

# A Dataset for N-ary Relation Extraction of Drug Combinations

Anonymous ACL submission

## Abstract

Combination therapies have become the standard of care for diseases such as cancer, tuberculosis, malaria and HIV. However, the combinatorial set of available multi-drug treatments creates a challenge in identifying effective combination therapies available in a situation. To assist medical professionals in identifying beneficial drug-combinations, we construct an expert-annotated dataset for extracting information about the efficacy of drug combinations from the scientific literature. Beyond its practical utility, the dataset also presents a unique NLP challenge, as the first relation extraction dataset consisting of variable-length relations. Furthermore, the relations in this dataset predominantly require language understanding beyond the sentence level, adding to the challenge of this task. We provide a promising baseline model and identify clear areas for further improvement. We release our dataset and code<sup>1</sup> publicly to encourage the NLP community to participate in this task.

## 1 Introduction

“So far, many monotherapies have been tested, but have been shown to have limited efficacy against COVID-19. By contrast, **combinational** therapies are emerging as a useful tool to treat SARS-CoV-2 infection.” (Ianevski et al., 2021).

Indeed, combining two or more drugs together has proven to be useful for treatments of various medical conditions, including cancer (DeVita et al., 1975; Carew et al., 2008; Shuhendler et al., 2010), AIDS (Bartlett et al., 2006), malaria (Eastman and Fidock, 2009), tuberculosis (Bhusal et al., 2005), hypertension (Rochlani et al., 2017) and COVID-19 (Ianevski et al., 2020).

In this work, we examine the clinically significant and challenging NLP task of extracting known

drug combinations from the scientific literature. We present an expert-annotated dataset and baseline models for this new task. Our dataset contains 1600 manually annotated abstracts, each mentioning between 2 and 15 drugs. 840 of these abstracts describe one or more positive drug combinations, varying in size from 2 to 11 drugs. The remaining 760 abstracts either contain mentions of drugs not used in combination, or discuss combinations of drugs that do not give a combined positive effect.

From a clinical perspective, solving the drug combination identification task will assist researchers in suggesting and validating complex treatment plans. For example, when searching for effective treatments for cancer, knowing which drugs interact synergistically with the first line treatment allows researchers to suggest new treatment plans that can subsequently be validated in-vivo and become a standard protocol (Wasserman et al., 2001; Katzir et al., 2019; Ianevski et al., 2020; Niezni et al., 2021).

From an NLP perspective, the drug combination identification task and dataset pushes the boundaries of relation extraction (RE) research, by introducing a relation extraction task with several challenging characteristics:

**Variable-length n-ary relations** Most work on relation extraction is centered on *binary relations* (e.g. Li et al. (2016), see full listing in §5), or on *n-ary relations with a fixed n* (e.g. Peng et al. (2017)). In contrast, the drug combination task involves *variable-length n-ary relations*: different passages discuss combinations of different numbers of drugs. For each subset of drugs mentioned in a passage, the model must predict if they are used together in a combination therapy and whether this drug combination is effective.

**No type hints** As noted by Rosenman et al. (2020) and Sabo et al. (2021), in many relation extraction benchmarks (Han et al., 2018; Sabo et al., 2021; Zhang et al., 2017), the argument types serve as

<sup>1</sup>Dataset and code can be found at <https://anonymous.4open.science/r/drug-synergy-models--C8B7/README.md>

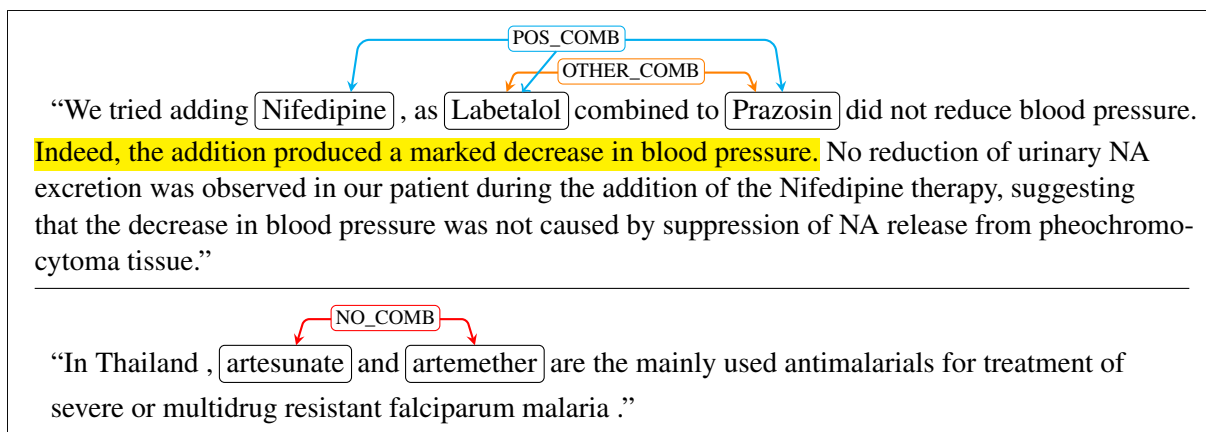


Figure 1: Examples of our label scheme. The top example contains two relations: a binary `OTHER_COMB` relation and a ternary `POS_COMB` relation. The evidence required to annotate the latter relation is found in a different sentence (highlighted). In the bottom example, each drug is described as a separate treatment rather than a combination therapy.

080 an effective clue. However, argument types do not  
 081 apply naturally to the drug combination task, in  
 082 which all possible relation arguments are entities  
 083 of the same type (drugs) and we need to identify  
 084 specific subsets of them.

085 **Long range dependencies** The information de-  
 086 scribing the efficacy of a combination is often  
 087 spread-out across multiple sentences. Indeed, our  
 088 annotators reported that for 67% of the instances,  
 089 the label could not be determined based on a single  
 090 sentence, requiring reasoning with a larger textual  
 091 context. Interestingly, our experiments show that  
 092 our models *are not* helped by the availability of  
 093 longer context, showing the limitations of current  
 094 standard modeling approaches. This suggests our  
 095 dataset can be a test-bed for models that attempt to  
 096 incorporate longer context.

097 **Challenging inferences** As we show in our qualita-  
 098 tive analysis (§4.2), instances in this dataset require  
 099 processing a range of phenomena, including coordi-  
 100 nation, numerical reasoning, and world knowledge.

101 We hope that by releasing this dataset we will  
 102 encourage NLP researchers to engage in this impor-  
 103 tant clinical task, while also pushing the boundaries  
 104 of relation extraction.

## 105 2 The Drug Combinations Dataset

106 A set of drugs in a biomedical abstract are classi-  
 107 fied to one of the following labels:

108 **Positive combination (`POS_COMB`):** the sen-  
 109 tence indicates the drugs are used in combination,  
 110 and the passage suggests that the combination has  
 111 additive, synergistic, or otherwise beneficial effects

112 which warrant further study.

113 **Non-positive combination (`OTHER_COMB`):**  
 114 the sentence indicates the drugs are used in com-  
 115 bination, but there is no evidence in the passage  
 116 that the effect is positive (it is either negative or  
 117 undetermined).<sup>2</sup>

118 **Not a combination (`NO_COMB`):** the sentence  
 119 does not state that the given drugs are used in com-  
 120 bination, even if a combination is indicated some-  
 121 where else in the wider context. An example is  
 122 given in the lower half of Figure 1, where each of  
 123 the drugs Artesunate and Artemether is given in  
 124 isolation, and no combination is reported.

