

Learning from Missing Relations: Contrastive Learning with Commonsense Knowledge Graphs for Commonsense Inference

Anonymous ACL submission

Abstract

Commonsense inference poses a unique challenge to reason and generate the physical, social, and causal conditions of a given event. Existing approaches to commonsense inference utilize commonsense transformers, which are large-scale language models that learn commonsense knowledge graphs. However, they suffer from a lack of coverage and expressive diversity of the graphs, resulting in a degradation of the representation quality. In this paper, we focus on addressing missing relations in commonsense knowledge graphs, and propose a novel contrastive learning framework called SOLAR¹. Our framework contrasts sets of semantically similar and dissimilar events, learning richer inferential knowledge compared to existing approaches. Empirical results demonstrate the efficacy of SOLAR in commonsense inference of diverse commonsense knowledge graphs. Specifically, SOLAR outperforms the state-of-the-art commonsense transformer on commonsense inference with ConceptNet by 1.84% on average among 8 automatic evaluation metrics. In-depth analysis of SOLAR sheds light on the effects of the missing relations utilized in learning commonsense knowledge graphs.

1 Introduction

Commonsense inference, reasoning of unobserved conditions from an observed event, is an important but challenging task in natural language processing (NLP) (Rashkin et al., 2018; Bosselet et al., 2019; Yuan et al., 2020; Hwang et al., 2021). This is easy for humans, but still out of the reach of current artificial intelligence systems. Commonsense inference aims to generate textual descriptions of the inference results, which is more in line with the

¹Code available at https://anonymous.4open.science/r/solar-commonsense_inference-37E7

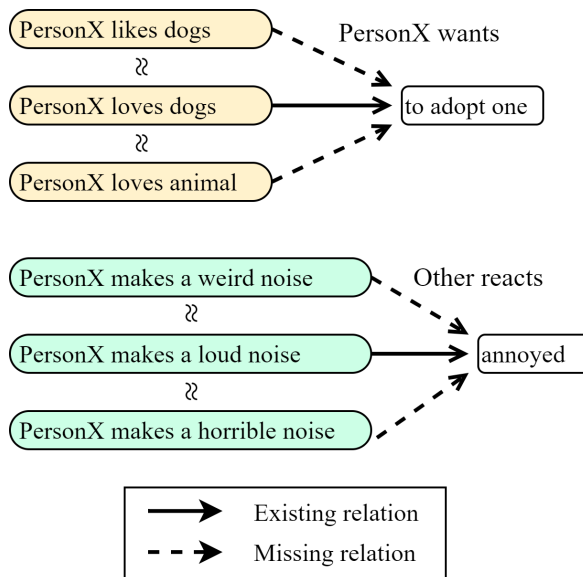


Figure 1: Illustration of missing relations of semantically similar events in commonsense KGs.

process of humans reasoning based on their knowledge. For a given event “X walks into a hospital”, the causal conditions (e.g., what to do before and after the event), physical conditions (e.g., capability and location of entities), and social conditions (the intention and reaction of X) of the event are to be inferred.

Recent studies on commonsense inference have adopted commonsense transformers (Bosselet et al., 2019), which are large-scale language models trained on commonsense knowledge graphs (KGs) like ATOMIC (Sap et al., 2019) and ConceptNet (Speer et al., 2017). Such models are grounded on the hypothesis that language models can memorize facts in their parameters during training (Petroni et al., 2019; Roberts et al., 2020). Despite these efforts, commonsense transformer models still suffer from two main obstacles inherent in commonsense KGs: (1) *lack of coverage* and (2) *expressive diversity* of the graphs. First, commonsense KGs lack the coverage required to be applicable for diverse

059	situations in the real world (Li et al., 2016; Saito et al., 2018). In ATOMIC, even with the possibility	110
060	of far more commonsense properties being relevant, any single event has only 2.2 commonsense	111
061	properties directly related on average. Second, with the non-canonical and free-form text representation	112
062	for the nodes in commonsense KGs, semantically identical or similar expressions of events are represented	113
063	as distinct nodes (Malaviya et al., 2020). For example, “PersonX is fond of dogs” and “PersonX	114
064	likes dogs” are semantically identical, but represented as distinct nodes. The expressive diversity	115
065	makes commonsense KGs substantially sparser than conventional KGs. Owing to the lack of coverage	116
066	and expressive diversity, significant amount of relations are missing in commonsense KGs.	117
067		118
068		119
069		120
070		121
071		122
072		123
073		124
074		125
075		126
076		127
077		128
078		129
079		130
080		131
081		132
082		133
083		134
084		135
085		136
086		137
087		138
088		139
089		140
090		141
091		142
092		143
093		144
094		145
095		146
096		147
097		148
098		149
099		150
100		151
101		152
102		153
103		154
104		155
105		156
106		
107		
108		
109		

