

Glider: Global and Local Instruction-Driven Expert Router

Anonymous ACL submission

Abstract

The development of performant pre-trained models has driven the advancement of routing-based expert models tailored to specific tasks. However, these methods often favor generalization over performance on held-in tasks. This limitation adversely impacts practical applicability, as real-world deployments require robust performance across both known and novel tasks. We observe that current token-level routing mechanisms neglect the global semantic context of the input task. To address this, we propose a novel method, **Global and Local Instruction Driven Expert Router** (GLIDER) that proposes a multi-scale routing mechanism, encompassing a semantic global router and a learned local router. The global router leverages recent LLMs' semantic reasoning capabilities to generate task-specific instructions from the input query, guiding expert selection across all layers. This global guidance is complemented by a local router that facilitates token-level routing decisions within each module, enabling finer control and enhanced performance on unseen and challenging tasks. Our experiments using T5-based expert models for T0 and FLAN tasks demonstrate that GLIDER achieves substantially improved held-in performance while maintaining strong generalization on held-out tasks. Additionally, we perform ablations experiments to dive deeper into the components of GLIDER and plot routing distributions to show that GLIDER can effectively retrieve the correct expert for held-in tasks while also demonstrating compositional capabilities for held-out tasks. Our experiments highlight the importance of our multi-scale routing that leverages LLM-driven semantic reasoning for MoErging methods.

1 Introduction

The emergence of highly capable large language models (LLMs) has marked an increased attention in downstream task specialization. This spe-

cialization often leverages parameter-efficient fine-tuning (PEFT) techniques, such as LoRA (Hu et al., 2021), which introduce minimal trainable parameters (“adapters”) to adapt pre-trained LLMs for specific tasks. The compact size of these specialized PEFT modules enables easy sharing, which has led to the distribution of an evergrowing number of adapters on various platforms.

This proliferation of expert models, *i.e.* specialized adapters, has led to the development of methods for re-using such experts to improve performance or generalization (Muqeeth et al., 2024; Ostapenko et al., 2024; Huang et al., 2024a). Central to these approaches are routing mechanisms that adaptively select relevant experts for a particular task or query. These routing methods have been referred to as “Model MoErging” (Yadav et al., 2024) since they frequently share methodologies and ideas with mixture-of-experts (MoE) models (Shazeer et al., 2017; Fedus et al., 2022; Du et al., 2022) and model merging (Yadav et al., 2023b,a; Ilharco et al., 2022). However, MoE methods train experts jointly from scratch (Gupta et al., 2022) while MoErging utilizes a decentralized, community-sourced pool of pre-trained experts. Furthermore, it departs from traditional model merging techniques by dynamically and adaptively combining these experts, optimizing performance at the query or task level. MoErging methods offer three key advantages: (1) They support decentralized model development by reusing and routing among independently trained experts, reducing reliance on centralized resources. (2) They facilitate modular capability expansion and “transparency” in updates as they either add or modify specialized expert models. (3) They allow for compositional generalization by recombining fine-grained skills from various experts, extending the system’s abilities to new unseen tasks beyond the capabilities of the individual expert models.

Most MoErging (Chronopoulou et al., 2023;

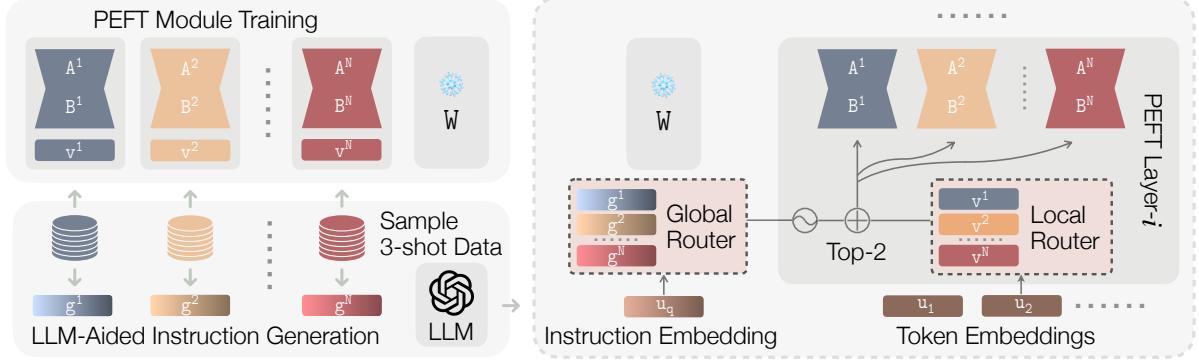


Figure 1: Overview of our method. **Contributor** (left): Each contributor utilizes local data to train several components: the PEFT module (comprising A_i and B_i), task vectors (v_i), and global routing vectors (g_i). For the latter, an LLM is employed to generate semantically-informed instructions based on 3 randomly selected examples, which are then embedded into g_i . **Aggregator** (right): The aggregator utilizes local and global task vectors to construct local routers $[\bar{v}^1; \dots; \bar{v}^N]$ and a global router $[g^1; \dots; g^N]$, respectively. For each query, the global router uses an LLM-generated instruction embedding to produce the global routing score. This score is then scaled and combined with the local routing score, enabling fine-grained control over expert selection.

084 Muqeeth et al., 2024; Zhao et al., 2024b) methods prioritize either known or unseen tasks, limiting real-world applicability where both are critical. 085 Real-world queries often span domains and defy clean categorization into predefined task boundaries. 086 For instance, translating and analyzing text requires collaboration between multiple experts rather than selecting a single specialized model. 087 Current approaches struggle in such scenarios. 088 Phatgoose demonstrates this tradeoff, excelling on 089 unseen tasks but underperforming on known ones. 090

091 We hypothesize that this gap arises from the model’s token-level routing mechanism. We show that for the held-in tasks, the independent routing decisions at each layer, based solely on individual token embeddings, lack sufficient global context to retrieve the correct expert for all tokens at every module. This leads to suboptimal routing, which may propagate noise through the network, further hindering accurate expert utilization in deeper layers. This highlights a critical limitation of token-level approaches to handling held-in tasks, which hence falls short of the goal of building a routing system that seamlessly handles arbitrary queries. We believe that adding a global routing mechanism based on semantic task information can aid the token-level router for the correct retrieval of held-in tasks. Hence, we ask the question.

(Q) Can we leverage LLMs to generate semantics-aware task instructions to guide routing mechanism to facilitate both specialization and generalization?

This paper addresses the challenges by inves-

tigating the potential of leveraging the inherent reasoning and generalization capabilities of LLMs to guide the routing process in an MoE-like model composed of specialized LoRA modules. We introduce, **Global and Local Instruction Driven Expert Router (GLIDER)** that hinges on a multi-scale routing mechanism that contains both local and global routers to select top-2 expert models as shown in Figure 1. The global router leverages LLM-generated, semantics-aware instructions (see Appendix B.2) for each input query to score expert models. This high-level guidance is then complemented by a learned local router, which makes token-level routing decisions at each module, enabling fine-grained control and improving performance on the challenging held-out tasks. Through this framework, we highlight the crucial role of LLM reasoning in unlocking the compositional generalization capabilities of MoE models.

To test the effectiveness of our GLIDER method, we follow Phatgoose (Muqeeth et al., 2024) and use T5 models (Raffel et al., 2020) to create expert models for T0 held-in (Sanh et al., 2022) and FLAN tasks (Longpre et al., 2023) and test performance on T0 held-in & held-out (Sanh et al., 2022) and big-bench lite (BIG-bench authors, 2023) & hard tasks (Suzgun et al., 2022). Our key contributions and findings are:

- We introduce GLIDER, which employs LLM-guided multi-scale global and local attention. Our experiments show that GLIDER outperforms previous methods, significantly improving performance on held-in tasks (e.g. 6.6% over Phatgoose on T0 held-in) while also en-

- 148 hancing zero-shot held-out compositional gen-
149 eralization (e.g. 0.9% on T0 held-out).
- 150 • We find that without LLM assistance, MoE
151 models underperform individual specialized
152 models on held-in tasks by 8.2%. Incorpor-
153 ating semantic-aware instructions enables
154 GLIDER to achieve comparable performance,
155 demonstrating the LLM’s capacity to effec-
156 tively infer task identity and guide module
157 selection without explicit task labels.
- 158 • GLIDER also maintains strong performance on
159 held-out tasks, showcasing its adaptability and
160 generalization capabilities. Our work high-
161 lights the critical role of LLMs in enhancing
162 MoE models’ compositional generalization,
163 advancing the development of more robust
164 and versatile AI systems capable of handling
165 both familiar and novel tasks.

2 Related Works

The abundance of specialized expert models has spurred the development of techniques to leverage “experts” models for enhanced performance and generalization. [Yadav et al. \(2024\)](#) called such techniques as “MoErging”¹ methods which rely on adaptive routing mechanisms to select relevant experts for specific tasks or queries. These methods can be broadly classified into four categories based on the design of their routing mechanisms.

Embedding-Based Routing: This category encompasses methods that derive routing decisions from learned embeddings of expert training data. These methods typically compare a query embedding against the learned expert embeddings to determine the optimal routing path. Examples include AdapterSoup ([Chronopoulou et al., 2023](#)), Retrieval of Experts ([Jang et al., 2023](#)), LoRaRetriever ([Zhao et al., 2024b](#)), Mo’LoRA ([Maxine, 2023](#)), the embedding-based approach of Airoboros ([Durbin, 2024](#)), and Dynamic Adapter Merging ([Cheng et al., 2024](#)).

Classifier-Based Routing: This category consists of methods that train a router to function as a classifier. This router is trained to predict the optimal routing path based on features extracted from expert datasets or unseen data. Representative methods in this category include Zooter ([Lu et al., 2023](#)), Branch-Train-Mix ([Sukhbaatar et al., 2024](#)),

¹See e.g. https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard

Routing with Benchmark Datasets ([Shnitzer et al., 2023](#)), Routoo ([Mohammadshahi et al., 2024](#)), and RoutELLM ([Ong et al., 2024](#)). The key distinction between embedding-based and classifier-based routing lies in the router’s architecture and training methodology. While embedding-based routing often employs a nearest neighbor approach, classifier-based routing typically relies on logistic regression or analogous classification techniques.

Task-Specific Routing: This category focuses on methods tailored to enhance performance on specific target tasks. These methods learn a task-specific routing distribution over the target dataset to optimize performance for the given task. Methods include LoraHub ([Huang et al., 2023](#)), LoRA-Flow ([Wang et al., 2024](#)), AdapterFusion ([Pfeiffer et al., 2021](#)), π -Tuning ([Wu et al., 2023](#)), Co-LLM ([Shen et al., 2024](#)), Weight-Ensembling MoE ([Tang et al., 2024](#)), MoLE ([Wu et al., 2024](#)), MeteoRA ([Xu et al., 2024](#)), PEMT ([Lin et al., 2024](#)), MixDA ([Diao et al., 2023](#)), and Twin-Merging ([Lu et al., 2024](#)).

Routerless Methods: This final category encompasses methods that do not rely on an explicitly trained router. Instead, these methods often employ alternative mechanisms, such as heuristics or rule-based systems, for routing decisions. Examples include Arrow ([Ostapenko et al., 2024](#)), PHAT-GOOSE ([Muqeeth et al., 2024](#)), the “ask an LLM” routing of Airoboros ([Durbin, 2024](#)) and LlamaIndex ([Liu, 2024](#)). Phatgoose and Arrow use only local routers, in contrast, GLIDER uses both local and global guidance for routing.

