Symbolic Knowledge Distillation: from General Language Models to Commonsense Models

Peter West^{†‡*} Chandra Bhagavatula[‡] Jack Hessel[‡] Jena D. Hwang[‡] Liwei Jiang^{†‡} Ronan Le Bras[‡] Ximing Lu^{†‡} Sean Welleck^{†‡} Yejin Choi ^{†‡*} [†]Paul G. Allen School of Computer Science & Engineering, University of Washington [‡]Allen Institute for Artificial Intelligence

Abstract

The common practice for training commonsense models has gone *from-human-to-corpus-to-machine*: humans author commonsense knowledge graphs in order to train commonsense models. In this work, we investigate an alternative, *from-machine-to-corpus-to-machine*: general language models author these commonsense knowledge graphs to train commonsense models.

Our study leads to a new framework, **Symbolic Knowledge Distillation**. As with prior art in Knowledge Distillation (Hinton et al., 2015), our approach uses larger models to teach smaller models. A key difference is that we distill knowledge symbolically—as text—in addition to the resulting neural model. We distill only one aspect—the commonsense of a general language model teacher, allowing the student to be a different type of model, a commonsense model. Altogether, we show that careful prompt engineering and a separately trained critic model allow us to selectively distill high-quality causal commonsense from GPT-3, a general language model.

Empirical results demonstrate that, for the first time, a *human-authored* commonsense knowledge graph is surpassed by our *automatically distilled* variant in all three criteria: quantity, quality, and diversity. In addition, it results in a neural commonsense model that surpasses the teacher model's commonsense capabilities despite its 100x smaller size. We apply this to the ATOMIC resource, and will share our new symbolic knowledge graph and commonsense models¹.

1 Introduction

Prior works have suggested that pre-trained language models possess limited understanding of commonsense knowledge (Merrill et al., 2021; Talmor et al., 2021; Davis and Marcus, 2017) despite

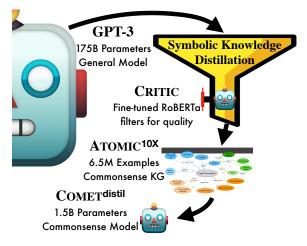


Figure 1: **Symbolic knowledge distillation** extracts the commonsense from the large, general language model GPT-3, into 2 forms: a large commonsense knowledge graph **ATOMIC**^{10x}, and a compact commonsense model **COMET**^{DIS}_{TIL}. The quality of this knowledge can be controlled and improved by adding a **critic** model, making GPT-3 a stronger teacher.

otherwise stellar performance on leaderboards. As a result, symbolic commonsense knowledge graphs (Speer et al., 2017; Sap et al., 2019; Hwang et al., 2021) and corresponding neural representations (Bosselut et al., 2019; Hwang et al., 2021; Zhang et al., 2020b) have supplemented past models with commonsense capabilities. This has enabled diverse downstream applications, including interactive learning through a conversational interface (Arabshahi et al., 2021), persona- and affect-aware conversation models (Kearns et al., 2020), figurative language understanding (Chakrabarty et al., 2020, 2021), story telling (Ammanabrolu et al., 2021a) and fantasy games (Ammanabrolu et al., 2021b).

The common practice for commonsense knowledge graph construction sees humans spell out as many pieces of knowledge as possible. This pipeline goes *from-human-to-corpus-to-machine*, with commonsense models trained from human-

¹We will share this following the anonymity period. We have permission from OpenAI to release GPT-3 generations

authored knowledge graphs. Yet, high-quality, human-authored knowledge is expensive to scale, limiting coverage; this motivates an alternative: from-machine-to-corpus-to-machine. Prior efforts toward automatic commonsense knowledge graphs have resulted in considerably lower quality than human-written data (Hwang et al., 2021; Zhang et al., 2020b), which in turn leads to less reliable neural models (Hwang et al., 2021). Broad literature consistently shows machine-authored knowledge graphs underperform human-authored graphs (Etzioni et al., 2011; Mitchell et al., 2015; Bollacker et al., 2008).

In this work, we propose Symbolic knowledge **distillation**, a new conceptual framework towards high-quality automatic knowledge graphs for commonsense, leveraging state-of-the-art models and novel methodology. Most prior art for automatic knowledge graph construction extracts knowledge from raw text (Bhakthavatsalam et al., 2020; Zhang et al., 2020a; Zhou et al., 2020; Zhang et al., 2020b; Li et al., 2020). In contrast, our approach is motivated by knowledge distillation (Hinton et al., 2015) wherein a larger teacher model transfers knowledge to a compact student model (§2.1). Our method differs from prior knowledge distillation in key ways: we distill a symbolic knowledge graph (i.e., generated text) in addition to a neural model, and we distill only a selective aspect of the teacher model. This selectively allows the student model to be of a different type (commonsense model), compared to the teacher (general language model), enriching the scope of distillation. An added benefit is that knowledge distilled as text is human readable: it can be understood and evaluated.

A general language model—GPT-3 in our case—is an imperfect commonsense teacher on its own, and the ability to evaluate distilled knowledge is useful in improving it. We empirically demonstrate that, by training a separate critic model to judge symbolic generation quality, a more precise teacher can be defined. Knowledge from this critical teacher is higher quality—even exceeding human-authored knowledge. Yet even before training a critic, our study makes the unexpected finding that the student model surpasses the commonsense of GPT-3, our knowledge source.

To test symbolic knowledge distillation against the *human–to–corpus–to–machine* paradigm, we compare with ATOMIC $_{20}^{20}$ (Hwang et al., 2021), which is a human-authored commonsense knowl-

edge graph. We find that **ATOMIC**^{10x}, our machine-generated corpus, exceeds the human generated corpus in *scale*, *accuracy*, and *diversity* with respect to 7 commonsense inference types that we focus on in this study. The resulting commonsense model, **COMET**^{DIS}_{TIL}, not only surpasses the human-trained equivalent COMET²⁰₂₀, but is also smaller, more efficient, and produces commonsense at a higher accuracy than its own teacher–GPT-3.

Symbolic knowledge distillation offers a promising new role for general language models, as commonsense knowledge sources, and humans, as small-scale evaluators to train critic models rather than authors of commonsense knowledge. Our work demonstrates that humans and LMs can be effective collaborators for curating commonsense knowledge graphs and training efficient and performant commonsense models.

2 Overview and Key Findings

Throughout our work, we describe the *machine*—*to*—*corpus*—*to*—*machine* methodology of symbolic knowledge distillation. We first go *machine*—*to*—*corpus* (§3), by decoding from GPT-3, then improve our knowledge with a specialized critic model (§4), and finally distill this knowledge into an efficient commonsense model (§5), going *corpus*—*to*—*machine*. Throughout this process, we evaluate against a human knowledge source, comparing our automatic knowledge graph ATOMIC ¹⁰x and commonsense model COMET ^{DIS} to the human-authored ATOMIC ²⁰20 and resulting model COMET ²⁰20 (Hwang et al., 2021).

2.1 Symbolic Knowledge Distillation

Our proposed methodology parallels knowledge distillation (Hinton et al., 2015), a method for compressing a large or complicated teacher distribution P_t into a smaller/simpler student distribution P_s . Key to knowledge distillation² is the notion of minimizing the cross-entropy between P_t and P_s :

$$H(P_t, P_s) = -\sum_{y \in Y} P_t(y) \log P_s(y) \qquad (1)$$

Knowledge is transferred to the student by encouraging it to match teacher predictions. Hinton et al. (2015) apply this to conditional classification: for

²In its simplest case, with temperature set to 1.0

X starts running	xEffect	gets in shape	X sings a song	HinderedBy but not if	X can't remember the lyrics
X and Y engage in an argument	xWant so, X wants	to avoid Y	X is not well	xReact so, X feels	lonely
X learns to type fast	xNeed X needed	to have taken typing lessons	X takes care of a monkey	xAttr X is seen as	kind
X steals his grandfather's sword	xEffect so, X	is punished by his grandfather	X butts in	HinderedBy but not if	X is too shy to speak up
X takes up new employment	xIntent because X wants	to be self sufficient	X waits for the storm to break	xEffect	is safe from the storm

Figure 2: Example **automatically generated** ATOMIC triples from our ATOMIC^{10x} commonsense knowledge graph. Each example includes a generated **event**, **relation** (with natural language interpretation), and generated **inference**.

each training input, P_t and P_s are model predictions over label set Y. Typically Y is a tractable set, over which this sum can reasonably be calculated.

