

# Analysing Temporality in General-Domain Entailment Graphs

Anonymous ACL submission

## Abstract

Entailment Graphs based on open relation extraction run the risk of learning spurious entailments (e.g. *win against*  $\models$  *lose to*) from antonymous predications that are observed with the same entities referring to different times. Previous research has demonstrated the potential of using temporality as a signal to avoid learning these entailments in the sports domain. We investigate whether this extends to the general news domain. Our method introduces a temporal window that is set dynamically for each eventuality using a temporally-informed language model. We evaluate our models on a sports-specific dataset, and ANT – a novel general-domain dataset based on WordNet antonym pairs. We find that whilst it may be useful to reinterpret the Distributional Inclusion Hypothesis to include time for the sports news domain, this does not apply to the general news domain.

## 1 Introduction

The ability to recognise textual entailment and paraphrase is essential to many NLP applications, including open-domain question answering over unstructured data. This setting frequently poses the challenge that the answer to the question is not explicitly stated in the text, and can only be inferred using entailment rules and/or paraphrases. For example, the question might ask “Did Arsenal play Man United last night?” and the post-match report states “Arsenal beat Man United 1-0”. A system that can recognise that *beat*  $\models$  *play* will be able to provide the correct answer (“yes”).

*Entailment Graphs* (Berant et al., 2011, 2015; Hosseini et al., 2018), learned using unsupervised methods applied over large text corpora, have been proposed as a means to support answering such questions. Entailment Graphs comprise nodes representing predicates, and edges representing the entailment relation between them. They can be

learned using the Distributional Inclusion Hypothesis (DIH) (Dagan et al., 1999; Geffet and Dagan, 2005), which can be used to learn entailments between predicates (e.g. *win*  $\models$  *play*) based on their context — co-occurring argument pairs (e.g. (*Man United*, *Arsenal*)).

However, the argument pair-based, atemporal formulation of the DIH does not support the class of predicate pairs that are antonyms and occur frequently within a window of time with the same argument pairs, such as *winning* and *losing* (Guillou et al., 2020). For example, sports teams often play against each other multiple times in a season, likely with different outcomes, so that predicates such as *win against* and *lose to* are both likely to apply to the same sports team argument pairs (e.g. Arsenal, Man United). Consequently, current state-of-the-art methods for learning Entailment Graphs are prone to learning erroneous entailment relations between these pairs of highly correlated predicates (e.g. *to win*  $\models$  *to lose*).

Guillou et al. (2020) propose an algorithm that circumvents this issue by considering argument pair occurrences only when the eventualities temporally overlap. This effectively reinterprets the DIH’s context set as containing both argument pairs and time. They refine Entailment Graph induction for the sports news domain. We extend their work by applying the method to the general news domain, and propose setting different size temporal comparison windows for the different predicates contained in the general domain. We dynamically assign a different window size for each eventuality in the corpus using a temporally-aware language model (Zhou et al., 2020) that predicts the expected duration of the eventuality. We evaluate the Entailment Graphs on the *Sports* dataset of Guillou et al. (2020), and ANT – a novel dataset derived from WordNet (Miller, 1995) antonyms.

We find that refining the DIH’s context to include time is beneficial for the sports news domain,

but that this does not extend to the general news domain. We do, however, identify predicates in legal news as another possible area in which temporal information may be useful for learning Entailment Graphs. Our contributions are: 1) ANT, a novel general-domain entailment dataset based on WordNet antonyms, 2) a novel method for dynamically assigning temporal windows to predicates based on contextualised predictions of eventuality duration, 3) a comparison of the temporal method across domains, revealing that the method works for specific predicate pairs, and 4) an analysis of why the method is successful under specific conditions.

## 2 Background

### 2.1 Entailment Graphs

Entailment Graph Induction uses directional distributional similarity measures to determine whether an entailment relation holds between two predicates  $p$  and  $q$ . Successful measures include the purely directional Weed’s precision score (Weeds and Weir, 2003), and the Balanced Inclusion score (BInc) (Szpektor and Dagan, 2008), which combines both symmetric and directional measures. We use these two scores as our baselines. Both scores are based on the Distributional Inclusion Hypothesis (DIH), which states that  $p$  entails  $q$  if the set of contexts in which  $p$  can be used is *included in* the context set of  $q$  (Dagan et al., 1999; Geffet and Dagan, 2005). When applying the hypothesis to predicates, the context set has mostly been taken to refer to argument pairs (e.g. by Berant et al. (2011) and Hosseini et al. (2018)).

Entailment Graphs have been built for a range of domains, including health (Levy et al., 2014), news (Hosseini et al., 2018), and commonsense (Yu et al., 2020). By focusing on the news domain we are able to leverage two sources of temporal information: the publication dates of the articles and the rich set of temporal expressions within them. Previous work has also considered a number of options for representing nodes in the graphs: typed binary predicates (Berant et al., 2011; Hosseini et al., 2018), Open-IE propositions (Levy et al., 2014), and eventualities (Yu et al., 2020). We use typed predicates, following Hosseini et al. (2018).

### 2.2 Temporality and Entailment Graphs

Guillou et al. (2020) incorporated temporal information into the graph learning framework of Hosseini et al. (2018), extending the local entailment

score computation method to incorporate temporal filtering of eventualities. This reinterprets the DIH to include time in the context set of any given predicate. Unlike in Hosseini et al. (2018) where all eventualities of pairs of predicates that share the same arguments are considered for comparison, Guillou et al. (2020) aim to compare only those eventualities of predicates with shared arguments for which the underlying eventualities are temporally close to each other. The strength of this method is its ability to separate out instances of recurring eventualities, e.g. sports matches that occur between the same pair of teams. Following promising results for the sports domain, we extend the method to the general news domain.

### 2.3 Evaluating Entailment Graphs

Entailment Graphs are typically evaluated using datasets comprised of premise-hypothesis sentence pairs with labels denoting the entailment relation that holds between them. Dataset construction has been framed as a number of manual annotation tasks, e.g. image captioning (Bowman et al., 2015), question answering (Levy and Dagan, 2016), and fact verification (Schmitt and Schütze, 2019).