3 Problem Statement

In our work, we aim to build a routing mechanism capable of performing well on diverse queries from various tasks, including both seen and unseen tasks. For each query/token and module, this routing mechanism dynamically selects a model from a large pool of specialized expert models to achieve high performance. To facilitate modular development, we adopt a *contributor-aggregator* framework ([Yadav et al., 2024](#)) where individual contributors create specialized expert models from a generalist model for their respective tasks and distribute these models to others for public usage. The aggregator builds a routing mechanism over the expert models that shared by the contributor to direct queries to the most relevant experts. Follow-

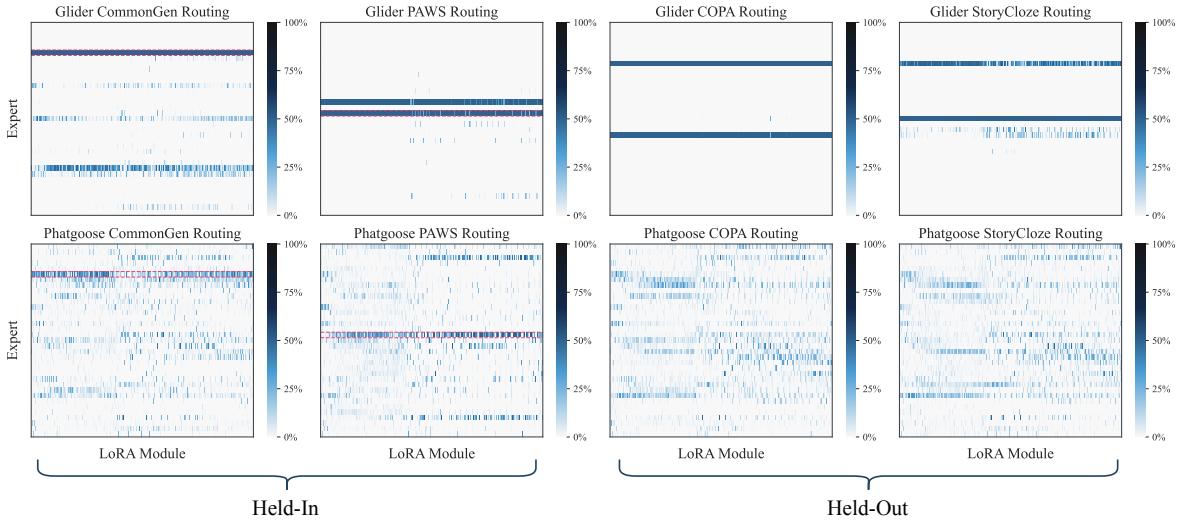


Figure 2: We present routing heatmaps for GLIDER and Phatgoose on two held-in and two held-out tasks. For held-in tasks, oracle experts are marked with red dashed lines. GLIDER selects oracle experts more frequently than Phatgoose for held-in tasks, leading to improvements of 3.3% on CommonGen and 6.5% on PAWS. For held-out tasks, GLIDER also tends to select the most relevant experts across most LoRA modules, resulting in improvements of 2.2% on COPA and 5.8% on StoryCloze.

ing recent works (Muqeeth et al., 2024; Ostapenko et al., 2024), we use parameter-efficient finetuning (PEFT) (Liu et al., 2022; Sung et al., 2022; Poth et al., 2023) methods like LoRA (Hu et al., 2022) to train the expert models. Since PEFT typically has lower computational and communication costs than full-model finetuning (Hu et al., 2022; Liu et al., 2022), the use of PEFT makes it easier to participate and contribute. PEFT methods introduce modules throughout the model – for example, LoRA (Hu et al., 2022) introduces a low-rank update at every linear layer in the model. We refer to each of these updates as a *module*. Subsequently, the trained expert models and additional information are shared with the aggregators. The aggregator’s job is to collect these expert models and the additional information and design the post-hoc routing mechanism. This mechanism will effectively direct incoming queries to the most appropriate expert model for each token and at each module to ensure optimal performance on both seen and unseen tasks. This approach allows for the seamless integration of new capabilities by adding expert models to the existing pool. Next, we formally define our contributor-aggregator framework.

Let us assume that there are N contributors, $\{c_1, c_2, \dots, c_N\}$, and each contributor c_i has access to a task-specific datasets \mathcal{D}_i . Each contributor, c_i , follows the predefined training protocol \mathcal{T} provided by the aggregator. The training protocol (\mathcal{T}) takes in a base model (θ_{base}) and a dataset (\mathcal{D}_i). It returns the expert model parameters (ϕ_i) along with any additional information (Ψ_i) that needs to

be shared with the aggregators, for example, the gate vectors described in Section 4.1. Specifically, $\{\phi_i, \Psi_i\} \leftarrow \mathcal{T}(\theta_{\text{base}}, \mathcal{D}_i)$. All contributors share this information with the aggregator, which creates a pool of models containing $\{(\phi_i, \Psi_i)\}_{i=1}^N$. The aggregators (\mathcal{A}) then uses these expert models and the auxiliary information to create a routing mechanism $\mathcal{R}(\cdot)$ that takes the user query q as the input and return routing path describing how the information is routed through the given set of expert models. Formally, $\mathcal{R}(\cdot) \leftarrow \mathcal{A}(\{(\phi_i, \Psi_i)\}_{i=1}^N)$. The function $\mathcal{R}(\cdot)$ describe the full path of input query by making various choices about 1) expert input granularity, choosing to route per-token, per-query, or per-task, 2) expert depth granularity, opting for either per-module or model-level routing, and 3) selecting between sparse or dense routing. Finally, the aggregator uses the routing mechanism to answer incoming queries.

4 Methodology

To recap, our goal is to build a MoErging method that dynamically routes queries to a diverse pool of specialized expert models, addressing the challenge of effectively handling queries from various tasks and ensuring both held-in and held-out performance. Our proposed method, **Global and Local Instruction Driven Expert Router** (GLIDER), leverages a combination of local and global routing vectors to achieve this goal. Specifically, contributors train task-specific routing vectors, while an LLM generates global semantic task instructions, which are then converted to global instruction routing

vectors. During inference, these local and global routing vectors are combined to perform top-k discrete routing, directing queries to the most suitable expert model. This process is visualized in Figure 1 and described in detail below.

4.1 Expert Training Protocol

Our expert training protocol \mathcal{T} takes as input the base model parameters, θ_{base} , and a dataset d and performs three steps to obtain the required output. First, we train the LoRA experts (ϕ) and then the local routing vectors (l) while keeping the LoRA experts fixed. Finally, we train the global routing vector (g) by using an LLM and an embedding model. Formally, in our case, $\phi, \Psi = \{l, g\} \leftarrow \mathcal{T}(\theta_{\text{base}}, d)$ which are then shared with the aggregators to create the routing mechanism. We described these steps in detail below.

PEFT Training of Expert Model. GLIDER is compatible with expert models trained using parameter-efficient finetuning methods (e.g. LoRA (Hu et al., 2022), Adapters (Houlsby et al., 2019)) that introduce small trainable modules throughout the model. We focus on PEFT experts because they typically have lower computational and communication costs than full-model finetuning (Yadav et al., 2023a), making it easier to train and share expert models. Following Phatgoose (Muqeeth et al., 2024), this work specifically focuses on LoRA (Hu et al., 2022) due to its widespread use. LoRA introduces a *module* comprising the trainable matrices $B \in \mathbb{R}^{d \times r}$ and $A \in \mathbb{R}^{r \times n}$ in parallel to each linear layer with parameters $W \in \mathbb{R}^{d \times n}$. Given the i^{th} input token activation u_i , LoRA modifies the output of the linear layer from Wu_i to $Wu_i + \frac{\alpha}{r} \cdot BAu_i$ where α is a constant and usually is set to 1. During training, the matrices A and B are trainable, while the original linear layer W is kept frozen. We denote the final trained expert parameters with $\phi = \{(A_1, B_1), \dots, (A_m, B_m)\}$, where m is the number of modules in the model.

Training Local Routing Vectors. Following Phatgoose (Muqeeth et al., 2024), after training the PEFT modules on their dataset, a local router is introduced before each PEFT module. This router, employing a shared vector across all queries and tokens, dynamically determines the utilization of the PEFT module based on the input token activations. The router is trained for a small number of steps using the same dataset and objective as

the PEFT module while keeping the expert PEFT parameters fixed. This process effectively learns to associate the token activation patterns with the learned expert model. For LoRA, the local router, represented by a trainable vector $v \in \mathbb{R}^d$, controls the contribution of the PEFT module to the final output. This results in a modified linear layer of the form $Wu_i + \frac{\alpha}{r} \cdot BAu_i \cdot \text{sigmoid}(v^\top u_i)$, where α, W, B , and A are frozen, and the local router v is learned. We denote the final local routing vectors as $l = \{v_1, \dots, v_m\}$ where m is the number of modules in the model.

Creating LLM-Aided Global Routing Vector. The local routing vectors capture the intricate relationships between token activations and expert models, enabling efficient query routing in cases where no dedicated expert is available. Conversely, for queries corresponding to held-in tasks, direct retrieval of the relevant expert model is preferred to process the full query. For this purpose, we create a global routing vector that utilizes an LLM to generate a semantically-informed instruction, termed as task description, which effectively captures the essence of the kind of queries the expert can handle. We prompt an LLM with three randomly selected in-context examples to generate this task description. We used the gpt-4-turbo model along with the prompt provided in Appendix B. The resulting task description is then embedded using an off-the-shelf embedding model, specifically the nomic-embed-text-v1.5 model, to produce a global routing vector for the task. We denote the global routing vector as $g \in \mathbb{R}^{d_g}$.

4.2 GLIDER: Inference Expert Aggregation

Following training, all contributors share their expert models along with the auxiliary information comprising of the local and global routing vectors, $\{\phi^t, l^t, g^t\}_{t=1}^N$, where t indexes the input tokens with the aggregators. The GLIDER method subsequently leverages this information to perform inference on arbitrary queries.

Local Router. Before each input module m , a separate local router weight $L_m \in \mathbb{R}^{N \times d}$ is inserted to make local per-token, per-module routing decisions. For a given module m and expert model c , we have $\bar{v}_m^c = \frac{v_m^c - \mu(v_m^c)}{\sigma(v_m^c)}$, where $\mu(\cdot)$ and $\sigma(\cdot)$ denote the mean and standard deviation respectively. Next, we obtain the local router for module m by stacking these standardised local routing vectors as $L_m = [\bar{v}_m^1; \dots; \bar{v}_m^N] \in \mathbb{R}^{N \times d}$. Next, for each token i with

activation u_i coming into module m , we standardise it to obtain $\bar{u}_i = \frac{u_i - \mu(u_i)}{\sigma(u_i)}$. We then compute the local affinity scores, $s_m^{loc} \in \mathbb{R}^N$ between the local router L_m and u_i as $s_m^{loc} = \text{cos-sim}(L_m, u_i)$.

Global Router. The global router aims to capture task semantics to retrieve relevant experts for any given input query. We create the global router weight $G \in \mathbb{R}^{N \times d_g}$ by stacking the global routing vectors from all the expert models as $G = [g^1; \dots; g^N]$. This router is not a part of the base model and is added before the model to independently process the full query. Given an input query u along with three few-shot input-output pairs of similar queries, we prompt an LLM (gpt-4-turbo) using the template provided in Appendix B to obtain a task description for the query. We then embed this task description using the same embedding model (nomic-embed-text-v1.5) to obtain the vector $q_u \in \mathbb{R}^{d_g}$. We then compute the global affinity score, $s^{glob} \in \mathbb{R}^N$, by computing the cosine similarity as $s^{glob} = \text{cos-sim}(G, q_u)$.

