | Attributes can belong to one of the following nine | 901 |
| 846 | 7(3):e33677. | types: | 902 |
| | | 1. Age: Numeric or non-numeric information | 903 |
| 847 | Aryeh Tiktinsky, Vijay Viswanathan, Danna Niezni, | about the age of the population under study. | 904 |
| 848 | Dana Meron Azagury, Yosi Shamy, Hillel Taub- | 2. Sex: Reported sex of the population under | 905 |
| 849 | Tabib, Tom Hope, and Yoav Goldberg. 2022. A | study. | 906 |
| 850 | dataset for n-ary relation extraction of drug combi- | 3. Size: Size of the population sample under study. | 907 |
| 851 | nations. In <i>Proceedings of the 2022 Conference of</i> | 4. Condition: Medical conditions prevalent in the | 908 |
| 852 | <i>the North American Chapter of the Association for</i> | study population, including diseases, symptoms, | 909 |
| 853 | <i>Computational Linguistics: Human Language Tech-</i> | prior medical history and procedures, etc. | 910 |
| 854 | <i>nologies</i> , pages 3190–3203. | 5. Demographic: Additional demographic infor- | 911 |
| | | mation reported about the population such as | 912 |
| 855 | Somin Wadhwa, Silvio Amir, and Byron Wallace. 2023. | location, race, etc. | 913 |
| 856 | Revisiting relation extraction in the era of large lan- | 6. Route: Description of the way an intervention | 914 |
| 857 | guage models . In <i>Proceedings of the 61st Annual</i> | is administered (e.g., a chemical may be admin- | 915 |
| 858 | <i>Meeting of the Association for Computational Lin-</i> | istered orally, topically, intravenously, etc.). | 916 |
| 859 | <i>guistics (Volume 1: Long Papers)</i> , pages 15566– | 7. Dosage: Quantity of administration for the in- | 917 |
| 860 | 15589, Toronto, Canada. Association for Computa- | tervention being studied. This is not necessarily | 918 |
| 861 | tional Linguistics. | limited to chemical/drug interventions (e.g., for | 919 |
| | | an intervention like educational sessions, num- | 920 |
| 862 | Qingyun Wang, Doug Downey, Heng Ji, and Tom Hope. | ber of sessions is considered “dosage”). | 921 |
| 863 | 2023. Learning to generate novel scientific directions | | |
| 864 | with contextualized literature-based discovery. <i>arXiv</i> | | |
| 865 | <i>preprint arXiv:2305.14259</i> . | | |
| | | | |
| 866 | Yijia Zhang, Qingyu Chen, Zhihao Yang, Hongfei Lin, | | |
| 867 | and Zhiyong Lu. 2019. Biwordvec, improving | | |
| 868 | biomedical word embeddings with subword infor- | | |
| 869 | mation and mesh. <i>Scientific data</i> , 6(1):52. | | |