Evaluating entailments that involve temporality has received less attention. The *FraCas* test suite (Cooper et al., 1996) contains only a small number of temporal examples that are based on entailments between predicates. *TEA* (Kober et al., 2019), which comprises sentence pairs in which temporally ordered predications have varying tense and aspect, does not include non-entailments that can be learned through the temporal separation of eventualities (e.g. outcome predicates *win - lose*). The *Sports* dataset (Guillou et al., 2020) of entailment pairs between paraphrases of the predicates *play*, *win*, *lose*, and *tie*, was developed to address this gap. However, its narrow focus on sports makes it unsuitable for evaluating graphs for the general news domain. This motivates the construction of our general-domain ANT dataset.

### 2.4 Antonym Detection

Related to our work is the field of antonym detection, in which antonyms are distinguished from other semantic relations such as synonymy. We focus on the related but distinct task of Recognizing Textual Entailment (RTE) in the presence of antonymy, which can be seen as a more challenging version of the typical RTE setup. Antonymy detection is evaluated using various datasets, notably

181 the relation classification-style EVALution dataset  
182 (Santus et al., 2015) and PPDB-based dataset of  
183 Rajana et al. (2017), and the multiple-choice GRE  
184 question dataset (Mohammad et al., 2013). To compare  
185 our work to previous Entailment Graph models, we  
186 instead opt for the RTE paradigm, focusing on  
187 sentences containing binary predications. We note  
188 that the labels in ANT can easily be remapped for  
189 the evaluation of antonym detection systems.

### 190 3 Method

#### 191 3.1 Relation Extraction

192 We start by extracting relation triples from a corpus  
193 of news articles. We use MONTEE (Bijl de Vroe et al.,  
194 2021), an open-domain system that uses the RotatingCCG  
195 parser (Stanojević and Steedman, 2019) and extracts  
196 relations consisting of predicates and their arguments  
197 by traversing the resulting CCG dependency graph. For  
198 each sentence we extract all potential *binary relations*  
199 of the form **arg1-predicate-arg2** (e.g. **Arsenal-beat-**  
200 **Man United**)<sup>1</sup>. Arguments, which may be either  
201 Named Entities or general entities (all other nouns  
202 and noun phrases), are mapped to their fine-grained  
203 FIGER types (Ling and Weld, 2012) (e.g. *person*,  
204 *disease*, etc.).

205 We extended MONTEE to add temporal intervals  
206 to binary relations where there is a path in the  
207 dependency graph between the predicate and a temporal  
208 expression in the text. The temporal intervals consist  
209 of the start and end date of the eventuality, and  
210 are derived using SUTime (Chang and Manning, 2012)  
211 – a tool for automatically identifying and resolving  
212 temporal expressions (such as “Monday 7th March  
213 2022”) found in the text, to a calendar date range.  
214 Expressions such as “yesterday” are resolved  
215 relative to the article’s publication date.

#### 217 3.2 Graph Learning with Temporal Filtering

218 To learn Entailment Graphs we use the temporal  
219 filtering method of Guillou et al. (2020). The  
220 temporal filtering method extends the graph learning  
221 framework of Hosseini et al. (2018) by adding a  
222 method to filter the counts of predicate  $p$  according  
223 to whether each eventuality’s time interval overlaps  
224 with any of  $q$ ’s. That is, an eventuality in  $p$   
225 is retained (and counted) if it is temporally close  
226 enough to any eventuality in  $q$ . The goal of this  
227 process is to separate out different instances of

228 recurring eventualities involving the same argu-  
229 ment pairs. For example, to separate out two sports  
230 matches between the same pair of teams which  
231 occur on different dates, and which likely have  
232 different outcomes.

233 The input to the method is the set of typed binary  
234 relations paired with their time intervals. The output  
235 is a set of graphs, one for each pair of FIGER  
236 types found in the set of binary relations. To study  
237 the role of the temporal signal in isolation, we use  
238 only *locally* learned entailment scores, leaving an  
239 investigation of the interaction between temporality  
240 and globalisation with soft constraints (Hosseini  
241 et al., 2018) to future work. See Guillou et al.  
242 (2020) for further details of the temporal filtering  
243 method.

#### 244 3.3 Dynamic Temporal Window

245 Although a uniform temporal window is suitable  
246 for sports matches, which are typically concluded  
247 within a single day, it may be less suitable for other  
248 eventualities. Following the recommendation of  
249 Guillou et al. (2020) we develop a novel method  
250 that applies a dynamic window on a per-predicate  
251 basis to reflect that different eventualities remain  
252 relevant for different lengths of time. For example,  
253 the window around information stating that a person  
254 *is president* should be larger than a report of a  
255 person *visiting a location*.

256 We incorporate a temporally-aware language  
257 model, *TacoLM* (Zhou et al., 2020), and use it  
258 as the basis for per-predicate dynamic windowing.  
259 *TacoLM* predicts the expected duration of an eventuality  
260 using the context provided by the sentence in which  
261 the eventuality mention occurs. For each eventuality  
262 in a sentence it assigns a duration label from the set  
263  $\{seconds, minutes, hours, days, weeks, months, years, decades, centuries\}$ .  
264 In a small number of cases *TacoLM* is unable to make  
265 a prediction, indicated by the *no\_prediction* label<sup>2</sup>.

266 In the uniform window model, each eventuality  
267  $e$  is assigned a temporal interval  $e_t = [t_{start} - w, t_{end} + w]$ ,  
268 where  $t_{start}$  and  $t_{end}$  are predicted using  $SUTime(e)$ ,  
269 and  $w$  is the model’s fixed window size. In the  
270 dynamic window model, we instead assign  $e_t = [t_{start} - map(TLM(e)), t_{end} + map(TLM(e))]$ .  
271 Here  $map(TLM(e))$  is *TacoLM*’s prediction mapped to a concrete duration  
272 value:  $\{seconds, minutes, hours, days\} \mapsto 5$ ,  
273  
274  
275

<sup>1</sup>As we are not concerned with the intersection of temporality and modality, we do not tag relations for modality

<sup>2</sup>249,262 [0.61%] eventuality mentions in the NewsSpike corpus

276 *weeks*  $\mapsto$  15, *months*  $\mapsto$  30, *years*  $\mapsto$  365,  
277 *decades*  $\mapsto$  3,650, *centuries*  $\mapsto$  36,500. That  
278 is, for shorter durations we maintain a uniform win-  
279 drow of 5 days, extending it only for eventualities  
280 with longer durations.

### 281 3.4 Similarity Measures

282 We compute both a symmetric and a directional  
283 temporally-informed similarity measure to learn  
284 entailments, making use of the temporally filtered  
285 counts and PMI scores described in Section 3.2.  
286 We adapted BInc (Szpektor and Dagan, 2008) and  
287 Weeds’ precision (Weeds and Weir, 2003).

288 We compute **Temporal Weed’s precision** using  
289 the temporally-filtered counts. For the BInc-based  
290 measure **Temporal BInc (Ratio PMI)**, we use the  
291 temporal PMI scores. We scale the PMI scores used  
292 to compute BInc. The temporally filtered  $PMI_t$   
293  $= PMI \cdot (c_t/c)$ , i.e. the original PMI multiplied  
294 by the ratio of filtered counts ( $c_t$ ) to regular counts  
295 ( $c$ ). We also include the count-based temporal and  
296 atemporal versions of BInc.

## 297 4 Evaluation

298 We evaluate the Entailment Graphs using two dif-  
299 ferent entailment datasets. 1) the *Sports* dataset  
300 (Guillou et al., 2020) which contains 1,312 entail-  
301 ment pairs, focusing on events that occur between  
302 two sports teams. 2) ANT, a novel dataset based on  
303 WordNet antonym pairs, which addresses the need  
304 for a *general-domain*, RTE-style dataset containing  
305 antonyms.

### 306 4.1 ANT Dataset Construction Overview

307 ANT<sup>3</sup> contains entailment pair examples of the  
308 form *premise*, *hypothesis*, *label*. The premise and  
309 hypothesis take the form of natural English sen-  
310 tences containing a subject, predicate, and object.  
311 The label denotes one of four types of entailment  
312 relation: 1) *Antonym*: non-entailments between  
313 antonymous predicates (e.g. *acquit* - *convict*), 2)  
314 *Directional Entailments*: an antonymous predi-  
315 cate and a related third predicate (e.g. *acquit*  $\models$   
316 *indict*), 3) *Directional Non-Entailments*: the re-  
317 verse of each Directional Entailment (e.g. *indict*  
318  $\not\models$  *acquit*), and 4) *Paraphrases* of each predicate  
319 in the antonym pair (e.g. *acquit* - *absolve*). For  
320 a standard entailment evaluation setup, we map:  
321 (*Antonyms*, *Dir.Non-Entailments*)  $\mapsto$  0 and

<sup>3</sup><https://anonymous-link.com>

(*Paraphrases*, *Dir.Entailments*)  $\mapsto$  1. Our re-  
322 leased dataset contains the original four labels as  
323 these may be useful in future research. 324

325 Dataset construction was semi-automatic. The  
326 manual steps were carried out by two expert anno-  
327 tators: one native, and one fluent English speaker<sup>4</sup>.  
328 Our dataset generation method uses the entailment  
329 relations between manually annotated predicate  
330 *clusters* to generate entailment pairs. By ensuring  
331 that most of the annotation occurs at the *predicate*  
332 level, rather than the *predicate-pair* or *sentence-*  
333 *pair* level, we are able to generate thousands of  
334 high quality entailment pairs from hundreds of an-  
335 notated predicates. This is in contrast with the con-  
336 struction processes of the Levy (Levy and Dagan,  
337 2016) and SherLliC (Schmitt and Schütze, 2019)  
338 datasets, which involved generating large numbers  
339 of candidate entailment pairs of varying quality,  
340 prior to manual annotation by crowd-source work-  
341 ers. Our method also avoids the issue of selection  
342 bias present in Zeichner et al. (2012) and SherLliC,  
343 that arises from using a similarity measure to auto-  
344 matically pre-select candidate entailments.

### 345 4.2 Antonym Pair Selection

346 We started by automatically collecting a list of  
347 477 lemmatised verb antonym pairs from Word-  
348 Net (Miller, 1995) and propose these as possible  
349 conflicting predicate pairs. Although WordNet’s  
350 antonym set is not large, the high quality of its  
351 annotations makes it a reliable starting point.

352 We excluded antonym pairs that express a type  
353 of temporal entailment (e.g. *fall asleep* and *wake*  
354 *up*), as these appear to express a more complicated  
355 relationship than simple antonymy. While these  
356 predicate pairs are antonymous when interpreted  
357 as simultaneous eventualities, they also entail each  
358 other at some temporal distance (e.g. you cannot  
359 *fall asleep* and *wake up* at the same time, but you  
360 need to *fall asleep* before you can *wake up*). If one  
361 of the two human annotators marked the antonym  
362 pair as having a possible temporal entailment be-  
363 tween the predicates, we removed it from the set.  
364 This step resulted in 283 remaining antonym pairs.

365 We also removed pairs that were highly spe-  
366 cific (e.g. *dehydrogenate-hydrogenate*) as these  
367 are likely to be infrequent in the general domain,  
368 pairs resulting from simple alternation of prepo-  
369 sitions or morphemes (*scale up-scale down*; *de-*  
370 *ceive-undecieve*), and duplicate pairs in the British

<sup>4</sup>Both annotators were authors of this paper

371 spelling.<sup>5</sup> We were left with 114 antonym pairs.

### 372 4.3 Entailment Cluster Construction

373 For each antonym pair, we identified possible para-  
374 phrases and third predicates that are entailed by  
375 both. We used the online Merriam-Webster The-  
376 saurus (Merriam-Webster, 2021), which includes  
377 both (near) synonyms and antonyms, and the Relat-  
378 edwords website (RelatedWords, 2021) – an online  
379 tool for finding related words beyond synonyms,  
380 which combines a number of NLP resources in-  
381 cluding word embedding spaces, ConceptNet and  
382 WordNet. This helped us find less typical para-  
383 phrases and often suggested entailed predicates.

384 For each antonym pair we created an *entailment*  
385 *cluster*  $\mathcal{C} = (A_1, A_2, E)$ , where  $A_1$  and  $A_2$  are the  
386 sets of predicates containing the *first* and *second*  
387 predicate in the seed antonym pair respectively,  
388 plus their paraphrases, and  $E$  is a set of predicates  
389 entailed by all the elements in  $\cup(A_1, A_2)$ .

390 Each cluster was then manually annotated with  
391 a set of argument type pairs (distinct from the  
392 FIGER types for Named Entities), which were  
393 later used for instantiating simple sentences. For  
394 example, the cluster for the antonym seed pair  
395 *refresh-tire* receives a set containing a single ar-  
396 gument type pair, *activity-generic\_person*. We  
397 allowed predicates with a specific word sense  
398 to be assigned a specific set of types. For ex-  
399 ample, for the *enjoy-suffer through* pair, the en-  
400 tailed predicate *see* is assigned the set containing  
401 just the type *generic\_person-entertainment\_watch*,  
402 to avoid it being paired with arguments from  
403 the *entertainment\_read* type. This also enabled  
404 us to specify argument order, allowing a pred-  
405 icate pair like *refresh(activity-generic\_person) -*  
406 *do(generic\_person-activity)*.

### 407 4.4 Entailment Pair Generation

408 The aim of the generation step is to automatically  
409 convert the entailment clusters into the dataset for-  
410 mat required for evaluation: premise, hypothesis,  
411 and a label denoting the type of entailment relation  
412 that holds between them.

413 To generate entailment pairs we take the cross  
414 product of different sets in the cluster. *Directional*  
415 *Entailments* are generated by  $\cup(A_1 \times E, A_2 \times E)$ ,  
416 *Antonyms* by  $\cup(A_1 \times A_2, A_2 \times A_1)$ , *Directional*

<sup>5</sup>We prefer American English spellings (e.g. *colonize*) over British English spellings (*colonise*) as the training corpus contains mostly American English news articles.

*Non-Entailments* by  $\cup(E \times A_1, E \times A_2)$  and *Para-*  
417 *phrases* by  $\cup(A_1 \times A_1, A_2 \times A_2)$ , excluding du-  
418 plicate predicates. We exclude an entailment pair  
419 if no intersection is found in the sets of its argu-  
420 ment types, or if it already occurs as part of another  
421 antonym pair’s cluster. 422

423 To generate a sentence for a predicate we need  
424 to populate its subject and object arguments. We  
425 therefore manually created argument strings for  
426 each argument type, ensuring they combine effec-  
427 tively with all predicates in the cluster. For exam-  
428 ple, the argument type *politician* maps to arguments  
429 like *Hillary Clinton*, used to instantiate sentences  
430 for predicates like *govern*. We used the Related-  
431 edwords website (RelatedWords, 2021) for inspira-  
432 tion. We then sampled an argument type pair from  
433 the intersection of those that apply for both pred-  
434 icates in the entailment pair. For each argument  
435 type we sampled non-identical argument strings.  
436 This produces an entailment example of the form  
437 (*arg1, predicate1, arg2. arg1, predicate2, arg2.*  
438 *label*). For example, (*The school, admitted, Jean.*  
439 *The school, evaluated, Jean. 1*) represents the direc-  
440 tional entailment *admit*  $\models$  *evaluate*. Finally, both  
441 annotators made a complete pass over the dataset,  
442 working separately to identify errors and correct  
443 the clusters accordingly. For example, they iden-  
444 tified unforeseen predicate-argument mismatches  
445 stemming from word sense ambiguity. Whilst this  
446 refinement method may be repeated indefinitely,  
447 we found that after a single pass the quality of the  
448 generated sentence pairs was very high.

449 The test subset<sup>6</sup> of ANT (based on 100 Word-  
450 Net antonym pairs) contains 6,300 entailment  
451 pairs: 1,800 Antonyms, 1,465 Directional Entail-  
452 ments, 1,465 Directional Non-Entailments, and  
453 1,570 Paraphrases. For the purpose of evalua-  
454 tion we used the following data subsets: 1) **Base:**  
455 *Antonyms* and *Directional Entailments*, and 2) **Di-**  
456 **rectional:** *Directional Entailments* and *Directional*  
457 *Non-Entailments*.

### 458 4.5 Error Analysis

459 To verify the dataset’s quality we conducted an er-  
460 ror analysis on 200 examples, with 50 examples  
461 per label sampled randomly from the test set. We  
462 found 82.5% (165/200) examples to be *correct*,  
463 confirming that the dataset is of high quality. Of  
464 the 35 *incorrect* examples we labelled five as a

<sup>6</sup>ANT also contains a small development set (based on 14 antonym pairs) for use with supervised learning techniques

465 syntactic error, 18 as a semantic error, and 12 as  
 466 unnatural/disfluent. The syntactic errors were at-  
 467 tributed to wrong verb tense or a missing auxiliary  
 468 verb in the predication. Sometimes semantic errors  
 469 resulted from the introduction of subtle meaning  
 470 change, such as for the directional non-entailment  
 471 “Morgan **changed** the server” - “Morgan upgraded  
 472 the server” (here **changed** might be interpreted as  
 473 **replaced**). They also arose due to predicate pairs  
 474 that were overlooked in cluster construction, e.g.  
 475 **look down on** is an antonym of **like** but not neces-  
 476 sarily a paraphrase of **dislike** - you can **dislike** (a  
 477 person) without **looking down on** (them). *Unnatu-*  
 478 *ral* sentences were often the result of odd argument-  
 479 predicate combinations, e.g. “Gale **expended** gas”.

## 5 Experimental Setup

481 We used the NewsSpike corpus of multi-source  
 482 news text (Zhang and Weld, 2013), comprising  
 483 approximately 0.5M articles collected over a pe-  
 484 riod of 6 weeks. Using MONTEE, we extracted  
 485 40,669,470 binary relation triples from NewsSpike.  
 486 Of these 8,107,944 (19.94%) binary relations are  
 487 extracted with a temporal interval resolved by SU-  
 488 Time (Chang and Manning, 2012) from a temporal  
 489 expression in the text. We use the SUTime tem-  
 490 poral interval if it is available and back off to the  
 491 document publication date if not.

492 We used the *entGraph*<sup>7</sup> framework with the ex-  
 493 tension of temporal filtering by Guillou et al. (2020)  
 494 to train the Entailment Graphs. See Appendix A for  
 495 hardware requirements and parameter settings. We  
 496 conducted two main experiments. In Section 6.1  
 497 we first show results comparing the uniform and  
 498 dynamic window strategies on the ANT dataset (us-  
 499 ing the default 5 day window suggested by Guillou  
 500 et al. (2020)). We then compare performance of  
 501 the temporal method on the general domain ANT  
 502 dataset compared to *Sports* (Section 6.2)<sup>8</sup>.

## 6 Results

### 6.1 Uniform vs. Dynamic Windowing

505 Results of the temporal scores with a uniform and  
 506 a dynamic window, alongside the atemporal scores,  
 507 are presented in Table 1. All methods here were  
 508 evaluated on the ANT dataset. We find that the dy-  
 509 namic windows consistently achieve higher Area  
 510 Under the precision-recall Curve (AUC) scores

<sup>7</sup><https://github.com/mjhosseini/entGraph>

<sup>8</sup>We also include AUC scores on the commonly used Levy/Holt dataset for the interested reader (Appendix B).

Window Method	ANT Base		ANT Dir.	
	Uni.	Dyn.	Uni.	Dyn.
<b>Similarity measures:</b>				
Weed’s Pr (Count)	<b>0.181</b>	<b>0.181</b>	<b>0.199</b>	<b>0.199</b>
T. Weed’s Pr (Count)	0.164	0.180	0.177	0.198
BInc (PMI)	0.161	0.161	0.178	0.178
T. BInc (Ratio PMI)	0.144	0.161	0.157	0.178
BInc (Count)	0.159	0.159	0.167	0.167
T. BInc (Count)	0.144	0.160	0.148	0.167

Table 1: Dynamic vs. Uniform window. AUC scores for the **Base** and **Directional** subsets of the *Sports* and ANT datasets.

Data subset	Sports		ANT	
	Base	Dir.	Base	Dir.
<b>Recall &lt; threshold</b>	0.75	0.75	0.3	0.3
<b>Similarity measure:</b>				
Weed’s Pr (Count)	0.440	0.460	<b>0.181</b>	<b>0.199</b>
T. Weed’s Pr (Count)	0.455	<b>0.472</b>	0.164	0.177
BInc (PMI)	0.471	0.432	0.161	0.178
T. BInc (Ratio PMI)	<b>0.495</b>	0.437	0.144	0.157
BInc (Count)	0.462	0.419	0.159	0.167
T. BInc (Count)	0.481	0.430	0.144	0.148

Table 2: Sports vs. General domain. AUC scores for the **Base** and **Directional** subsets of the ANT dataset, with a uniform window.

511 than uniform windows. This seems initially hopeful,  
 512 because it suggests that the duration informa-  
 513 tion from the TacoLM model is useful, and that a  
 514 dynamic window is more effective. However, we  
 515 note that all dynamic scores are very close to the  
 516 atemporal scores. These effects hold for both the  
 517 Base and Directional subsets of ANT.

518 We recognise this result may be caused by spuri-  
 519 ous overlaps — the larger windows result in fewer  
 520 occasions of temporal filtering, making the dy-  
 521 namic scores almost equal to the atemporal scores.  
 522 For now we will focus on uniform scores for the  
 523 remaining results and analysis, as it is likely more  
 524 representative of the temporal signal. Although  
 525 a window set dynamically per eventuality is theo-  
 526 retically appealing, we leave it to future research  
 527 further develop this idea.

### 6.2 Sports vs. General News Domain

529 Table 2 contains AUC scores for the Base and Di-  
 530 rectional subsets of the *Sports*<sup>9</sup> and ANT datasets.

<sup>9</sup>Note that the results are different to those reported by Guillou et al. (2020) due to updated relation extraction and a

		True	False	$\delta(T - F)$
Sports	% Scaled	31.5	35.8	-4.2
	% Overlap	72.8	65.8	7.1
ANT	% Scaled	53.0	51.8	1.2
	% Overlap	50.4	50.4	0.0

Table 3: Analysing the difference in effect of temporal filtering between the Sports and ANT base datasets.

Performance of the temporal measures is consistently higher than their atemporal counterparts for the *Sports* Base and Directional subsets. We tested statistical significance of the difference on the *Sports* dataset using bootstrap resampling (10K samples) (Efron and Tibshirani, 1985; Koehn, 2004). For **Base**, p-values are .074, .065 and .091 for Weed’s, BInc (PMI) and BInc (Count), respectively, while for **Dir** they are .14, .28 and .11, so that on **Base** the difference is significant for all scores at the  $< 0.1$  level. For the subsets of ANT, however, the performance of temporal measures is consistently lower than that of the atemporal counterparts. This difference suggests that the atemporal formulation of the DIH by Dagan et al. (1999) and Geffet and Dagan (2005) is appropriate for the general domain, while the temporal formulation is more applicable in the sports domain.

Figure 1 contains the precision-recall curves for the *Sports* and ANT datasets. To provide a fair comparison between scores that have different recall ranges, we compute AUC under a recall threshold, chosen separately for each dataset (see values in Table 2). For *Sports* we observe higher precision for the temporal measures compared with their atemporal counterparts at lower recall ranges. For the ANT dataset, recall is very low. This is due to the absence of many of the entailment pairs in the Entailment Graphs. Furthermore, in contrast to *Sports*, the temporal curves are consistently below the atemporal curves, even at the low recall level. It seems that the temporal distributions for eventualities in the more general domain are such that temporal filtering has a counterproductive effect. We analyse this in Section 7.

## 7 Analysis and Discussion

Table 3 contains statistics of temporal separation for the Base subset of the *Sports* and ANT datasets.

correction to the method for applying temporal windows. We consider the results in this paper to be the definitive results on the sports domain.

*% Scaled\_down* is the percentage of PMI scores (for each co-occurrence of triples) that are scaled down by the temporal method. *% Overlap* is the percentage of eventuality comparisons ( $e_p, e_q$ ) that result in a temporal overlap. When the method is effective, we expect *% Scaled* to be higher for false predicate pairs than true predicate pairs (as scores of antonymous predicate pairs should be scaled down). Scaling should be inversely related to the average *Overlap*, which should be higher for true predicate pairs than false ones.

We indeed find that *% Scaled* is higher for false predicates pairs in the *Sports* dataset, whereas there is a small difference in the wrong direction for the ANT dataset. This helps explain the differences observed in the Base dataset precision-recall graphs A (*Sports* dataset) and C (ANT dataset) in Figure 1. Furthermore, *% Overlap* has the expected correlation, showing that our method works for the temporal distribution of the sports domain data, but not the general-domain data. That is, it can be applied successfully when argument co-occurrences of the antonym pairs are temporally disjoint more often than the argument co-occurrences of entailing predicate pairs. In our training corpus, this distribution holds for sports predicate pairs but not for general domain predicate pairs. We expect that the performance decrease for the general domain is due to the scaling and overlap being (fairly) uniformly distributed over True and False predicate pairs — in that case the method simply reduces the amount of data available, while providing no added accuracy.

Breaking down the *% Scaled* statistic per predicate pair in the ANT dataset, we do find antonyms for which many scores are scaled down, indicating that there may still be specific predicates in the general domain where temporality is a useful signal. For example, the antonymous predicate pairs that are scaled most include *violate-respect*, *convict-acquit*, *allow-prohibit* and *(thing) kills (person)-(person) survives (thing)*, suggesting that predicates in legal news are worth exploring. Examples found in the corpus also support this idea for other predicate pairs. We find “Cameron, who ... , leaves London today ... .” and “Cameron will instead stay in London ... .”, referring to dates a month apart. The atemporal baseline models use this data to erroneously support that *leave*  $\models$  *stay in*, whereas our method successfully disentangles the evidence. Future research could further investigate which other subdomains or predicates stand to benefit from tem-

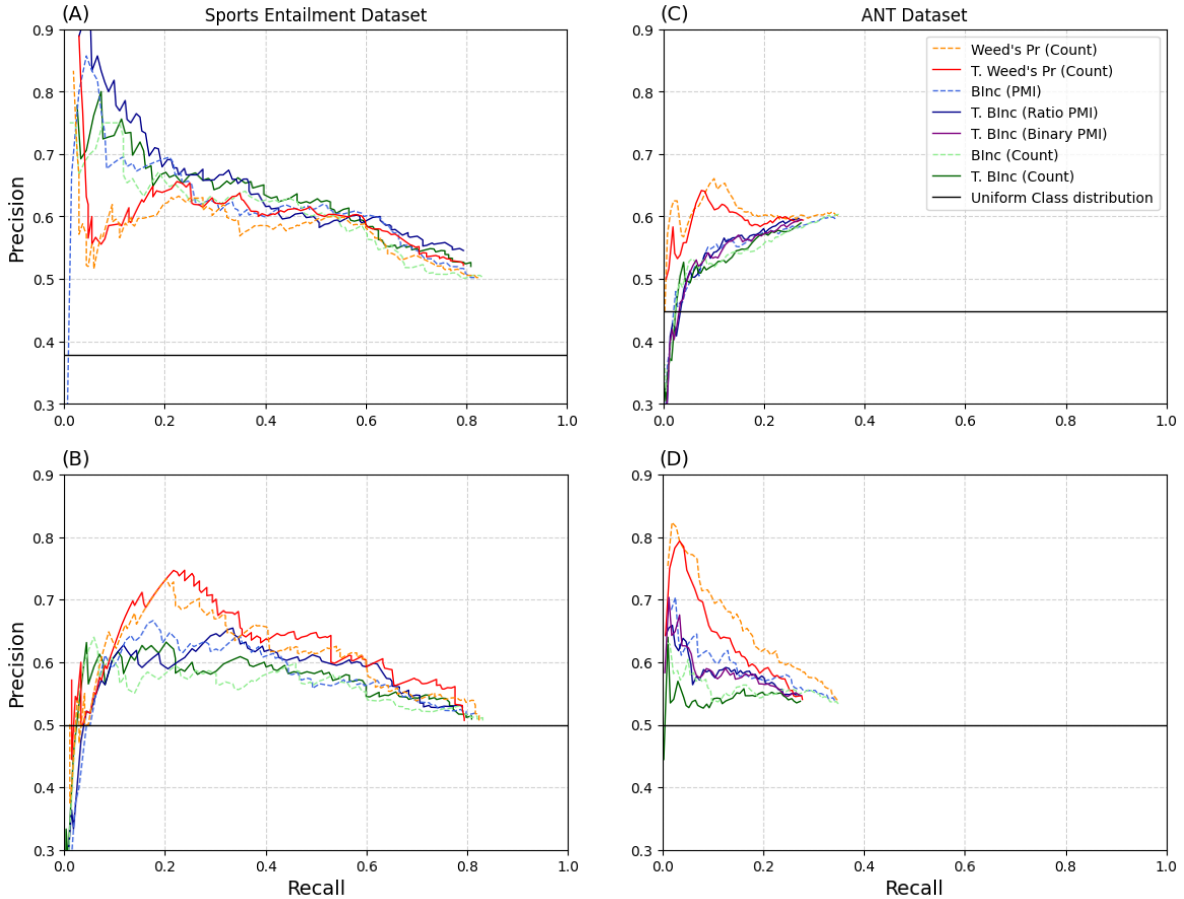


Figure 1: Precision-recall plots for the *Sports* dataset (A) base and (B) directional subsets, and the ANT dataset (C) base and (D) directional subsets

poral information. This could inform models that are able to decide whether to apply temporal filtering for particular predicate pairs.

Although dynamic windowing was unsuccessful in this iteration, further improvements to the algorithm in terms of robustness to noise, and an evaluation that includes predicates with a greater variety of durations, may yet prove the usefulness of the idea. Another direction is to explore how Entailment Graphs can be used to learn temporal entailments (such as *wake up - go to sleep*), which were excluded from the ANT dataset. Combining this with recent work on Multivalent Entailment Graphs will be essential, as many of the entailment edges may be multivalent (e.g. “A kills B”  $\models$  “B is dead”, see also McKenna et al. (2021)). We might also consider the interaction of temporality and modality, since the temporal signal should be more able to separate antonymous data when it does not include binary relations that are stated as occurring with some degree of uncertainty (see also Guillou et al. (2021)).

## 8 Conclusion

We applied a temporal Entailment Graph induction method (Guillou et al., 2020) to the general news domain. We evaluated performance using the *Sports* dataset and the ANT dataset. ANT is designed with underlying categories of (non-) entailment in mind, and uses a semi-automatic construction method that still results in high-quality annotation. Results show that a reformulation of the Distributional Inclusion Hypothesis to incorporate time is beneficial for the sports domain, while the atemporal formulation of the DIH is appropriate for the general domain (although there may be specific predicates for which the temporal formulation is effective). Our analysis suggests this is due to (temporal) distributional properties of the predicates. The method works when there are false predicate pairs for which the scores are scaled down (due to co-occurrences with the same argument pairs at different times) and true predicate pairs that maintain a high degree of temporal overlap.



## References

- Jonathan Berant, Noga Alon, Ido Dagan, and Jacob Goldberger. 2015. [Efficient global learning of entailment graphs](#). *Computational Linguistics*, 41(2):221–263.
- Jonathan Berant, Ido Dagan, and Jacob Goldberger. 2011. [Global learning of typed entailment rules](#). In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 610–619, Portland, Oregon, USA. Association for Computational Linguistics.
- Sander Bijl de Vroe, Liane Guillou, Miloš Stanojević, Nick McKenna, and Mark Steedman. 2021. [Modality and negation in event extraction](#). In *Proceedings of the 4th Workshop on Challenges and Applications of Automated Extraction of Socio-political Events from Text (CASE 2021)*, pages 31–42, Online. Association for Computational Linguistics.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. [A large annotated corpus for learning natural language inference](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.
- Angel X. Chang and Christopher Manning. 2012. [SU-Time: A library for recognizing and normalizing time expressions](#). In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC’12)*, pages 3735–3740, Istanbul, Turkey. European Language Resources Association (ELRA).
- Robin Cooper, Dick Crouch, Jan Van Eijck, Chris Fox, Josef Van Genabith, Jan Jaspars, Hans Kamp, David Milward, Manfred Pinkal, Massimo Poesio, Steve Pulman, Ted Briscoe, Holger Maier, and Karsten Konrad. 1996. Using the framework. FraCaS: A framework for computational semantics.
- Ido Dagan, Lillian Lee, and Fernando C. N. Pereira. 1999. Similarity-based models of word cooccurrence probabilities. *Machine Learning*, 34:43–69.
- Bradley Efron and Robert Tibshirani. 1985. The bootstrap method for assessing statistical accuracy. *Behaviormetrika*, 12(17):1–35.
- Maayan Geffet and Ido Dagan. 2005. [The distributional inclusion hypotheses and lexical entailment](#). In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL’05)*, pages 107–114, Ann Arbor, Michigan. Association for Computational Linguistics.
- Liane Guillou, Sander Bijl de Vroe, Mohammad Javad Hosseini, Mark Johnson, and Mark Steedman. 2020. [Incorporating temporal information in entailment graph mining](#). In *Proceedings of the Graph-based Methods for Natural Language Processing (TextGraphs)*, pages 60–71, Barcelona, Spain (Online). Association for Computational Linguistics.
- Liane Guillou, Sander Bijl de Vroe, Mark Johnson, and Mark Steedman. 2021. [Blindness to modality helps entailment graph mining](#). In *Proceedings of the Second Workshop on Insights from Negative Results in NLP*, pages 110–116, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Xavier Holt. 2018. Probabilistic models of relational implication. Master’s thesis, Macquarie University.
- Mohammad Javad Hosseini, Nathanael Chambers, Siva Reddy, Xavier R. Holt, Shay B. Cohen, Mark Johnson, and Mark Steedman. 2018. [Learning typed entailment graphs with global soft constraints](#). *Transactions of the Association for Computational Linguistics*, 6:703–717.
- Mohammad Javad Hosseini, Shay B. Cohen, Mark Johnson, and Mark Steedman. 2019. [Duality of link prediction and entailment graph induction](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4736–4746, Florence, Italy. Association for Computational Linguistics.
- Mohammad Javad Hosseini, Shay B. Cohen, Mark Johnson, and Mark Steedman. 2021. [Open-domain contextual link prediction and its complementarity with entailment graphs](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2790–2802, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Thomas Kober, Sander Bijl de Vroe, and Mark Steedman. 2019. [Temporal and aspectual entailment](#). In *Proceedings of the 13th International Conference on Computational Semantics - Long Papers*, pages 103–119, Gothenburg, Sweden. Association for Computational Linguistics.
- Philipp Koehn. 2004. [Statistical significance tests for machine translation evaluation](#). In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 388–395, Barcelona, Spain. Association for Computational Linguistics.
- Omer Levy and Ido Dagan. 2016. [Annotating relation inference in context via question answering](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 249–255, Berlin, Germany. Association for Computational Linguistics.
- Omer Levy, Ido Dagan, and Jacob Goldberger. 2014. [Focused entailment graphs for open IE propositions](#). In *Proceedings of the Eighteenth Conference on Computational Natural Language Learning*, pages 87–97, Ann Arbor, Michigan. Association for Computational Linguistics.

773	Xiao Ling and Daniel S. Weld. 2012. Fine-grained entity recognition. In <i>Proceedings of the Twenty-Sixth AAAI Conference on Artificial Intelligence, AAAI'12</i> , page 94–100. AAAI Press.	
774		
775		
776		
777	Nick McKenna, Liane Guillou, Mohammad Javad Hosseini, Sander Bijl de Vroe, Mark Johnson, and Mark Steedman. 2021. Multivalent entailment graphs for question answering. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 10758–10768, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	
778		
779		
780		
781		
782		
783		
784		
785	Merriam-Webster. 2021. Merriam-webster online thesaurus. <a href="https://www.merriam-webster.com/thesaurus">https://www.merriam-webster.com/thesaurus</a> . Accessed: 2021-12-16.	
786		
787		
788	George A. Miller. 1995. Wordnet: A lexical database for english. <i>Commun. ACM</i> , 38(11):39–41.	
789		
790	Saif M. Mohammad, Bonnie J. Dorr, Graeme Hirst, and Peter D. Turney. 2013. Computing lexical contrast. <i>Computational Linguistics</i> , 39(3):555–590.	
791		
792		
793	Sneha Rajana, Chris Callison-Burch, Marianna Apidianaki, and Vered Shwartz. 2017. Learning antonyms with paraphrases and a morphology-aware neural network. In <i>Proceedings of the 6th Joint Conference on Lexical and Computational Semantics (*SEM 2017)</i> , pages 12–21, Vancouver, Canada. Association for Computational Linguistics.	
794		
795		
796		
797		
798		
799		
800	RelatedWords. 2021. Relatedwords.org website. <a href="https://www.relatedwords.org/">https://www.relatedwords.org/</a> . Accessed: 2021-12-16.	
801		
802		
803	Enrico Santus, Frances Yung, Alessandro Lenci, and Chu-Ren Huang. 2015. EVALution 1.0: an evolving semantic dataset for training and evaluation of distributional semantic models. In <i>Proceedings of the 4th Workshop on Linked Data in Linguistics: Resources and Applications</i> , pages 64–69, Beijing, China. Association for Computational Linguistics.	
804		
805		
806		
807		
808		
809		
810	Martin Schmitt and Hinrich Schütze. 2019. SherLiC: A typed event-focused lexical inference benchmark for evaluating natural language inference. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 902–914, Florence, Italy. Association for Computational Linguistics.	
811		
812		
813		
814		
815		
816	Miloš Stanojević and Mark Steedman. 2019. CCG parsing algorithm with incremental tree rotation. 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 228–239, Minneapolis, Minnesota. Association for Computational Linguistics.	
817		
818		
819		
820		
821		
822		
823		
824	Idan Szepktor and Ido Dagan. 2008. Learning entailment rules for unary templates. In <i>Proceedings of the 22nd International Conference on Computational Linguistics (Coling 2008)</i> , pages 849–856, Manchester, UK. Coling 2008 Organizing Committee.	
825		
826		
827		
828		
	Julie Weeds and David Weir. 2003. A general framework for distributional similarity. In <i>Proceedings of the 2003 Conference on Empirical Methods in Natural Language Processing</i> , pages 81–88.	829 830 831 832
	Changlong Yu, Hongming Zhang, Yangqiu Song, Wilfred Ng, and Lifeng Shang. 2020. Enriching large-scale eventuality knowledge graph with entailment relations. In <i>Conference on Automated Knowledge Base Construction, AKBC 2020, Virtual, June 22-24, 2020</i> .	833 834 835 836 837 838
	Naomi Zeichner, Jonathan Berant, and Ido Dagan. 2012. Crowdsourcing inference-rule evaluation. In <i>Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 156–160, Jeju Island, Korea. Association for Computational Linguistics.	839 840 841 842 843 844
	Congle Zhang and Daniel S. Weld. 2013. Harvesting parallel news streams to generate paraphrases of event relations. In <i>Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing</i> , pages 1776–1786, Seattle, Washington, USA. Association for Computational Linguistics.	845 846 847 848 849 850
	Ben Zhou, Qiang Ning, Daniel Khashabi, and Dan Roth. 2020. Temporal common sense acquisition with minimal supervision. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7579–7589, Online. Association for Computational Linguistics.	851 852 853 854 855 856

## 857 A Experimental Settings / Requirements

858 With the following exceptions we used MONTEE’s  
 859 default settings to extract binary relations. We  
 860 enabled the SUTime component to ensure that  
 861 each binary relation with a predicate that could  
 862 be linked to a time expression was assigned a  
 863 time interval derived from SUTime [includeTempo-  
 864 ral=True]. These time intervals were used when  
 865 computing the temporal similarity measures but  
 866 ignored during the computation of the atemporal  
 867 measures. We disabled unary relation extraction  
 868 [writeUnaryRels=False], and restricted binary rela-  
 869 tions to only those that include at least one named  
 870 entity [acceptGGBinary=False].

871 We used the *entGraph* framework of Hosseini  
 872 et al. (2018) to construct Entailment Graphs. We  
 873 raised the threshold values for infrequent predicates  
 874 [minPredForArgPair=4] and argument pairs [mi-  
 875 nArgPairForPred=4] for all type-pair graphs (with  
 876 the exception of the very large *thing-thing* graph  
 877 for which we used settings of 6 and 6 respectively),  
 878 and used the default values for all other parameters.

879 All of the experiments were conducted on a sin-  
 880 gle server which has two Intel Xeon E5-2697 v4  
 881 2.3GHz CPUs (each with 18 cores) and 330GB  
 882 RAM. The computational cost of training a single  
 883 Entailment Graph is approximately one day and  
 884 160GB RAM. Evaluation of both the Levy/Holt  
 885 and ANT datasets using the *entGraph* evaluation  
 886 scripts takes approximately 6 hours per graph.

## 887 B Results on the Levy/Holt Dataset

888 Previous work on Entailment Graphs has reported  
 889 performance on the general-domain Levy/Holt  
 890 (Levy and Dagan, 2016; Holt, 2018) dataset of  
 891 18,407 entailment pairs (Hosseini et al., 2018, 2019,  
 892 2021; McKenna et al., 2021; Guillou et al., 2021).  
 893 Although not designed for evaluating performance  
 894 on the task of temporally separating eventualities,  
 895 we also include results on the Levy/Holt dataset for  
 896 the interested reader. We use the same dev/test split  
 897 proposed by (Hosseini et al., 2018): 5,486 pairs for  
 898 dev and 12,921 pairs for test.

899 AUC scores for the uniform widow experiment  
 900 (with a 5-day window) are provided in Table 4. We  
 901 observe a similar pattern for the Levy/Holt dataset  
 902 as we did for ANT – performance of the temporal  
 903 measures is lower than their atemporal counterparts.  
 904 This further supports the claim in Section 6 that  
 905 the temporal distributions for the eventualities in  
 906 the general domain are not well suited to temporal

Data subset	Dev		Test	
	All	Dir.	All	Dir.
Recall < threshold	0.45	0.5	0.45	0.5
<b>Similarity measure:</b>				
Weed’s Pr (Count)	0.215	<b>0.217</b>	0.207	<b>0.220</b>
T. Weed’s Pr (Count)	0.215	0.183	0.193	0.194
BInc (PMI)	0.221	0.203	<b>0.212</b>	0.203
T. BInc (Ratio PMI)	<b>0.224</b>	0.164	0.196	0.176
BInc (Count)	0.217	0.208	0.205	0.201
T. BInc (Count)	0.218	0.173	0.191	0.174

Table 4: AUC scores for the Levy/Holt datasets: **All** examples and **Directional** only examples for the dev and test sets. Uniform window.

907 filtering compared to sports domain predicates. In  
 908 the case of Levy/Holt, the result may also be due to  
 909 the lack of antonym pairs, which may benefit more  
 910 from the temporal signal than other non-entailment  
 911 categories.