Combining Global and Local Router. At each module m , we have the global and local affinity scores s^{glob} and s_m^{loc} respectively. Following Phatgoose (Muqeeth et al., 2024), we scale the local scores with a factor of $1/\sqrt{N}$. However, the global router’s main goal is to retrieve the correct expert for the held-in tasks. Therefore, we first check if the expert with the highest global affinity score ($\max(s^{glob})$) is above a threshold (p). If such experts exist, then we set a high α to enforce retrieval and vice versa. Hence, we propose to scale the global scores with α , where $\alpha = \gamma \cdot \mathbb{I}_{\{\max(s^{glob}) - p > 0\}} + \beta$, where p is the cosine similarity threshold, and γ and β are scaling hyperparameters. Using our ablation experiments in Section 5.4, we set $p = 0.8$, $\gamma = 100$ and $\beta = 3$. We then obtain the final affinity score $s \in \mathbb{R}^N = \alpha \cdot s^{glob} + s_m^{loc}/\sqrt{N}$. Then GLIDER selects the top- k experts after performing softmax over the final affinity score s as $\mathcal{E}_{top} = \text{top-}k(\text{softmax}(s))$. Finally, the output of the module for token activation u_i is computed as $W u_i + \sum_{k \in \mathcal{E}_{top}} s_k \cdot B_k A_k u_i$.

5 Experiments

5.1 Setting

Dataset. Our experiments utilize the multitask prompted training setup (**T0-HI**) introduced by Sanh et al. (2021), which has become a standard

benchmark for evaluating held-in performance as well as generalization to unseen tasks (Chung et al., 2022; Longpre et al., 2023; Jang et al., 2023; Zhou et al., 2022). Phatgoose (Muqeeth et al., 2024) shows how local routing can be used for generalization to unseen domain, hence, following them, we employ LM-adapted T5.1.1 XL (Lester et al., 2021) as our base model which is a 3B parameter variant of T5 (Raffel et al., 2020) further trained on the C4 dataset using a standard language modeling objective. For held-out evaluations, we follow Phatgoose (Muqeeth et al., 2024) and use three held-out benchmark collections. We use the T0 held-out (**T0-HO**) datasets used in Sanh et al. (2021) and the two subsets of BIG-bench (**BIG-bench authors**, 2023). Specifically, we use BIG-bench Hard (**BBH**) (Suzgun et al., 2022), consisting of 23 challenging datasets, and BIG-bench Lite (**BBL**) (**BIG-bench authors**, 2023), a lightweight 24-dataset proxy for the full benchmark. Similar to Muqeeth et al. (2024), we exclude certain BIG-bench datasets due to tokenization incompatibility with the T5 tokenizer.

Expert Creation. To create the pool of expert module for routing, we follow Muqeeth et al. (2024) and use two distinct dataset collections: ① T0 Held-In (Sanh et al., 2021) consisting of the 36 held-in prompted datasets for tasks from the T0 training procedure. ② The “FLAN Collection” (Longpre et al., 2023) which significantly expands the T0 tasks by incorporating prompted datasets from SuperGLUE (Wang et al., 2019a), Super Natural Instructions (Wang et al., 2022b), dialogue datasets, and Chain-of-Thought datasets (Wei et al., 2022b). Following Muqeeth et al. (2024), we create 166 specialized models from the FLAN Collection. For each dataset in these collections, we train Low-Rank Adapters (LoRAs) (Hu et al., 2021) modules resulting in pools of 36 and 166 expert models for T0 Held-In and FLAN, respectively. Similar to Phatgoose, we use a rank of $r = 16$ and train for 1000 steps using the AdamW optimizer (Loshchilov and Hutter, 2017) with a learning rate of 5×10^{-3} and a warmup ratio of 0.06. After training the LoRA module, we freeze it and train the local routing vectors for an additional 100 steps with the same hyperparameters. Finally, following prior work (Shazeer et al., 2016; Du et al., 2022; Lepikhin et al., 2020), GLIDER performs top- k routing with $k = 2$.

Table 1: Performance evaluated on the T0 set and FLAN set. We present the performance on both held-in tasks (*i.e.* T0-HI) and held-out tasks (*i.e.* T0-HO, BBH, and BBL). We compare the following methods: (1) performance upper bound, *i.e.* Oracle Expert; (2) zero-shot baselines, *i.e.* Multi-Task Fine-Tuning, Expert Merging, Arrow, and Phatgoose; (3) few-shot baselines, *i.e.* LoRA Hub and GLIDER. We mark the best performance besides the upper bound (*i.e.*, Oracle Expert) in **bold**.

Method	T0				FLAN	
	T0-HI	T0-HO	BBH	BBL	BBH	BBL
Oracle Expert	69.60	51.60	34.90	36.60	38.90	45.40
Multi-Task Fine-Tuning	55.90	51.60	34.90	36.60	38.90	45.40
Expert Merging	30.73	45.40	35.30	36.00	34.60	34.00
Arrow	39.84	55.10	33.60	34.50	30.60	29.60
Phatgoose	61.42	56.90	34.90	37.30	35.60	35.20
LORA Hub	31.90	46.85	31.35	31.18	34.50	30.54
GLIDER	68.04	57.78	35.29	37.46	35.07	35.52

508 5.2 Baselines

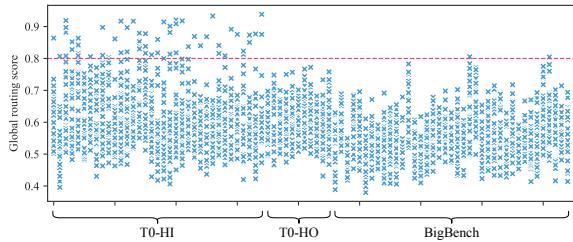
509 **Expert Merging.** Model Merging (Yadav et al.,
510 2023b; Choshen et al., 2022) involves averaging
511 the parameters of multiple models or modules to
512 create a single aggregate model. We merge by
513 multiplying the LoRA matrices and then taking
514 an unweighted average of all the experts within
515 the pool. It is important to note that this merging
516 strategy requires homogeneous expert module ar-
517 chitectures; in contrast, GLIDER can accomodate
518 heterogeneous expert modules.

519 **Arrow.** Following Ostapenko et al. (2024), we
520 employ a routing mechanism where gating vectors
521 are derived from LoRA expert modules. Specif-
522 ically, the first right singular vector of the outer
523 product of each module’s LoRA update (BA) serves as
524 its gating vector. Input routing is determined by a
525 probability distribution based on the absolute dot
526 product between the input representation and each
527 gating vector. We utilize top- k routing with $k = 2$.

528 **Phatgoose.** Phatgoose (Muqeeth et al., 2024)
529 first learn the LoRA modules for each, followed by
530 learning a sigmoid gating vector similar to our local
531 router. During inference, they make routing deci-
532 sions for each token independently for all modules.
533 Specifically, they first standardize the input token
534 activations and gating vectors from all experts and
535 then perform similarity-based top-2 routing.

536 **LoRA Hub.** LoraHub (Huang et al., 2023)
537 method performs gradient-free optimization us-
538 ing few-shot task samples to learn mixing coef-
539 ficients for different expert models while keeping
540 them fixed. Once the coefficients are learned, they
541 merge the experts with the learned weight and route
542 through the merged expert.

543 **Multi-task Fine-Tuning.** Multitask training is a



544 Figure 3: Global routing scores for tasks in the T0 set.
545 The red horizontal line indicates our design threshold
546 of 0.8. Each column represents an evaluated task from
547 T0-HI, T0-HO, BigBench using T0 held-in experts. All
548 global routing scores for each task are plotted, corre-
549 sponding to the 35 experts in total.

550 proven method for enhancing zero-shot generalization
551 (Sanh et al., 2021; Wei et al., 2022a) but is in-
552 feasible given our problem setting and data access
553 limitations. We include it as a baseline using pub-
554 licly available models. Specifically, we utilize the
555 T0-3B model (Sanh et al., 2021) for the T0 Held-In
556 datasets, given its training on a matching dataset
557 collection. For FLAN, a directly comparable pub-
558 licly available model is unavailable; therefore, we
559 report FLAN-T5 XL results trained on a different,
560 undisclosed dataset mixture, while acknowledging
561 the limitations of this indirect comparison.

562 **Oracle.** Following (Jang et al., 2023) and
563 (Muqeeth et al., 2024), we employ an Oracle rout-
564 ing scheme as a performance upper bound. This
565 scheme selects the expert exhibiting optimal per-
566 formance on a given evaluation dataset, thus repre-
567 senting a non-zero-shot approach.

568 5.3 Main Results

569 Table 1 presents the comparison results among our
570 GLIDER and six baselines on both held-in and held-
571 out settings. We report the average performance
572 across all tasks for each setting, please see Ap-
573 pendix D for each tasks metric. To further illus-
574 trate the performance, we also include the results of
575 Oracle Expert, which has extra access to the task
576 identities of expert modules and evaluated datasets
577 and can be regarded as an *upper bound*.

578 **T0 Setting.** In the T0 task set, the following ob-
579 servations can be drawn: ① For the held-in tasks,
580 *i.e.* T0-HI, GLIDER significantly outperforms other
581 baselines and almost matches the performance of
582 Oracle Expert upper bound. ② For T0-HO and
583 BBL tasks, GLIDER achieves the best performance
584 among all the methods, including Oracle Expert up-
585 per bound. ③ GLIDER has negligible lower perfor-
586 mance, *i.e.* 0.01%, compared to the Expert Merg-
587 ing baseline in BBH but outperforms it by around

Table 2: Ablation on the instruction coefficient α . We mark the best performance in **bold** and the performance corresponding to the selected α by GLIDER in blue.

α	T0			
	T0-HI	T0-HO	BBH	BBL
0	61.42	56.90	34.90	37.30
1	62.20	57.04	35.05	37.79
3	63.40	57.78	35.29	37.46
10	65.52	57.98	34.80	37.04
100	68.04	53.22	31.73	34.97
1000	66.88	52.91	30.71	34.31
3000	66.69	52.37	30.03	33.24

12% on T0-HO and 1.5% on BBL. Besides Expert Merging, GLIDER outperforms all other methods on BBH, including the Oracle Expert upper bound.

5.4 Ablation Study and Further Investigation

Ablation on the global routing scale α . To illustrate how the specialization and generalization abilities change as we scale the coefficient α of the global routing score, we conduct the ablation study of α ranging $\{1, 3, 10, 100, 1000, 3000\}$. As shown in Table 2, we present experimental results of the T0 task set on both held-in and held-out tasks. For held-in tasks, *i.e.* T0-HI, GLIDER can select the optimal α to scale the global routing score. For held-out tasks, *i.e.* {T0-HO, BBH, BBL}, GLIDER produce either the optimal α (for BBH) or the sub-optimal α with slightly lower performance to the optimal ones (for T0-HO and BBL). Lastly, note that Phatgoose correspond to the setting where there is no global semantics used, *i.e.*, $\alpha = 0$.

Ablation on the routing strategy. There exists a trade-off between performance and efficiency when using different top-k routing strategies (Rachamandran and Le, 2019). To investigate the impact of routing strategy in GLIDER, we evaluate top-k routing of k in $\{1, 2, 3\}$. Moreover, we further evaluate the top-p routing (Huang et al., 2024c; Zeng et al., 2024) of p in $\{25\%, 50\%, 75\%\}$, where each token selects experts with higher routing probabilities until the cumulative probability exceeds threshold p . As shown in Table 3, we can draw the following conclusions: (1) For top-k routing, $k = 2$ shows comparable or better performance than $k = 3$, particularly for T0-HO and BBH, while offering improved efficiency. (2) For top-p routing, higher p values consistently yield better performance at the cost of efficiency. Therefore, we use top-2 routing in GLIDER by default.