For distilling the knowledge of generative models, we can think of an unconditional language model (LM e.g. GPT-3) as P_t . This makes Y the set of all strings, over which LMs define probability. Unfortunately Y is an exponential set, intractable to sum over in Eq 1. Kim and Rush (2016) address this problem by simply taking the mode of P_t over Y, truncating most of the teacher distribution to the most likely sequence and discarding information.

Instead, we consider a sampling-based interpretation of the same objective:

$$H(P_t, P_s) = \underset{y \sim P_t(y)}{\mathbb{E}} \left[-\log P_s(y) \right]$$
 (2)

which exactly equals the cross-entropy of Eq 1, at the limit under pure sampling from P_t .³

Yet distilling all knowledge from the teacher may not be desirable—our work is specifically focused on distlling commonsense knowledge from GPT-3. The ideal teacher P_t is a commonsense expert, but GPT-3 can approximate such a teacher, off-the-shelf, via prompting. This ability to select information is one explicit benefit of the sampling-based interpretation of Eq 2: while Eq 1 uses continuous logits over existing data, sampling gives discrete control over transferred information, by selecting which samples are elicited and used. For the general language model GPT-3, We encourage domain/quality with prompting, and sample truncation (Holtzman et al., 2020). We call this the loose teacher P_t^L —knowledge is generated and

transferred from GPT-3, but without critical assessment of correctness (§3).

In fact, sampling knowledge in Eq 2 offers even more control, as generations can be individually interpreted and judged. Given an indicator function A(x) for which knowledge x is correct, we can define a stronger teacher model. Using a Product of Experts (Hinton, 2002) between the loose teacher P_t^L and and the critic A(x), we define a critical teacher:

$$P_t(x) \propto P_t^L(x|p) \cdot A(x)$$
 (3)

In practice, A(x) is a textual classifier learned on human judgements, 1 for knowledge predicted to be correct and 0 otherwise. Thus, the critic gives control over the correctness and confidence of the knowledge that is transferred (§4).

2.2 Key Findings

Applying symbolic knowledge distillation in practice results in promising and surprising findings:

- 1. Learning symbolic knowledge from language models can be framed as a symbolic extension to knowledge distillation. In §2.1, we describe learning commonsense as a symbolic extension to knowledge distillation, with GPT-3 a knowledge source. We elaborate on this process with positive results in §3,4, and 5.
- 2. Symbolic knowledge distillation constructs a high quality knowledge graph at scale. Our method naturally yields a machine-generated commonsense knowledge graph, which can achieve impressive quality (§4), beyond that of human-authored data. An effective critic which filters incorrect generated knowledge is key.
- **3.** A critical teacher results in a higher quality student. In §4, we show that making the teacher

³A useful consequence of this framing is that access to the full model distribution is not required. Our experiments (§3) use GPT-3, for which the distribution is **not available**, thus our method is applicable while knowledge distillation is not.

more critical results in higher quality knowledge, even as it reduces the scale of knowledge transferred. This demonstrates that *quality* matters, not just *quantity*, as higher quality knowledge results in a higher quality commonsense model in §5 despite smaller scale data.

- **4.** Critical teacher or not, a student can outperform the knowledge source. In §5, we show the unexpected result that all student models exceed the quality of GPT-3, the knowledge source.
- 5. Machines can win over humans for automatic knowledge graph construction. In §4 and §5, we show that machine generated knowledge and the resulting commonsense model can outperform their equivalents that use a human knowledge source. Our symbolic knowledge exceeds humans at scale, quality, and diversity. The resulting commonsense model achieves the most accurate commonsense KG completions.

3 Machine-to-Corpus Verbalization

Symbolic knowledge distillation begins by going *machine–to–corpus*, i.e. generating many commonsense facts, which results in a commonsense knowledge graph. §2.1 frames this as sampling to estimate the knowledge distillation objective–a student commonsense model learns from the generations of a teacher (GPT-3).

We start with a *loose teacher*, transferring knowledge by prompted generation with truncated sampling alone—this is in contrast to the *critical teacher* (§4) which explicitly judges and filters the generated samples. The loose teacher uses few-shot prompting as in Brown et al. (2020). We use a few-shot template:

```
\begin{array}{l} \texttt{<TASK-PROMPT>} \\ \texttt{<EX}_1 - \texttt{INP>} \texttt{<EX}_1 - \texttt{OUT>} \\ \dots \\ \texttt{<EX}_{N-1} - \texttt{INP>} \texttt{<EX}_{N-1} - \texttt{OUT>} \\ \texttt{<EX}_N - \texttt{INP>} \end{array}
```

where $\langle EX_i-INP \rangle/\langle EX_i-OUT \rangle$ are humanauthored, natural language ATOMIC entries, and $\langle TASK-PROMPT \rangle$ is a description of the problem. Given such a prompt, GPT-3 generates the *missing piece*, output $\langle EX_N-OUT \rangle$ for input $\langle EX_N-INP \rangle$, following the pattern of earlier examples (1 to N-1). We find important aspects for producing high-quality commonsense knowledge:

• Examples should be numbered. e.g.

- $\langle \text{EX}_5 \text{INP} \rangle$ might begin with "5)" to indicate it is the 5th example.
- The format of $\langle EX_i INP \rangle$ and $\langle EX_i OUT \rangle$ should linguistically imply the relationship between them. See below for examples.
- <TASK-PROMPT> can be used to give extra specification to complicated problems.

3.1 Data: ATOMIC

We demonstrate symbolic knowledge distillation on the ATOMIC *if-then* resource (Sap et al., 2019). This follows an event-relation-inference (triple) format. The corpus links *events* (e.g. *X attacks Y*) to relations, e.g. **HinderedBy** which describes what might hinder an event. For a relation/event, the goal is to generate a resulting inference, e.g. *X attacks Y* **HinderedBy** *X is restrained*.

Of the 23 relations from the most recent version—ATOMIC₂₀²⁰—we limit our investigation to 7 relations that correspond to *causal* commonsense knowledge: **xAttr** (how X is perceived after *event*), **xReact** (how X reacts in response to *event*), **xEffect** (what X does after *event*), **xIntent** (X's intent in *event*), **xWant** (what X wants after *event*), **xNeed** (what X needed for *event* to take place) and **HinderedBy**. We describe how **verbalization** is applied to ATOMIC data in 2 steps: generating underlying events (heads), then full examples (inference given event).

3.2 Event Generation

Events are context-free premises in ATOMIC involving PersonX (and sometimes a second PersonY) in various scenarios. These events form heads in knowledge graph triples. We generate events by filling in the elements of our template:

```
    Event: X overcomes evil with good
    Event: X does not learn from Y
    Event: X looks at flowers
    11.
```

The format is simple, as events are generated *unconditionally*. We use 100 high-quality events from the ATOMIC $_{20}^{20}$ corpus for our prompt, selected to avoid grammatical or logical errors, and minimize semantic overlap. We randomly sample 10 of these seed events for each generation batch, resulting in randomized prompts. We use nucleus sampling (p=0.9) (Holtzman et al., 2020), and presence/frequency penalties of 0.5 from the GPT-3 interface. We generate 165K unique events using

the 175B-parameter Davinci model⁴ from Brown et al. (2020) (human-authored ATOMIC $_{20}^{20}$ contains only 6.2K events).

3.3 Inference Generation

Generating ATOMIC inferences requires reasoning about events and relations together. We design verbalization templates fo reach relation, with iterative design and small-scale verification by the authors⁵ e.g. we prompt the **xNeed** relation as follows:

What needs to be true for this event to take place?

. . .

Event <i>: X goes jogging

Prerequisites: For this to

happen, X needed to wear running

shoes

...

Event <N>: X looks at flowers
Prerequisites: For this to
happen,

The language of this template implies the relation-specific task, both "Prerequisites:" and beginning with "for this to happen" suggest the **xNeed** relation. As well, we include an xNeed-specific <TASK-PROMPT>. We use 10 few-shot examples for each prompt.⁶

For each event/relation (165K X 7) we generate 10 inferences with the Curie GPT-3 model⁷ and earlier hyperparameters. Removing duplicate and degenerate (e.g. fewer than 3 characters) generations yields 6.46M ATOMIC-style data triples (examples in Figure 2). We call this ATOMIC^{10x}, as it contains an order of magnitude more triples than ATOMIC²⁰ for the 7 relations we study.