125 Our primary interest is to identify sets of drugs  
 126 that are positive combinations.

### 127 2.1 Relevant Context Size for Classifying 128 Drug Combinations

129 When formulating the extraction task and design-  
 130 ing our data collection methodology, we first an-  
 131 alyzed the locality of the phenomenon: to what  
 132 extent are drug combinations are expressed in a  
 133 single sentence, or is a larger context is needed?  
 134 We sampled 275 abstracts that contained known  
 135 drug combinations according to DrugComboDB.<sup>3</sup>  
 136 Analysis showed that 51% of these abstracts men-  
 137 tioned attempted drug combinations. In 97% of the

<sup>2</sup>We also experimented with another label for combinations that are discouraged (antagonistic, harmful or not effective). The agreement for this label was low, leading us to keep it as a subset of `OTHER_COMB`.

<sup>3</sup>We used `Syner&Antag_voting.csv` taken from <http://drugcombdb.denglab.org/download/> and ranked according to the *Voting* metric.

abstracts containing drug combinations, all participating drugs in the attempted combination could be located within a single sentence in the abstract (for an example, see the OTHER\_COMB relation in Figure 1). However, establishing the efficacy of the combination frequently required a larger context (such as the context accompanying the POS\_COMB relation in Figure 1).

## 2.2 Task Definition

We define each instance in the Drug Combination Extraction (DCE) task to consist of a sentence, drug mentions within the sentence, and an enclosing context (e.g. paragraph or abstract).

The output of the task is a set of relations, each consisting of a set of participating drug spans and a relation label (POS\_COMB or OTHER\_COMB). Each subset of drug mentions not included in the output set is implicitly considered to have relation label NO\_COMB.

More formally, DCE is the task of labeling an instance  $X = \{C, i, D\}$  with a set of relation instances  $R$ , where  $C = (S_1, \dots, S_n)$  is an ordered list of context sentences (e.g. all the sentences in an abstract or paragraph),  $1 \leq i \leq n$  is an index of a target sentence  $S_i = (w_1, \dots, w_{n(i)})$  with  $n(i)$  words, and  $D = \{(d_{1start}, d_{1end}), \dots, (d_{mstart}, d_{mend})\}$  is a set of  $m \geq 2$  spans of drug mentions in  $S$ . The output is a set  $R = \{(c_i, y_i)\}$  where  $c_i \in \mathcal{P}(D)$  is a drug combination from  $\mathcal{P}(D)$ , the set of all possible drug combinations, and  $y_i \in \{\text{POS\_COMB}, \text{OTHER\_COMB}\}$  is a combination label.

## 2.3 Evaluation Metric

We consider two settings: “Exact Match”, a strict version which considers identifying exact drug combinations, and “Partial Match”, a more relaxed version which assigns partial credits to correctly identified subsets.

We use standard precision, recall and F1 metrics for both settings. For the partial-match case, we replace the binary 0 or 1 score for a given combination with a refined score:  $shared\_drugs/total\_drugs$ . If there are multiple partial matches with gold relations, we take the one with maximum overlap. We compute **recall** as  $identified\_relations/all\_gold\_relations$ , and **precision** as  $correct\_relations/identified\_relations$ .

We consider two metrics, the averaged Positive Combination F1 score which compares

POS\_COMB to the rest, and the averaged Any Combination F1 score which counts correct predictions for any combination label (POS or OTHER) as opposed to NO\_COMB. The latter is an easier task, but still valuable for identifying drug combinations irrespective of their efficacy.

## 2.4 Collecting Data for Annotation

To collect data for annotation we curated a list of 2411 drugs from DrugBank<sup>4</sup> and sampled from PubMed a set of sentences which mention 2 or more drugs. Analysis of the first 50 sentences from this sample showed that only 8/50 of the sentences included mentions of drug combinations. This meant that annotating the full sample will be costly, and will result in a dataset that’s highly skewed toward relatively trivial NO\_COMB instances.

We therefore repeated this experiment, sampling sentences whose PubMed abstract included a trigger phrase.<sup>5</sup> 48% of 50 sampled sentences included mentions of drug combinations. Evaluating the coverage of the trigger list against a new sample of abstracts with known drug combinations showed that 90% of these new abstracts included one of the trigger words. This suggests our trigger list is useful for fetching label-balanced data, without prohibitively restricting coverage and diversity.

Accordingly, we collected the majority of instances for annotation, 90%, using a basic search for sentences that contain at least two different drugs and whose abstract contains one of the trigger phrases. To overcome the lexical restrictions imposed by our trigger list, we sample the remaining 10% of instances using distant supervision: fetching sentences containing pairs of drugs known to be synergistic according DrugComboDB, but whose abstract does not include one of our trigger phrases. All data collecting queries were performed using the SPIKE Extractive Search tool (Shlain et al., 2020; Taub-Tabib et al., 2020). The process is illustrated at the top of Figure 2.

<sup>4</sup>Curation included downloading a premade drug list from DrugBank’s website, while removing non pharmacological intervention such as Vitamins and Supplements. The later we got from the FDA orange book.

<sup>5</sup>Triggers were selected by manually identifying words and phrases which frequently appear in abstracts mentioning drug combinations. These are phrases like “combination”, “followed by”, “prior to”, etc. (see full list in Appendix A.3). The triggers are recall oriented, so while a presence of a trigger increases the chances that an abstract mentions a drug combination, it is definitely not clearly indicative. Importantly, since we’re dealing with a wide context, the presence of a trigger in an abstract which includes multiple drugs does not mean the trigger is related to the drugs.

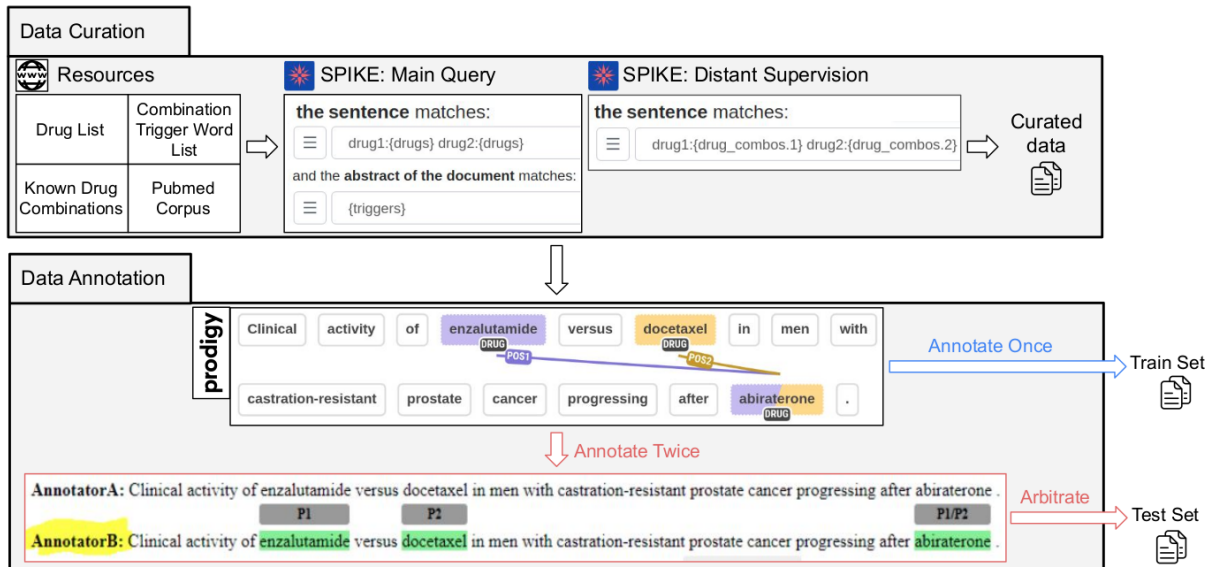


Figure 2: Illustration of the data construction process. First we construct the required knowledge resources. Then, we collect data using SPIKE –an extractive search tool– over the PubMed database. The train and test sets were annotated using Prodigy over the curated data. For test data, we collected two annotations for each sample, and then had a domain expert resolve annotation disagreements.