Investigation on the threshold design of global scores. As in Section 4, we compute the scale

Table 3: Ablation on the routing strategy. GLIDER employs top-2 routing. We mark the best performance among top-k and top-p routing in **bold**, respectively.

Method	T0			
	T0-HI	T0-HO	BBH	BBL
Top-1	67.96	56.07	33.91	35.82
Top-2	68.04	57.78	35.39	37.46
Top-3	68.06	57.52	35.08	38.55
Top-25%	67.98	56.53	34.10	36.32
Top-50%	67.95	57.25	35.07	37.49
Top-75%	68.02	57.86	35.38	38.65

α for global scores using the formula $\alpha = \gamma * \mathbb{I}_{\{\max(\text{sglob}) - 0.8 > 0\}} + \beta$, where we establish a threshold of 0.8 to differentiate evaluated tasks. Figure 3 presents the global routing scores for each task in the T0 set to motivate the rationale behind this design. For all held-in tasks (*i.e.*, T0-HI), at least one expert (typically the oracle expert trained on the evaluated task) achieves global routing scores exceeding 0.8. Consequently, GLIDER applies a higher $\alpha = 100$, enabling effective identification of tasks corresponding to a specifically trained expert and enhancing retrieval of this oracle expert. For nearly all held-out tasks (*i.e.*, T0-HO and BigBench), no global routing score surpasses 0.8, prompting GLIDER to utilize a lower $\alpha = 3$. Two exceptions among the held-out tasks are bbq_lite_json and strange_stories in BigBench, where one score marginally exceeds 0.8 in each case. For these two, GLIDER employs the higher $\alpha = 100$, resulting in performance improvements of 1.3% and 2.9% respectively over $\alpha = 3$, thus showing the effectiveness of our design.

6 Conclusion

This paper introduces GLIDER, a novel multi-scale routing mechanism that incorporates both global semantic and local token-level routers. By leveraging the semantic reasoning capabilities of LLMs for global expert selection and refining these choices with a learned local router, GLIDER addresses the limitations of existing methods that often perform poorly on held-in tasks. Our empirical evaluation on T0 and FLAN benchmarks, using T5-based experts, demonstrates that GLIDER achieves substantial improvements in held-in task performance while maintaining competitive generalization on held-out tasks. These findings suggest that incorporating global semantic task context into routing mechanisms is crucial for building robust and practically useful routing-based systems.

660 7 Limitation

661 The main limitation of GLIDER lies in its heavy
662 dependence on large language models (specifically
663 GPT-4) for generating semantic task descriptions.
664 This reliance introduces potential accessibility bar-
665 riers due to API costs. Furthermore, investigating
666 the application of GLIDER to other modalities be-
667 yond language tasks, such as vision or multi-modal
668 expert models, could unlock new capabilities for
669 specialized model routing.

670 References

671 [Wikiquote, russian proverbs.](#)

672 Abubakar Abid, Maheen Farooqi, and James Zou. 2021.
673 [Persistent anti-Muslim bias in large language models.](#)
674 *arXiv preprint*.

675 Alessandro Achille, Michael Lam, Rahul Tewari,
676 Avinash Ravichandran, Subhransu Maji, Charless C
677 Fowlkes, Stefano Soatto, and Pietro Perona. 2019.
678 Task2vec: Task embedding for meta-learning. In
679 *Proceedings of the IEEE/CVF international conference*
680 *on computer vision*, pages 6430–6439.

681 Joshua Ackerman and George Cybenko. 2020. [A sur-
682vey of neural networks and formal languages.](#) *arXiv
683 preprint*.

684 Aishwarya Agrawal, Jiasen Lu, Stanislaw Antol, Mar-
685 garet Mitchell, C. Lawrence Zitnick, Dhruv Batra,
686 and Devi Parikh. 2015. [VQA: Visual question an-
687 swering.](#) *arXiv preprint*.

688 Akshay Agrawal, Shane Barratt, and Stephen Boyd.
689 2020. [Learning convex optimization models.](#) *arXiv
690 preprint*.

691 Scott Alexander. 2020. [A very unlikely chess game.](#)

692 Mohammad Aliannejadi, Hamed Zamani, Fabio
693 Crestani, and W. Bruce Croft. 2019. [Asking clar-
694 ifying questions in open-domain information-seeking
695 conversations.](#) In *Proceedings of the 42nd Interna-
696 tional ACM SIGIR Conference on Research and De-
697 velopment in Information Retrieval*, New York, NY,
698 USA. Association for Computing Machinery.

699 Miltiadis Allamanis, Earl T. Barr, Premkumar Devanbu,
700 and Charles Sutton. 2018. [A survey of machine learn-
701 ing for big code and naturalness.](#) *ACM Comput. Surv.*,
702 51(4).

703 Miltiadis Allamanis, Pankajan Chanthirasegaran, Push-
704 meet Kohli, and Charles Sutton. 2016. [Learning
705 continuous semantic representations of symbolic ex-
706 pressions.](#) *arXiv preprint*.

707 Uri Alon, Shaked Brody, Omer Levy, and Eran Yahav.
708 2018. [code2seq: Generating sequences from struc-
709 tured representations of code.](#) *arXiv preprint*.

710 Uri Alon, Roy Sadaka, Omer Levy, and Eran Yahav.
711 2020. [Structural language models of code.](#) In *Pro-
712 ceedings of the 37th International Conference on
713 Machine Learning*, volume 119 of *Proceedings of
714 Machine Learning Research*, pages 245–256. PMLR.

715 Miriam Amin and Manuel Burghardt. 2020. [A survey
716 on approaches to computational humor generation.](#)
717 In *Proceedings of the 4th Joint SIGHUM Workshop
718 on Computational Linguistics for Cultural Heritage,
719 Social Sciences, Humanities and Literature*, pages
720 29–41, Online. International Committee on Compu-
721 tational Linguistics.

722 Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik
723 Koncel-Kedziorski, Yejin Choi, and Hannaneh Ha-
724 jishirzi. 2019. [MathQA: Towards interpretable math
725 word problem solving with operation-based for-
726 malisms.](#) In *Proceedings of the 2019 Conference
727 of the North American Chapter of the Association for
728 Computational Linguistics: Human Language Tech-
729 nologies, Volume 1 (Long and Short Papers)*, pages
730 2357–2367, Minneapolis, Minnesota. Association for
731 Computational Linguistics.

732 Prithviraj Ammanabrolu, William Broniec, Alex
733 Mueller, Jeremy Paul, and Mark O. Riedl. 2019. [To-
734 ward automated quest generation in text-adventure
735 games.](#) *arXiv preprint*.

736 Prithviraj Ammanabrolu, Wesley Cheung, Dan Tu,
737 William Broniec, and Mark O. Riedl. 2020. [Bringing
738 stories alive: Generating interactive fiction worlds.](#)
739 *arXiv preprint*.

740 Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Chris-
741 tiano, John Schulman, and Dan Mané. 2016. [Con-
742 crete problems in AI safety.](#) *arXiv preprint*.

743 Brandon Amos and J. Zico Kolter. 2017. [Optnet: Dif-
744 ferentiable optimization as a layer in neural networks.](#)
745 *arXiv preprint*.

746 Issa Annamoradnejad and Gohar Zoghi. 2020. [Col-
747 BER: Using BERT sentence embedding for humor
748 detection.](#) *arXiv preprint*.

749 Mikel Artetxe, Sebastian Ruder, and Dani Yogatama.
750 2020. [On the cross-lingual transferability of mono-
751 lingual representations.](#) In *Proceedings of the 58th
752 Annual Meeting of the Association for Computational
753 Linguistics*, pages 4623–4637, Online. Association
754 for Computational Linguistics.

755 Shima Asaadi, Saif Mohammad, and Svetlana Kir-
756 itchenko. 2019. [Big BiRD: A large, fine-grained,
757 bigram relatedness dataset for examining semantic
758 composition.](#) In *Proceedings of the 2019 Conference
759 of the North American Chapter of the Association for
760 Computational Linguistics: Human Language Tech-
761 nologies, Volume 1 (Long and Short Papers)*, pages
762 505–516, Minneapolis, Minnesota. Association for
763 Computational Linguistics.

764 Pepa Atanasova, Jakob Grue Simonsen, Christina Li-
765 oma, and Isabelle Augenstein. 2020. [Generating fact](#)

766	checking explanations.	In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7352–7364, Online. Association for Computational Linguistics.	819
767			820
768			821
769			
770	Salvatore Attardo.	2017. Humor in language . In <i>Oxford Research Encyclopedia of Linguistics</i> . Oxford University Press.	822
771			823
772			824
773	Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives.	2007. Dbpedia: A nucleus for a web of open data . In <i>The Semantic Web</i> .	825
774			
775			
776			
777	R H. Baayen, R Piepenbrock, and L Gulikers.	1995. Celex2 ldc96l14 .	826
778			827
779	Yasaman Bahri, Ethan Dyer, Jared Kaplan, Jaehoon Lee, and Utkarsh Sharma.	2021. Explaining neural scaling laws .	828
780			829
781			830
782	Anton Bakhtin, Sam Gross, Myle Ott, Yuntian Deng, Marc’Aurelio Ranzato, and Arthur Szlam.	2019. Real or fake? Learning to discriminate machine from human generated text . <i>arXiv preprint</i> .	831
783			832
784			
785			
786	Matej Balog, Alexander L. Gaunt, Marc Brockschmidt, Sebastian Nowozin, and Daniel Tarlow.	2016. Deepcoder: Learning to write programs . <i>arXiv preprint</i> .	833
787			834
788			835
789	Satanjeev Banerjee and Ted Pedersen.	2003. Extended gloss overlaps as a measure of semantic relatedness . In <i>IJCAI’03: Proceedings of the 18th International Joint Conference on Artificial Intelligence</i> , page 805–810, San Francisco. Morgan Kaufmann.	836
790			837
791			838
792			839
793			
794	Roy Bar-Haim, Ido Dagan, Bill Dolan, Lisa Ferro, Danilo Giampiccolo, Bernardo Magnini, and Idan Szpektor.	2006. The second pascal recognising textual entailment challenge. In <i>Proceedings of the second PASCAL challenges workshop on recognising textual entailment</i> , volume 6, pages 6–4. Venice.	840
795			841
796			842
797			843
798			844
799			845
800	Oren Barkan and Noam Koenigstein.	2016. ITEM2VEC: Neural item embedding for collaborative filtering . In <i>2016 IEEE 26th International Workshop on Machine Learning for Signal Processing (MLSP)</i> , pages 1–6, Piscataway, NJ. Institute of Electrical and Electronics Engineers.	846
801			847
802			848
803			
804			
805			
806	Solon Barocas and Andrew D. Selbst.	2016. Big data’s disparate impact . <i>California Law Review</i> , 104(3):671–732.	849
807			850
808			
809	Max Bartolo, Alastair Roberts, Johannes Welbl, Sebastian Riedel, and Pontus Stenetorp.	2020. Beat the ai: Investigating adversarial human annotation for reading comprehension. <i>Transactions of the Association for Computational Linguistics</i> , 8:662–678.	851
810			852
811			853
812			854
813			855
814	Sumit Basu and Janara Christensen.	2013. Teaching classification boundaries to humans . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 27, pages 109–115, Menlo Park, CA. Association for the Advancement of Artificial Intelligence.	856
815			857
816			
817			
818			
819	Jason Baumgartner, Savvas Zannettou, Brian Keegan, Megan Squire, and Jeremy Blackburn.	2020. The Pushshift Reddit dataset . <i>arXiv preprint</i> .	861
820			862
821			863
822	Nihat Bayat and Gökhan Çetinkaya.	2020. The relationship between inference skills and reading comprehension . <i>TED EĞİTİM VE BİLİM (Education and Science)</i> , 45(203):177–190.	864
823			865
824			866
825			
826	Mayur J. Bency, Ahmed H. Qureshi, and Michael C. Yip.	2019. Neural path planning: Fixed time, near-optimal path generation via oracle imitation . In <i>2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)</i> , pages 3965–3972, Piscataway, NJ. Institute of Electrical and Electronics Engineers.	867
827			868
828			869
829			870
830			871
831			872
832			
833	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 <i>Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, FAccT ’21</i> , page 610–623, New York, NY, USA. Association for Computing Machinery.	873
834			874
835			875
836			876
837			877
838			878
839			
840	Emily M. Bender and Alexander Koller.	2020. Climbing towards NLU: On meaning, form, and understanding in the age of data . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5185–5198, Online. Association for Computational Linguistics.	879
841			880
842			881
843			882
844			883
845			
846	Luisa Bentivogli, Peter Clark, Ido Dagan, and Danilo Giampiccolo.	2009. The fifth pascal recognizing textual entailment challenge. In <i>TAC</i> .	884
847			885
848			886
849	Jean Berko.	1958. The child’s learning of english morphology . < <i>i</i> >WORD</ <i>i</i> >, 14(2-3):150–177.	887
850			
851	Tarek R. Besold, Artur d’Avila Garcez, Sebastian Bader, Howard Bowman, Pedro Domingos, Pascal Hitzler, Kai-Uwe Kuehnberger, Luis C. Lamb, Daniel Lowd, Priscila Machado Vieira Lima, Leo de Penning, Gadi Pinkas, Hoifung Poon, and Gerson Zaverucha.	2017. Neural-symbolic learning and reasoning: A survey and interpretation . <i>arXiv preprint</i> .	888
852			889
853			890
854			891
855			892
856			893
857			
858	Gregor Betz, Christian Voigt, and Kyle Richardson.	2020. Critical thinking for language models . <i>arXiv preprint</i> .	894
859			895
860			896
861	Chandra Bhagavatula, Ronan Le Bras, Chaitanya Malaviya, Keisuke Sakaguchi, Ari Holtzman, Hannah Rashkin, Doug Downey, Scott Wen-tau Yih, and Yejin Choi.	2019. Abductive commonsense reasoning . <i>arXiv preprint</i> .	897
862			898
863			899
864			900
865			
866	Sumithra Bhakthavatsalam, Daniel Khashabi, Tushar Khot, Bhavana Dalvi Mishra, Kyle Richardson, Ashish Sabharwal, Carissa Schoenick, Oyvind Tafjord, and Peter Clark.	2021. Think you have solved direct-answer question answering? Try ARCD-DA, the direct-answer AI2 reasoning challenge . <i>arXiv preprint</i> .	901
867			902
868			903
869			904
870			905
871			906
872			