| | | | |
|-----|---|--|------|
| 922 | 8. Strength: Strength of chemical/drug interventions administered. | as “X is usually treated using Y...”, do not annotate Y unless Y was one of the treatments actually given to a population in the current study. | 971 |
| 923 | | | 972 |
| 924 | 9. Duration: Interval of time over which an intervention was administered. | | 973 |
| 925 | | | |
| 926 | A.3 Relation Types | C Dataset Construction Details | 974 |
| 927 | Our schema allows for both binary and n-ary relations (with variable n), to capture four types of structure: | Characteristics of sampled abstracts: Since the EBM-NLP corpus sampled randomized clinical trials from PubMed with an emphasis on cardiovascular diseases, cancer, and autism, the clinical trials portion of our dataset also heavily features these topics. On the other hand, for case reports, comparing MeSH term distributions across all reports (2M abstracts) with case reports containing numeric information (the 900k we sample from), we see a massive reduction (> 30%) in terms associated with the following topics: surgery and post-surgery care, dentistry, ophthalmology, prostheses and rehab, patient care and nursing, some mental disorders and circulatory diseases/issues. Hence, we expect these topics to be relatively undersampled in our pool of case reports. | 975 |
| 928 | | | 976 |
| 929 | 1. AttributeOf: N-ary relations linking population and intervention entities with their associated attributes. | | 977 |
| 930 | | | 978 |
| 931 | 2. Subpopulation: N-ary relations capturing parent-child relationships between population and subpopulation entities. | | 979 |
| 932 | | | 980 |
| 933 | 3. InterventionOf: Binary relations linking population and subpopulations entities with the intervention(s) tested on them. | | 981 |
| 934 | | | 982 |
| 935 | 4. Result: N-ary relations capturing all numeric or non-numeric outcome results and comparisons reported by linking together the population, subpopulation, intervention, measurement, numericfinding and/or qualifier and temporal entities involved in each result/comparison. | | 983 |
| 936 | | | 984 |
| 937 | All n-ary relations can contain multiple entities of a single type. For example, a result relation can involve multiple interventions or populations. The only cardinality constraints imposed are that every result relation should focus on a <i>single</i> measurement entity and always contain <i>at least one</i> population/intervention entity. | | 985 |
| 938 | | | 986 |
| 939 | | | 987 |
| 940 | | | 988 |
| 941 | | | 989 |
| 942 | | | 990 |
| 943 | | | 991 |
| 944 | | | 992 |
| 945 | | | 993 |
| 946 | | | 994 |
| 947 | | | 995 |
| 948 | | | 996 |
| 949 | | | 997 |
| 950 | | | 998 |
| 951 | | | 999 |
| 952 | | | 1000 |
| 953 | B Additional Annotation Rules | Annotation Pilots: During pilots, we also conducted one or more disagreement discussion sessions per pilot round. These discussions were helpful in providing annotators the opportunity to highlight important spans/relations being missed by the schema, which led to the addition of the subpopulation entity, demographic attribute, and subpopulationof and treatmentof relations. Despite the introduction of some new elements, inter-annotator agreement continued to increase steadily over the pilot rounds, as shown in Table 9 before plateauing at the end of round 3. | 1001 |
| 954 | While using this annotation schema to annotate clinical knowledge, we also keep in mind the following rules: | | 1002 |
| 955 | | | 1003 |
| 956 | • For every entity/attribute span, only annotate its first occurrence in the text, unless there is a more descriptive span later. We follow this rule to avoid conducting an additional coreference annotation step to link all spans referring to the same entity. | Full-Scale Annotation: During rounds 1-3 of full-scale annotation, annotators were provided batches of 25 abstracts each. As their familiarity with the annotation schema and ability to handle ambiguous cases improved, we provided larger batches of 100 abstracts each during rounds 4-6. After each round, agreement was assessed and disagreement discussions were conducted to discuss ambiguous cases, if needed, which ensured that agreement was maintained across rounds as seen from Table 9. Tables 10, 11 and 12 present final agreement scores per entity type, attribute type and relation type respectively. From these tables, we can see that Subpopulation and Intervention entities are the trickiest to annotate, leading to lower agreement on SubpopulationOf and InterventionOf relation types due to error cascading (i.e., if entity annotations don’t match, relation annotations are unlikely to | 1004 |
| 957 | | | 1005 |
| 958 | • Ignore misspellings and include all associated modifiers and abbreviations while annotating spans | | 1006 |
| 959 | | | 1007 |
| 960 | • Do not annotate generic or high-level spans (e.g., genetic disorder), or generic terms (e.g., complications, deficiency, disease, syndrome, gene, drug, protein, nucleotide, etc.). | | 1008 |
| 961 | | | 1009 |
| 962 | • Do not annotate background occurrences of entities. For example, if a treatment Y is mentioned | | 1010 |
| 963 | | | 1011 |
| 964 | | | 1012 |
| 965 | | | 1013 |
| 966 | | | 1014 |
| 967 | | | 1015 |
| 968 | | | 1016 |
| 969 | | | 1017 |
| 970 | | | 1018 |
| | | | 1019 |
| | | | 1020 |

| Round | Entity F1 | | Attribute F1 | | Relation F1 | |
|---------|---------------|---------------|---------------|---------------|---------------|---------------|
| | Exact | Partial | Exact | Partial | Exact | Partial |
| Pilot 1 | 0.6240 | 0.7579 | 0.7215 | 0.8163 | 0.2193 | 0.6379 |
| Pilot 2 | 0.7206 | 0.8818 | 0.6923 | 0.7385 | 0.4997 | 0.7878 |
| Pilot 3 | 0.6449 | 0.7900 | 0.5370 | 0.6852 | 0.4449 | 0.7960 |
| Batch 1 | 0.5130 | 0.7318 | 0.7611 | 0.8496 | 0.3899 | 0.6979 |
| Batch 2 | 0.6094 | 0.7900 | 0.6216 | 0.8508 | 0.6397 | 0.9137 |
| Batch 3 | 0.5312 | 0.7797 | 0.6364 | 0.8182 | 0.3121 | 0.7595 |
| Batch 4 | 0.5714 | 0.7817 | 0.7347 | 0.7755 | 0.5399 | 0.7343 |
| Batch 5 | 0.5643 | 0.6929 | 0.4717 | 0.6762 | 0.3382 | 0.6766 |
| Batch 6 | 0.6358 | 0.7930 | 0.5417 | 0.7582 | 0.3122 | 0.6890 |
| Overall | 0.5764 | 0.7578 | 0.6174 | 0.7801 | 0.4209 | 0.7414 |