## 2.5 The Annotation Process

Seven graduate students in biomedical engineering took part in the annotation task. The students all completed a course in combination therapies for cancer and were supervised by a principled researcher with expertise in this field.

We provided the participants with annotation guidelines which specified how the annotation process should be carried out (see Appendix A.1) and conducted an initial meeting where we reviewed the guidelines with the group and discussed some of the examples together.

Each of the participants had access to a separate instance of the Prodigy annotation tool (Montani and Honnibal, 2018), pre-loaded with the candidate annotation instances. Once a session starts, the instances (containing of a sentence with marked drug entities, and its context) appear in a sequential manner, with no time limit. For each instance we instructed the annotators to mark all subsets of drugs that participated in a combination, and for each subset to indicate its label (POS\_COMB or OTHER\_COMB). Moreover, we instructed them to indicate whether the context was needed in order to determine the positive efficacy of the relation.

Despite the considerable time required for expert annotation, we collected annotations for 1634 passages. Among these, 272 were assigned to at least two annotators. After further arbitration by the lead

| Metric                       | Partial Match | Exact Match |
|------------------------------|---------------|-------------|
| Avg. Any Combination F1      | 88.9          | 86.1        |
| Avg. Positive Combination F1 | 83.4          | 79.6        |

Table 1: Agreement scores using our adaptation of F1 score to allow for partial-match.

researcher, these were used for the test set. The process is illustrated in the bottom part of Figure 2.

## 2.6 Inter-annotator Agreement

During the course of the task we calculated inter-annotator agreement multiple times to identify cases of disagreement and provide feedback to annotators. Each time, a set of 25 instances were randomly selected and assigned to all annotators. Agreement was calculated based on a pairwise F1 measure (with some modifications as described in §2.3) and averaged over all pairs of annotators (see discussion of alternative metrics in Appendix A.2).

Final agreement numbers, in Table 1, are satisfactory (Aroyo and Welty, 2013; Araki et al., 2018).

## 2.7 Resulting Dataset

The dataset consists of 1634 annotated abstracts,<sup>6</sup> split into 1362 train and 272 test instances. These

<sup>6</sup>This is a similar size to existing human-labeled biomedical relation extraction datasets, such as BioCreative V CDR (Li et al., 2016), which has 1500 abstracts annotated, BioCreative VI (Krallinger et al., 2017), which has 2432 abstracts, and DDI (Herrero-Zazo et al., 2013), which has 714 abstracts.



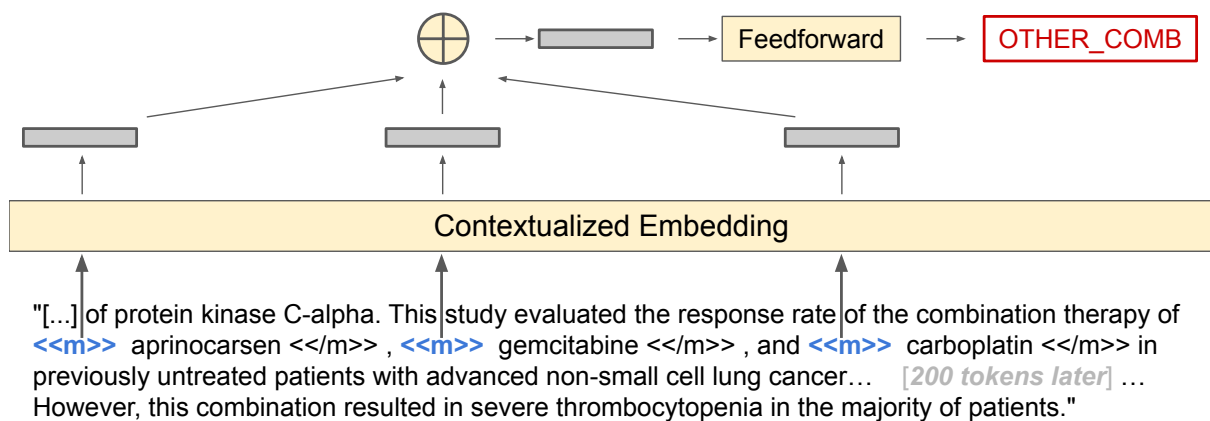


Figure 3: Our baseline architecture, adapted from the PURE model (Zhong and Chen, 2021)

include 1248 relations; 838 are POS\_COMB and 410 are OTHER\_COMB (with the same label ratio in the train and test sets). 591 sentences contain no drug combination, 877 contain one relation (either POS\_COMB or OTHER\_COMB), and 166 contain two or more different combinations. Among annotated relations, 900 are binary, 226 are 3-ary, 69 are 4-ary, and 53 are 5-ary or more.

For each instance in the resulting dataset we include the context-required indication provided by the annotators. In 835 out of 1248 relations the annotator marked the context as needed which is 67% of the time, showing the importance of the context in the DCE task.

### 3 Experiments

#### 3.1 Baseline Model Architecture

We establish a baseline model to measure the difficulty of our dataset and reveal areas for improvement. For our underlying baseline model architecture, we adopt the PURE architecture from Zhong and Chen (2021), which is state-of-the-art on several relation classification benchmarks, including the SciERC binary scientific RE dataset (Luan et al., 2018). The PURE architecture, designed for 2-ary and 3-ary relation extraction, consists of three components. First, special “entity marker” tokens are inserted around all entities in a candidate relation. Next, these marker tokens are encoded with a contextualized embedding model. Finally, the entity marker embeddings are concatenated and fed to a feedforward layer for prediction.

Unlike the original PURE architecture, we consider the more challenging case of extracting relations of variable arity. To support this setting, we *average* the entity marker tokens in a relation rather

than concatenate. The final baseline model architecture is shown in Figure 3. For the contextual embedding component of this architecture, we experiment with four different pretrained scientific language understanding models (SciBERT (Beltagy et al., 2019), BlueBERT (Peng et al., 2019), PubmedBERT (Gu et al., 2020), and BioBERT (Lee et al., 2020)). During training, we only finetune the final \*BERT layer. We train each model architecture for 10 epochs on a single NVIDIA Tesla T4 GPU with 15GB of GPU memory, which takes roughly 7 hours to train for each model.

To our knowledge, there are no other models designed for variable-length  $N$ -ary relation extraction, so we consider no other baselines.

#### 3.2 Domain-Adaptive Pretraining

Our baseline model architecture relies heavily on a pretrained contextual embedding model to provide discriminative features to the relation classifier. Gururangan et al. (2020) showed that continued domain-adaptive pretraining almost always leads to significantly improved downstream task performance. Following this paradigm, we performed continued domain-adaptive pretraining (“DAPT”) on our contextual embedding models.