873	Satwik Bhattacharya, Kabir Ahuja, and Navin Goyal.	929
874	2020a. On the ability and limitations of transformers to recognize formal languages. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 7096–7116, Online. Association for Computational Linguistics.	930
875		931
876		932
877		933
878		
879	Satwik Bhattacharya, Kabir Ahuja, and Navin Goyal.	934
880	2020b. On the practical ability of recurrent neural networks to recognize hierarchical languages. In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 1481–1494, Barcelona, Spain (Online). International Committee on Computational Linguistics.	935
881		936
882		937
883		938
884		939
885		
886	Surya Bhupatiraju, Rishabh Singh, Abdel-rahman Mohamed, and Pushmeet Kohli. 2017. Deep API programmer: Learning to program with APIs. <i>arXiv preprint</i> .	940
887		941
888		942
889		943
890	Alan W. Biermann. 1978. The inference of regular LISP programs from examples. <i>IEEE Transactions on Systems, Man, and Cybernetics</i> , 8(8):585–600.	944
891		945
892		946
893	BIG-bench authors. 2023. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. <i>Transactions on Machine Learning Research</i> .	947
894		948
895		949
896		950
897	Joachim Bingel and Anders Søgaard. 2017. Identifying beneficial task relations for multi-task learning in deep neural networks. <i>arXiv preprint arXiv:1702.08303</i> .	951
898		952
899		
900		
901	Julia Birke and Anoop Sarkar. 2006. A clustering approach for nearly unsupervised recognition of nonliteral language. In <i>11th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 329–336, Trento, Italy. Association for Computational Linguistics.	953
902		954
903		
904		
905		
906		
907	Sebastian Bischoff, Niklas Deckers, Marcel Schliebs, Ben Thies, Matthias Hagen, Efstathios Stamatatos, Benno Stein, and Martin Potthast. 2020. The importance of suppressing domain style in authorship analysis. <i>arXiv preprint</i> .	955
908		956
909		957
910		958
911		
912	Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. 2020. PIQA: reasoning about physical commonsense in natural language. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 34, pages 7423–7439, Menlo Park, CA. Association for the Advancement of Artificial Intelligence.	959
913		960
914		961
915		962
916		963
917		964
918		
919	Yuri Bizzoni and Shalom Lappin. 2018. Predicting human metaphor paraphrase judgments with deep neural networks. In <i>Proceedings of the Workshop on Figurative Language Processing</i> , pages 45–55, New Orleans, Louisiana. Association for Computational Linguistics.	965
920		966
921		967
922		
923		
924		
925	Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Pieler, USVSN Sai Prashanth, Shivanshu Purohit,	968
926		969
927		970
928		971
929	Laria Reynolds, Jonathan Tow, Ben Wang, and Samuel Weinbach. 2022. GPT-NeoX-20B: An open-source autoregressive language model. In <i>Proceedings of the ACL Workshop on Challenges & Perspectives in Creating Large Language Models</i> .	972
930		973
931		974
932		975
933		
934	Vladislav Blinov, Valeria Bolotova-Baranova, and Pavel Braslavski. 2019. Large dataset and language model fun-tuning for humor recognition. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4027–4032, Florence, Italy. Association for Computational Linguistics.	976
935		977
936		978
937		979
938		980
939		981
940		982

983	Samuel R. Bowman and George Dahl. 2021. What will it take to fix benchmarking in natural language understanding? In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 4843–4855, Online. Association for Computational Linguistics.	1040
984		1041
985		1042
986		1043
987		
988		
989		
990	Bozhidar Bozhanov and Ivan Derzhanski. 2013. Rosetta stone linguistic problems. In <i>Proceedings of the Fourth Workshop on Teaching NLP and CL</i> , pages 1–8, Sofia, Bulgaria. Association for Computational Linguistics.	1044
991		1045
992		1046
993		
994		
995	Matko Bošnjak, Tim Rocktäschel, Jason Naradowsky, and Sebastian Riedel. 2017. Programming with a differentiable Forth interpreter. In <i>Proceedings of the 34th International Conference on Machine Learning</i> , volume 70, pages 547–556.	1047
996		1048
997		1049
998		1050
999		1051
1000	Gwern Branwen. 2020. GPT-3 creative fiction. <i>Gwern.net</i> .	1052
1001		1053
1002	Luke Breitfeller, Emily Ahn, David Jurgens, and Yulia Tsvetkov. 2019. Finding microaggressions in the wild: A case for locating elusive phenomena in social media posts. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 1664–1674, Hong Kong, China. Association for Computational Linguistics.	1054
1003		
1004		
1005		
1006		
1007		
1008		
1009		
1010		
1011	Ralf Brown. 2014. Non-linear mapping for improved identification of 1300+ languages. In <i>Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 627–632, Doha, Qatar. Association for Computational Linguistics.	1055
1012		1056
1013		1057
1014		
1015		
1016		
1017	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. 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, Chris 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. In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 1877–1901. Curran Associates, Inc.	1058
1018		1059
1019		1060
1020		1061
1021		
1022		
1023		
1024		
1025		
1026		
1027		
1028		
1029		
1030		
1031	Thomas Bugnyar, Stephan A. Reber, and Cameron Buckner. 2016. Ravens attribute visual access to unseen competitors. <i>Nature Communications</i> , 7:article 10506.	1066
1032		1067
1033		1068
1034		1069
1035	Franck Burlot, Yves Scherrer, Vinit Ravishankar, Ondřej Bojar, Stig-Arne Grönroos, Maarit Koponen, Tommi Nieminen, and François Yvon. 2018. The WMT’18 morpheme test suites for English-Czech, English-German, English-Finnish	1070
1036		1071
1037		
1038		
1039		
	and Turkish-English. In <i>Proceedings of the Third Conference on Machine Translation: Shared Task Papers</i> , pages 546–560, Belgium, Brussels. Association for Computational Linguistics.	1072
		1073
		1074
	Corrado Böhm. 1964. On a family of Turing machines and the related programming language. <i>ICC Bulletin</i> , 3:187–194.	1075
		1076
		1077
	Lucas Caccia, Edoardo Ponti, Zhan Su, Matheus Pereira, Nicolas Le Roux, and Alessandro Sordoni. 2023. Multi-head adapter routing for cross-task generalization. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	1078
		1079
	Kate Cain and Jane V. Oakhill. 1999. Inference making ability and its relation to comprehension failure. <i>Reading and Writing</i> , 11(5–6):489–503.	1080
		1081
	Maya Cakmak and Andrea L. Thomaz. 2014. Eliciting good teaching from humans for machine learners. <i>Artificial Intelligence</i> , 217:198–215.	1082
		1083
	Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. <i>Science</i> , 356(6334):183–186.	1084
		1085
	Josep Call and Michael Tomasello. 2008. Does the chimpanzee have a theory of mind? 30 years later. <i>Trends in Cognitive Sciences</i> , 12:187–192.	1086
		1087
	Tracy Canfield. 2010. Machine translation of Klingon.	1088
		1089
	Nathanael Chambers. 2012. Labeling documents with timestamps: Learning from their time expressions. In <i>Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 98–106, Jeju Island, Korea. Association for Computational Linguistics.	1090
		1091
	Sharath Chandra Guntuku, Mingyang Li, Louis Tay, and Lyle H. Ungar. 2019. Studying cultural differences in emoji usage across the East and the West. In <i>Proceedings of the International AAAI Conference on Web and Social Media</i> , volume 13, pages 226–235, Menlo Park, CA. Association for the Advancement of Artificial Intelligence.	1092
		1093
	Nick Chater and Paul Vitányi. 2003. Simplicity: A unifying principle in cognitive science? <i>Trends in Cognitive Sciences</i> , 7:19–22.	1094
		1095
	Antonio Chella, Arianna Pipitone, Alain Morin, and Famira Racy. 2020. Developing self-awareness in robots via inner speech. <i>Frontiers in Robotics and AI</i> , 7.	1096
		1097
	Howard Chen, Alane Suhr, Dipendra Misra, Noah Snavely, and Yoav Artzi. 2019a. Touchdown: Natural language navigation and spatial reasoning in visual street environments. In <i>2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 12530–12539, Piscataway, NJ. Institute of Electrical and Electronics Engineers.	1098
		1099
		1100