Table 9: Evolution of inter-annotator agreement during pilots and full-scale annotation rounds

| Type | Exact F1 | Partial F1 |
|----------------|----------|------------|
| Population | 0.4333 | 0.8665 |
| Subpopulation | 0.4299 | 0.6168 |
| Intervention | 0.4333 | 0.5781 |
| Measurement | 0.5230 | 0.7554 |
| Temporal | 0.6230 | 0.6885 |
| NumericFinding | 0.7063 | 0.8812 |
| Qualifier | 0.6911 | 0.7749 |

Table 10: Inter-annotator agreement per entity type

| Type | Exact F1 | Partial F1 |
|-------------|----------|------------|
| Age | 0.8500 | 0.9756 |
| Sex | 0.9231 | 0.9231 |
| Size | 0.6462 | 0.7385 |
| Condition | 0.5091 | 0.7429 |
| Demographic | 0.6667 | 0.8000 |
| Route | 0.8000 | 0.8000 |
| Dosage | 0.6923 | 0.9630 |
| Strength | - | - |
| Duration | 0.0800 | 0.4800 |

Table 11: Inter-annotator agreement per attribute type. Note that the agreement sample did not include any strength entities.

| Type | Exact F1 | Partial F1 |
|-----------------|----------|------------|
| AttributeOf | 0.7654 | 0.7654 |
| InterventionOf | 0.3797 | 0.3797 |
| SubpopulationOf | 0.1633 | 0.5185 |
| Result | 0.2561 | 0.7994 |

Table 12: Inter-annotator agreement per relation type

match either).

D Inter-Annotator Agreement

Table 9 shows the evolution in inter-annotator agreement over our initial pilot rounds, as well as the level of inter-annotator agreement maintained during each round of the full-scale annotation process. We see a large increase in relation agreement from pilot 1 to pilot 2, and consistent agreement scores across all tasks in all rounds thereafter. Tables 10, 11 and 12 present inter-annotator agreement breakdown according to entity, attribute and relation types in our schema.

E Hyperparameter Details

Extractive Models:

- **OneIE:** We use an overall learning rate and weight decay of $1e - 3$, and a learning rate and weight decay of $1e - 5$ for the BERT component, a batch size of 10, and gradient clipping value of 5.0. The model is trained for 60 epochs with a 5-epoch warmup phase.
- **PURE:** We use a context window size of 300 words, overall learning rate of $1e - 5$, task learning rate of $5e - 4$, batch size of 16, and train for 100 epochs.
- **LocLabel:** We use a learning rate of $5e - 6$, warmup rate of 0.1, weight decay of 0.01, gradient clipping value of 1.0, batch size of 6 and train for 35 epochs. LocLabel also requires word vectors, for which we use the 200-dimensional Pubmed-trained word2vec embeddings (BioWordVec) released by Zhang et al.

(2019), which are available at <https://github.com/ncbi-nlp/BioWordVec>.

- **W2NER:** We use an overall learning rate of $1e - 3$ and a learning rate of $5e - 6$ for the BERT component, no weight decay, warmup factor fo 0.1, gradient clipping value of 5.0, batch size of 8, and train for 10 epochs.

Generative Models: All models are trained for 10 epochs with a learning rate of $1e - 5$, input context length of 1024, output length of 128, and a batch size of 2.

GPT3.5/GPT4: We test the 16k and 8k context length versions of GPT3.5 and GPT4 respectively since our extraction tasks are abstract-level and require longer input contexts. We use the June 2023 versions of both models due to their *function calling* capabilities, which leverage a structured JSON output format to improve information extraction capabilities. All experiments are run with a temperature of 0 and max output length of 512 tokens.

F Computing Infrastructure

All LLM experiments are carried out on NVIDIA RTX A6000 GPUs with 48 GB RAM. Each finetuning run (FLAN-T5, BioGPT, BioMedLM) requires two GPUs with runtimes ranging from 9-17 hours depending on task size and model size. We use the DeepSpeed integration from Huggingface, with ZeRO-3 optimization, for multi-GPU training.