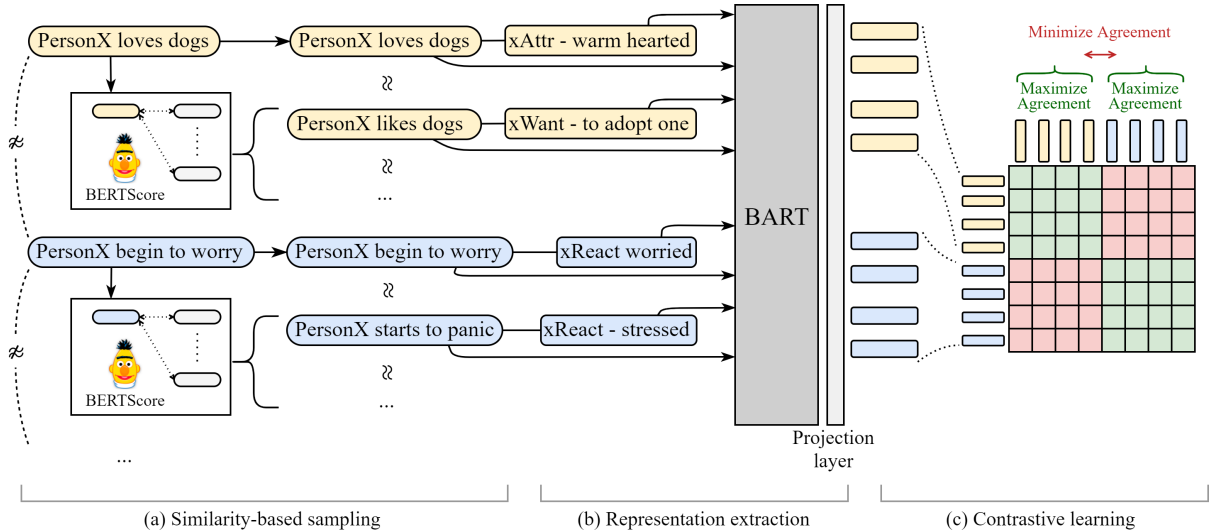


Figure 2: Illustration of contrastive learning of commonsense tuples. (a) Based on adversarially sampled root subjects, semantically similar subjects are sampled. (b) Subjects and relation-object pairs connected to them are projected to separate hidden representations through a generative language model and a projection layer. (c) Hidden representations obtained from the same root subject are considered as positive pairs, and those obtained from other root subjects are considered as negative pairs for contrastive learning.

KGs must be designed to contain knowledge that is not already expressible by language models. Gabriel et al. (2021) focused on discourse-level commonsense inference, and Yuan et al. (2020) proposed a language model architecture for logically consistent commonsense reasoning. Previous studies have proposed training language models on existing tuples in commonsense KGs for commonsense inference. In our work, we focus on addressing the missing relations of commonsense KGs for better commonsense inference.

2.2 Contrastive Learning

Contrastive learning has shown promising performances in computer vision (Henaff, 2020; He et al., 2020). SimCLR (Chen et al., 2020b) introduced a simple but powerful contrastive learning approach and showed a competitive performance with supervised learning approaches. Contrastive learning is also widely used in natural language processing, where a model obtains unsupervised representations by learning to predict positive or negative pairs. Mikolov et al. (2013) proposed an efficient method for learning word representations by classifying whether given words appear in the same context or not. Furthermore, contrastive learning has been adopted to improve the representations of pre-trained language models. Reimers and Gurevych (2019); Zhang et al. (2020b); Yan et al. (2021) introduced contrastive learning frameworks for enhanc-

ing the sentence representations. Lee et al. (2020) proposed a contrastive learning method to mitigate the exposure bias problem. Inspired by these studies, we propose a novel contrastive learning framework for commonsense representation learning with commonsense KGs. With our proposed framework, the model learns inferential knowledge from both existing and missing relations.

3 Methodology

In this section, we describe the model architecture and training procedure of the proposed framework.

3.1 Notation

We define $G = (V, E)$ as the commonsense knowledge graph that consists of a set of nodes V and a set of edges E . Following the notation from COMET (Bosselut et al., 2019), we denote each knowledge tuple from the knowledge graph as $\{s, r, o\}$, where s is the phrase subject, r is the relation, and o is the phrase object of the tuple. Here, s and o are natural language sequences, and r is a single special token (e.g., $\langle x\text{Intent} \rangle$). Note that $s, o \in V$ and $\{s, r, o\} \in E$. We define S as the set of all existing subjects from the knowledge graph, and it follows that $S \subset V$. Finally, we denote the generative language model to be trained as $f(\cdot)$ and a projection layer at the top of the model as $g(\cdot)$.

Algorithm 1 Set Construction Algorithm.

Input: root subjects S_{root} , number of root subjects N , edges E , set size $2m$, threshold δ , BERTScore function $b(\cdot, \cdot)$, base model $f(\cdot)$, projection layer $g(\cdot)$

for $s_i \in S_{root}$ **do**
 Initialize G_i as \emptyset
 for $j \in \{1, \dots, m\}$ **do**
 if $j = 1$ **then**
 $s_j^i \leftarrow s_i$
 else
 repeat ▷ Sample similar subject
 $s_j^i \leftarrow \text{sample}(S)$
 until $b(f(s_j^i), f(s_i)) > \delta$
 end if
 get tuple $\{s_j^i, r_j^i, o_j^i\} \in E$ containing s_j^i
 $z_{2j-1}^i \leftarrow g(f(s_j^i))$
 $z_{2j}^i \leftarrow g(f(r_j^i \oplus o_j^i))$
 $G_i \leftarrow G_i \cup \{z_{2j-1}^i, z_{2j}^i\}$
 end for
end for
return G_1, G_2, \dots, G_N

3.2 Commonsense Representation Learning

To improve commonsense representations of the language model prior to learning commonsense inference, we first proceed with commonsense representation learning through contrastive learning of commonsense tuples and commonsense reconstruction.

Contrastive learning of commonsense tuples.

Inspired by our key observation that semantically identical or similar events can have same relations, we propose a novel commonsense representation learning method based on contrastive learning.

The overall procedure of the proposed method is depicted in Figure 2. First, we obtain a set of N root subjects $S_{root} = \{s_1, s_2, \dots, s_N\}$ through adversarial sampling on S . The adversarial sampling procedure is designed such that pairwise semantic similarity of subjects in S_{root} lies between minimum similarity α and maximum similarity β . Here, we use BERTScore (Zhang et al., 2020a) between phrase subjects as the semantic similarity metric.

We then obtain positive and negative pairs by constructing N sets G_1, G_2, \dots, G_N containing hidden representations, where each G_i corresponds to a root subject $s_i \in S_{root}$. For an arbitrary element $s_i \in S_{root}$, we first sample m tuples

$\{s_j, r_j, o_j\}$ ($j = 1, 2, \dots, m$) from E that contain subjects s_j semantically similar to s_i . Each s_j and $r_j \oplus o_j$ is projected to hidden representations $z_{2j-1}^i = g(f(s_j))$ and $z_{2j}^i = g(f(r_j \oplus o_j))$, and added to G_i . Here, \oplus denotes concatenation of two sequences. Repeating for m times, the constructed set G_i contains $2m$ hidden representations derived from subjects that are semantically similar to the root subject s_i , and the relation-object pairs connected to them. Algorithm 1 summarizes the construction procedure.

We consider samples from the same set as positive pairs, and those from different sets are negative pairs in contrastive learning. We use NT-Logistic (the normalized temperature-scaled logistic) objective function (Chen et al., 2020b) as our training objective to maximize the agreement between positive pairs while minimizing the agreement between negative pairs. The formal expression of our objective function is given by the following equations:

$$l_i^{pos} = -\frac{\sum_{p,q=1}^{2m} \log \sigma(z_p^{i,T} z_q^i / \tau)}{2m}, \quad (1)$$

$$l_i^{neg} = -\frac{\sum_{i < j \leq N} \sum_{p,q=1}^{2m} \log \sigma(-z_p^{i,T} z_q^j / \tau)}{m(N-1)}, \quad (2)$$

$$L_{cont} = \frac{1}{N} \sum_{i=1}^N (l_i^{pos} + l_i^{neg}), \quad (3)$$

where l_i^{pos} is the loss function over positive pairs in set G_i , and l_i^{neg} is the loss function over negative pairs among set G_i and the other sets. In addition, τ denotes the temperature parameter for temperature scaling. The model is trained to minimize the final objective L_{cont} , which is the mean of l_i^{pos} and l_i^{neg} for all $i = 1, 2, \dots, N$.

Commonsense reconstruction. To further improve the representation of a single tuple, we propose a commonsense reconstruction task inspired by Lewis et al. (2020), in which the model learns to reconstruct noisy tuples into their original form. More specifically, we noise a commonsense tuple $\{s, r, o\}$ by randomly choosing one of the three elements, masking the span of the chosen element, and shuffling the order of the tuple. The model is trained to reconstruct the original tuple from the noisy tuple. This task complements the contrastive learning method by training the model to better understand the commonsense tuple itself. The objective of the commonsense reconstruction task is

285 to minimize L_{recon} computed by cross-entropy be- 329
286 tween the decoder output and the original tuple. 330

287 The model learns commonsense representations 331
288 through multitask learning on the two aforemen- 332
289 tioned tasks simultaneously. Therefore, the final 333
290 objective function of our framework is to minimize 334
291 the combined loss: 335

$$292 \quad L_{rep} = \omega L_{cont} + (1 - \omega)L_{recon}. \quad (4) \quad 336$$

293 3.3 Fine-tuning on Commonsense KGs 337

294 After learning commonsense representations, we 340
295 remove the projection head and fine-tune the model 341
296 with commonsense KGs to learn commonsense 342
297 inference. The model learns to generate a phrase 343
298 object o given a concatenation of phrase subject s 344
299 and relation r . The objective function of the task is 345
300 as follows: 346

$$301 \quad L_{infer} = - \sum_{i=0}^{|E|} \log P_{\theta}(o_i | s_i, r_i) \quad (5) \quad 347$$

302 3.4 Language Model Architecture 348

303 While SOLAR is agnostic to its generative lan- 349
304 guage model architecture, for our experiments, 350
305 we use BART (Lewis et al., 2020) with its pre- 351
306 trained parameters as our base generative language 352
307 model. BART is a transformer-based sequence- 353
308 to-sequence language model with a bidirectional 354
309 encoder and a left-to-right autoregressive decoder. 355
310 For commonsense representation learning (Section 356
311 3.2), we add a projection layer that maps the BART 357
312 decoder output representations to a space where 358
313 contrastive loss is applied. The projection head 359
314 is then removed for fine-tuning on commonsense 360
315 KGs (Section 3.3). 361

316 4 Experiments 362

317 In this section, we demonstrate the efficacy of our 363
318 framework by comparing the commonsense infer- 364
319 ence performances of SOLAR with those of the 365
320 state-of-the-art commonsense transformers. 366

321 4.1 Dataset 367

322 Commonsense KGs are widely used for evaluat- 370
323 ing the commonsense inference capability by mea- 371
324 suring the plausibility of the generated inferences 372
325 given unobserved events or entities. Hwang et al. 373
326 (2021) developed an adversarial splitting method 374
327 for dividing training, validation, and test sets that 375
328 prevent overlapping subjects of knowledge tuples 376

329 between the sets. We utilize the splitting method 330
331 to evaluate the inference capability of the model 332
333 for unseen events or entities. We use three com- 334
335 monsense KGs in our experiments: ConceptNet 336
337 (Speer et al., 2017), ATOMIC (Sap et al., 2019), 338
339 and ATOMIC₂₀²⁰ (Hwang et al., 2021). 340

ConceptNet is a general commonsense knowledge 341
342 graph. We use a subset of the graph provided by 343
344 Li et al. (2016), which involves 36 relations and 345
346 300K tuples. The subset is divided into 265K, 5K, 347
348 and 30K tuples for training, validation, and testing 349
350 respectively. 351

ATOMIC is a social commonsense knowledge 352
353 graph that involves 9 relations with 877K tuples. 354
355 The split of ATOMIC includes 710K, 80K, and 356
357 87K tuples for training, validation, and testing, re- 358
359 spectively. 360

ATOMIC₂₀²⁰ is a recently proposed large-scale com- 361
362 monsense knowledge graph, which involves 23 363
364 commonsense dimensions and contains 1.33M tu- 364
365 ples. It includes physical-entity, social-interaction, 365
366 and event-centered commonsense. ATOMIC₂₀²⁰ is 366
367 split into 1.08M, 10K, and 15K tuples for training, 367
368 validation, and testing, respectively. 368
369