We acquired in-domain pretraining data using the same procedure used to collect data for annotation: running a SPIKE query against PubMed to find abstracts containing multiple drug names and a “trigger phrase” (from the list in Appendix A.3). This query resulted in 190K unique abstracts. We do not include any paragraphs from our annotated dataset. We then perform domain-adaptive training against this dataset using the Huggingface Transformers library. We train for 10 epochs using a learning rate of  $5e-4$ ,

| Model       | Positive Combination F1  |                          | Any Combination F1       |                          |
|-------------|--------------------------|--------------------------|--------------------------|--------------------------|
|             | Exact Match              | Partial Match            | Exact Match              | Partial Match            |
| Human-Level | 79.6                     | 83.4                     | 86.1                     | 88.9                     |
| Rule-based  | 31.8                     | 45.6                     | 39.1                     | 57.4                     |
| SciBERT     | 44.6 ( $\pm$ 4.6)        | 55.0 ( $\pm$ 5.9)        | 50.2 ( $\pm$ 1.9)        | 63.6 ( $\pm$ 2.7)        |
| w/ DAPT     | 54.8 ( $\pm$ 3.2)        | 63.6 ( $\pm$ 2.0)        | 61.8 ( $\pm$ 2.7)        | 72.8 ( $\pm$ 2.1)        |
| BlueBERT    | 41.2 ( $\pm$ 4.8)        | 51.7 ( $\pm$ 6.0)        | 47.3 ( $\pm$ 4.2)        | 59.9 ( $\pm$ 6.2)        |
| w/ DAPT     | 56.6 ( $\pm$ 2.3)        | 63.5 ( $\pm$ 3.1)        | 64.2 ( $\pm$ 2.6)        | 74.7 ( $\pm$ 2.7)        |
| PubmedBERT  | 50.7 ( $\pm$ 5.5)        | 59.6 ( $\pm$ 5.8)        | 55.9 ( $\pm$ 3.2)        | 66.7 ( $\pm$ 3.8)        |
| w/ DAPT     | <b>61.8</b> ( $\pm$ 5.1) | <b>67.7</b> ( $\pm$ 4.8) | <b>69.4</b> ( $\pm$ 1.7) | <b>77.5</b> ( $\pm$ 2.2) |
| BioBERT     | 45.4 ( $\pm$ 3.7)        | 55.8 ( $\pm$ 2.2)        | 46.7 ( $\pm$ 3.6)        | 58.3 ( $\pm$ 5.1)        |
| w/ DAPT     | 56.0 ( $\pm$ 6.5)        | 63.5 ( $\pm$ 7.5)        | 65.6 ( $\pm$ 1.8)        | 75.7 ( $\pm$ 2.2)        |

Table 2: Comparing different foundation models (with and without continued domain-adaptive pretraining) on Exact-Match and Partial-Match relation extraction metrics. Mean score from 4 different random seeds is reported, and standard deviation is computed across seeds.

345 finetuning all \*BERT layers and using the same  
346 optimization parameters specified by Gururangan  
347 et al. (2020). This pretraining took roughly 8 hours  
348 using four 15GB NVIDIA Tesla T4 GPUs.

### 3.3 Relation Prediction

350 To apply the model to drug combination extraction,  
351 we reduce the RE task to an RC task by consider-  
352 ing all subsets of drug combinations in a sentence,  
353 treating each one as a separate classification input,  
354 and combining the predictions.

355 This poses two challenges: there may be a large  
356 number of candidate relations for a given document,  
357 and each relation is classified independently despite  
358 the combinatorial structure. To handle these issues,  
359 we use a greedy heuristic of choosing the smallest  
360 set of disjoint relations whose union covers as many  
361 drug entities as possible in the sentence. We do  
362 this iteratively: at each step, we choose the largest  
363 predicted relation that does not contain any drugs  
364 found in the relations chosen at previous iterations.

365 This greedy heuristic favors large (high arity)  
366 relations. Nonetheless, we empirically find this  
367 method is helpful for extracting high-precision  
368 drug combinations from our model architecture.

### 3.4 Rule-based baseline

370 To further validate that the trigger words do not in-  
371 troduce bias to our task, we consider an additional  
372 baseline based on the following rule: if a trigger  
373 word is found in the same sentence with multiple  
374 drugs, this set of drugs is tagged as POS\_COMB.

## 4 Results

### 4.1 Effect of Pretrained LMs and Domain-Adaptive Pretraining

375 We show results of our baseline model architec-  
376 tures in Table 2. For each model, we report the  
377 mean and standard deviation of each metric over  
378 four identical models trained with different seeds.<sup>7</sup>  
379 Among the four base scientific language under-  
380 standing models in our experiments, we observe  
381 PubmedBERT to be the strongest on every metric.  
382 We additionally find that domain-adaptive pretrain-  
383 ing provides significantly improvements for every  
384 base model, consistently giving 5-10 points of im-  
385 provement on Positive Combination F1 score. The  
386 value of domain-adaptive pretraining supports our  
387 observation that encoding domain knowledge is  
388 critical to solving this new task.

389 The rule-based approach underperformed all  
390 learned models (30 F1 points under our strongest  
391 model, PubmedBERT-DAPT). This shows this task  
392 cannot be reduced to keyword identification.

### 4.2 Qualitative Error Analysis

393 We identify classes of challenges that make this  
394 task difficult, both in terms of human annotation  
395 and machine prediction.

396 **Coordination Ambiguity:** A known linguistic  
397 challenge is the ambiguity that stems from vague  
398 coordination. In cases where explicit combination  
399 words (e.g. combination, plus, together with, etc)  
400 are not used, it may be unclear whether two drugs  
401 are being used together or separately. For example  
402 in “*These findings may help clinicians identify pa-  
403 tients for whom acamprosate and naltrexone may*”  
404  
405  
406  
407

<sup>7</sup>Seeds used are 2021, 2022, 2023, and 2024

| Model                       | Positive Combination F1 |                   | Any Combination F1 |                   |
|-----------------------------|-------------------------|-------------------|--------------------|-------------------|
|                             | Exact Match             | Partial Match     | Exact Match        | Partial Match     |
| No Extra-sentential Context | 63.4 ( $\pm$ 0.6)       | 68.5 ( $\pm$ 1.1) | 69.7 ( $\pm$ 1.3)  | 76.8 ( $\pm$ 1.7) |
| 1 Sentence of Context       | 63.9 ( $\pm$ 2.3)       | 69.4 ( $\pm$ 3.5) | 71.9 ( $\pm$ 1.1)  | 78.6 ( $\pm$ 1.8) |
| 2 Sentences of Context      | 61.9 ( $\pm$ 9.0)       | 67.6 ( $\pm$ 9.2) | 70.1 ( $\pm$ 2.3)  | 77.9 ( $\pm$ 1.3) |
| 3 Sentences of Context      | 65.2 ( $\pm$ 2.3)       | 72.4 ( $\pm$ 1.3) | 70.8 ( $\pm$ 1.7)  | 78.7 ( $\pm$ 1.2) |

Table 3: The effect of extra-sentential context on model performance.  $n$  sentences are included on each side of the relation-bearing sentence. Mean and standard deviation of each metric are reported over 4 different random seeds.

408 *be most beneficial*” it is unclear if *acamprosate* and  
409 *naltrexone* are being described in combination or  
410 as independent treatments, leading to either a POS  
411 label for the former or NO\_COMB for the latter.

412 **Numerical and Relative Reasoning:** In some  
413 cases, the effect of a treatment is described in rela-  
414 tive or numerical terms, rather than an absolute  
415 claim. Consider the example, “*The infection rate*  
416 *in the control group was 3.5% and in the treated*  
417 *group 0.5%.*”. Here, the reader must compare the  
418 control vs experimental groups and deduce that the  
419 experimental outcome is positive, because the treat-  
420 ment yields a lower infection rate.

421 **Domain Knowledge:** Similarly, classifying rela-  
422 tions in this dataset may require an understand-  
423 ing of domain knowledge. In “*Growth inhibition*  
424 *and apoptosis were significantly higher in BxPC-3,*  
425 *HPAC, and PANC-1 cells treated with celecoxib*  
426 *and erlotinib than cells treated with either cele-*  
427 *coxib or erlotinib*”, one must understand that hav-  
428 ing higher values of *Growth inhibition and apopto-*  
429 *sis* in specific cells is a positive outcome, in order  
430 to classify this combination as positive.

431 **Context related Complications:** The following  
432 are kinds of complications found when the evi-  
433 dence lies in the wider part of the context.