3.4 Evaluating a Generated Commonsense Knowledge Graph

Machine generation enables a large scale of unique generations at a much lower cost than human-authored knowledge (Table 1), but what kind of examples are produced by GPT-3, and how does

it differ from knowledge produced by humans? In this section, we conduct an in-depth analysis to answer these questions.

Lexical Differences: Diversity and Uniqueness Recent work finds that machine generations can be repetitive and lack diversity (Welleck et al., 2020; Holtzman et al., 2020); one way generated knowledge may differ from human-authored is less creative word choice, diversity, or more repetition.

To test this, we begin with lexical diversity (i.e. unique words used, Table 2). While there is variation by relation, the diversity of $ATOMIC^{10x}$ actually exceeds $ATOMIC^{20}_{20}$ here, 5.2M unique words to 1.5M. In addition, it contains significantly more strictly unique generated inferences (Table 2, unique tails).

BLEU Soft Uniqueness. Exact match (above) fails to capture the notion of *similar* text. Following the intuition of self-BLEU (Zhu et al., 2018), we define *soft uniqueness* to describe diversity of generations in a corpus. An inference x is softly-unique if:

$$BLEU_2(C, x) < 0.5$$

where C is the set of inferences for a given input (in our case, event + relation), and 0.5 is an empirical threshold. To find soft-uniqueness of a corpus, we iteratively remove examples until all are softly unique, i.e. low mutual lexical overlap; higher diversity means more such examples (thus a larger softly unique corpus is preferable). Softly-unique corpus sizes are given in Table 4 ("Size (div)"). ATOMIC^{10x} has a smaller *fraction* of softly-unique examples than ATOMIC²⁰₂₀, yet it contains many more such examples. ATOMIC^{10x} contains 4.38M such examples (full size 6.5M) vs. ATOMIC²⁰₂₀, which has 560K (full size 600K).

Model-based Diversity Measurement. Lexical notions of diversity reward differences in surface form, which may not always reflect diversity of *information*, only format. Thus, we next study information-theoretic measures for diversity. Intuitively, diverse information should be less predictable, or higher entropy. With GPT-2 XL models finetuned on ATOMIC²⁰ and ATOMIC^{10x} (§5) we estimate **entropy**—roughly, how difficult it is for a model to capture the corpus information (Table 3). This is 4 times higher for ATOMIC^{10x}, suggesting more content from a modeling perspective. We also estimate **cross-entropy**—how

⁴the largest available version of GPT-3

⁵See Appendix D for full prompts.

⁶We also replace anonymous names ("X") with sampled generic names as this improved quality, See Appendix D. Once generation is complete, we substitute in generic markers ("X") for the final dataset.

⁷for the largest, Davinci, 12M generations is computationally/monetarily intractable.

Relation	$ATOMIC_{20}^{20}$	$ATOMIC^{10x}$
HinderedBy	77,616	1,028,092
xNeed	100,995	760,232
xWant	109,098	730,223
xIntent	54,839	965,921
xReact	62,424	1,033,123
xAttr	113,096	884,318
xEffect	90,868	1,054,391
Total Count	608,936	6,456,300
Est Total Cost	~\$40,000	~\$6,000
Est Cost Per Triple	~\$0.06	~\$0.001

Table 1: Number of unique triples with the given relation, $|(\cdot, \text{relation}, \cdot)|$. The estimated cost for ATOMIC^{10x} comes at a fraction of a conservative estimation for ATOMIC²⁰ crowdsourcing costs.

	Ler	igth		ique ns (K)		que s (K)
	${f A}_{20}^{20}$	A ^{10x}	\mathbf{A}_{20}^{20}	A ^{10x}	${f A}_{20}^{20}$	A ^{10x}
xWant	4.69	5.16	322	784	69	152
xAttr	1.42	2.73	15	21	11	8
xEffect	3.92	4.66	216	864	55	185
xIntent	4.59	5.92	136	800	30	135
xNeed	4.51	5.97	289	1378	64	231
xReact	4.03	1.77	48	5	12	2
HinderedBy	7.93	7.49	522	1775	290	874
Events	5.20	5.32	109	881	6.2	165

Table 2: Average length, total unique tokens and total unique examples (in K, i.e. 1000s) by relation type and in events (bottom row) from ATOMIC_{20}^{20} (A_{20}^{20}) and ATOMIC^{10x} (A^{10x}).

Entropy	Cross Entropy	KL Divergence
$H(D_1) = 1.27 \mid$	$H(D_1, D_2) = 9.31$	$ D_{KL}(D1 D2) = 8.04$
$H(D_2) = 7.80 \mid$	$H(D_2, D_1) = 41.48$	$ D_{KL}(D_2 D_1) = 33.68 $

Table 3: Entropy, cross-entropy, and divergence of $\mathrm{ATOMIC}_{20}^{20}\left(D_{1}\right)$ and $\mathrm{ATOMIC}_{20}^{\mathbf{10x}}\left(D_{2}\right)$.

well a model trained on one corpus describes the other. From ATOMIC^{10x} to ATOMIC^{20}_{20} , this is 9.31, only 2 points higher than its entropy suggesting ATOMIC^{20}_{20} is describable with information from ATOMIC^{10x} . In reverse, this is 41.48 suggesting much of ATOMIC^{10x} is not captured by ATOMIC^{20}_{20} – ATOMIC^{10x} is surprising given only information from ATOMIC^{20}_{20} .

Human Evaluation of Quality. Perhaps most importantly, we study the *quality* of knowledge in each corpus. We conduct human evaluation with Amazon Mechanical Turk. 3 annotators rate each triple resulting in "accepted", "rejected" or "no

Corpus	Accept	Reject	N/A	Size	Size (div)
$ATOMIC_{20}^{20}$	86.8	11.3	1.9	0.6M	0.56M
ATOMIC ^{10x}	78.5	18.7	2.8	6.5M	4.38M
	88.4	9.5	2.1	5.1M	3.68M
(critic _{low})	91.5	6.8	1.7	4.4M	3.25M
	95.3	3.8	1.0	3.0M	2.33M
(critichigh)	96.4	2.7	0.8	2.5M	2.00M
+ GPT-J	72.0	27.6	0.4	-	-
+ T5-11B LM	71.7	26.9	1.4	-	-

Table 4: Attributes of ATOMIC^{10x} and ATOMIC^{10x} (row 2) including the critic model (§4, rows 3 - 6) with various filtering cutoffs. Accept and Reject are by majority human vote unless any mark N/A. Size is in unique examples⁹. The highest precision corpus is ATOMIC^{10x} with (critic_{high}), but multiple versions surpass ATOMIC²⁰. We also include alternate models (GPT-J and T5-11B) as the loose teacher.

judgement". We evaluate 3000 examples⁸ from ATOMIC^{10x}, and 1000 from ATOMIC²⁰ (Table 4). We find Fleiss' kappa (Fleiss, 1971) of 40.8 indicating moderate agreement (Landis and Koch, 1977), and 90.5% accuracy agreement. We require workers meet an Amazon Mechanical Turk qualification for annotation quality based on past commonsense evaluations. We compensate workers \$0.17 per task, which we estimate require 30 seconds. Further details and task template are in appendix §A.

For the *loose teacher*, consider the top row of ATOMIC^{10x} in Table 4 (other rows add the critic §4). ATOMIC^{10x} exceeds ATOMIC²⁰ in scale, but is somewhat less acceptable by human raters—by roughly 8 percentage points. Yet, the larger scale of ATOMIC^{10x} implies a significantly higher *number* of accurate examples. Increasing the proportion of these is the main objective of the critic (§4).

How do Knowledge Sources Compare? To understand the robustness of our approach, we assess other language models as the knowledge source (i.e. loose teacher): GPT-J (Wang and Komatsuzaki, 2021) and T5-11B adapted for language modelling (Lester et al., 2021). We substitute both for GPT-3 as in §3.2,3.3, generating a small-scale corpus to evaluate. We conduct human evaluation on 1000 examples as above (Table 4). Both models attain roughly 72% accuracy, 6 points below GPT-3 (78.5). This suggests strong potential, but higher quality from GPT-3. We explore this further in Appendix B.

⁸this ensures at least 1000 after filtering by the critic §4)

4 Making the Teacher More Critical

Symbolic knowledge distillation requires a strong teacher model to maximize the quality of the generated knowledge graph and resulting student model (§5). While the *loose teacher* (GPT-3 alone) results in a viable commonsense knowledge graph, evaluation shows this isn't a perfect commonsense teacher. Thus, we multiply in a critic model, to filter lower-quality knowledge, correcting the teacher (§2.1). With modest supervision (a small-scale human evaluation) we train a classifier to predict and discriminate unacceptable examples. We multiply this with the loose teacher §3, creating a critical teacher product of experts. In practice this means filtering ATOMIC^{10x} to create new corpora that are higher quality, yet still larger scale than humanauthored ATOMIC $_{20}^{20}$.

Training a knowledge critic We gather a training set of *correct vs. incorrect* human judgments on a randomly-sampled set of 10K entries of ATOMIC^{10x}, as in §3.4 but with one annotation per example. We take a (random) train/dev/test split of 8k/1k/1k. While this step requires human annotation, humans take on the role of high-level supervisors here—critiquing a small number of generations rather than authoring the entire knowledge graph as in previous work. Indeed, the cost/complexity of this step is similar to a typical human evaluation, making it far cheaper/easier than eliciting human-authored knowledge in past work.

We train binary classifiers (critics) for human acceptability using RoBERTa-Large (Liu et al., 2019). We find pretraining on MNLI results in the best model in terms of precision and recall, and we suggest this technique for future studies. We give more detail in Appendix C, including baselines. Our best model vastly improves the accuracy of ATOMIC 10x (Table 4), demonstrating that a small amount of human supervision can consistently help to correct GPT-3's mistakes.

Size-accuracy trade-off Using our critic to filter knowledge results in a natural trade-off between size and accuracy. We test several cutoffs for ATOMIC^{10x}, i.e. confidence at which the critic rejects examples. We report human-measured accuracy (Accept/Reject column Table 4) following §3.4. We compare the loose

	Random	Inf	Event	EMAP	Full
AP	79.3	81.9	86.2	87.1	94.0

Table 5: Average Precision for ablated critic models. The critic not only filters *awkward phrasings* which can be identified by either the event (**Event**) or inference (**Inf**) in isolation (EMAP only identifies these), but also *logical misalignments*, which require modeling interactions between event/inference, i.e. the full critic (**Full**).

teacher (unfiltered) to critical teachers. Discarding 20% of instances that the critic judges as least acceptable (reducing corpus size from 6.5M to 5.1M), ATOMIC 10x 's accuracy rises $78.5 \rightarrow 88.4$; human-authored ATOMIC $^{20}_{20}$ contains 600K entries at 86.8% accuracy. Reducing to total size to 2.5M examples (38% of full size), we attain 96.4% accuracy, nearly 10 points above ATOMIC $^{20}_{20}$ while still 4X larger.

What gets filtered out? We qualitatively identify two types of filtered triples: 1) logical misalignments, events/inferences joined in an inconsistent manner. Recognizing these requires understanding events-inference interactions, e.g., X cannot find his shirt as a result X is wearing a shirt; 2) awkward phrasings, in which events/inferences are individually incoherent e.g. PersonX has a fire in the bath—resulting triples are invalid as the event is implausible.

To understand what is filtered, we ablate the critic (Table 5): our full model is compared to a random predictor, event-only model, and inference-only model. We also compare to an EMAP (Hessel and Lee, 2020) version, i.e. an ensemble of event and inference-only, without interactions between event/inference (needed for *logical misalignments*).

We find GPT-3 produces both independent awkwardly-phrased events/inferences (filtered by X-only models) and logical misalignments. The classifier, trained on validated knowledge triples, helps in both cases. The EMAP of our full model (identifies only awkward phrasings) achieves 87% AP, and our full model (which additionally identifies logical misalignments) improves to 94% AP.

Does filtering hurt diversity? One concern is that the critic may keep only similar "safe" examples, lacking novelty. We repeat our diversity analysis (§3.4) for critical corpora (Table 4, "Size (div)", higher=better). As we filter, we surprisingly observe proportionally *more* diverse examples: full

 $^{^9}$ Size of ATOMIC $_{20}^{20}$ is given as the number of comparable datapoints, i.e. those with the same relations as ATOMIC 10x .

CKG Completion Model	Train Corpus Acc	Accept	Reject	N/A
GPT2-XL zero-shot	-	45.1	50.3	4.6
GPT-3	-	73.3	24.1	2.6
COMET ²⁰ ₂₀	86.8	81.5	16.3	2.2
COMET _{TIL} +critic _{low} +critic _{high}	78.5	78.4	19.2	2.4
	91.5	82.9	14.9	2.2
	96.4	87.5	10.2	2.3

Table 6: Model performance on knowledge base completion, measured by human judgement. Inferences are generated on held-out events from ATOMIC_{20}^{20} . Models besides GPT-3 use GPT-2 XL architecture. COMET_{TIL}^{DIS} with a strong critic (+critic_high) achieves the highest acceptance rate overall–87.5.

ATOMIC^{10x} has a diverse subset 68% of its size; rising to 80% with the most extreme filtering. One possibility is that GPT-3 gravitates towards common sentence structures for inconsistent knowledge. These would be recognizable to the critic, and removing them would increase both quality and diversity. This surprising result warrants further study.

5 Corpus-to-Machine: Distillation

The final step of symbolic knowledge distillation trains a compact model on the generated natural language knowledge graph. Our base model is GPT2-XL trained on all of ATOMIC^{10x}: we denote this model by COMET^{DIS}_{TIL}. We additionally train the model on critical versions of ATOMIC^{10x}–crit_{low} denotes training on the corpus achieving 91.5% accuracy, and crit_{high} on the 96.4% accuracy corpus. Models are trained for 1 epoch, with default parameters using the Huggingface Transformers library (Wolf et al., 2019).

5.1 Evaluating a Symbolically Distilled Model

Evaluation follows past work (Hwang et al., 2021; Bosselut et al., 2019; Sap et al., 2019) testing the ability of models to do knowledge base completion, i.e. generating inferences for test events, specifically from the ATOMIC $_{20}^{20}$ test set. We use human evaluation $_{20}^{10}$ following Section 3.4, on 1000 inputs (event + relation), with results in Table 6. We compare to the GPT2-XL-based COMET $_{20}^{20}$ model trained on human-generated ATOMIC $_{20}^{20}$, and GPT-

3 using the same generation method as §3–in effect, comparing the student COMET_{TIL} to the *loose teacher* GPT-3. We omit the *critical teacher* (GPT-3 + critic), which is not assured to produce an inference for each input, as the critic may reject all tails for some inputs. We also compare to zero-shot GPT2-XL (Radford et al., 2019) using the same methodology (Table 6).

How does COMET_{TIL} compare to GPT-3? In knowledge distillation, the student model often deteriorates in performance (Hinton et al., 2015; Kim and Rush, 2016) compared to its teacher. Comparing our base teacher–GPT-3–to the simplest version of COMET_{TIL} (top-row COMET_{TIL} of Table 6) surprisingly shows the student *surpasses* GPT-3, the model that generates its training data¹¹. We posit that the superior performance of COMET_{TIL} may have to do with mistakes of GPT-3 being filtered by verbalization and training of GPT-2, and possibly the focus of COMET_{TIL} on one commonsense domain while GPT-3 covers a more general domain. We leave further study of this effect for future work.

How does COMET DIS compare to human knowledge? While COMET Without the critic is slightly outperformed by COMET DIS in terms of accuracy, this reverses with the critic. For both cutoffs tested, COMET SUPPLIES SUPPLIES COMET DIS SUPPLIES COMET DIS SUPPLIES OF THE PROPERTY OF THE PROPER

Usefulness of COMET_{TIL}^{DIS} For on-demand inference, where a single high quality inference for some input event/relation is required, COMET $_{\rm TIL}^{\rm DIS}$ is the **best available model**: the most performant version surpasses COMET $_{20}^{20}$ by 5 points and GPT-3 by over 10. The critical teacher (GPT-3 + critic) yields a more accurate *corpus*, but may filter all inferences for an input, giving no output.

Limits and Future Work The success of symbolic knowledge distillation is a first step—demonstrating superior performance to human authoring on the commonsense relations tested here. No aspect of our approach is specific to these relations, yet further work is needed to explore the feasibility of generation for other aspects of commonsense and knowledge, beyond these relations, to concepts like physical or temporal commonsense.

¹⁰We find Fleiss' kappa (Fleiss, 1971) of 47.1 for acceptance, indicating moderate agreement. (Landis and Koch, 1977), and accuracy agreement of 88.7%.

¹¹The slight difference in acceptability for GPT-3 from Table 4 is likely due to variance in raters between rounds of evaluation, and a different distribution of events–Table 4 uses generated events while Table 6 uses events from ATOMIC $_{20}^{20}$.

6 Related Work

Commonsense Knowledge Graphs (CKG) CKGs provide knowledge for commonsense reasoning. Some are manually constructed, e.g. ATOMIC (Sap et al., 2019; Hwang et al., 2021). ConceptNet (Speer et al., 2017) contains taxonomy and physical commonsense, authored by humans or compiled from such sources. Some CKGs are automatically constructed: TransOMCS (Zhang et al., 2020a) extracts 18.48M tuples from syntactic parses and CausalBank (Li et al., 2020) extracts 314M cause-effect pairs by pattern-matching. In contrast, we *generate* commonsense.

Extracting Knowledge from LMs Past work uses models for automatic knowledge graph completion (Bosselut et al., 2019; Hwang et al., 2021; Li et al., 2020). Yet, models are trained on *existing* resources; ATOMIC^{10x} is generated without these. Other works mine factual/commonsense knowledge directly from off-the-shelf LMs (Petroni et al., 2019; Davison et al., 2019; Xiong et al., 2020), but not resulting in the quality at scale of ATOMIC^{10x}.

Knowledge Distillation Other works use knowledge distillation (Hinton et al., 2015) for generation. (Sanh et al., 2019) follow a label smoothing formulation, while Kim and Rush (2016) follow a similar formulation to us (§2.1), but use the mode of the teacher distribution rather than sampling. Our work is unique in distilling *specific* information (commonsense) from a general language model.

Data Generation While manual dataset creation is expensive and complex (Schwartz et al., 2017; Agrawal et al., 2018; Tsuchiya, 2018; Bras et al., 2020), crowdsourcing is the most popular method for goal-oriented, high quality/coverage datasets.

Past automatic data mainly use extractive approaches, e.g. syntactic parsing (Zhang et al., 2020a) or pattern matching (Li et al., 2020) from unstructured text (Lehmann et al., 2015; Buck et al., 2014). These scale, but are noisy and limited in format—ATOMIC knowledge will not appear simply in natural text. Some works explore automatic data synthesis/expansion by finetuning LMs on existing labeled data (Anaby-Tavor et al., 2020; Papanikolaou and Pierleoni, 2020; Kumar et al., 2020; Yang et al., 2020), but are limited by data quality.

7 Conclusions

We introduce symbolic knowledge distillation, a *machine–to–corpus–to–machine* pipeline for commonsense that does not require human-authored knowledge–instead, using machine generation. Knowledge is transferred from a large, general model to a compact commonsense model, through a commonsense corpus–yielding a commonsense knowledge graph and model. Our resulting symbolic knowledge graph has greater scale, diversity, and quality than human authoring. symbolic knowledge distillation offers an alternative to human-authored knowledge in commonsense research.

Acknowledgments

This work was funded in part by the Natural Sciences and Engineering Research Council of Canada (NSERC) (funding reference number 401233309), DARPA MCS program through NIWC Pacific (N66001-19-2-4031), and the Allen Institute for AI.

Ethical Considerations

One aspect of our work with the potential for ethical pitfalls is large-scale generation from pretrained language models, in constructing ATOMIC^{10x}. Recent work (Bender et al., 2021) has highlighted the risks of models trained on massive text resources, as GPT-3 (Brown et al., 2020) is, which we use for generation. Indeed, open generations from pretrained language models can often contain harmful, biased, or offensive aspects. We argue here that this risk is largely mitigated in our work, mainly due to the narrow and constrained nature of our generations. The goal of our work is characterising simple and generic anonymous situations, specifically in terms of commonsense causes and effects. We ensure generations are focused on these topics through careful prompting, which we found to be quite effective at keeping these generations ontopic. As such, the potential for harmful generation is very low; indeed, in a manual inspection of 100 generated examples, we found none that were significant harmful, besides one that contained adult content.

A related concern is the potential for large models and training sets to make automated oppression or exploitation possible, for instance in surveillance or generating fake news. As above, we argue that the generic, commonsense nature of our data and models makes this concern less relevant here. Our data does not contain any information directly related to these harmful domains (e.g. social media or fake news generation). While our data may assist machines in understanding basic situations, this is unlikely to be useful for harmful models given the simplicity of our data and still-flawed commonsense capabilities of even the most advanced models.

Finally, we note that we ensure fair and generous compensation for all human evaluators we hire through Amazon Mechanical Turk. Based on our estimates of time required per task, we ensure that the effective pay rate is at least \$15 per hour.

References

- Aishwarya Agrawal, Dhruv Batra, Devi Parikh, and Aniruddha Kembhavi. 2018. Don't just assume; look and answer: Overcoming priors for visual question answering. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 4971–4980.
- Prithviraj Ammanabrolu, Wesley Cheung, William Broniec, and Mark O. Riedl. 2021a. Automated storytelling via causal, commonsense plot ordering. In *AAAI*.
- Prithviraj Ammanabrolu, Jack Urbanek, Margaret Li, Arthur D. Szlam, Tim Rocktaschel, and Jason Weston. 2021b. How to motivate your dragon: Teaching goal-driven agents to speak and act in fantasy worlds. In *NAACL*.
- Ateret Anaby-Tavor, Boaz Carmeli, Esther Goldbraich, Amir Kantor, George Kour, Segev Shlomov, Naama Tepper, and Naama Zwerdling. 2020. Do not have enough data? deep learning to the rescue! *Proceedings of the AAAI Conference on Artificial Intelligence*, 34:7383–7390.
- Forough Arabshahi, Jennifer Lee, Antoine Bosselut, Yejin Choi, and Tom. Mitchell. 2021. Conversational multi-hop reasoning with neural commonsense knowledge and symbolic logic rules. In *EMNLP*.
- Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, page 610–623, New York, NY, USA. Association for Computing Machinery.
- Sumithra Bhakthavatsalam, Chloe Anastasiades, and Peter E. Clark. 2020. Genericskb: A knowledge base of generic statements. *ArXiv*, abs/2005.00660.

- Kurt D. Bollacker, Colin Evans, Praveen K. Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In *SIGMOD Conference*.
- Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, A. Çelikyilmaz, and Yejin Choi. 2019. Comet: Commonsense transformers for automatic knowledge graph construction. In *ACL*.
- Ronan Le Bras, Swabha Swayamdipta, Chandra Bhagavatula, Rowan Zellers, Matthew E. Peters, Ashish Sabharwal, and Yejin Choi. 2020. Adversarial filters of dataset biases. In *ICML*.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners.
- Christian Buck, Kenneth Heafield, and Bas Van Ooyen. 2014. N-gram counts and language models from the common crawl. In *LREC*, volume 2, page 4. Citeseer
- Tuhin Chakrabarty, Debanjan Ghosh, Smaranda Muresan, and Nanyun Peng. 2020. R³: Reverse, retrieve, and rank for sarcasm generation with commonsense knowledge. In *ACL*.
- Tuhin Chakrabarty, Xurui Zhang, Smaranda Muresan, and Nanyun Peng. 2021. MERMAID: Metaphor generation with symbolism and discriminative decoding. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4250–4261, Online. Association for Computational Linguistics.
- Ernest Davis and Gary Marcus. 2017. Causal generative models are just a start. *Behavioral and Brain Sciences*, 40.
- Joe Davison, Joshua Feldman, and Alexander M Rush. 2019. Commonsense knowledge mining from pretrained models. 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 1173–1178.
- Oren Etzioni, Anthony Fader, Janara Christensen, Stephen Soderland, et al. 2011. Open information extraction: The second generation. In *Twenty-Second International Joint Conference on Artificial Intelligence*.

- Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological bulletin*, 76(5):378.
- Jack Hessel and Lillian Lee. 2020. Does my multimodal model learn cross-modal interactions? it's harder to tell than you might think! In *EMNLP*.
- Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean. 2015. Distilling the knowledge in a neural network. In NIPS Deep Learning and Representation Learning Workshop.
- Geoffrey E Hinton. 2002. Training products of experts by minimizing contrastive divergence. *Neural computation*, 14(8):1771–1800.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net.
- 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. AAAI.
- William R. Kearns, Neha Kaura, Myra Divina, Cuong Viet Vo, Dong Si, Teresa M. Ward, and Weichao Yuwen. 2020. A wizard-of-oz interface and persona-based methodology for collecting health counseling dialog. *Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems*.
- Yoon Kim and Alexander M. Rush. 2016. Sequence-level knowledge distillation. In *EMNLP*.
- Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980.
- Varun Kumar, Ashutosh Choudhary, and Eunah Cho. 2020. Data augmentation using pre-trained transformer models. In *Proceedings of the 2nd Workshop on Life-long Learning for Spoken Language Systems*, pages 18–26, Suzhou, China. Association for Computational Linguistics.
- J Richard Landis and Gary G Koch. 1977. The measurement of observer agreement for categorical data. *biometrics*, pages 159–174.
- Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, D. Kontokostas, Pablo N. Mendes, Sebastian Hellmann, M. Morsey, Patrick van Kleef, S. Auer, and C. Bizer. 2015. Dbpedia - a large-scale, multilingual knowledge base extracted from wikipedia. *Semantic Web*, 6:167–195.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. In *EMNLP*.

- Zhongyang Li, Xiao Ding, Ting Liu, J. Edward Hu, and Benjamin Van Durme. 2020. Guided generation of cause and effect. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- William Merrill, Yoav Goldberg, Roy Schwartz, and Noah A. Smith. 2021. Provable limitations of acquiring meaning from ungrounded form: What will future language models understand? *Transactions of the Association for Computational Linguistics*, 9:1047–1060.
- Tom Michael Mitchell, William W. Cohen, Estevam R. Hruschka, Partha P. Talukdar, Bo Yang, Justin Betteridge, Andrew Carlson, Bhavana Dalvi, Matt Gardner, Bryan Kisiel, Jayant Krishnamurthy, N. Lao, Kathryn Mazaitis, Thahir Mohamed, Ndapandula Nakashole, Emmanouil Antonios Platanios, Alan Ritter, Mehdi Samadi, Burr Settles, Richard C. Wang, D. Wijaya, Abhinav Gupta, Xinlei Chen, Abulhair Saparov, Malcolm Greaves, and Joel Welling. 2015. Never-ending learning. *Communications of the ACM*, 61:103 115.
- Yannis Papanikolaou and A. Pierleoni. 2020. Dare: Data augmented relation extraction with gpt-2. *ArXiv*, abs/2004.13845.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? 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.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *ArXiv*, abs/1910.01108.
- 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 ifthen reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3027–3035.
- Roy Schwartz, Maarten Sap, Ioannis Konstas, Leila Zilles, Yejin Choi, and Noah A. Smith. 2017. The effect of different writing tasks on linguistic style: A

- case study of the ROC story cloze task. In *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)*, pages 15–25, Vancouver, Canada. Association for Computational Linguistics.
- Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31.
- Alon Talmor, Ori Yoran, Ronan Le Bras, Chandra Bhagavatula, Yoav Goldberg, Yejin Choi, and Jonathan Berant. 2021. Commonsenseqa 2.0: Exposing the limits of ai through gamification.
- Masatoshi Tsuchiya. 2018. Performance impact caused by hidden bias of training data for recognizing textual entailment. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. https://github.com/ kingoflolz/mesh-transformer-jax.
- Sean Welleck, Ilia Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston. 2020. Neural text generation with unlikelihood training. In International Conference on Learning Representations
- Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122. Association for Computational Linguistics.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, R'emi Louf, Morgan Funtowicz, and Jamie Brew. 2019. Huggingface's transformers: State-of-the-art natural language processing. *ArXiv*, abs/1910.03771.
- Wenhan Xiong, Jingfei Du, William Yang Wang, and Veselin Stoyanov. 2020. Pretrained encyclopedia: Weakly supervised knowledge-pretrained language model. In *International Conference on Learning Representations*.
- Yiben Yang, Chaitanya Malaviya, Jared Fernandez, Swabha Swayamdipta, Ronan Le Bras, Ji-Ping Wang, Chandra Bhagavatula, Yejin Choi, and Doug Downey. 2020. Generative data augmentation for commonsense reasoning. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1008–1025, Online. Association for Computational Linguistics.

- Hongming Zhang, Daniel Khashabi, Y. Song, and D. Roth. 2020a. TransOMCS: From linguistic graphs to commonsense knowledge. In *IJCAI*.
- Hongming Zhang, Xin Liu, Haojie Pan, Y. Song, and C. Leung. 2020b. Aser: A large-scale eventuality knowledge graph. *Proceedings of The Web Conference* 2020.
- Ben Zhou, Qiang Ning, Daniel Khashabi, and Dan Roth. 2020. Temporal common sense acquisition with minimal supervision. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7579–7589, Online. Association for Computational Linguistics.
- Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. 2018. Texygen: A benchmarking platform for text generation models. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, pages 1097–1100.

A Human Evaluation Details

We conduct human evaluations on Amazon Mechanical Turk using the template of Figures 4,5. Workers are presented with ATOMIC-style triples, replacing relations with natural language templates (e.g. *HinderedBy* becomes "can be hindered by"). 3 annotators rate each triple, with options for acceptability: "always/often", "sometimes/likely", "farfetched/never", "invalid", or "too unfamiliar to judge". The first two are considered "accepted", the second two "rejected" and the final is "no judgement". For reporting acceptance rates, and training a critic model, we only distinguish between "accepted" and not "accepted".

Workers are compensated \$0.17 per task (i.e. completing all questions in the evaluation template Figures 4,5). We estimate an upper bound of 30s to complete a single task, which gives an hourly rate of \$20.4. Workers are selected based on an Amazon Mechanical Turk qualification, specifically filtering for workers with high accuracy on past knowledge base triple evaluations. We follow the same setup for all evaluations, besides number of annotators. This setup is shown to result in consistent and reliable annotations, with an inter-annotator agreement given by Fleiss' kappa (Fleiss, 1971) of 40.8 when evaluating with 3 annotators, in §3.4.

B Using Alternate Models as Knowledge Sources

One natural question that arises from the strong performance of symbolic knowledge distillation is whether other sources of knowledge (i.e. language models) would similarly benefit from this method. In this section, we particularly measure the capacity of other language models to serve as the "loose teacher" which generated the base knowledge of the resulting corpus.

We expand our study beyond GPT-3 here (the model used in our work), to include 2 contemporary large language models, GPT-J (Wang and Komatsuzaki, 2021) and T5-11B (Lester et al., 2021) finetuned for language modelling. For knowledge generation (verbalization) we follow the same procedure as $\S 3$ along with simple adjustments to improve quality. We are investigating the effect of the critic on knowledge precision here, so we also include ATOMIC $^{20}_{20}$ to probe the usefulness of automatic filtering for human-authored knowledge.

For each knowledge source, we follow the human evaluation setup in §3.4 to obtain quality an-

notations of 2000 examples, with 1 annotation per example. This follows a similar setup to §4-indeed, we are replicating the earlier critic experiments but at a smaller scale (2000 annotations vs. 10000) to allow for more knowledge sources. For each knowledge source, we randomly split into sizes of 1400/300/300 for train, dev, and test sets. We follow §4 to train a critic model for each knowledge source.

We plot different thresholds (% of corpus filtered) against the resulting precision (proportion of corpus that is judged to be "valid" knowledge) in Figure 3, and give numbers at various sizes in Table 7. One striking aspect is that a critic model can raise the precision of any of these knowledge sources to approximately 90% while retaining 30% of the original corpus size. While this discards a significant portion of the original generated knowledge, it raises the exciting prospect of using more cost-effective models at a large scale to generate strong commonsense corpora like ATOMIC^{10x}. GPT-J and T5-11B can both be run locally by researchers, unlike GPT-3 which uses a pay-pergeneration API. Thus, one can imagine producing a large and high-quality corpus like ATOMIC^{10x} at a lower cost by instead generating a larger volume of knowledge from such an accessible model, and simply filtering to a greater extent.

Another interesting aspect is how the various knowledge sources diverge. Under little to no critical filtering (i.e. corpus size = 1.0), the precision of various knowledge sources is widely spread. Before applying a critic, quality of knowledge source is very important. Indeed, precision is ordered by cost of generation: human $ATOMIC_{20}^{20}$ has the highest precision while being the most expensive, followed by GPT-3 (used here) which is pay-pergeneration, and finally the two publicly available models. Another point of divergence is for extreme filtering (at approximately 20% of the original corpus size. All knowledge sources but GPT-3 plateau at approximately 90% accuracy, while GPT-3 rises towards 100%. Indeed, this supports our use of GPT-3 in this work, as a high-quality automatic knowledge source.

C Critic Model

We train binary classifiers (critics) for human acceptability using RoBERTa-Large (Liu et al., 2019), fine-tuning all parameters, along with a 2-layer MLP on the [CLF] representation. We conduct

]	Precisi	on at C	orpus	Size (%	·)		
Knowledge Source	100	90	80	70	60	50	40	30	20	10
$\overline{ATOMIC_{20}^{20}}$	84.0	86.3	87.9	89.0	88.3	88.7	91.7	90.0	90.0	90.0
GPT-J	71.7	76.7	81.7	83.8	86.7	88.0	88.3	87.8	93.3	90.0
T5-11B	64.7	66.7	70.8	74.8	79.4	84.7	89.2	92.2	91.7	93.3
GPT-3 curie	79.3	81.5	85.0	86.2	88.3	90.7	91.7	90.0	98.3	100.0

Table 7: Knowledge precision at various corpus sizes (from 100% to 10%) based on filtering by the critic model. Precision is calculated by human annotation of valid or invalid knowledge. We consider 4 knowledge sources, as described in Appendix B. This corresponds to the data plotted in Figure 3.

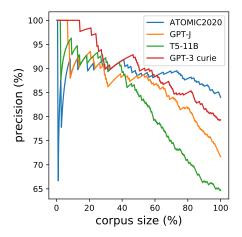


Figure 3: Precision resulting from the critic step from §4, with various thresholds. We include corpora generated by GPT-3 (ATOMIC^{10x}), GPT-J, T5-11B, and humans (ATOMIC²⁰). Without filtering (corpus size = 1.0), different corpora have a variety of precisions. As more examples are filtered by the critic, precision rises significantly demonstrating the strong value of the critic step.

a small grid search on the validation set finding batch size 128, dropout .1, and Adam (Kingma and Ba, 2015) learning rate 5e-6 to be effective. We use early stopping and decay learning rate on validation performance plateauing, to maximize R@80% on the validation set. We find RoBERTa pretrained on MNLI (Williams et al., 2018) effective, outperforming other options. As well, we substitute randomly-sampled names in for person designations "X"/"Y". We include as a baseline an unsupervised filtration metric inspired by (Davison et al., 2019): they propose a model estimate of PMI to score mined commonsense triples. In our case, we use Negative Log-Likelihood (NLL) and token-mean-NLL from GPT-3 itself.

The validation precision/recall of our best performing model, the baselines, and the in-optimal hyperparameter configurations are given in Figure 6. Once fixing our model, we applied it to the test set (also in Fig 6), verifying that it generalizes to ATOMIC^{10x} entries. Overall, our trained critic model is more effective than the baselines in identifying high and low quality teacher generations at all levels of precision and recall. This result demonstrates that a small amount of human supervision can consistently help to correct GPT-3's mistakes.

D ATOMIC^{10x} Generation Prompts

We include example prompts for all generations we do, from Table 8 to 15. Note that elements of generation prompts are randomized for each batch. For event generation, the few-shot examples and order are randomly sampled from a seed set of 100 high-quality examples from ATOMIC $_{20}^{20}$ in each batch. For inference generation, the natural names used for PersonX and PersonY are randomly sampled from a small predefined set of names.

Instructions (click to expand/collapse)

(WARNING: This HIT may contain adult content. Worker discretion is advised.)

Thanks for participating in this HIT!

If the data is good, it's good. If bad, then bad. Please annotate as you see not worrying about how many of each label you find yourself assigning! If you understand the words but the Phrases or the complete assertation makes poor sense, please mark as INVALID. Thank you!

You will evaluate how often assertions are true. Each assertion is comprised of 3 parts: Phrase A, Relation,

Phrase A, Phrase B Short phrases. May describe objects, object properties, events, actions, etc.

How A relates to B.

For each assertion, determine how true it is:

always/often Always or quite often true.

sometimes/likely Sometimes is true or true for some people. -or- Likely true. False or farfetched, at best. -or- Unlikely to be true. farfetched/never This assertion makes no sense (i.e., "what does this even mean?!"). too unfamiliar to judge Cannot make a fair evaluation. Unfamiliar with one or both of the phrase.

If you see "nothing in particular" for *Phrase B*, assess Phrase B in context:

- Sometimes certain actions can simply be responded to by doing nothing!
- Other times, doing nothing in particular is simply a weird or unlikely reaction to something.
- $\bullet \ \ \text{See examples under tricky relations tagged with} \\ \boxed{ \textbf{nothing in particular example} }$

Please report any prejudiced or inappropriate language:

- Profane or offensive content (NSFW, R-rated material etc)
- Prejudiced assumptions or derogatory language that <u>villainizes</u> people. HOWEVER, please note, not all negative content is derogatory especially if Phrase B is intrinsically what Phrase A means. For example:

- criminals are characterized by committing crime is OK.

 → This isn't necessarily villianizing people since "criminal" means "a person who has commited a crime".

 homeless are characterized by being lazy is prejudiced.

 → There are many reason a person is rendered homeless. This is a gratuitous prejudice about homelessness.
- Material that people may find disturbing, off-putting, or improper

A couple NOTES:

- Please be **forgiving** of *spelling or grammatical errors*
- If the terms are too obscure or you don't know the truth of the fact at the top of your head, it is okay to mark is "too unfamiliar to judge". If you can answer (e.g., based on likelihood), please provide a response

Tricky Relations (click to expand/collpase)

Examples (click to expand/collapse)

Figure 4: Page 1 of template used for human evaluation.

☐ This fact is true☐ I would count to		farfetched/never	invalid	too unfamiliar to judge	
🗆 I would count t	but outdated				
	his as an inappropria	ate, prejudiced or offe	nsive mat	erial	
PersonX aske	d PersonY out or	a date, can be h	indered	by, PersonX is still dat	ing Sarah
How often doe	es the assertion	hold true?			
always/often	sometimes/likely	farfetched/never	invalid	too unfamiliar to judge	
This fact is true	but outdated				
I would count to	his as an inappropria	ate, prejudiced or offe	nsive mat	erial	
PersonX fails	to go home as a	result, PersonX,	is groun	ded	
	es the assertion		8		
	sometimes/likely	farfetched/never	invalid	too unfamiliar to judge	
☐ This fact is true				,	
		es, as a result, Pe			
PersonX make	es her own clothers	es, as a result, Pe	rsonX fe	eels, <i>artistic</i>	
PersonX make How often doe always/often	es her own clothe es the assertion sometimes/likely	es, as a result, Pe			
PersonX make How often doe always/often This fact is true	es her own clothes the assertion sometimes/likely	es, as a result, Pe hold true? farfetched/never	rsonX fe	teels, <i>artistic</i> too unfamiliar to judge	
PersonX make How often doe always/often This fact is true	es her own clothes the assertion sometimes/likely	es, as a result, Pe	rsonX fe	teels, <i>artistic</i> too unfamiliar to judge	
PersonX make How often doe always/often This fact is true	es her own clother es the assertion sometimes/likely but outdated this as an inappropria	es, as a result, Pe hold true? farfetched/never	rsonX fe invalid	teels, <i>artistic</i> too unfamiliar to judge	d by the music
PersonX make How often doe always/often This fact is true I would count to	es her own clother es the assertion sometimes/likely but outdated this as an inappropria	es, as a result, Pe hold true? farfetched/never ate, prejudiced or offe	rsonX fe invalid	tels, <i>artistic</i> too unfamiliar to judge	d by the music
PersonX make How often doe always/often This fact is true I would count th PersonX notic How often doe	es her own clothers the assertion as sometimes/likely but outdated this as an inappropriates Persony's res	es, as a result, Pe hold true? farfetched/never ate, prejudiced or offe	rsonX fe invalid	tels, <i>artistic</i> too unfamiliar to judge	d by the music
PersonX make How often doe always/often This fact is true I would count th PersonX notic How often doe	es her own clothes the assertion as sometimes/likely but outdated this as an inappropriates Persony's respect the assertion as sometimes/likely	es, as a result, Pe hold true? farfetched/never ate, prejudiced or offe ponse, can be hir hold true?	invalid ensive mat	teels, <i>artistic</i> too unfamiliar to judge erial by, <i>PersonX is distracte</i>	d by the music

Figure 5: Page 2 of template used for human evaluation.

You must ACCEPT the HIT before you can submit the results.

1. Event: PersonX unwraps PersonY's hands

2. Event: PersonX overcomes evil with good

3. Event: PersonX is fed up with the present situation

4. Event: PersonX breaks PersonX's back

5. Event: PersonX calls no one

6. Event: PersonX never gets angry

7. Event: PersonX does not learn from PersonY

8. Event: PersonX refuses to touch PersonY's hands

9. Event: PersonX looks at flowers

10. Event: PersonX unloads an atomic bomb

11. Event:

Table 8: Prompt for head generation.

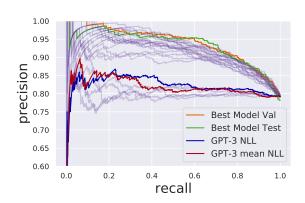


Figure 6: Precision vs. recall of our critic model on the human labelled validation set. The best trained models are labelled, and other hyper-parameter settings are shown as faded lines. We also include generation negative log-likelihood (nll) and token-wise mean nll as cutoff measures—these perform much worse than the supervised model.

Next, how are people seen in each situation? Examples:

Situation 1: Devin bullies Jean.

Devin is seen as dominant.

Situation 2: Jamie moves to another city.

Jamie is seen as adventurous.

Situation 3: Sydney changes Ryan's mind.

Sydney is seen as influential.

Situation 4: Lindsay writes a story.

Lindsay is seen as creative.

Situation 5: Rowan covers Pat's expenses.

Rowan is seen as wealthy.

Situation 6: Lee takes time off.

Lee is seen as carefree.

Situation 7: Riley advises Noel.

Riley is seen as informed.

Situation 8: Adrian bursts into tears.

Adrian is seen as depressed.

Situation 9: Hunter deals with problems.

Hunter is seen as responsible.

Situation 10: Sam follows Charlie.

Sam is seen as suspicious.

Situation 11: Alex makes Chris wait.

Alex is seen as

Table 9: Prompt for generating xAttr.

Next, what do situations make people do? Examples:

Situation 1: Devin gets a divorce.

As a result, Devin dates someone new.

Situation 2: Jamie lifts weights.

As a result, Jamie has sore muscles.

Situation 3: Sydney takes Ryan to a bar.

As a result, Sydney gets drunk.

Situation 4: Lindsay decides to hire a tutor.

As a result, Lindsay gets better grades.

Situation 5: Rowan buys Pat drinks.

As a result, Rowan is thanked by Pat.

Situation 6: Lee hears bad news.

As a result, Lee begins to cry.

Situation 7: Riley buys a chocolate bar.

As a result, Riley gets change.

Situation 8: Adrian does a lot of work.

As a result, Adrian gets mental fatigue.

Situation 9: Hunter attends a concert.

As a result, Hunter hears a new song.

Situation 10: Sam gets the job done.

As a result, Sam gets more responsibilities.

Situation 11: Alex makes Chris wait.

As a result, Alex

Table 10: Prompt for generating xEffect.

For each situation, describe the intent. Examples:

Situation 1: Devin gets the newspaper.

Devin intends to read the newspaper.

Situation 2: Jamie works all night.

Jamie intends to meet a deadline.

Situation 3: Sydney destroys Ryan.

Sydney intends to punish Ryan.

Situation 4: Lindsay clears her mind.

Lindsay intends to be ready for a new task.

Situation 5: Rowan wants to start a business.

Rowan intends to be self sufficient.

Situation 6: Lee ensures Ali's safety.

Lee intends to be helpful.

Situation 7: Riley buys lottery tickets.

Riley intends to become rich.

Situation 8: Alex makes Chris wait.

Alex intends

Table 11: Prompt for generating xIntent.

Next, we will discuss what people need for certain situations. Examples:

- 1. Before Devin makes many new friends, Devin has to spend time with people.
- 2. Before Jamie gets a date, Jamie has to ask someone out.
- 3. Before Sydney changes Ryan's mind, Sydney has to think of an argument.
- 4. Before Lindsay gets a job offer, Lindsay has to apply.
- 5. Before Rowan takes a quick nap, Rowan has to lie down.
- 6. Before Lee tries to kiss Ali, Lee has to approach Ali.
- 7. Before Riley rides Noel's skateboard, Riley has to borrow it.
- 8. Before Adrian eats the food, Adrian has to prepare a meal.
- 9. Before Hunter watches Netflix, Hunter has to turn on the TV.
- 10. Before Sam has a baby shower, Sam has to invite some friends.
- 11. Before Alex makes Chris wait, Alex has

Table 12: Prompt for generating xNeed.

Next, how do people feel in each situation? Examples:

Situation 1: Devin lives with Jean's family.

Devin feels loved.

Situation 2: Jamie expects to win.

Jamie feels excited.

Situation 3: Sydney comes home late.

Sydney feels tired.

Situation 4: Lindsay sees dolphins.

Lindsay feels joyful.

Situation 5: Rowan causes Pat anxiety.

Rowan feels guilty.

Situation 6: Lee goes broke.

Lee feels embarrassed.

Situation 7: Riley has a drink.

Riley feels refreshed.

Situation 8: Adrian has a heart condition.

Adrian feels scared about their health.

Situation 9: Hunter shaves Avery's hair.

Hunter feels helpful.

Situation 10: Sam loses all of Charlie's money.

Sam feels horrible.

Situation 11: Alex makes Chris wait.

Alex feels

Table 13: Prompt for generating xReact.

Next, what do people want in each situation? Examples:

Situation 1: Devin mows the lawn.

Devin wants to take a shower.

Situation 2: Jamie is going to a party.

Jamie wants to take an Uber home.

Situation 3: Sydney bleeds a lot.

Sydney wants to go to the ER.

Situation 4: Lindsay works as a cashier.

Lindsay wants to find a better job.

Situation 5: Rowan gets dirty.

Rowan wants to do a load of laundry.

Situation 6: Lee stays up all night studying.

Lee wants to rest.

Situation 7: Riley gets Noel's autograph.

Riley wants to tell some friends.

Situation 8: Adrian sees Taylor's point.

Adrian wants to agree with Taylor.

Situation 9: Hunter leaves Avery's bike.

Hunter wants to keep the bike safe.

Situation 10: Sam wants a tattoo.

Sam wants to find a tattoo design.

Situation 11: Alex makes Chris wait.

Alex wants

Table 14: Prompt for generating xWant.

Next, what can hinder each situation? Examples:

Situation 1: Devin makes a doctor's appointment,

This is hindered if Devin can't find the phone to call the doctor.

Situation 2: Jamie rubs Wyatt's forehead,

This is hindered if Jamie is afraid to touch Wyatt.

Situation 3: Sydney eats peanut butter,

This is hindered if Sydney is allergic to peanuts.

Situation 4: Lindsay looks perfect,

This is hindered if Lindsay can't find any makeup.

Situation 5: Rowan goes on a run,

This is hindered if Rowan injures her knees.

Situation 6: Lee takes Ali to the emergency room,

This is hindered if Ali has no health insurance to pay for medical care.

Situation 7: Riley spends time with Noel's family,

This is hindered if Noel's family doesn't like spending time with Riley.

Situation 8: Adrian moves from place to place,

This is hindered if Adrian can't afford to move.

Situation 9: Hunter protests the government,

This is hindered if Hunter is arrested.

Situation 10: Sam has a huge fight,

This is hindered if Sam does not like confrontation.

Situation 11: Alex makes Chris wait,

This is hindered if

Table 15: Prompt for generating HinderedBy.