369 4.2 Experimental Settings 370

Baseline We use COMET (Bosselut et al., 2019), 371
372 the state-of-the-art commonsense transformers in 372
373 commonsense inference, as the baseline. We use 373
374 the public HuggingFace (Wolf et al., 2019) imple- 374
375 mentation of pre-trained BART (Lewis et al., 2020) 375
376 as a language model and train it using SOLAR and 376
377 COMET for comparison. BART-base has 6 trans- 377
378 former layers for encoder and decoder each with 378
379 a hidden size of 768, whereas BART-large has 12 379
380 transformer layers for encoder and decoder each 380
381 with a hidden size of 1024. For fine-tuning, we em- 381
382 pirically choose the best number of epochs, learn- 382
383 ing rate, and batch size among $\{1, 3, 5, 7, 11, 13\}$, 383
384 $\{1e-4, 1e-5, 1e-6\}$, and $\{16, 32, 64, 128\}$, respec- 384
385 tively, and use the Adam optimizer with $\beta_1 = 0.9$, 385
386 $\beta_2 = 0.999$. 386

Training details of SOLAR. In contrastive learn- 370
371 ing of commonsense tuples, we extract $n \in$ 371
372 $\{4, 8, 16, 32\}$ root subjects while maintaining the 372
373 similarity (%) between subjects with a minimum of 373
374 $\alpha \in \{40, 50, 60\}$ and a maximum of $\beta \in \{70, 80\}$. 374
375 We then sample $m \in \{8, 16, 32\}$ semantically sim- 375
376 ilar subjects based on previously extracted subjects. 376
377 We set the temperature parameter τ to 0.1. 377

In reconstructive learning tasks, we corrupt tu- 378

		BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE-L	CIDEr	BERTScore
ConceptNet	COMET-base	15.60	10.26	6.88	4.84	11.79	16.61	33.41	53.18
	SOLAR-base	17.12	11.55	8.10	5.79	12.90	18.25	38.91	53.86
ATOMIC	COMET-base	53.03	33.97	23.13	16.90	34.05	56.07	74.63	64.57
	SOLAR-base	53.59	34.51	23.89	17.82	34.42	56.60	75.24	64.78
ATOMIC ₂₀	COMET-base	44.99	26.95	17.44	11.77	31.20	48.33	59.48	63.11
	SOLAR-base	45.42	27.62	18.15	12.47	31.59	48.84	61.12	63.27

Table 1: Evaluation results (%) of commonsense inference with base models.

		BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE-L	CIDEr	BERTScore
ConceptNet	COMET-large	17.88	11.35	7.13	4.00	13.47	19.36	37.72	54.07
	SOLAR-large	19.28	12.73	8.57	5.62	14.69	20.89	43.15	54.71
ATOMIC	COMET-large	54.05	34.92	24.04	17.62	35.06	56.93	75.46	64.84
	SOLAR-large	54.31	35.77	25.41	19.45	35.30	57.11	76.33	64.91
ATOMIC ₂₀	COMET-large	46.08	28.23	18.70	12.86	32.22	49.44	62.13	63.52
	SOLAR-large	46.51	28.99	19.52	13.73	32.53	49.76	63.24	63.58

Table 2: Evaluation results (%) of commonsense inference with large models.

	Cont.	Recon.	BLEU-3	CIDEr
SOLAR-base	✓	✓	18.27	61.15
	✓	✗	18.02	61.02
	✗	✓	17.89	60.90
	✗	✗	17.43	59.48

Table 3: Ablation study of commonsense representation learning methods on ATOMIC₂₀

379 ples by masking the span of each tuple elements
380 and randomly shuffling the order. The span length
381 is drawn from a Poisson distribution ($\lambda = 3$). SO-
382 LAR learns commonsense representation through
383 multi-task approach, and we set the task weight as
384 $\omega = 0.8$. In addition, we optimize the model using
385 the RecAdam (Chen et al., 2020a) optimizer to pre-
386 vent catastrophic forgetting during commonsense
387 representation learning. We set the hyperparam-
388 eters of the optimizer to $k = 0.001$ and $t_0 = 1000$.
389 After representation learning, we set the same hy-
390 perparameters as the baseline. We report the best
391 results among possible hyperparameter settings.

392 **Metrics.** To measure the commonsense inference
393 capability of SOLAR, we use common evaluation
394 metrics in the text generation: BLEU (Papineni
395 et al., 2002), ROUGE (Lin, 2004), CIDEr (Vedan-
396 tam et al., 2015) and BERTScore (Zhang et al.,
397 2020a).

398 4.3 Results

399 **Overall performance.** We evaluate SOLAR and
400 COMET on three commonsense KGs and report

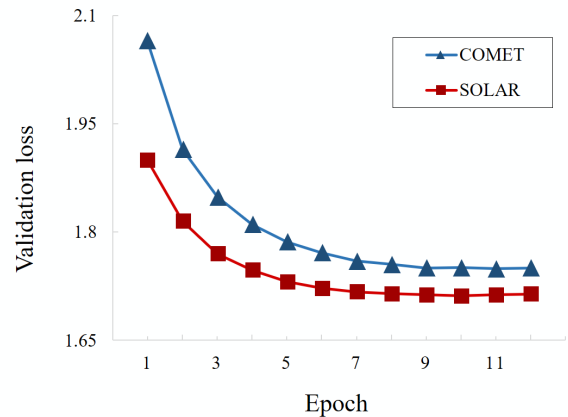


Figure 3: Validation loss of COMET-large and SOLAR-large on ATOMIC₂₀

401 the automatic evaluation results of generated in-
402 ferences. In our result tables, we denote model
403 names in form of (framework)-(BART model con-
404 figuration). For example, SOLAR and COMET
405 with BART-base are denoted by SOLAR-base and
406 COMET-base, respectively.

407 Table 1 shows that SOLAR-base outperforms
408 COMET-base for all KGs. By averaging over all
409 metrics, SOLAR-base improves the performance
410 of COMET-base on ConceptNet, ATOMIC, and
411 ATOMIC₂₀ by 1.74%, 0.57%, and 0.65%, respec-
412 tively. Experiments on large model configura-
413 tions establish the new state-of-the-art results on
414 commonsense inference with KGs. Table 2 shows
415 that SOLAR-large outperforms COMET-large, the
416 previous state-of-the-art, for all KGs and evalua-

Subject	Relation	Ground truth	COMET	SOLAR
PersonX is always busy	xReact	exhausted	busy	tired
sugar cube	ObjectUse	eat as food	mix with sugar	swetten coffee
PersonX gives PersonY a cup	HinderedBy	PersonY is not thirsty	PersonX is allergic to water	PersonX doesn't have a cup
PersonX likes the movie	HinderedBy	They were too busy texting	PersonX is allergic to the movie	The movie is too boring

Table 4: Examples of commonsense inference from COMET and SOLAR in ATOMIC₂₀²⁰.

metrics. We observe 1.84%, 0.70%, and 0.58% average performance improvement on ConceptNet, ATOMIC, and ATOMIC₂₀²⁰ respectively. Furthermore, SOLAR-base performs comparably to COMET-large on ATOMIC and ATOMIC₂₀²⁰, and performs better on ConceptNet, despite using only one-third of parameters. This shows the parameter-efficiency of our approach compared to COMET.

Analysis on commonsense inference. We provide further analysis on commonsense inference results of SOLAR and COMET. Figure 3 shows the validation loss curve for COMET-large and SOLAR-large. It is clearly observed that SOLAR gives smaller loss than COMET on validation sets, which indicates that SOLAR generalizes commonsense better than COMET. In addition, Table 4 shows examples of commonsense inference results by COMET and SOLAR. It can be observed that SOLAR generates plausible inferences with novel expressions, whereas COMET extracts words from the subject phrase to generate inferences, leading to trivial or wrong results. Another observation is that COMET is vulnerable to the annotation bias in KGs. For example, in ATOMIC₂₀²⁰, the word “allergic” frequently appears with relation “HinderedBy”, and COMET is biased to generate wrong inferences like “allergic to the moive”. In contrast, SOLAR makes better inference results without such bias.

Ablation study. We conduct an ablation study to measure the effectiveness of each component of our proposed framework. Table 3 shows that learning on both tasks performs better than learning on only one of the two tasks. We observe that contrastive learning of commonsense tuples is the key to our performance improvement that SOLAR achieves, and the reconstruction task also plays a role in the

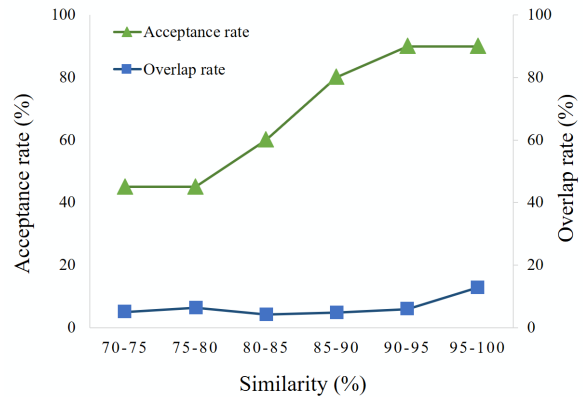


Figure 4: Acceptance and overlap rates of generated missing relations. Similarity is measured by BERTScore.

framework.

Acceptance of missing relations. We conduct a qualitative analysis of missing relations generated through our approach. Table 5 shows examples of tuple pairs and their similarity values measured by BERTScore. In the first row, “PersonX throws a huge party” and “PersonX throws a big party” are semantically similar, and each relation-object can be shared with the subject of the other (e.g., PersonX throws a huge party - oEffect - smile). In contrast, as in the last example, tuple pairs with a low similarity between subjects cannot share relation-object with one another. From these examples, we observe that tuple pairs with higher similarity between subjects generate more plausible tuples when their relation-object pair are shared, consistent with our intuition.

We further provide a quantitative analysis by measuring the acceptance rate of missing relations generated through our approach and comparing it with the overlap rate. Overlap rate is the probability of a missing relation already existing in the

Similarity (%)	Subject	Relation – object	Plausible
95.8	PersonX throws a huge party PersonX throws a big party	oReact-important oEffect-smile	✓
95.3	handgun pistol	AtLocation-army AtLocation-pants	✓
90.3	protective clothing safety gear	ObjectUse-keep them safe ObjectUse-protect from injury	✓
87.0	trash bags trashbins	ObjectUse-put things in ObjectUse-get rid of garbage	✓
82.0	PersonX takes PersonY to see a doctor PersonX takes PersonY to the vet	oEffect-get checked by doctor xWant-get dog checked	✗
70.1	PersonX hugs PersonY back PersonX screams at PersonY	oReact-loved and needed oEffect-sweats in terror	✗

Table 5: Qualitative analysis on examples of similarity-based tuple extraction from ATOMIC₂₀²⁰. Similarity is measured by BERTScore between the subjects of tuples. Humans evaluate whether the tuples are plausible after the relation-objects are replaced by that of each other.

Method	BLEU-3	CIDEr	BERTScore
Baseline	17.44	59.48	63.11
Fine-tuning	17.38	59.11	63.08
Contrastive Learning	18.15	61.12	63.27

Table 6: Evaluation results of methods for learning from missing relations.

graph. To measure the acceptance rate of missing relations, we randomly sample 20 missing relations per similarity interval (total 120 samples) and ask human annotators to determine their plausibility. Three workers annotated each missing relation as *accept* if it is plausible or *reject* otherwise, and we used majority voting as the final annotation. Figure 4 shows the acceptance rate of the missing relations regarding semantic similarity of subjects. It shows that the acceptance rate of missing relation is proportional to the similarity, and if the tuples have a similarity of greater than 90%, then 90% of the missing tuples are then valid. In contrast, when the similarity dropped below 85%, the acceptance rate decreased drastically. The blue line in Figure 4 represents the overlap rate according to the similarity. For tuple pairs of high similarity exceeding 90%, the overlap rate is significantly lower (< 20%) than the acceptance rate, which shows that novel missing relations can be effectively identified through our method.

Methods for learning from missing relations.

We investigate the effectiveness of our method for learning from missing relations. We compare our contrastive learning method with a fine-tuning method where missing relations are directly

added to a commonsense KG and subsequently learned. We use missing relations generated on subjects with exceeding 90% similarity. Table 6 shows that our proposed contrastive learning method shows best performance, while fine-tuning method is worse than the baseline. We speculate that direct fine-tuning is vulnerable to unacceptable relations, while our proposed contrastive learning framework is robust to them. These results indicate that directly learning from missing tuples harm the commonsense inference capability of the model. We speculate that our approach can handle noise or incorrect missing relations by implicitly learning from missing relations.

5 Conclusion

We have presented a novel contrastive learning framework of commonsense transformers, called SOLAR, to effectively learn from missing relations in commonsense KGs. Moreover, we have developed a new construction scheme for positive and negative sets of examples based on similarities in language model representations. By utilizing our carefully designed methods, SOLAR effectively learns both existing and missing relations of events, alleviating the difficulties in learning commonsense KGs. Our empirical evaluations of diverse commonsense KGs demonstrate the efficacy of SOLAR in commonsense inference. In particular, SOLAR consistently outperforms the state-of-the-art commonsense transformers across all the evaluation metrics and commonsense KGs.