1093	Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. 2020. Generative pretraining from pixels . In <i>Proceedings of the 37th International Conference on Machine Learning</i> , volume 119 of <i>Proceedings of Machine Learning Research</i> , pages 1691–1703. PMLR.	1146
1094		1147
1095		1148
1096		1149
1097		1150
1098		
1099	Peng-Yu Chen and Von-Wun Soo. 2018. Humor recognition using deep learning . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)</i> , pages 113–117, New Orleans, Louisiana. Association for Computational Linguistics.	1151
1100		1152
1101		1153
1102		
1103		
1104		
1105		
1106	Ricson Chen. 2020. Transformers play chess .	1154
1107		1155
1108		1156
1109		1157
1110		1158
1111		
1112		
1113		
1114	Feng Cheng, Ziyang Wang, Yi-Lin Sung, Yan-Bo Lin, Mohit Bansal, and Gedas Bertasius. 2024. DAM: Dynamic adapter merging for continual video qa learning. <i>arXiv preprint arXiv:2403.08755</i> .	1159
1115		1160
1116		1161
1117		1162
1118		1163
1119		
1120		
1121		
1122		
1123	Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? Try ARC, the AI2 reasoning challenge . <i>arXiv preprint</i> .	1164
1124		1165
1125		1166
1126		
1127	Peter Clark, Oyvind Tafjord, and Kyle Richardson. 2020. Transformers as soft reasoners over language . <i>arXiv preprint</i> .	1167
1128		1168
1129		1169
1130	Robert J. Clark and Robert R. Jackson. 1994. Self recognition in a jumping spider: Portia labiata females discriminate between their own draglines and those of conspecifics . <i>Ethology Ecology & Evolution</i> , 6(3):371–375.	1170
1131		1171
1132		
1133		
1134		
1135		
1136	Lidia Contreras-Ochando, Cèsar Ferri, José Hernández-Orallo, Fernando Martínez-Plumed, María José Ramírez-Quintana, and Susumu Katayama. 2020. Automated data transformation with inductive programming and dynamic background knowledge . In <i>Machine Learning and Knowledge Discovery in Databases</i> , pages 735–751, Cham. Springer.	1172
1137		1173
1138		1174
1139		1175
1140		1176
1141		1177
1142		1178
1143	Lidia Contreras-Ochando, César Ferri, José Hernández-Orallo, Fernando Martínez-Plumed, Marfa José Ramírez-Quintana, and Susumu Katayama. 2018. General-purpose declarative inductive programming with domain-specific background knowledge for data wrangling automation . <i>arXiv preprint</i> .	1179
1144		1180
1145		1181
1146		1182
1147		1183
1148		1184
1149		
1150		
1151		
1152		
1153		
1154		
1155		
1156		
1157		
1158		
1159		
1160		
1161		
1162		
1163		
1164		
1165		
1166		
1167		
1168		
1169		
1170		
1171		
1172		
1173		
1174		
1175		
1176		
1177		
1178		
1179		
1180		
1181		
1182		
1183		
1184		
1185		
1186		
1187		
1188		
1189		
1190		
1191		
1192		
1193		
1194		
1195		
1196		
1197		
1198		
1199		
1200		

1201	Andrew Cropper and Stephen H. Muggleton. 2016.	Marie-Catherine de Marneffe, Mandy Simons, and Judith Tonhauser. 2019. <i>The CommitmentBank: Investigating projection in naturally occurring discourse</i> . <i>Proceedings of Sinn und Bedeutung</i> , 23(2):107–124.	1254
1202	Metagol system.		1255
1203	Andrew Cropper, Alireza Tamaddoni-Nezhad, and Stephen H. Muggleton. 2016. <i>Meta-interpretive learning of data transformation programs</i> . In <i>Inductive Logic Programming</i> , pages 46–59, Cham.		1256
1204	Springer.		1257
1205			
1206	Joe Cruse. 2015. <i>Emoji usage in TV conversation</i> . <i>Twitter blog</i> .	Scott Deerwester, Susan T. Dumais, George W. Furnas, Thomas K. Landauer, and Richard Harshman. 1990. <i>Indexing by latent semantic analysis</i> . <i>Journal of the American Society for Information Science</i> , 41(6):391–407.	1258
1207			1259
1208		Judith Degen, Robert D. Hawkins, Caroline Graf, Elisa Kreiss, and Noah D. Goodman. 2020. <i>When redundancy is useful: A Bayesian approach to “overinformative” referring expressions</i> . <i>Psychological Review</i> , 127:591–621.	1260
1209			1261
1210	Marc-Alexandre Côté, Ákos Kádár, Xingdi Yuan, Ben Kybartas, Tavian Barnes, Emery Fine, James Moore,		1262
1211	Ruo Yu Tao, Matthew Hausknecht, Layla El Asri, Mahmoud Adada, Wendy Tay, and Adam Trischler.		1263
1212	2018. <i>TextWorld: A learning environment for text-based games</i> . <i>arXiv preprint</i> .		1264
1213			1265
1214	Ido Dagan, Oren Glickman, and Bernardo Magnini. 2005. The pascal recognising textual entailment challenge. In <i>Machine Learning Challenges Workshop</i> , pages 177–190. Springer.		1266
1215			1267
1216	Jim Daley. 2021. <i>White Chicago cops use force more often than Black officers</i> . <i>Scientific American</i> .	Sunipa Dev, Tao Li, Jeff M. Phillips, and Vivek Sriku-	1268
1217		mar. 2020. <i>On measuring and mitigating biased in-</i>	1269
1218		<i>ferences of word embeddings</i> . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , vol-	1270
1219		ume 34, pages 7659–7666, Menlo Park, CA. Associa-	1271
1220		tion for the Advancement of Artificial Intelligence.	1272
1221			1273
1222	Sahith Dambekodi, Spencer Frazier, Prithviraj Am-	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	1274
1223	manabrolu, and Mark O. Riedl. 2020. <i>Playing text-</i>	Kristina Toutanova. 2018. <i>BERT: Pre-training of</i>	1275
1224	<i>based games with common sense</i> . <i>arXiv preprint</i> .	<i>deep bidirectional transformers for language under-</i>	1276
1225		<i>standing</i> . <i>arXiv preprint</i> .	1277
1226	Pradeep Dasigi, Nelson F. Liu, Ana Marasovic, Noah A.	Jacob Devlin, Jonathan Uesato, Surya Bhupatiraju,	1278
1227	Smith, and Matt Gardner. 2019. <i>Quoref: A read-</i>	Rishabh Singh, Abdel-rahman Mohamed, and Push-	1279
1228	<i>ing comprehension dataset with questions requir-</i>	meet Kohli. 2017. <i>RobustFill: Neural program learn-</i>	1280
1229	<i>ing coreferential reasoning</i> . In <i>Proceedings of the</i>	<i>ing under noisy I/O</i> . In <i>Proceedings of the 34th</i>	1281
1230	<i>2019 Conference on Empirical Methods in Natu-</i>	<i>International Conference on Machine Learning</i> , vol-	1282
1231	<i>ral Language Processing and the 9th International</i>	ume 70, pages 990–998, New York, NY, USA. Associa-	1283
1232	<i>Joint Conference on Natural Language Processing</i>	tion for Computing Machinery.	1284
1233	(EMNLP-IJCNLP).		
1234	Wayne Davis. 2019. Implicature. In Edward N. Zalta,	Bhuwan Dhingra, Kathryn Mazaitis, and William W.	1285
1235	editor, <i>The Stanford Encyclopedia of Philosophy</i> , Fall	Cohen. 2017. <i>Quasar: Datasets for question answer-</i>	1286
1236	2019 edition. Metaphysics Research Lab, Stanford	<i>ing by search and reading</i> . <i>arXiv preprint</i> .	1287
1237	University.		
1238	Marie-Catherine de Marneffe, Christopher D. Man-	Kaustubh Dhole, Gurdeep Singh, Priyadarshini P. Pai,	1288
1239	nning, and Christopher Potts. 2012. <i>Did it happen?</i>	and Sukanta Mondal. 2014. <i>Sequence-based predic-</i>	1289
1240	<i>The pragmatic complexity of veridicality assessment</i> .	<i>tion of protein–protein interaction sites with l1-logreg</i>	1290
1241	<i>Computational Linguistics</i> , 38(2):301–333.	<i>classifier</i> . <i>Journal of Theoretical Biology</i> , 348:47–	1291
1242		54.	1292
1243			
1244	Marie-Catherine de Marneffe, Anna N. Rafferty, and	Shizhe Diao, Tianyang Xu, Ruijia Xu, Jiawei Wang,	1293
1245	Christopher D. Manning. 2008. <i>Finding contradic-</i>	and T. Zhang. 2023. <i>Mixture-of-domain-adapters:</i>	1294
1246	<i>tions in text</i> . In <i>Proceedings of ACL-08: HLT</i> , pages	<i>Decoupling and injecting domain knowledge to pre-</i>	1295
1247	1039–1047, Columbus, Ohio. Association for Com-	<i>trained language models’ memories</i> . In <i>Annual Meet-</i>	1296
1248	putational Linguistics	<i>ing of the Association for Computational Linguistics</i> .	1297
1249			
1250	Marie-Catherine de Marneffe, Mandy Simons, and Ju-	Emily Dinan, Angela Fan, Adina Williams, Jack Ur-	1298
1251	dith Tonhauser. 2017. <i>The commitmentbank: Investi-</i>	banek, Douwe Kiela, and Jason Weston. 2019. <i>Queens</i>	1299
1252	<i>gating projection in naturally occurring discourse</i> .	<i>are powerful too: Mitigating gender bias in dia-</i>	1300
1253	In <i>Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers</i> , pages 48–54, Valencia, Spain. Association for Computational Lin-	<i>logue generation</i> . <i>arXiv preprint</i> .	1301
1254	guistics.	William B Dolan and Chris Brockett. 2005. Automati-	1302
1255		cally constructing a corpus of sentential paraphrases.	1303
1256		In <i>Proceedings of the Third International Workshop on Paraphrasing (IWP2005)</i> .	1304
1257			1305
1258		Tiansi Dong, Chengjiang Li, Christian Bauckhage,	1306
1259		Juanzi Li, Stefan Wrobel, and Armin B. Cre-	1307
1260		mers. 2020. <i>Learning syllogism with Euler neural-</i>	1308
1261		<i>networks</i> . <i>arXiv preprint</i> .	1309

1310	Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. 2022. Glam: Efficient scaling of language models with mixture-of-experts. In <i>International Conference on Machine Learning</i> , pages 5547–5569. PMLR.	1367
1311	Avia Efrat, Uri Shaham, Dan Kilman, and Omer Levy. 2021. <i>Cryptonite: A cryptic crossword benchmark for extreme ambiguity in language</i> . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 4186–4192, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	1368
1312		1369
1313		1370
1314		1371
1315		1372
1316		1373
1317	Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. <i>DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs</i> . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 2368–2378, Minneapolis, Minnesota. Association for Computational Linguistics.	1374
1318	Liat Ein Dor, Alon Halfon, Yoav Kantor, Ran Levy, Yosi Mass, Ruty Rinott, Eyal Shnarch, and Noam Slonim. 2018. <i>Semantic relatedness of Wikipedia concepts – benchmark data and a working solution</i> . In <i>Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)</i> , Miyazaki, Japan. European Language Resources Association.	1375
1319		1376
1320		1377
1321		1378
1322		1379
1323		1380
1324		1381
1325		
1326	Liam Dugan, Daphne Ippolito, Arun Kirubarajan, and Chris Callison-Burch. 2020. <i>RoFT: A tool for evaluating human detection of machine-generated text</i> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 189–196, Online. Association for Computational Linguistics.	1382
1327		1383
1328		1384
1329		1385
1330		1386
1331		1387
1332		1388
1333	Ben Eisner, Tim Rocktäschel, Isabelle Augenstein, Matko Bošnjak, and Sebastian Riedel. 2016. <i>emoji2vec: Learning emoji representations from their description</i> . In <i>Proceedings of The Fourth International Workshop on Natural Language Processing for Social Media</i> , pages 48–54, Austin, TX, USA. Association for Computational Linguistics.	1389
1334		1390
1335		1391
1336		1392
1337		
1338	Ahmed El-Kishky, Frank Xu, Aston Zhang, and Jiawei Han. 2019. <i>Parsimonious morpheme segmentation with an application to enriching word embeddings</i> . <i>arXiv preprint</i> .	1393
1339		1394
1340	Ran El-Yaniv and David Yanay. 2013. <i>Semantic sort: A supervised approach to personalized semantic relatedness</i> .	1395
1341		
1342		
1343	Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhilasha Ravichander, Eduard Hovy, Hinrich Schütze, and Yoav Goldberg. 2021. <i>Measuring and improving consistency in pretrained language models</i> . <i>arXiv preprint</i> .	1396
1344		1397
1345		1398
1346		1399
1347		1400
1348		
1349		
1350	Kevin Ellis and Sumit Gulwani. 2017. <i>Learning to learn programs from examples: Going beyond program structure</i> . In <i>Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17</i> , pages 1638–1645.	1401
1351		1402
1352		1403
1353		1404
1354		1405
1355		
1356	Kevin Ellis, Catherine Wong, Maxwell Nye, Mathias Sable-Meyer, Luc Cary, Lucas Morales, Luke Hewitt, Armando Solar-Lezama, and Joshua B. Tenenbaum. 2020. <i>Dreamcoder: Growing generalizable, interpretable knowledge with wake-sleep Bayesian program learning</i> . <i>arXiv preprint</i> .	1406
1357		1407
1358		1408
1359		1409
1360		1410
1361		1411
1362		
1363		
1364		
1365		
1366		
1367	Richard Evans, Jose Hernandez-Orallo, Johannes Welbl, Pushmeet Kohli, and Marek Sergot. 2019. <i>Making sense of sensory input</i> . <i>arXiv preprint</i> .	1412
1368		1413
1369		1414
1370		
1371	Richard Evans, David Saxton, David Amos, Pushmeet Kohli, and Edward Grefenstette. 2018. <i>Can neural networks understand logical entailment?</i> <i>arXiv preprint</i> .	1415
1372		1416
1373		1417
1374		1418
1375		
1376	Matan Eyal, Tal Baumel, and Michael Elhadad. 2019. <i>Question answering as an automatic evaluation metric for news article summarization</i> . In <i>Proceedings</i>	1419
1377		1420
1378		1421

1422	<i>of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 3938–3948, Minneapolis, Minnesota. Association for Computational Linguistics.	Jerry A. Fodor. 1975. <i>The Language of Thought</i> . Harvard University Press, Cambridge, MA.	1477
1423			1478
1424			1479
1425			1480
1426			1481
1427	Alexander R. Fabbri, Irene Li, Tianwei She, Suyi Li, and Dragomir R. Radev. 2019. Multi-news: A large-scale multi-document summarization dataset and abstractive hierarchical model . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> .	Mark Forsyth. 2014. <i>The Elements of Eloquence: Secrets of the Perfect Turn of Phrase</i> . Berkley, New York.	1482
1428			1483
1429			1484
1430			1485
1431			1486
1432			1487
1433	Felix Faltings, Michel Galley, Gerold Hintz, Chris Brockett, Chris Quirk, Jianfeng Gao, and Bill Dolan. 2020. Text editing by command . <i>arXiv preprint</i> .	Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. 2020. Linear mode connectivity and the lottery ticket hypothesis. In <i>International Conference on Machine Learning</i> , pages 3259–3269. PMLR.	1488
1434			1489
1435			1490
1436	Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 889–898, Melbourne, Australia. Association for Computational Linguistics.	Lea Frermann, Shay B. Cohen, and Mirella Lapata. 2018. Whodunnit? Crime drama as a case for natural language understanding . <i>Transactions of the Association for Computational Linguistics</i> , 6:1–15.	1491
1437			1492
1438			1493
1439			1494
1440			1495
1441			1496
1442	Xiaochao Fan, Hongfei Lin, Liang Yang, Yufeng Diao, Chen Shen, Yonghe Chu, and Yanbo Zou. 2020. Humor detection via an internal and external neural network . <i>Neurocomputing</i> , 394:105–111.	Kathryn J. Friedlander and Philip A. Fine. 2018. “The penny drops”: Investigating insight through the medium of cryptic crosswords. <i>Frontiers in Psychology</i> , 9.	1497
1443			1498
1444			1499
1445			1500
1446	William Fedus, Barret Zoph, and Noam Shazeer. 2022. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. <i>Journal of Machine Learning Research</i> , 23(120).	Martins Frolovs. 2019. Teaching GPT-2 transformer a sense of humor: How to fine-tune large transformer models on a single GPU in PyTorch . <i>Towards Data Science, Medium</i> .	1501
1447			1502
1448			1503
1449			1504
1450	Bjarke Felbo, Alan Mislove, Anders Søgaard, Iyad Rahwan, and Sune Lehmann. 2017. Using millions of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and sarcasm . In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> . Association for Computational Linguistics.	Saadie Gabriel, Asli Celikyilmaz, Rahul Jha, Yejin Choi, and Jianfeng Gao. 2020. Go figure: A meta evaluation of factuality in summarization . <i>arXiv preprint</i> .	1505
1451			1506
1452			1507
1453			1508
1454			1509
1455			1510
1456			1511
1457	Christiane Fellbaum, editor. 1998. <i>WordNet: An Electronic Lexical Database</i> . MIT Press, Cambridge, MA.	Evgeniy Gabrilovich and Shaul Markovitch. 2007. Computing semantic relatedness using Wikipedia-based explicit semantic analysis . In <i>IJCAI’07: Proceedings of the 20th International Joint Conference on Artificial Intelligence</i> , page 1606–1611, San Francisco. Morgan Kaufmann.	1512
1458			1513
1459			1514
1460	John K. Feser, Swarat Chaudhuri, and Isil Dillig. 2015. Synthesizing data structure transformations from input-output examples . In <i>Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI ’15</i> , page 229–239, New York, NY, USA. Association for Computing Machinery.	Ge Gao, Eunsol Choi, Yejin Choi, and Luke Zettlemoyer. 2018. Neural metaphor detection in context . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 607–613, Brussels, Belgium. Association for Computational Linguistics.	1515
1461			1516
1462			1517
1463			1518
1464			1519
1465			1520
1466			1521
1467	Susan T. Fiske. 1993. Controlling other people: The impact of power on stereotyping . <i>American Psychologist</i> , 48:621–628.	Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2021. The pile: An 800GB dataset of diverse text for language modeling . <i>arXiv preprint</i> .	1522
1468			1523
1469			1524
1470	Dawn P. Flanagan and Shauna G. Dixon. 2014. The Cattell-Horn-Carroll theory of cognitive abilities . In <i>Encyclopedia of Special Education</i> . John Wiley & Sons, Ltd.	Artur d’Avila Garcez and Luis C. Lamb. 2020. Neurosymbolic AI: The 3rd wave . <i>arXiv preprint</i> .	1525
1471			1526
1472			1527
1473			1528
1474	Pierre Flener and Ute Schmid. 2008. An introduction to inductive programming . <i>Artificial Intelligence Review</i> , 29:45–62.	Matt Gardner, Yoav Artzi, Victoria Basmova, Jonathan Berant, Ben Bogin, Sihao Chen, Pradeep Dasigi, Dheeru Dua, Yanai Elazar, Ananth Gottumukkala, Nitish Gupta, Hanna Hajishirzi, Gabriel Ilharco, Daniel Khashabi, Kevin Lin, Jiangming Liu, Nelson F. Liu, Phoebe Mulcaire, Qiang Ning, Sameer	1529
1475			1530
1476			1531

1531	Singh, Noah A. Smith, Sanjay Subramanian, Reut Tsarfaty, Eric Wallace, Ally Zhang, and Ben Zhou. 2020. Evaluating models' local decision boundaries via contrast sets . <i>arXiv preprint</i> .	Separating retrievability from inferential soundness. <i>Cognitive Psychology</i> , 25(4):524–575.	1590 1591
1532			
1533			
1534			
1535	Albert Gatt, Anja Belz, and Eric Kow. 2009. The TUNA-REG challenge 2009: Overview and evaluation results . In <i>Proceedings of the 12th European Workshop on Natural Language Generation (ENLG 2009)</i> , pages 174–182, Athens, Greece. Association for Computational Linguistics.	Elizabeth Jasmi George and Radhika Mamidi. 2020. Conversational implicatures in English dialogue: Annotated dataset . <i>Procedia Computer Science</i> , 171:2316–2323. Special issue: Third International Conference on Computing and Network Communications (CoCoNet'19).	1592 1593 1594 1595 1596 1597
1536			
1537			
1538			
1539			
1540			
1541	Alexander L. Gaunt, Marc Brockschmidt, Rishabh Singh, Nate Kushman, Pushmeet Kohli, Jonathan Taylor, and Daniel Tarlow. 2016. TerpreT: A probabilistic programming language for program induction . <i>arXiv preprint</i> .	Mor Geva, Ankit Gupta, and Jonathan Berant. 2020a. Injecting numerical reasoning skills into language models . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 946–958, Online. Association for Computational Linguistics.	1598 1599 1600 1601 1602 1603
1542			
1543			
1544			
1545			
1546	Manas Gaur, Amanuel Alambo, Joy Prakash Sain, Ugur Kursuncu, Krishnaprasad Thirunarayan, Ramakanth Kavuluru, Amit Sheth, Randy Welton, and Jyotishman Pathak. 2019. Knowledge-aware assessment of severity of suicide risk for early intervention . In <i>The World Wide Web Conference, WWW '19</i> , page 514–525, New York, NY, USA. Association for Computing Machinery.	Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. Did Aristotle Use a Laptop? A Question Answering Benchmark with Implicit Reasoning Strategies . <i>Transactions of the Association for Computational Linguistics</i> , 9:346–361.	1604 1605 1606 1607 1608 1609
1547			
1548			
1549			
1550			
1551			
1552			
1553			
1554	Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. 2020. RealToxicityPrompts: Evaluating neural toxic degeneration in language models . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 3356–3369, Online. Association for Computational Linguistics.	Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2020b. Transformer feed-forward layers are key-value memories . <i>arXiv preprint</i> .	1610 1611 1612
1555			
1556			
1557			
1558			
1559			
1560			
1561	Sebastian Gehrmann, Tosin Adewumi, Karmanyai Aggarwal, Pawan Sasanka Ammanamanchi, Anuoluwapo Aremu, Antoine Bosselut, Khyathi Raghavi Chandu, Miruna-Adriana Clinciu, Dipanjan Das, Kaustubh Dhole, Wanyu Du, Esin Durmus, Ondřej Dušek, Chris Chinenye Emezue, Varun Gangal, Cristina Barbacea, Tatsunori Hashimoto, Yufang Hou, Yacine Jernite, Harsh Jhamtani, Yangfeng Ji, Shailza Jolly, Mihir Kale, Dhruv Kumar, Faisal Ladhak, Aman Madaan, Mounica Maddela, Khyati Mahajan, Saad Mahamood, Bodhisattwa Prasad Majumder, Pedro Henrique Martins, Angelina McMillan-Major, Simon Mille, Emiel van Miltenburg, Moin Nadeem, Shashi Narayan, Vitaly Nikolaev, Andre Niyongabo Rubungo, Salomey Osei, Ankur Parikh, Laura Perez-Beltrachini, Nirajan Ramesh Rao, Vikas Raunak, Juan Diego Rodriguez, Sashank Santhanam, João Sedoc, Thibault Sellam, Samira Shaikh, Anastasia Shimorina, Marco Antonio Sobrevilla Cabezudo, Hendrik Strobelt, Nishant Subramani, Wei Xu, Diyi Yang, Akhila Yerukola, and Jiawei Zhou. 2021. The GEM benchmark: Natural language generation, its evaluation and metrics . In <i>Proceedings of the 1st Workshop on Natural Language Generation, Evaluation, and Metrics (GEM 2021)</i> , pages 96–120, Online. Association for Computational Linguistics.	Bilal Ghanem, Jihen Karoui, Farah Benamara, Paolo Rosso, and Véronique Moriceau. 2020. Irony detection in a multilingual context . In Joemon M. Jose, Emine Yilmaz, João Magalhães, Pablo Castells, Nicola Ferro, Mário J. Silva, and Flávio Martins, editors, <i>Advances in Information Retrieval. ECIR 2020. Lecture Notes in Computer Science</i> , volume 12036. Springer, Cham.	1613 1614 1615 1616 1617 1618 1619 1620
1562			
1563			
1564			
1565			
1566			
1567			
1568			
1569			
1570			
1571			
1572			
1573			
1574			
1575			
1576			
1577			
1578			
1579			
1580			
1581			
1582			
1583			
1584			
1585			
1586			
1587			
1588	Dedre Gentner, Mary Jo Rattermann, and Kenneth D. Forbus. 1993. The roles of similarity in transfer :	Debanjan Ghosh, Avijit Vajpayee, and Smaranda Muresan. 2020. A report on the 2020 sarcasm detection shared task . <i>arXiv preprint</i> .	1621 1622 1623
1589			
1624	Sayan Ghosh and Shashank Srivastava. 2021. ePiC: Employing proverbs in context as a benchmark for abstract language understanding . <i>CoRR</i> , arXiv:2109.06838.		1624 1625 1626 1627
1625			
1626			
1627			
1628	Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and Bill Dolan. 2007. The third pascal recognizing textual entailment challenge. In <i>Proceedings of the ACL-PASCAL workshop on textual entailment and paraphrasing</i> , pages 1–9. Association for Computational Linguistics.		1628 1629 1630 1631 1632 1633
1629			
1630			
1631			
1632			
1633			
1634	Edward Gibson, Richard Futrell, Julian Jara-Ettinger, Kyle Mahowald, Leon Bergen, Sivalogeswaran Ratnasingam, Mitchell Gibson, Steven T. Piantadosi, and Bevil R. Conway. 2017. Color naming across languages reflects color use . <i>Proceedings of the National Academy of Sciences</i> , 114(40):10785–10790.		1634 1635 1636 1637 1638 1639
1635			
1636			
1637			
1638			
1639			
1640	Matthew L. Ginsberg. 2014. Dr.Fill: Crosswords and an implemented solver for singly weighted CSPs . <i>CoRR</i> , arXiv:1401.4597.		1640 1641 1642
1641			
1642			

1643	Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. 2019. Samsum corpus: A human-annotated dialogue dataset for abstractive summarization . In <i>Proceedings of the 2nd Workshop on New Frontiers in Summarization</i> .	1698
1644		1699
1645		1700
1646		
1647		
1648	Yoav Goldberg. 2019. Assessing BERT’s syntactic abilities . <i>arXiv preprint</i> .	1701
1649		1702
1650	Arthur S. Goldberger. 1972. Structural equation methods in the social sciences . <i>Econometrica</i> , 40(6):979–1001.	1703
1651		1704
1652		1705
1653	Hila Gonen and Yoav Goldberg. 2019. Lipstick on a pig: Debiasing methods cover up systematic gender biases in word embeddings but do not remove them . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 609–614, Minneapolis, Minnesota. Association for Computational Linguistics.	1706
1654		1707
1655		1708
1656		1709
1657		
1658		
1659		
1660		
1661		
1662	Roberto González-Ibáñez, Smaranda Muresan, and Nina Wacholder. 2011. Identifying sarcasm in Twitter: A closer look . In <i>Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies</i> , pages 581–586, Portland, Oregon, USA. Association for Computational Linguistics.	1710
1663		1711
1664		
1665		
1666		
1667		
1668		
1669	Noah D. Goodman and Michael C. Frank. 2016. Pragmatic language interpretation as probabilistic inference . <i>Trends in Cognitive Sciences</i> , 20:818–829.	1712
1670		1713
1671		1714
1672	Karthik Gopalakrishnan, Behnam Hedayatnia, Qinlang Chen, Anna Gottardi, Sanjeev Kwatra, Anu Venkatesh, Raefer Gabriel, and Dilek Hakkani-Tür. 2019. Topical-chat: Towards knowledge-grounded open-domain conversations . In <i>Proc. Interspeech 2019</i> , pages 1891–1895.	1715
1673		1716
1674		1717
1675		
1676		
1677		
1678	Karthik Gopalakrishnan, Behnam Hedayatnia, Longshaokan Wang, Yang Liu, and Dilek Hakkani-Tür. 2020. Are neural open-domain dialog systems robust to speech recognition errors in the dialog history? An empirical study . In <i>Proc. Interspeech 2020</i> , pages 911–915.	1718
1679		1719
1680		1720
1681		1721
1682		1722
1683		1723
1684	Andrew S. Gordon. 2010. Choice of plausible alternatives (COPA) .	1724
1685		1725
1686		1726
1687		
1688		
1689	David Graff, Junbo Kong, Ke Chen, and Kazuaki Maeda. 2003. English gigaword . <i>Linguistic Data Consortium, Philadelphia</i> , 4(1):34.	1727
1690		1728
1691	Alex Graves, Greg Wayne, and Ivo Danihelka. 2014. Neural Turing machines . <i>arXiv preprint</i> .	1729
1692		
1693		
1694		
1695		
1696		
1697		
1698	Kavukcuoglu, and Demis Hassabis. 2016. Hybrid computing using a neural network with dynamic external memory . <i>Nature</i> , 538:471–476.	1730
1699		1731
1700		1732
1701	C. Cordell Green, Richard J. Waldinger, David R. Barstow, Robert Elschlager, Douglas B. Lenat, Brian P. McCune, David E. Shaw, and Louis I. Steinberg. 1974. Progress report on program-understanding systems (AIM-240) .	1733
1702		1734
1703		1735
1704		1736
1705		
1706	Cordell Green. 1981. Application of theorem proving to problem solving . In Bonnie Lynn Webber and Nils J. Nilsson, editors, <i>Readings in Artificial Intelligence</i> , pages 202–222. Morgan Kaufmann.	1737
1707		1738
1708		1739
1709		
1710	H. Paul Grice and Peter F. Strawson. 1956. In defense of a dogma . <i>The Philosophical Review</i> , 65(2):141–158.	1740
1711		1741
1712	Aditya Grover, Eric Wang, Aaron Zweig, and Stefano Ermon. 2019. Stochastic optimization of sorting networks via continuous relaxations . <i>arXiv preprint</i> .	1742
1713		1743
1714		1744
1715	Jiatao Gu, Hany Hassan, Jacob Devlin, and Victor O. K. Li. 2018. Universal neural machine translation for extremely low resource languages . <i>arXiv preprint</i> .	1745
1716		1746
1717		1747
1718	Kristina Gulordava, Piotr Bojanowski, Edouard Grave, Tal Linzen, and Marco Baroni. 2018. Colorless green recurrent networks dream hierarchically . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 1195–1205, New Orleans, Louisiana. Association for Computational Linguistics.	1748
1719		1749
1720		1750
1721		1751
1722		1752
1723		1753
1724		1754
1725		1755
1726		
1727	Sumit Gulwani. 2011. Automating string processing in spreadsheets using input-output examples . <i>SIGPLAN Not.</i> , 46(1):317–330.	1756
1728		1757
1729		1758
1730	Sumit Gulwani, William R. Harris, and Rishabh Singh. 2012. Spreadsheet data manipulation using examples . <i>Commun. ACM</i> , 55(8):97–105.	1759
1731		1760
1732		1761
1733	Sumit Gulwani, José Hernández-Orallo, Emanuel Kitzelmann, Stephen H. Muggleton, Ute Schmid, and Benjamin Zorn. 2015. Inductive programming meets the real world . <i>Commun. ACM</i> , 58(11):90–99.	1762
1734		1763
1735		1764
1736		1765
1737	Sumit Gulwani, Oleksandr Polozov, and Rishabh Singh. 2017a. Program synthesis . <i>Foundations and Trends in Programming Languages</i> , 4(1–2):1–119.	1766
1738		1767
1739		1768
1740	Sumit Gulwani, Oleksandr Polozov, and Rishabh Singh. 2017b. Program Synthesis . NOW, Boston.	1769
1741		1770
1742	Aditya Gupta, Jiacheng Xu, Shyam Upadhyay, Diyi Yang, and Manaal Faruqui. 2021. Disfl-QA: A benchmark dataset for understanding disfluencies in question answering . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 3309–3319, Online. Association for Computational Linguistics.	1771
1743		1772
1744		1773
1745		1774
1746		1775
1747		1776
1748		1777

1749	Shashank Gupta, Subhabrata Mukherjee, Krishan Subudhi, Eduardo Gonzalez, Damien Jose, Ahmed H Awadallah, and Jianfeng Gao. 2022. Sparsely activated mixture-of-experts are robust multi-task learners. <i>arXiv preprint arXiv:2204.07689</i> .	1804
1750		1805
1751		1806
1752		1807
1753		1808
1754	Isidor S. Gvarjalaže and Dzhuansher I. Mchedlishvili. 1971. <i>English Proverbs and Sayings</i> . Vysshaya shkola, Moscow.	1809
1755		1810
1756		1811
1757	Samuel Gyasi Obeng. 1996. The proverb as a mitigating and politeness strategy in Akan discourse . <i>Anthropological Linguistics</i> , 38(3):521–549.	1812
1758		1813
1759		1814
1760	Ivan Habernal, Henning Wachsmuth, Iryna Gurevych, and Benno Stein. 2018. The argument reasoning comprehension task: Identification and reconstruction of implicit warrants . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 1930–1940, New Orleans, Louisiana. Association for Computational Linguistics.	1815
1761		1816
1762		1817
1763		1818
1764		1819
1765		1820
1766		1821
1767		1822
1768		1823
1769	Michael Hahn. 2020. Theoretical limitations of self-attention in neural sequence models . <i>Transactions of the Association for Computational Linguistics</i> , 8:156–171.	1824
1770		1825
1771		1826
1772		1827
1773	Rowan Hall Maudslay, Hila Gonen, Ryan Cotterell, and Simone Teufel. 2019. It's all in the name: Mitigating gender bias with name-based counterfactual data substitution . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 5267–5275, Hong Kong, China. Association for Computational Linguistics.	1828
1774		1829
1775		1830
1776		1831
1777		1832
1778		1833
1779		1834
1780		1835
1781		1836
1782	Joseph Y. Halpern. 2016. <i>Actual causality</i> . MIT Press, Cambridge, MA.	1837
1783		1838
1784	Rujun Han, Xiang Ren, and Nanyun Peng. 2020. ECONET: Effective continual pretraining of language models for event temporal reasoning . <i>arXiv preprint</i> .	1839
1785		1840
1786		1841
1787		1842
1788	Maria Hanzén. 2007. When in Rome, do as the Romans do: Proverbs as a part of EFL teaching . Master’s thesis, Jönköping University, School of Education and Communication, Jönköping.	1843
1789		1844
1790		1845
1791		1846
1792		1847
1793		1848
1794		1849
1795	Yiding Hao, William Merrill, