# **Critical Thinking for Language Models**

Gregor Betz<sup>†</sup> and Christian Voigt<sup>†</sup> and Kyle Richardson<sup>‡</sup>

<sup>†</sup> Karlsruhe Institute of Technology, Karlsruhe, Germany

{gregor.betz, christian.voigt}@kit.edu

<sup>‡</sup> Allen Institute for AI, Seattle, WA, USA

{kyler}@allenai.org

## **Abstract**

This paper takes a first step towards a critical thinking curriculum for neural autoregressive language models. We introduce a synthetic corpus of deductively valid arguments, and generate artificial argumentative texts to train and evaluate GPT-2. Significant transfer learning effects can be observed: Training a model on three simple core schemes allows it to accurately complete conclusions of different, and more complex types of arguments, too. language models generalize the core argument schemes in a correct way. over, we obtain consistent and promising results for NLU benchmarks. In particular, pre-training on the argument schemes raises zero-shot accuracy on the GLUE diagnostics by up to 15 percentage points. The findings suggest that intermediary pretraining on texts that exemplify basic reasoning abilities (such as typically covered in critical thinking textbooks) might help language models to acquire a broad range of reasoning skills. The synthetic argumentative texts presented in this paper are a promising starting point for building such a "critical thinking curriculum for language models."

#### 1 Introduction

Pre-trained autoregressive language models (LM) such as GPT-2 and GPT-3 achieve, remarkably, competitive results in a variety of language modeling benchmarks without task-specific fine-tuning (Radford et al., 2019; Brown et al., 2020). Yet, it is also widely acknowledged that these models struggle with reasoning tasks, such as natural language inference (NLI) or textual entailment (Askell, 2020). Actually, that doesn't come as a surprise, given the tendency of humans to commit errors in reasoning (Kahneman, 2011; Sunstein and Hastie, 2015), their limited critical think-

ing skills (Paglieri, 2017), the resulting omnipresence of fallacies and biases in texts and the frequently low argumentative quality of online debates (Hansson, 2004; Guiaşu and Tindale, 2018; Cheng et al., 2017). Neural language models are known to pick up and reproduce normative biases (e.g., regarding gender or race) present in the dataset they are trained on (Gilburt, 2019), as well as other annotation artifacts (Gururangan et al., 2018); no wonder this happens with argumentative biases and reasoning flaws, too (Kassner and Schütze, 2020; Talmor et al., 2020). This diagnosis suggests that there is an obvious remedy for LMs' poor reasoning capability: make sure that the training corpus contains a sufficient amount of exemplary episodes of sound reasoning.

In this paper, we take a first step towards the creation of a "critical thinking curriculum" for neural language models. Critical thinking can be loosely defined as "reasonable reflective thinking that is focused on deciding what to believe or do." (Norris and Ennis, 1989) Generally speaking, our study exploits an analogy between teaching critical thinking to students and training language models so as to improve their reasoning skill. More specifically, we build on three key assumptions that are typically made in critical thinking courses and textbooks: First, there exist fundamental reasoning skills that are required for, or highly conducive to, a large variety of more specific and advanced critical thinking skills (e.g., Fisher, 2001, p. 7). Second, drawing deductive inferences is one such basic ability (e.g., Fisher, 2001, pp. 7–8). Third, reasoning skills are not (just) acquired by learning a theory of correct reasoning, but by studying lots of examples and doing "lots of good-quality exercises" (Lau and Chan, 2020), typically moving from simple to more difficult problems (e.g., Bowell and Kemp, 2014).

These insights from teaching critical thinking translate, with respect to our study, as follows.

First of all, we design and build 'lots of good-quality exercises': a synthetic corpus of deductively valid arguments which instantiate a variety of (syllogistic) argument schemes, and which are rendered as text paragraphs (Section 3). Next, we use our synthetic argument text corpus to train and to evaluate GPT-2 (Section 4). The training, which maximizes a causal language modeling objective, can be conceived of as a generic, intermediary pre-training in the spirit of STILTS (Phang et al., 2018).

Evaluating the models' ability to correctly complete conclusions of arguments, we observe strong transfer learning effects/generalization (Section 5): Just training the models on a few central core schemes (generalized modus ponens, contraposition and chain rule) allows them to accurately complete conclusions of different types of arguments, too (e.g., complex argumentative forms that involve dilemma and de Morgan). The language models appear to connect and generalize the core argument schemes in a correct way. In addition, the models are equally able to apply learned argument patterns beyond the training corpus' domain. Tests with a simple manually authored argument produce evidence that generic language modeling skill facilitates the successful generalization of learned argument patterns.

Moreover, we test the trained models on different reasoning benchmarks. Because we are particularly interested in transfer learning effects, we do so in a zero-shot set-up (i.e., evaluating our argumentation models on entirely unrelated NLU tasks, which follows recent work by Mitra et al. (2019); Shwartz et al. (2020); Ma et al. (2020)). We obtain consistent and promising results for the GLUE diagnostics (Wang et al., 2018) and SNLI (Bowman et al., 2015) benchmarks (Section 5), finding that training on core schemes clearly improves NLU skill. However, training on the argument corpus doesn't affect the performance with regard to the semantically more demanding Argument Reasoning Comprehension task (Habernal et al., 2018) or the critical thinking assessment compiled in LogiQA (Liu et al., 2020).

All these transfer learning effects observed strengthen the analogy between teaching critical thinking and training language models: A variety of reasoning skills are improved by generic, intermediary pre-training on high-quality texts that exemplify a basic reasoning skill, namely simple deductive argumentation. Obviously, drawing correct inferences is just one of the elementary skills typically covered in critical thinking courses (Fisher, 2001). Critical thinking involves more than deduction. And it would hence, by analogy, be unreasonable to expect that intermediary pretraining on the synthetic argument corpus suffices to turn language models into accomplished reasoners. However, we have shown that argumentative texts (with valid syllogistic arguments) are certainly a good starting point when building a more comprehensive dataset for initial or intermediary pre-training that might help language models to acquire a broad range of reasoning skills. Or, to put it differently, the synthetic argumentative texts might belong to the core of a "critical thinking curriculum for language models." In the final section, we advance some ideas for complementing the artificial argument corpus so as to further improve the performance of LMs with regard to different reasoning benchmarks.

#### 2 Related Work

To our knowledge, this paper is, together with Gontier et al. (2020), among the first to show that autoregressive language models like GPT-2 can learn to reason by training on a text corpus of correct natural language arguments. By contrast, previous work in this field, described below, has typically modeled natural language reasoning problems as classification tasks and trained neural systems to accomplish them. For example, Schick and Schütze (2020a,b), using pattern verbalizations, construct structured training data that is suitable for training a masked language model with classification head, and thusly achieve remarkable NLU performance. This paper explores the opposite route: We start with highly structured (synthetic) data, render it as unstructured, plain text and train a uni-directional language model on the synthetic text corpus.

Over and above the methodological novelty of our approach, we discuss, in the following, related reasoning benchmarks and explain what sets our synthetic argument corpus apart from this work.

Rule reasoning in natural language Various datasets have been developed for (deductive) rule reasoning in natural language. In these tasks, one or multiple rules, i.e. (generalized) conditionals, must be applied to a fact base in order to deductively infer a conclusion. Facts and conclusions

are represented by atomic statements. Rule application closely resembles the conclusion completion task for *generalized modus ponens* and *generalized modus tollens* schemes described below. However, we go beyond previous work in investigating the ability of language models to infer conclusions that have a more complex logicosemantic structure (e.g., existential or universal statements).

The question answering bAbI dataset (Weston et al., 2016) contains a task which involves applying very specific rules of the form "Xs are afraid of Ys" to an instance (for example: "Mice are afraid of cats. Jerry is a mouse. What is Jerry afraid of? *A:cats*"). Equally simple, one-step rule applications are tested in Richardson et al. (2020), and also contained in the QuaRTz dataset (Tafjord et al., 2019).

ROPES (Lin et al., 2019) is a reading comprehension task that involves applying background knowledge to a given situation (both being presented as paragraph long text). Correct answers can be inferred by one-step rule application; part of the challenge is to identify the relevant rule and fact in the text.

RuleTaker, arguably the most general system for natural rule reasoning in natural language so far, is a transformer model that has been fine-tuned to predict whether a conclusion can be inferred from a set of rules and facts, not all of which are necessarily required to draw the conclusion (Clark et al., 2020). Moreover, inferring the conclusion from the premise set might involve multiple inference steps. The authors show that the transformer model can be trained to perform this task nearly flawlessly and, moreover, to 'explain' its inferences by identifying relevant premises. They also observe substantial transfer learning effects.

PRover extends RuleTaker by a component for proof generation (Saha et al., 2020). Technically, the QA head of the RoBERTa language model (Liu et al., 2019) is complemented by two additional neural classifiers (for nodes and edges) that are used to to construct proof chains. Saha et al. (2020) show that PRover can construct valid proofs and outperforms RuleTaker in terms answer accuracy in a zero-shot setting.

Training on synthetic knowledge-graph data (such as "Paris CapitalOf France" and "France HasCapital Paris") *from scratch*, Kassner et al. (2020) find that BERT is able to correctly infer

novel facts. This confirms that language models can, in principle, learn basic conceptual rules, which, e.g., express that a relation is symmetric or that two terms are equivalent.

Benchmarks for enthymematic reasoning An 'enthymeme' is an argument whose premises are not explicitly stated, e.g.: "Jerry is a mouse. Therefore, Jerry is afraid of cats." The three tasks described below involve such reasoning with implicit assumptions, whereas our synthetic argument corpus doesn't: all premises are transparent and explicitly given.

Commensense Transformers (COMET) are autoregressive language models for generating commonsense knowledge graphs (Bosselut et al., 2019). Being trained on seed data, the models are able to meaningfully relate subject phrases to object phrases in terms of multiple binary relations (by doing the type of completion tasks we introduce in Section 4), and can thereby both reproduce and extend a given knowledge graph. In particular, this includes generating statements about causal relationships, which can be construed as enthymematic reasoning with commonsense background assumptions. For example, given the input "PersonX is re-elected. As a result, PersonX wants" the model generates as completions: "to get a raise", "to go to office", "to go home", "to make a speech", "to celebrate" - all of which are plausible fill-ins. The implicit commonsense premises that underlie this (entyhmematic) inference are principles such as "If someone has been re-elected, then they want to celebrate."

The Argument Reasoning Comprehension (ARC) dataset (Habernal et al., 2018) comprises simple informal arguments. Each argument contains two premises: whereas the first premise is explicitly stated, there are two alternative formulations of the second premise. The task consists in identifying which of these two alternative formulations is actually assumed in the argument. For example: "Miss America gives honors and education scholarships. And since [scholarships would give women a chance to study | scholarships would take women from the home], Miss America is good for women." ARC therefore assesses the ability to make implicit premises explicit. An adversarial ARC dataset that eliminates clues in the original benchmark is also available in Niven and Kao (2019).

CLUTRR is a task generator for relational rea-

soning on kinship graphs (Sinha et al., 2019). CLUTTR takes a set of (conceptual) rules about family relations as given and constructs settheoretic possible worlds (represented as graphs) which instantiate these rules. In such a possible (kinship) world, a target fact and a set of base facts are identified such that the base facts together with the rules deductively entail the target fact. The task consists in inferring the target fact from the base facts alone - the conceptual rules remain implicit. For example: "Kristin and her son Justin went to visit her mother Carol on a nice Sunday afternoon. They went out for a movie together and had a good time. Q: How is Carol related to Justin? A: Carol is the grandmother of Justin." So, CLUTRR assesses entyhmematic deductive reasoning with implicit conceptual rules. Gontier et al. (2020) have trained a generative Transformer language model on a synthetic text corpus (with each argumentative text containing a story, a proof chain and a conclusion from CLUTTR) and show that the language model does not only learn to draw the correct conclusion (given an argument with implicit commonsense premises), but also seems to acquire the ability to generate valid proof chains.

Critical thinking tasks LogiQA (Liu et al., 2020) is a collection of publicly available critical thinking questions, used by the National Civil Servants Examination of China to assess candidates' critical thinking and problem solving skills. LogiQA covers tasks of various types: different kinds of natural language inference problems as well as the identification of implicit premises or (practical) instrumental reasoning. Its scope is much broader than our highly specific and carefully designed argument corpus. The LogiQA tasks are shown to be hard for current AI systems, of which a fine-tuned transformer model performs best with an accuracy score of 35% – 50 percentage points below human performance.

## 3 An Artificial Argument Corpus

This section describes the construction of a synthetic corpus of natural language arguments used for training and evaluating GPT-2.<sup>1</sup>

The corpus is built around eight simple, deductively valid syllogistic argument schemes (top row

in Figure 1). These *base schemes* have been chosen because of their logical simplicity as well as their relevance in critical thinking and argument analysis (Feldman, 2014; Bowell and Kemp, 2014; Brun and Betz, 2016). Each of these eight base schemes is manually varied in specific ways to create further valid variants.

Negation variants of base schemes (second row in Figure 1) are created by substituting a subformula with its negation and/or by applying duplex negatio affirmat.

Complex predicates variants (third row in Figure 1) build on base schemes or their respective negation variants and are obtained by substituting atomic predicates with compound disjunctive or conjunctive ones.

*De Morgan* variants of base schemes (fourth row in Figure 1) are finally derived by applying de Morgan's law to the respective variants created before.

With 2-3 different versions for each of these variations of a base scheme (parameter "n" in Figure 1), we obtain, all in all, 71 distinct hand-crafted argument schemes. Obviously, some of these schemes can be derived from others. For example, generalized modus ponens and generalized contraposition (base schemes) entail a *negation variant* of generalized modus tollens. Likewise, generalized contraposition and hypothetical syllogism 1 entail a *(negation variant of)* hypothetical syllogism 2.

In view of their simplicity and prominence in natural language argumentation, three of the eight *base schemes* are marked as *core schemes*: generalized modus ponens, generalized contraposition, hypothetical syllogism 1.

Natural language instances of the argument schemes can be created by means of a first-orderlogic domain (with names and predicates) and natural language templates for the formal schemes. In order to obtain a large variety of realistic natural language arguments, we have devised

- a multi-stage templating process with
- · alternative templates at each stage and
- multiple domains.

As shown in Figure 2, this process can be split into five consecutive steps.

In *step 1*, the argument scheme, which serves as formal template for the natural language argument, is chosen.

<sup>&</sup>lt;sup>1</sup>The corpus as well as the source code used to generate it will be released at https://github.com/debatelab/aacorpus.

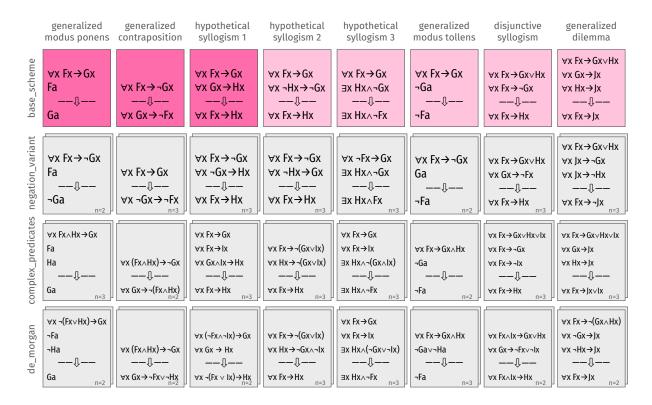


Figure 1: Syllogistic argument schemes used to create an artificial argument corpus.

In step 2, each sentence in the formal scheme (premises and conclusion) is individually replaced by a natural language pattern in accordance with a randomly chosen template. For example, the formula " $\forall xFx \rightarrow Gx$ " might be replaced by any of the following natural language sentence schemes:

- "Every F is a G."
- "Whoever is a F is also a G."
- "Being a G is necessary for being a F."
- "If someone is a F, then they are a G."\*

Some of these patterns are not used for training, but are reserved for generating an out-of-domain test dataset (e.g., the template marked with an asterisk in the above list).

In *step 3*, the entity- and property-placeholders in the resulting argument scheme are replaced argument-wise with names and predicates from a domain. We hence obtain an instance of the formal argument scheme as premise-conclusion list. Each domain provides hundreds of entity-names, which can be paired with different binary predicates to create thousands of different unary predicates. The following example predicates illustrate the domains used in this study:

• Female Relatives: sister of Anna, grand-daughter of Elsa, cousin of Sarah, ...

- *Male Relatives:* grandson of Ryan, nephew of Jim, cousin of Lee, ...
- Football Fans: supporter of Real Madrid CF, ex-fan of Sevilla FC, member of SSC Napoli, ...
- Personal Care: regular consumer of Dove shampoo, infrequent user of L'Oreal shampoo, loyal buyer of Redken shampoo, ...
- Chemical Ingredients: ingredient of Maypole Soap, ingredient of OASIS CREAM, ingredient of BB concealer, ...
- Dinosaurs\*: contemporary of Megalosaurus, predator of Iguanodon, ancestor of Allosaurus....
- Philosophers\*: teacher of Aeschines of Neapolis, pupil of Cratylus, reader of Democritus, ...

Domains marked with an asterisk are used for testing only, and not for training (see below and Section 4.2).

In *step 4*, the premises of the natural language argument are randomly re-ordered.

In *step 5*, the premise-conclusion list is packed into a text paragraph by adding an argument intro, framing the premises, and adding an inference indicator. Again, multiple templates are available for doing so, which yields a large variety of textual

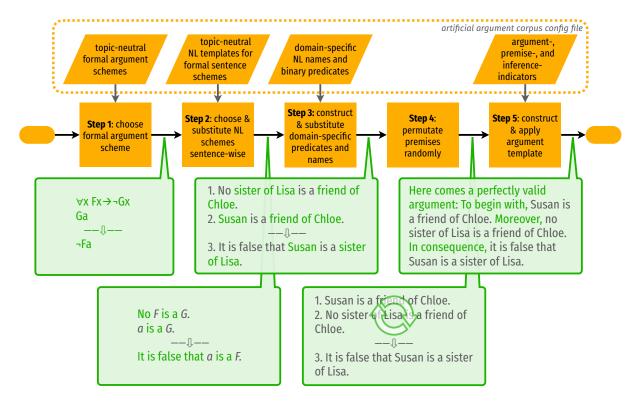


Figure 2: Pipeline for creating natural language instances of argument schemes with multiple templating.

renderings of an argument.

Following this pipeline, we generate natural language instances of each formal argument scheme, thus creating:

- 1. a training set of argumentative texts, based on the default domains and templates (TRAIN);
- 2. an evaluation set of argumentative texts, based on the default domains and templates, which are used for development (DEV);
- a test set of argumentative texts, based on the default domains and templates and used for final tests (TEST\_OUT-OF-SAMPLE);
- 4. a test set of argumentative texts, based on the domains and templates reserved for testing (TEST\_OUT-OF-DOMAIN).

This represents the artificial argument text corpus we use to train and evaluate GPT-2.

### 4 Experiments with GPT-2

We train and evaluate three compact versions of GPT-2 with 117M, 345M and 762M parameters respectively using the implementation from Wolf et al. (2019). We note that all of these models fall short of the full-scale model with 1542M parameters.<sup>2</sup>

#### 4.1 Training

From the training items in the Artificial Argument Corpus (TRAIN) we sample three types of differently-sized training sets as follows (see also the color pattern in Figure 1):

- TRAIN01: all training items which are instances of a *core scheme*, i.e. generalized modus ponens, generalized contraposition, hypothetical syllogism 1 (N=4.5K, 9K, 18K, 36K)
- TRAIN02: all training items which are instances of a base scheme (N=4.5K, 9K, 18K, 36K)
- TRAIN03: all training items in the corpus (N=4.5K, 9K, 18K, 36K)

In an attempt to avoid over-fitting, we blend the training arguments with snippets from Reuters news stories (Lewis et al., 2004) and the standardized Project Gutenberg Corpus (Gerlach and Font-Clos, 2018), trying a mixing ratio of 1:1 and thus doubling training size to N=9K, 18K, 36K, 72K. (We find that fine-tuning on the accordingly enhanced argument corpus still increases the model's perplexity on the Wiki103 dataset by a factor of 1.5 (see Appendix B), which suggests to mix a higher proportion of common texts into the train-

 $<sup>^2</sup>$ The fine-tuned models will be released through https://huggingface.co/models.

ing data in future work.) The three different versions of GPT-2 are fine-tuned (causal language modeling objective, using default training scripts by Wolf et al. (2019)) on each of the 12 enhanced training sets (hyper-parameters are detailed in Appendix A). This gives us 36 fine-tuned model versions plus the three BASE models to evaluate. Unless explicitly stated otherwise, we report results of 762M parameter model trained on 72K items.

# 4.2 Testing

## **Conclusion Completion on Artificial Argument**

Corpus To test whether language models can reason correctly, we assess their ability to accurately complete conclusions of arguments in the artificial argument corpus. Here, we make use of the fact that, by construction, the conclusion of every argument in the corpus ends with a predicate (a property-term such as "sister of Chloe" or "supporter of Tottenham Hotspurs"), which is potentially preceded by a negator. First of all, as shown in Table 1, we test whether the model is able to correctly fill in the final predicate (task split). The second, more difficult task consists in completing the final predicate plus, if present, the preceding negator (task extended). With a third, adverserial task we check how frequently the model wrongly adjoins the complement of the correct completion of the extended task (task inverted). Consider, for example, the following argument:

It is not always easy to see who is related to whom – and in which ways. The following argument pertains to this question: First premise: Every workmate of Brad is a classmate of James. Second premise: Every classmate of James is not a classmate of Theodore. So, necessarily, everyone who is a workmate of Brad is  $[not \ a]_E$  [classmate of Theodore.] $_S$ "

In the *split* task, we prompt the model with the argument, dropping  $[]_S$ , and check whether it generates "classmate of Theodore". In the *extended* task, we prompt the model with the argument, dropping  $[]_E[]_S$ , and check whether it generates "not a classmate of Theodore". Finally, in the *inverted* task, we prompt the model as before and check whether it generates "a classmate of Theodore".

Clearly, the higher the accuracy in the *split* and *extended* tasks, and the lower the accuracy in the

Task	Conclusion with cloze-style prompt	Comple- tion
split	Every $F$ is a $G$	G
	Some $F$ is not a $G$	G
	a is a $F$ or not a $G$	G
extended	Every $F$ is a $G$	a $G$
	Some $F$ is not a $G$	not a $G$
	a is a $F$ or not a $G$	not a $G$
inverted	Every F is a G	not a G
	Some $F$ is not a $G$	$\frac{1}{2}$ not a $G$
	a is a $F$ or not a $G$	$\frac{1}{2}$ not a $G$

Table 1: Three conclusion completion tasks

*inverted* task, the stronger the model's reasoning performance.

Based on the artificial argument corpus (see Section 3), we generate and distinguish three different test datasets, each of which comprises the three tasks described above, as follows:

- out of sample: contains items from TEST\_OUT-OF-SAMPLE, which share domain and natural language templates with the training data;
- paraphrased: a sample of 100 items, randomly drawn from TEST\_OUT-OF-SAMPLE, which have been manually reformulated so as to alter the premises' grammatical structure imposed by the natural language templates;
- *out of domain*: contains items from TEST\_OUT-OF-DOMAIN, which belong to different domains instantiate grammatical patterns other than the training data.

Technically, conclusion completions, in all tasks and tests, are generated by the language model with top-p nucleus sampling (p = 0.9).

Classification for NLU Benchmarks To investigate transfer learning effects, we evaluate the trained models on standard NLU benchmarks, such as GLUE AX and SNLI. These benchmark tasks are classification problems. In the following, we describe how we use the generative language models to perform such classification.

Using simple templates, we translate each benchmark entry into alternative prompts (e.g., context and question) and/or alternative completions (e.g., answers). Consider for example a

GLUE-style problem given by two sentences "The girl is eating a pizza." and "The girl is eating food" and the question whether one entails, contradicts, or is independent of the other. We can construct three prompts, corresponding to the three possible answers (entail / contradict / independent):

*Prompt1*: The girl is eating a pizza. Therefore,

*Prompt2*: The girl is eating a pizza. This rules out that

*Prompt3*: The girl is eating a pizza. This neither entails nor rules out that *Completion*: the girl is eating food.

In this case, the correct match is obviously *Prompt1–Completion*. The ability of a language model to discern that "The girl is eating pizza" entails (and does not contradict) "The girl is eating food" will be reflected in a comparatively low conditional perplexity of *Completion* given *Prompt1* and a correspondingly high conditional perplexity of *Completion* given *Prompt2* or *Prompt3*.

Let us describe this procedure in more general terms and consider a textual classification problem with categories  $k=1\dots N$ . To classify a given input X, one constructs n alternative prompts  $p_1,\dots p_n$  and m alternative completions  $c_1,\dots,c_m$   $(N=m\cdot n)$ , such that each pair  $(p_i,c_j)$  corresponds to a class k of the classification problem, i.e.,

$$L:(p_i,c_i)\mapsto \{1\ldots N\}.$$

In the above pizza example, we have N=n=3 and m=1. Moreover, let  $\mathrm{PP}_{\mathbb{L}}(c|p)$  refer to the conditional perplexity of the completion c given prompt p according to the language model  $\mathbb{L}$ . Rather than directly using this conditional perplexity as a prediction score (as for instance in Shwartz et al., 2020), which doesn't account for varying 'prima facie' or 'prior' perplexities of alternative completions, we consider the degree to which prompting the model  $\mathbb{L}$  with p changes the the perplexity of c, i.e.

$$\operatorname{relPP}_{\mathbb{L}}(c,p) := \frac{\operatorname{PP}_{\mathbb{L}}(c|p)}{\operatorname{PP}_{\mathbb{L}}(c)}.$$

In analogy to Bayesian confirmation theory, this might be termed a (perplexity-based) *relevance measure*, as opposed to a measure of absolute confirmation (cf. Carnap, 1950, pp. 346-48). We now

use relevance perplexity as a score function to predict the category of X:

$$\operatorname{category}(X) = L(\underset{(p_i, c_j)}{\operatorname{argmin}}(\operatorname{relPP}(c_j, p_i))).$$

### 5 Results

# **Conclusion Completion on Artificial Argument**

Corpus Does the (fine-tuned) GPT-2 model correctly complete conclusions of natural language arguments? Figure 3 displays the evaluation results in an aggregated way. Each subplot visualizes the accuracy of the models in the three completion tasks for a different test dataset (see Section 4.2), comparing the BASE model (points at the very left) with the fine-tuned models trained on TRAIN01, TRAIN02, and TRAIN03 (in this order from left to right). The task-specific accuracy values are distinguished by line color.

We may observe, first of all, that training on the argument corpus effectively improves conclusioncompletion-skill. In all three test datasets, the accuracy in the split and extended tasks increases as the models are trained on more and more argument schemes, far exceeding the base model's performance. Once the model has seen all schemes (TRAIN03), accuracy levels reach 100% for indomain and 70%-90% for out-of-domain tests. However, the TRAIN01 and TRAIN02 models do also generate more incorrect completions than the BASE model (inverted task). But the frequency of such incorrect completions increases much less than the frequency of correct ones (the gap between blue and gray curve widens), and it actually falls back to almost zero with the TRAIN03 model. Out-of-domain performance of the models (right-hand plot) is qualitatively similar and only slightly less strong than in-domain performance (left-hand and middle plot). The models trained on arguments from a given domain are able to effectively exercise the reasoning skill thus acquired in other domains, and have hence gained topicneutral, universal reasoning ability.

The strong performance of TRAIN01 models, averaged over all schemes, suggests that significant transfer learning occurs and that training on a few argument schemes positively affects performance on other schemes, too. To further investigate this issue, Table 2 contrasts (a) the models' accuracy on schemes they have not been trained on – averaged over TRAIN01 and TRAIN02 mod-

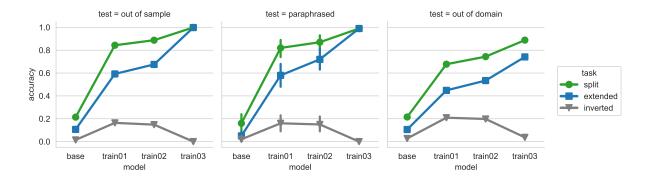


Figure 3: Accuracy of four model versions in three conclusion completion tasks and on different test datasets (out of sample, paraphrased, out of domain).

	BASE	(a) schemes not in training data (TR01-02)			(b) trained on schemes (TR01-03)		
Task		o-o-sample	paraphr.	o-o-domain	o-o-sample	paraphr.	o-o-domain
split	21.4	85.4	82.0	69.4	99.9	99.2	89.0
extended	10.7	60.3	59.3	45.8	99.9	99.2	76.2
inverted	1.5	16.9	18.0	22.1	0.0	0.0	3.2

Table 2: Accuracy of models in three conclusion completion tasks and on different test datasets (out of sample, paraphrased, out of domain). Columns report, separately, the performance (a) on schemes the model has not been trained on, and (b) on schemes that are covered by the model's training data.

els – with (b) their accuracy on schemes that are instantiated in their respective training corpus – averaged over TRAIN01, TRAIN02, and TRAIN023 models. The upshot is that trained models perform way more strongly than the base model not only on argument schemes they've been trained, but also on those schemes they haven't seen yet. We take this to be a promising result as it strengthens the analogy between teaching critical thinking and training language models: generic intermediary pre-training on high-quality texts that exemplify a specific, basic reasoning skill – namely, simple deductive argumentation – improves other, more complex reasoning skills.

Figure 4 gives further insights by differentiating evaluation results according to argument type. Its subplots are arranged in a grid that mirrors the organisation of argument schemes in Figure 1. Each subplot visualizes the ability of the models to correctly complete arguments of the corresponding scheme (given the out-of-sample test dataset). Accordingly, the left-hand plot in Figure 3 in effect averages all curves in Figure 4. Reported accuracy values that fall within gray background areas are attained by models which have seen the corresponding scheme during training. Vice versa, thick lines on white background visualize model

performance on unknown schemes. Figure 4 reveals, first of all, that even the BASE models (only pre-training, no fine-tuning) display a significant ability to correctly complete conclusions of some kinds of arguments. For example, GPT-2-762M achieves 50% accuracy (*split* task) in completing contrapositions, 30% accuracy in completing generalized modus ponens, and still 20% accuracy in completing disjunctive syllogism and dilemma arguments. These findings further corroborate the hypothesis that NLMs learn (basic) linguistic and reasoning skills "on the fly" by training on a large generic corpus (Radford et al., 2019).

In addition, the matrix plot (Figure 4) demonstrates that some types of arguments are much easier to master, given training on the core and possibly base schemes, than others. For instance, *complex\_predicates* variants of generalized modus ponens or *de\_morgan* variants of generalized modus tollens seem to be easily mastered by the TRAIN01 model. In contrast, even the TRAIN02 model, which has been fine-tuned on all eight base schemes, struggles with the *negation\_variants* of generalized modus ponens (generating substantially more incorrect than correct completions). All in all, the picture that emerges is plausible: Generalization towards novel types

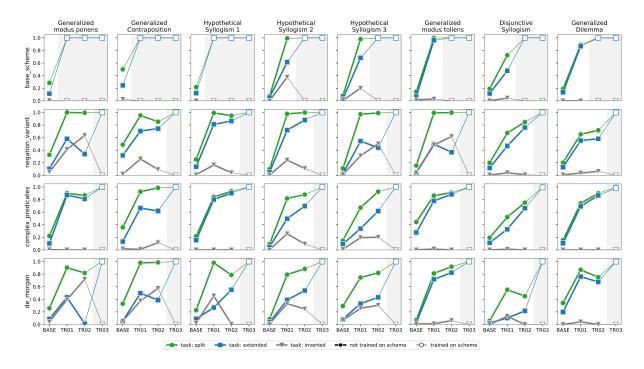


Figure 4: Accuracy of conclusion completions (three tasks) for instances of different argument schemes (see Figure 1) and four model versions.

of argument appears to be comparatively difficult whenever the new scheme involves negations (compare 2nd and 4th row in Figure 4 with 3rd row). This is consistent with the finding that some NLMs seemingly fail to understand simple negation (Kassner and Schütze, 2020; Talmor et al., 2020).

The results reported so far suggest that reasoning skills acquired on (a subset of) the artificial argument corpus generalize rather well – both to other domains and other types of arguments. We have further cross-checked these statistical findings by letting the models complete a conclusion of a simple manually authored argument:

[Hermes] Every philosopher is mortal. Hermes is not mortal. Therefore, Hermes...

This text differs syntactically and semantically from any argument possibly contained in the artificial argument corpus (where predicates have always the form "is/being a Y of X," and no domain covers philosophers or mortality). Obviously, it follows that Hermes "is not a philosopher." The argument instantiates *generalized modus tollens*, which is not a core scheme in TRAIN01. Can TRAIN01-models nonetheless fill out the unfinished argument in a sensible way?

Table 3 counts and compares the most frequent

		762M		117M
Completion		TR01	BASE	TR01
is not a philosopher.	*	100	2	2
is immortal.	=	0	12	0
is not a critic.	0	0	0	9
is mortal.	†	0	8	0
is not mortal.	=	0	6	0
is not Hermes.	†	0	2	0
does not exist.	0	0	2	0
is not God.	0	0	2	0
is not a friend of Eckhardt.	0	0	0	1
is not an expert of BSI Ar-	0	0	0	1
senal FC.				
is not a friend of Atalanta.	0	0	0	1
is not an infrequent user of	0	0	0	1
Neutrogena shampoo.				
others		0	66	85

Table 3: Absolute frequency of predicted completions for the hand-written [Hermes] query by three different models. Completions are – relative to the premises – entailed  $(\star)$ , redundant (=), contradictory  $(\dagger)$  or independent  $(\circ)$ .

completions generated by two TRAIN01 models (762M and 117M) and by the large untrained BASE model (762M). Exclusively the 762Mmodel trained on the core schemes reliably predicts the correct conclusion. The large BASE model rather repeats a premise or even generate a contradiction, whereas the small TRAIN01 model (117M) changes the topic. This is consistent with and illustrates our previous findings. Remarkably, although both the small and the large TRAIN01 models have been fine-tuned on precisely the same arguments, only the large model seems to correctly recognize the logical structure of the [Hermes] argument. Generic language modeling skill, it is suggested, facilitates the successful generalization of learned argument patterns beyond the templates used to create the synthetic training data.

To further understand transfer learning effects, we next examine whether intermediary pretraining on the artificial argument corpus improves zero-shot performance in other NLP reasoning tasks (i.e., without task-specific fine-tuning).

**GLUE AX** The GLUE datasets (Wang et al., 2018) represent standard benchmarks for natural language understanding (NLU). We evaluate our models' NLU skill in terms of accuracy on the curated GLUE diagnostics dataset (Figure 5).

Training on the artificial argument corpus substantially boosts accuracy on the GLUE diagnostics. Accuracy increases by at least 5 and up to 17 percentage points, depending on model size. Remarkably, training on the core scheme alone suffices to bring about these improvements.

This is a major finding and our clearest evidence so far that training on the AAC involves substantial transfer learning effects.

**SNLI** The SNLI dataset (Bowman et al., 2015) is another standard benchmark for NLI. Like the GLUE dataset, it consists in pairs of sentences which entail, contradict, or don't bear on each other. The assessment of our models with respect to SNLI data proceeds in close analogy to the GLUE benchmark.

The results, reported in Figure 5, are consistent with, albeit less definite than our previous findings for the GLUE benchmark: First and foremost, fine-tuning on all schemes (TRAIN03) improves the performance by up to 8 percentage points. Training on fewer schemes is slightly less effective. However, it is only the small and medium

sized model that profit from fine-tuning on the AAC; the SNLI performance of the 762M parameter model gets rather deteriorated. This might be due to a coincidentally strong performance of the corresponding BASE model (see Figure 7), or suggest that the large model, unlike the smaller ones, has already learned during pre-training whatever is of relevance for SNLI in the AAC. (Further experiments, preferably involving more model versions, are required to clarify this.)

Argument Reasoning Comprehension Task The Argument Reasoning Comprehension (ARC) task (Habernal et al., 2018) assesses the ability to identify a missing premise in an informally reconstructed and not necessarily deductively valid argument. It is a multiple-choice task where two alternative sentences are provided, one of which is the missing premise.

We design and apply specific templates to construct prompts and completions, and calculate relative perplexity as described in Section 4.2.

As shown in Figure 5, we find no evidence of transfer learning effects with respect to ARC.

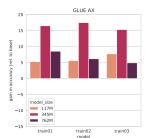
**LogiQA** LogiQA (Liu et al., 2020) is a collection of nearly 9,000 multiple-choice questions (four alternative answers each) used in critical thinking assessments. These questions span the whole range of critical thinking tasks.

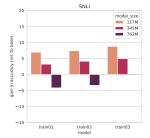
We design and apply specific templates to construct prompts and completions (one prompt and four completions per question), and use perplexity scores to predict classifications as described above (Section 4.2).

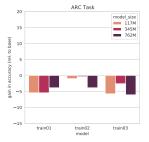
As can be seen from Figure 5, training on the artificial argument corpus has no effect whatsoever on the ability of the models to handle the critical thinking tasks collected in LogiQA.

### 6 Conclusion

This paper has taken a first step towards the creation of a critical thinking curriculum for neural language models. It presents a corpus of deductively valid, artificial arguments, and uses this artificial argument corpus to train and evaluate GPT-2. The observation of strong transfer learning effects/generalization is its main finding: Training a model on a few central core schemes allows it to accurately complete conclusions of different types of arguments, too. The language models seem to connect and to generalize the core argument







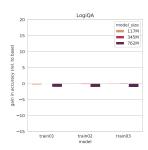


Figure 5: Gains in accuracy due to fine-tuning on the AAC (accuracy TRAIN model – accuracy BASE model) for differently sized models and different NLP benchmark tasks: the GLUE diagnostics data, the SNLI dataset, the argument reasoning comprehension (ARC) benchmark, and the LogiQA dataset.

schemes in a correct way. Moreover, the models are equally able to apply learned argument patterns beyond the domain they have been trained on, and there is evidence that generic language modeling skill facilitates the successful generalization of learned argument patterns. These findings are consistent with previous work on rule reasoning (Clark et al., 2020). They suggest that there exist (learning-wise) fundamental reasoning skills in the sense that generic intermediary pre-training on texts which exemplify these skills leads to spillover effects and can improve performance on a broad variety of reasoning tasks. The synthetic argumentative texts might be a good starting point for building such a "critical thinking curriculum for language models."

Moreover, the trained models have been tested on different reasoning benchmarks. We obtain clear and promising results for the GLUE and SNLI benchmarks. But training on the argument corpus doesn't affect the performance with regard to the semantically more demanding Argument Reasoning Comprehension task or the critical thinking assessment compiled in LogiQA.

Our work suggests different directions for advancing the approach adopted in this paper and further improving the general reasoning skill of neural language models:

• The syllogistic argument text corpus might be complemented with corpora of arguments that instantiate different kinds of correct schemes, e.g., propositional inference schemes, modal schemes, argument schemes for practical reasoning, complex argument schemes with intermediary conclusions or assumptions for the sake of the argument, etc. (Technically, we provide the infrastructure for doing so, as all this might be achieved

- through adjusting the argument corpus configuration file.)
- To succeed in NLI tasks, it doesn't suffice to understand 'what follows.' In addition, a system needs to be able to explicitly discern contradictions and *non sequiturs* (relations of logical independence). This suggests that the artificial argument corpus might be fruitfully supplemented with corpora of correctly identified aporetic clusters (Rescher, 1987) as well as corpora containing correctly diagnosed fallacies.
- In addition, the idea of curriculum learning for ML (Bengio et al., 2009) might be given a try. Accordingly, a critical thinking curriculum with basic exemplars of good reasoning would not only be used to fine-tune a pre-trained model, but would be employed as starting point for training a language model from scratch.

Natural language templating is a fundamental technique used throughout this paper: both in constructing the artificial argument corpus as well as in transforming the NLP benchmark datasets into text that can be processed by language models. The concrete templates applied have been designed in a trial-and-error process. It is far from clear that these represent optimal choices for effectively eliciting a language model's skills. Still, following (Jiang et al., 2020), it seems of great importance to gain a more systematic understanding of different templating strategies and their effects on metrics based on accuracy and perplexity.

In conclusion, designing a critical thinking curriculum for neural language models seems to be a promising and worthwhile research program to pursue.

# A Appendix: Training Parameters

We train the models on 8 GPUs for 2 epochs with batch size = 2, learning rate =  $5 \times 10^{-5}$ , gradient accumulation steps = 2, and default parameters of the HuggingFace implementation otherwise (Wolf et al., 2019).

# B Appendix: Performance Metrics for Differently Sized Training Sets

Figure 6 displays accuracy values on conclusion completion tasks for models trained on differently sized datasets.

Figure 7 reports perplexity and NLU accuracy metrics for models trained on differently sized datasets.

### References

- Amanda Askell. 2020. Gpt-3: Towards renaissance models. In *Daily Nous Blog: Philosophers On GPT-3*.
- Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. 2009. Curriculum learning. In *Proceedings of the 26th Annual International Conference on Machine Learning*, ICML '09, pages 41–48, New York, NY, USA. ACM.
- Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Çelikyilmaz, 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 (ACL)*.
- Tracey Bowell and Gary Kemp. 2014. *Critical Thinking: A Concise Guide*, 4th edition edition. Routledge, London.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics.
- 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.
- Georg Brun and Gregor Betz. 2016. Analysing practical argumentation. In Sven Ove Hansson and Gertrude Hirsch-Hadorn, editors, *The Argumentative Turn in Policy Analysis. Reasoning about Uncertainty*, pages 39–77. Springer, Cham.
- Rudolf Carnap. 1950. Logical Foundations of Probability. University of Chicago Press, Chicago.
- J. Cheng, M. Bernstein, C. Danescu-Niculescu-Mizil, and J. Leskovec. 2017. Anyone can become a troll: Causes of trolling behavior in online discussions. CSCW: Proceedings of the Conference on Computer-Supported Cooperative Work. Conference on Computer-Supported Cooperative Work, 2017, page 1217–1230.
- Peter Clark, Oyvind Tafjord, and Kyle Richardson. 2020. Transformers as soft reasoners over language. *arXiv preprint arXiv:2002.05867v2*.
- Richard Feldman. 2014. *Reason and Argument*. Pearson, Harlow.
- Alec Fisher. 2001. Critical Thinking: An Introduction. Cambridge University Press, Cambridge.
- Martin Gerlach and Francesc Font-Clos. 2018. A standardized project gutenberg corpus for statistical analysis of natural language and quantitative linguistics. *CoRR*, abs/1812.08092.
- Ben Gilburt. 2019. Examining gender bias in openai's gpt-2 language model. *hackernoon.com*.
- Nicolas Gontier, Koustuv Sinha, Siva Reddy, and Christopher Pal. 2020. Measuring systematic generalization in neural proof generation with transformers.

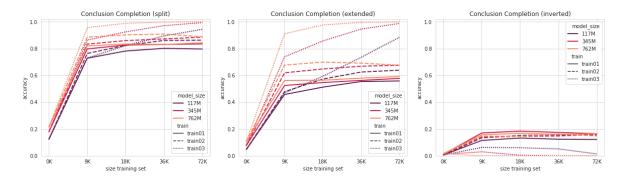


Figure 6: Accuracy on three conclusion completion tasks as a function of training corpus size.

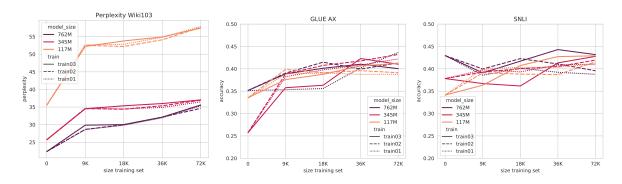


Figure 7: Perplexity and NLI metrics as a function of training corpus size.

Radu Cornel Guiaşu and Christopher W Tindale. 2018. Logical fallacies and invasion biology. *Biology & philosophy*, 33(5-6):34.

Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and Noah A Smith. 2018. Annotation artifacts in natural language inference data. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 107–112.

Ivan Habernal, Henning Wachsmuth, Iryna Gurevych, and Benno Stein. 2018. The argument reasoning comprehension task: Identification and reconstruction of implicit warrants. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long Papers), pages 1930–1940. Association for Computational Linguistics.

Sven Ove Hansson. 2004. Fallacies of risk. *Journal of Risk Research*, 7(3):353–360.

Zhengbao Jiang, Frank F. Xu, Jun Araki, and

Graham Neubig. 2020. How can we know what language models know? *Transactions of the Association for Computational Linguistics*, 8:423–438.

Daniel Kahneman. 2011. *Thinking, fast and slow*, 1st edition. Farrar, Straus and Giroux, New York.

Nora Kassner, Benno Krojer, and Hinrich Schütze. 2020. Are pretrained language models symbolic reasoners over knowledge?

Nora Kassner and Hinrich Schütze. 2020. Negated and misprimed probes for pretrained language models: Birds can talk, but cannot fly.

Joe Lau and Jonathan Chan. 2020. Critical thinking web. https://philosophy.hku.hk/think.

D. D. Lewis, Y. Yang, T. Rose, and F. Li. 2004. Rcv1: A new benchmark collection for text categorization research. *Journal of Machine Learning Research*, 5:361–397.

Kevin Lin, Oyvind Tafjord, Peter Clark, and Matt Gardner. 2019. Reasoning over paragraph effects in situations. *Proc. MRQA Workshop (EMNLP'19)*.

- Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. 2020. Logiqa: A challenge dataset for machine reading comprehension with logical reasoning. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJ-CAI 2020*, pages 3622–3628. ijcai.org.
- 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.
- Kaixin Ma, Filip Ilievski, Jonathan Francis, Yonatan Bisk, Eric Nyberg, and Alessandro Oltramari. 2020. Knowledge-driven self-supervision for zero-shot commonsense question answering. *arXiv preprint arXiv:2011.03863*.
- Arindam Mitra, Pratyay Banerjee, Kuntal Kumar Pal, Swaroop Mishra, and Chitta Baral. 2019. How additional knowledge can improve natural language commonsense question answering? *arXiv preprint arXiv:1909.08855*.
- Timothy Niven and Hung-Yu Kao. 2019. Probing neural network comprehension of natural language arguments. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4658–4664, Florence, Italy. Association for Computational Linguistics.
- SP Norris and RH Ennis. 1989. What is critical thinking. *The practitioner's guide to teaching thinking series: Evaluating critical thinking*, pages 1–26.
- Fabio Paglieri. 2017. A plea for ecological argument technologies. *Philosophy & Technology*, 30(2):209–238.
- Jason Phang, Thibault Févry, and Samuel R Bowman. 2018. Sentence encoders on stilts: Supplementary training on intermediate labeled-data tasks. arXiv preprint arXiv:1811.01088.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *Preprint*.

- Nicholas Rescher. 1987. Aporetic method in philosophy. *The Review of metaphysics*, 41(2):283–297.
- Kyle Richardson, Lawrence S. Moss, , and Ashish Sabharwal. 2020. Probing natural language inference models through semantic fragments. *AAAI'20*.
- Swarnadeep Saha, Sayan Ghosh, Shashank Srivastava, and Mohit Bansal. 2020. Prover: Proof generation for interpretable reasoning over rules.
- Timo Schick and Hinrich Schütze. 2020a. Exploiting cloze questions for few shot text classification and natural language inference.
- Timo Schick and Hinrich Schütze. 2020b. It's not just size that matters: Small language models are also few-shot learners.
- Vered Shwartz, Peter West, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2020. Unsupervised commonsense question answering with self-talk.
- Koustuv Sinha, Shagun Sodhani, Jin Dong, Joelle Pineau, and William L. Hamilton. 2019. Clutrr: A diagnostic benchmark for inductive reasoning from text. *arXiv* preprint arXiv:1908.06177v2.
- Cass R Sunstein and Reid Hastie. 2015. Wiser: getting beyond groupthink to make groups smarter. Harvard Business Review Press, Boston.
- Oyvind Tafjord, Matt Gardner, Kevin Lin, and Peter Clark. 2019. Quartz: An open-domain dataset of qualitative relationship questions. *EMNLP/IJCNLP*.
- Alon Talmor, Yanai Elazar, Yoav Goldberg, and Jonathan Berant. 2020. olmpics on what language model pre-training captures.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355.

- J. Weston, A. Bordes, S. Chopra, and T. Mikolov. 2016. Towards ai-complete question answering: A set of prerequisite toy tasks. *ICLR*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Huggingface's transformers: State-of-the-art natural language processing. *ArXiv*, pages arXiv–1910.