References

- 534 Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi. 2019. Comet: Commonsense transformers for automatic knowledge graph construction. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4762–4779.
- 540 Sanyuan Chen, Yutai Hou, Yiming Cui, Wanxiang Che, Ting Liu, and Xiangzhan Yu. 2020a. Recall and learn: Fine-tuning deep pretrained language models with less forgetting. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7870–7881.
- 546 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020b. A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning*, pages 1597–1607. PMLR.
- 551 Saadia Gabriel, Chandra Bhagavatula, Vered Shwartz, Ronan Le Bras, Maxwell Forbes, and Yejin Choi. 2021. Paragraph-level commonsense transformers with recurrent memory. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 12857–12865.
- 557 Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738.
- 562 Olivier Henaff. 2020. Data-efficient image recognition with contrastive predictive coding. In *International Conference on Machine Learning*, pages 4182–4192. PMLR.
- 566 Jena D Hwang, Chandra Bhagavatula, Ronan Le Bras, Jeff Da, Keisuke Sakaguchi, Antoine Bosselut, and Yejin Choi. 2021. (comet-) atomic 2020: On symbolic and neural commonsense knowledge graphs. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 6384–6392.
- 572 Seanie Lee, Dong Bok Lee, and Sung Ju Hwang. 2020. Contrastive learning with adversarial perturbations for conditional text generation. In *International Conference on Learning Representations*.
- 576 Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2020. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880.
- 584 Xiang Li, Aynaz Taheri, Lifu Tu, and Kevin Gimpel. 2016. Commonsense knowledge base completion. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1445–1455.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Chaitanya Malaviya, Chandra Bhagavatula, Antoine Bosselut, and Yejin Choi. 2020. Commonsense knowledge base completion with structural and semantic context. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 2925–2933.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*, pages 3111–3119.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? 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 2463–2473.
- Hannah Rashkin, Maarten Sap, Emily Allaway, Noah A Smith, and Yejin Choi. 2018. Event2mind: Commonsense inference on events, intents, and reactions. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 463–473.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. 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 3982–3992.
- Adam Roberts, Colin Raffel, and Noam Shazeer. 2020. How much knowledge can you pack into the parameters of a language model? In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5418–5426.
- Itsumi Saito, Kyosuke Nishida, Hisako Asano, and Junji Tomita. 2018. Commonsense knowledge base completion and generation. In *Proceedings of the 22nd Conference on Computational Natural Language Learning*, pages 141–150.
- Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A Smith, and Yejin Choi. 2019. Atomic: An atlas of machine commonsense for if-then reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3027–3035.

645 Robyn Speer, Joshua Chin, and Catherine Havasi. 2017.
646 Conceptnet 5.5: An open multilingual graph of gen-
647 eral knowledge. In *Proceedings of the Thirty-First*
648 *AAAI Conference on Artificial Intelligence*, pages
649 4444–4451.

650 Ramakrishna Vedantam, C Lawrence Zitnick, and Devi
651 Parikh. 2015. Cider: Consensus-based image de-
652 scription evaluation. In *Proceedings of the IEEE*
653 *Conference on Computer Vision and Pattern Recog-
654 nition*, pages 4566–4575.

655 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien
656 Chaumond, Clement Delangue, Anthony Moi, Pier-
657 ric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz,
658 et al. 2019. Huggingface’s transformers: State-of-
659 the-art natural language processing. *arXiv preprint*
660 *arXiv:1910.03771*.

661 Yuanmeng Yan, Rumei Li, Sirui Wang, Fuzheng Zhang,
662 Wei Wu, and Weiran Xu. 2021. ConSERT: A con-
663 trastive framework for self-supervised sentence repre-
664 sentation transfer. In *Proceedings of the 59th Annual*
665 *Meeting of the Association for Computational Lin-
666 guistics and the 11th International Joint Conference*
667 *on Natural Language Processing (Volume 1: Long*
668 *Papers)*, pages 5065–5075.

669 Chenxi Yuan, Chun Yuan, Yang Bai, and Ziran Li. 2020.
670 Logic enhanced commonsense inference with chain
671 transformer. In *Proceedings of the 29th ACM Inter-
672 national Conference on Information & Knowledge*
673 *Management*, pages 1763–1772.

674 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Wein-
675 berger, and Yoav Artzi. 2020a. Bertscore: Evaluating
676 text generation with bert. In *International Confer-
677 ence on Learning Representations*.

678 Yan Zhang, Ruidan He, Zuozhu Liu, Kwan Hui Lim,
679 and Lidong Bing. 2020b. An unsupervised sentence
680 embedding method by mutual information maximiza-
681 tion. In *Proceedings of the 2020 Conference on*
682 *Empirical Methods in Natural Language Processing*
683 *(EMNLP)*, pages 1601–1610.