434 Coreference: Anaphoric or coreferential reasoning  
435 is sometimes needed to understand the efficacy of  
436 the combination e.g. “*it was demonstrated that*  
437 *they could be combined with acceptable toxicity.*”.

438 Contradicting Evidence: the reader often must in-  
439 fer a conclusion given opposing claims within a  
440 given abstract. This can happen as combinations  
441 can be referred as e.g. *toxic but effective*.

442 Long Distance: The target sentence can be far—up  
443 to 41 sentences apart—from the evidence sentence,  
444 making it difficult for even humans to process.

### 4.3 Quantitative Error Analysis 445

446 To probe this task, we analyze the performance  
447 of our strongest model—the one using a Pubmed-  
448 BERT base model tuned with domain-adaptive  
449 pretraining—along different partitions of test data.  
450 We trained with four random seeds and perform  
451 comparisons using a paired multi-bootstrap hypoth-  
452 esis test where bootstrap samples are generated by  
453 sampling hierarchically over the four random seeds  
454 and the subsets of the test set (Sellam et al., 2021).  
455 We use 1000 bootstrap samples in each test.

#### 4.3.1 Do models leverage context effectively? 456

457 Each relation in our dataset consists of entities con-  
458 tained within a single sentence, but labeling the  
459 relation frequently requires extra-sentential con-  
460 text to make a decision. In our dataset, annota-  
461 tors record whether or not each relation requires  
462 paragraph-level context to label, reporting that 67%  
463 of drug combinations required context to annotate.

464 To understand the extent that models make use  
465 of paragraph-level context, we trained and evalu-  
466 ated our PubmedBERT-based model using varying  
467 amounts of extra-sentential context around the sen-  
468 tence containing drug entities. In Table 3, we see  
469 that adding context provides nearly identical perfor-  
470 mance to training a model with no extra-sentential  
471 context at all, with differences rarely exceeding one  
472 standard deviation of F1 score.

473 However, we see increased variability in “Pos-  
474 itive Combination F1” performance when extra-  
475 sentential context is used. To explain this, recall  
476 from §2.1 that determining the *efficacy* of a drug  
477 combination often requires paragraph-level context  
478 for annotators, while identifying *any combination*  
479 usually requires no context. From qualitative analy-  
480 sis of attention maps, we observed that our models  
481 are not able to consistently identify the salient parts  
482 of paragraph-level context, potentially causing in-  
483 stability with larger inputs.

484 These results suggest ample room for improve-  
485 ment in extracting document-level evidence. This

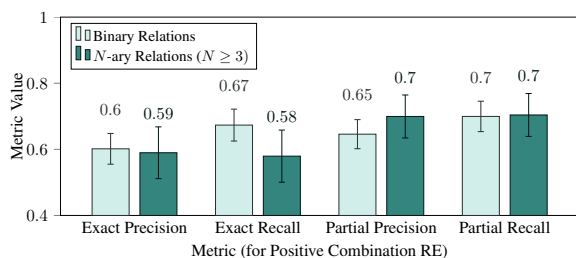


Figure 4: Comparing models performance on binary vs higher-order  $N$ -ary relations, averaged over 4 seeds of the PubmedBERT-DAPT model. No consistent significant differences were observed;  $p$ -values for these comparisons are 0.456, 0.149, 0.240, and 0.276.

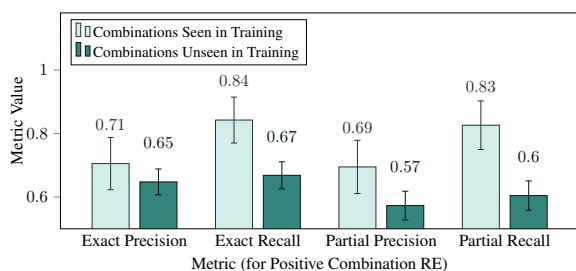


Figure 5: Comparing relation extraction on test set drug combinations that are observed in the training set or not, using the PubmedBERT-DAPT model. Paired multi-bootstrap test  $p$ -values for these four comparisons are 0.262, 0.025, 0.103, and 0.009, respectively.

486 makes our dataset a potentially useful benchmark  
487 for document-level language understanding.

### 4.3.2 Binary vs. higher-arity relations

488 Given that our dataset is the first relation extrac-  
489 tion dataset with variable-arity relations, do higher-  
490 order relations pose a particular challenge for our  
491 models? To answer this, we separate all predicted  
492 and ground truth relations for the test set into bi-  
493 nary relations and higher-arity relations. We then  
494 report precision among each subset of predicted  
495 relations and recall among each subset of ground  
496 truth relations. We perform this experiment across  
497 four different model seeds, and report results in  
498 aggregate using a paired multi-bootstrap procedure.  
499 In the results in Figure 4, we see no consistent sig-  
500 nificant differences between models of different  
501 arities, suggesting that our technique of computing  
502 relation representations by averaging entity repre-  
503 sentations scales well to higher-order relations.  
504

### 4.3.3 Generalizing to new drug combinations

505 How well can relation extraction models classify  
506 drug combinations not seen during training? Sim-  
507 ilar to the setup in §4.3.2, we divide all predicted  
508 and ground truth relations for the test set into the

510 set of drug combinations which are also annotated  
511 in our training set, and the set that have not been. In  
512 our dataset, over 80% of annotated test set relations  
513 are not found in the training set.

514 In Figure 5, performance is consistently better  
515 for relations observed in the training set than for  
516 unseen relations, by a margin of 10-15 points. Re-  
517 call, in particular, is significantly worse for rela-  
518 tions unseen during training (at 95% confidence),  
519 and precision is potentially also worse. Consider-  
520 ing that unseen drug combinations are practically  
521 more valuable than already-known combinations,  
522 improving generalization to new combinations is a  
523 critical area of improvement for this task.

## 5 Related Work

524 The DDI dataset (Herrero-Zazo et al., 2013) is the  
525 only work to our knowledge that annotates drug  
526 interactions for text mining. However, it funda-  
527 mentally differs from our dataset in the type of  
528 annotations provided: the DDI annotates the type  
529 of discourse context in which a drug combination is  
530 mentioned, without providing explicit information  
531 about combination efficacy. In contrast, our dataset  
532 is focused on semantically classifying the efficacy  
533 of drug combinations as stated in text.  
534

535 Other RE datasets exist in the biomedical field  
536 (Peng et al., 2017; Li et al., 2016; Wu et al., 2019;  
537 Krallinger et al., 2017), but do not focus on drug  
538 combinations. Similarly, several RE datasets tackle  
539 the  $N$ -arity problem in the scientific domain (Peng  
540 et al., 2017; Jain et al., 2020; Kardas et al., 2020;  
541 Hou et al., 2019), and in the non-scientific domain  
542 (Akimoto et al., 2019; Nguyen et al., 2016), how-  
543 ever, **all of them consider a fixed choice of  $N$ .**

## 6 Conclusions

544 We present a new resource for drug combination  
545 and efficacy identification. We establish base-  
546 line models that achieve promising results but re-  
547 veal clear areas for improvement. Beyond the  
548 immediate, application-ready value of this task,  
549 this task poses unique relation extraction chal-  
550 lenges as the first dataset containing variable-  
551 arity relations. We also highlight challenges with  
552 document-level representation learning and incor-  
553 porating domain knowledge. We encourage oth-  
554 ers to participate in this task, and our dataset  
555 and modeling code are all available to the public  
556 at [https://anonymous.4open.science/](https://anonymous.4open.science/r/drug-synergy-models--C8B7)  
557 [r/drug-synergy-models--C8B7](https://anonymous.4open.science/r/drug-synergy-models--C8B7).  
558



559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
  
569  
570  
571  
572  
573  
574  
  
575  
576  
577  
  
578  
579  
580  
581  
582  
583  
  
584  
585  
586  
587  
588  
589  
590  
591  
  
592  
593  
594  
595  
596  
  
597  
598  
599  
600  
  
601  
602  
603  
  
604  
605  
606  
607  
  
608  
609  
610  
611  
  
612  
613  
614

## References

Kosuke Akimoto, Takuya Hiraoka, Kunihiko Sadamasa, and Mathias Niepert. 2019. [Cross-sentence n-ary relation extraction using lower-arity universal schemas](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6225–6231, Hong Kong, China. Association for Computational Linguistics.

Jun Araki, Lamana Mulaffer, Arun Pandian, Yukari Yamakawa, Kemal Oflazer, and Teruko Mitamura. 2018. Interoperable annotation of events and event relations across domains. In *Proceedings 14th Joint ACL-ISO Workshop on Interoperable Semantic Annotation*, pages 10–20.

Lora Aroyo and Chris Welty. 2013. Measuring crowd truth for medical relation extraction. In *2013 AAAI Fall Symposium Series*.

John A Bartlett, Michael J Fath, Ralph Demasi, Ashwaq Hermes, Joseph Quinn, Elsa Mondou, and Franck Rousseau. 2006. An updated systematic overview of triple combination therapy in antiretroviral-naive hiv-infected adults. *Aids*, 20(16):2051–2064.

Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. [SciBERT: A pretrained language model for scientific text](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3615–3620, Hong Kong, China. Association for Computational Linguistics.

Y Bhusal, CM Shiohira, and N Yamane. 2005. Determination of in vitro synergy when three antimicrobial agents are combined against mycobacterium tuberculosis. *International journal of antimicrobial agents*, 26(4):292–297.

Jennifer S Carew, Francis J Giles, and Steffan T Nawrocki. 2008. Histone deacetylase inhibitors: mechanisms of cell death and promise in combination cancer therapy. *Cancer letters*, 269(1):7–17.

Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46.

VT DeVita, RC Young, and GP Canellos. 1975. [Combination versus single agent chemotherapy: a review of the basis for selection of drug treatment of cancer](#). *Cancer*, 35(1):98–110.

Richard T Eastman and David A Fidock. 2009. Artemisinin-based combination therapies: a vital tool in efforts to eliminate malaria. *Nature Reviews Microbiology*, 7(12):864–874.

Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological bulletin*, 76(5):378.

Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2020. [Domain-specific language model pretraining for biomedical natural language processing](#). 615  
616  
617  
618  
619

Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020. [Don’t stop pretraining: Adapt language models to domains and tasks](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360, Online. Association for Computational Linguistics. 620  
621  
622  
623  
624  
625  
626  
627

Xu Han, Hao Zhu, Pengfei Yu, Ziyun Wang, Yuan Yao, Zhiyuan Liu, and Maosong Sun. 2018. [Fewrel: A large-scale supervised few-shot relation classification dataset with state-of-the-art evaluation](#). 628  
629  
630  
631

Andrew F Hayes and Klaus Krippendorff. 2007. Answering the call for a standard reliability measure for coding data. *Communication methods and measures*, 1(1):77–89. 632  
633  
634  
635

María Herrero-Zazo, Isabel Segura-Bedmar, Paloma Martínez, and Thierry Declerck. 2013. The DDI corpus: An annotated corpus with pharmacological substances and drug–drug interactions. *Journal of biomedical informatics*, 46(5):914–920. 636  
637  
638  
639  
640

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 *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5203–5213, Florence, Italy. Association for Computational Linguistics. 641  
642  
643  
644  
645  
646  
647  
648

Aleksandr Ianevski, Rouan Yao, Svetlana Biza, Eva Zusinaite, Andres Mannik, Gaily Kivi, Anu Planken, Kristiina Kurg, Eva-Maria Tombak, Mart Ustav, et al. 2020. Identification and tracking of antiviral drug combinations. *Viruses*, 12(10):1178. 649  
650  
651  
652  
653

Aleksandr Ianevski, Rouan Yao, Hilde Lysvand, Gunnveig Grødeland, Nicolas Legrand, Valentyn Oksenysh, Eva Zusinaite, Tanel Tenson, Magnar Bjørås, and Denis E. Kainov. 2021. [Nafamostat–interferon- combination suppresses sars-cov-2 infection in vitro and in vivo by cooperatively targeting host tmprss2](#). *Viruses*, 13(9). 654  
655  
656  
657  
658  
659  
660

Sarthak Jain, Madeleine van Zuylen, Hannaneh Hajishirzi, and Iz Beltagy. 2020. [Scirex: A challenge dataset for document-level information extraction](#). 661  
662  
663

Marcin Kardas, Piotr Czapla, Pontus Stenetorp, Sebastian Ruder, Sebastian Riedel, Ross Taylor, and Robert Stojnic. 2020. [AxCell: Automatic extraction of results from machine learning papers](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8580–8594, Online. Association for Computational Linguistics. 664  
665  
666  
667  
668  
669  
670  
671

|     |  |     |
|-----|--|-----|
| 672 | Itay Katzir, Murat Cokol, Bree B Aldridge, and Uri Alon. 2019. Prediction of ultra-high-order antibiotic combinations based on pairwise interactions. <i>PLoS computational biology</i> , 15(1):e1006774.  | 727 |
| 673 |  | 728 |
| 674 |  | 729 |
| 675 |  | 730 |
| 676 | Martin Krallinger, Obdulia Rabal, Saber Ahmad Akhondi, Martín Pérez Pérez, Jesus Santamaría, Gael Pérez Rodríguez, Georgios Tsatsaronis, Ander Intxaurreondo, José Antonio Baso López, Umesh K. Nandal, Erin M. van Buel, Anjana Chandrasekhar, Marleen Rodenburg, Astrid Læg Reid, Marius A. Doornenbal, Julen Oyarzábal, Anália Lourenço, and Alfonso Valencia. 2017. Overview of the BioCreative VI chemical-protein interaction track. | 731 |
| 677 |  | 732 |
| 678 |  | 733 |
| 679 |  | 734 |
| 680 |  | 735 |
| 681 |  | 736 |
| 682 |  | 737 |
| 683 |  | 738 |
| 684 |  | 739 |
| 685 | Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. <i>Bioinformatics</i> , 36(4):1234–1240.   | 740 |
| 686 |  | 741 |
| 687 |  | 742 |
| 688 |  | 743 |
| 689 |  | 744 |
| 690 |  | 745 |
| 691 | Jiao Li, Yueping Sun, Robin Johnson, Daniela Sciaky, Chih-Hsuan Wei, Robert Leaman, Allan Peter Davis, Carolyn Mattingly, Thomas Wieggers, and Zhiyong lu. 2016. BioCreative V CDR task corpus: a resource for chemical disease relation extraction. <i>Database</i> , 2016:baw068.  | 746 |
| 692 |  | 747 |
| 693 |  | 748 |
| 694 |  | 749 |
| 695 |  | 750 |
| 696 |  | 751 |
| 697 | 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.  | 752 |
| 698 |  | 753 |
| 699 |  | 754 |
| 700 |  | 755 |
| 701 |  | 756 |
| 702 |  | 757 |
| 703 |  | 758 |
| 704 | Ines Montani and Matthew Honnibal. 2018. Prodigy: A new annotation tool for radically efficient machine teaching. <i>Artificial Intelligence</i> , to appear.  | 759 |
| 705 |  | 760 |
| 706 |  | 761 |
| 707 | Kiem-Hieu Nguyen, Xavier Tannier, Olivier Ferret, and Romaric Besançon. 2016. A dataset for open event extraction in English. In <i>Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)</i> , pages 1939–1943, Portorož, Slovenia. European Language Resources Association (ELRA).  | 762 |
| 708 |  | 763 |
| 709 |  | 764 |
| 710 |  | 765 |
| 711 |  | 766 |
| 712 |  | 767 |
| 713 |  | 768 |
| 714 | Danna Niezni, Yakir Amrusi, Shaked Launer-Wachs, Yuval Harris, Hagit Sason, and Yosi Shamay. 2021. High complexity combination therapy planning. <i>in submission</i> .  | 769 |
| 715 |  | 770 |
| 716 |  | 771 |
| 717 |  | 772 |
| 718 | Nanyun Peng, Hoifung Poon, Chris Quirk, Kristina Toutanova, and Wen tau Yih. 2017. Cross-sentence n-ary relation extraction with graph lstms.  | 773 |
| 719 |  | 774 |
| 720 |  | 775 |
| 721 | Yifan Peng, Shankai Yan, and Zhiyong Lu. 2019. Transfer learning in biomedical natural language processing: An evaluation of BERT and ELMo on ten benchmarking datasets. In <i>Proceedings of the 2019 Workshop on Biomedical Natural Language Processing (BioNLP 2019)</i> , pages 58–65.   | 776 |
| 722 |  | 777 |
| 723 |  | 778 |
| 724 |  | 779 |
| 725 |  | 780 |
| 726 |  | 781 |
|     | Yogita Rochlani, Mohammed Hasan Khan, Maciej Banach, and Wilbert S Aronow. 2017. Are two drugs better than one? a review of combination therapies for hypertension. <i>Expert opinion on pharmacotherapy</i> , 18(4):377–386.  | 782 |
|     |  | 783 |
|     |  | 784 |
|     | Shachar Rosenman, Alon Jacovi, and Yoav Goldberg. 2020. Exposing shallow heuristics of relation extraction models with challenge data.   | 785 |
|     |  | 786 |
|     |  | 787 |
|     | Ofer Sabo, Yanai Elazar, Yoav Goldberg, and Ido Dagan. 2021. Revisiting few-shot relation classification: Evaluation data and classification schemes.  | 788 |
|     |  | 789 |
|     |  | 790 |
|     | Thibault Sellam, Steve Yadlowsky, Jason Wei, Naomi Saphra, Alexander D’Amour, Tal Linzen, Jasmijn Bastings, Iulia Turc, Jacob Eisenstein, Dipanjan Das, Ian Tenney, and Ellie Pavlick. 2021. The multiberts: Bert reproductions for robustness analysis. <i>ArXiv</i> , abs/2106.16163.  | 791 |
|     |  | 792 |
|     |  | 793 |
|     | Micah Shlain, Hillel Taub-Tabib, Shoval Sadde, and Yoav Goldberg. 2020. Syntactic search by example. In <i>ACL</i> .   | 794 |
|     |  | 795 |
|     |  | 796 |
|     | Adam J Shuhendler, Richard Y Cheung, Janet Manias, Allegra Connor, Andrew M Rauth, and Xiao Yu Wu. 2010. A novel doxorubicin-mitomycin c co-encapsulated nanoparticle formulation exhibits anti-cancer synergy in multidrug resistant human breast cancer cells. <i>Breast cancer research and treatment</i> , 119(2):255–269.   | 797 |
|     |  | 798 |
|     |  | 799 |
|     | Hillel Taub-Tabib, Micah Shlain, Shoval Sadde, Dan Lahav, Matan Eyal, Yaara Cohen, and Yoav Goldberg. 2020. Interactive extractive search over biomedical corpora. <i>arXiv preprint arXiv:2006.04148</i> .  | 800 |
|     |  | 801 |
|     |  | 802 |
|     | Ernesto Wasserman, William Sutherland, and Esteban Cvitkovic. 2001. Irinotecan plus oxaliplatin: a promising combination for advanced colorectal cancer. <i>Clinical colorectal cancer</i> , 1(3):149–153.   | 803 |
|     |  | 804 |
|     |  | 805 |
|     | Ye Wu, Ruibang Luo, Henry CM Leung, Hing-Fung Ting, and Tak-Wah Lam. 2019. Renet: A deep learning approach for extracting gene-disease associations from literature. In <i>International Conference on Research in Computational Molecular Biology</i> , pages 272–284. Springer.  | 806 |
|     |  | 807 |
|     |  | 808 |
|     | Yuhao Zhang, Victor Zhong, Danqi Chen, Gabor Angeli, and Christopher D. Manning. 2017. Position-aware attention and supervised data improve slot filling. In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 35–45, Copenhagen, Denmark. Association for Computational Linguistics.  | 809 |
|     |  | 810 |
|     |  | 811 |
|     | Zexuan Zhong and Danqi Chen. 2021. A frustratingly easy approach for entity and relation extraction. In <i>North American Association for Computational Linguistics (NAACL)</i> .  | 812 |
|     |  | 813 |
|     |  | 814 |
|     |  | 815 |
|     |  | 816 |
|     |  | 817 |

## A Appendices

780

### A.1 Annotation Guidelines

781

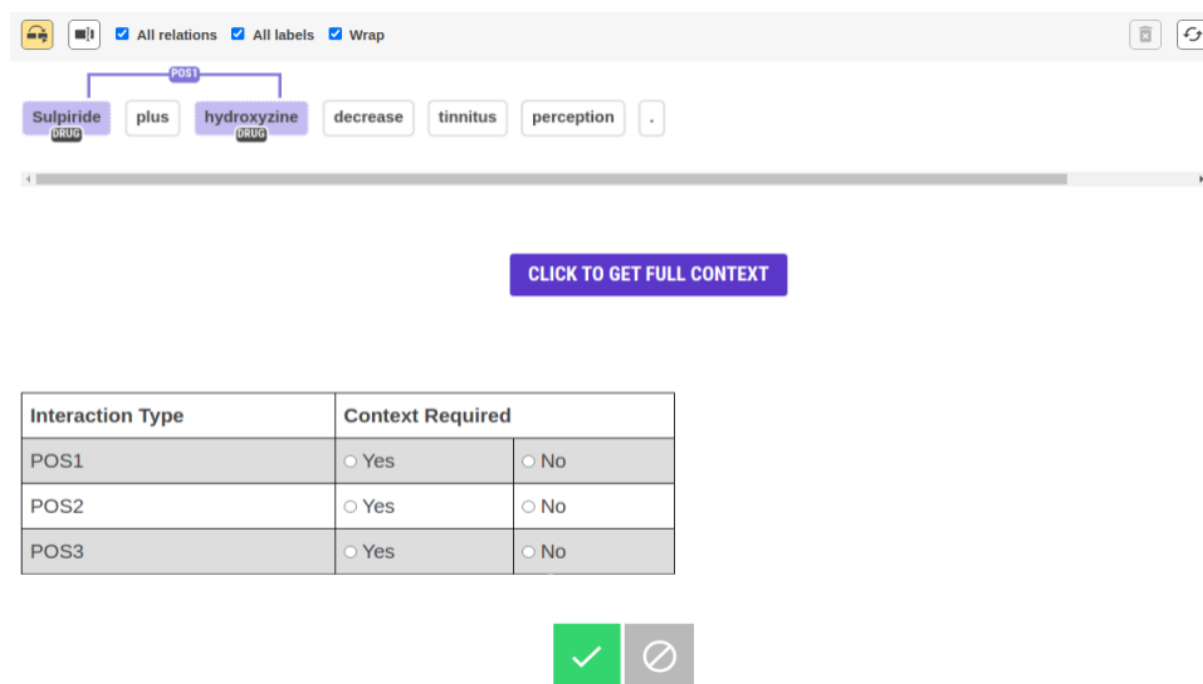


Figure 6: Annotation instance in the Prodigy environment. The screen is constructed of the sentence where they should mark relations, a button to show the full context and a selection per relation to indicate the necessity of the context.

For this task we recruited 7 annotators all studying for advanced degrees in biomedical engineering. The annotators were paid by their advisor, an amount that is standard for annotation projects in their country of residence. All participating annotators were provided with annotation guidelines. The guidelines specified how the annotation process should be carried out and provided definitions and examples for the different labels used. As the task progressed, the guidelines were also expanded to include discussion of frequently encountered issues.

782

783

784

785

786

787

For a given instance, such as presented in the top of Figure 6 the annotator needs to first recognize any missing drugs and mark them, and then label any interactions they find among the drugs. In case they need to consult a wider context they can press on a 'show more context' button and a text box with the wider context will appear. This context can be again hidden by clicking the same button if needed. Lastly, in the bottom of the sample page, we present a table with questions regarding the necessity of using the context.

788

789

790

791

792

Then the annotator should decide if they need to ignore the current sample or to complete the current instance and accept it, by pressing the accept and ignore buttons.

793

794

The annotators are instructed as follows. They should read the sentence carefully, and try to answer a two phase question to themselves. First, if the drugs are mentioned in any form of combination or they should be given separately. Second, if indeed the annotator recognized the drugs as a combination can they determine the efficacy of the combination by the sole sentence.

795

796

797

798

In case they can not determine the efficacy they are instructed to press on the 'get more context' button and read the entire context in order to determine what is the correct efficacy. If after reading the context they can still not determine the efficacy then the label of the interaction should be OTHER\_COMB (aside from negative label experimentation mentioned in Footnote 2). Otherwise it should be POS\_COMB. In case that they recognized that there is no combination between the drugs in the sentence then they should not use any label and simply accept the current instance. Then they should answer the context related questions for the POS\_COMB label in order to signal if the context was needed.

799

800

801

802

803

804

805

While reading the sentence if the annotators find unmarked drugs they can mark them before continuing

806

to the interaction-labeling phase and treat them the same as the other drugs, but, it is not required to mark a word as drug in order to use it in an interaction. If a drug is marked in a wrong manner they should try and fix it, e.g. the span of the drug is incorrect.

In order to achieve more consistent and accurate annotations, they are also instructed to annotate all the interactions that they can find in a given sentence. They should always use the *accept* button even if there are no interactions in the sentence. Only in cases where they want to skip a sentence (e.g. when there is an inherent problem with it) or leave it for a future discussion they should use the *ignore* button. An interaction can occur between more than two drugs, if so they should notice that they don't need each pair from this group to have a marked interaction, as long as they all connect to the same graph. e.g. "Drugs A, B and C are synergistic." connecting A to B and B to C is sufficient, no need to connect drug A to drug C. Each interaction should be marked with a different tag (POS\_COMB1, POS\_COMB2..., OTHER\_COMB1, OTHER\_COMB2...).

### A.2 Evaluation Metric Discussion

For measuring the agreement, we chose to use our adaptation of F1 score and not other common metrics such as Cohen's Kappa (Cohen, 1960) or one of its variations (e.g. Fleiss's Kappa (Fleiss, 1971) and Krippendorff's Alpha (Hayes and Krippendorff, 2007)). These metrics expect a setup where the *relation* candidates are already marked and the task is only to label them – a labeling task and not an extraction task. This causes two problems, one is that they inherently do not need to handle a partial match. So if for example there are three drugs in a sentence, the first annotator annotated a relation between drugs A and B, while a second annotator annotated the same relation between drugs A, B and C. So we will either underestimate or overestimate their agreement score if we considered this a mismatch or a match respectively. Moreover, their calculations depend on the *hypothetical agreement by chance* normalization factor, but this will not reflect the difficulty of random choosing in our setup as they ignore the size of the combinatorial set of relation candidates we can possibly have.

### A.3 Trigger List

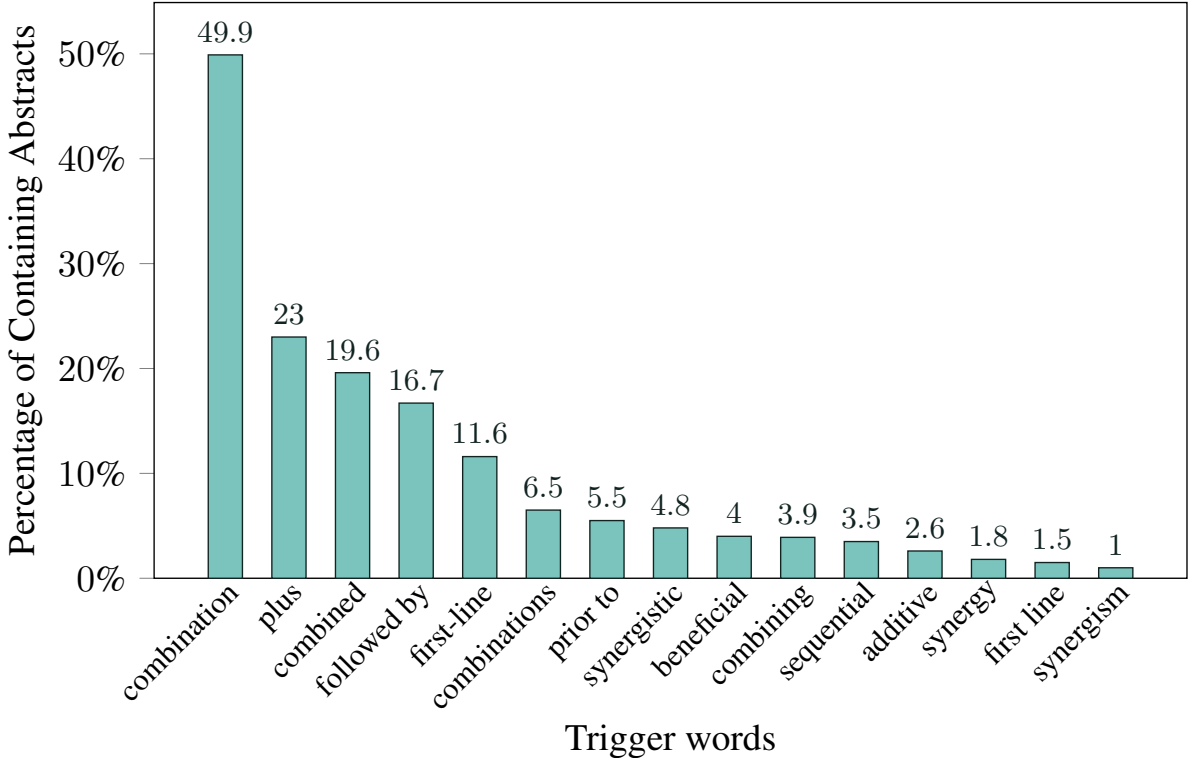


Figure 7: Abstracts percentage including each trigger word (1634 abstracts included; 44 words in the full word list; Words <1% were neglected from the figure).



In [Figure 7](#) we show the triggers that we used in the Spike queries. We show the percentage of abstracts that included each trigger (others under 1%: *conjunction, two-drug, first choice, additivity, combinational, synergetic, simultaneously with, supra-additive, five-drug, combinatory, over-additive, timed-sequential, co-blister, super-additive, synergisms, synergic, synergistical, less-than-additive, greater-than-additive, additivesynergistic, supraadditive, superadditive, overadditive, subadditive, first-choice, 2-drug, sub-additive, more-than-additive, 3-drug*).

832  
833  
834  
835  
836  
837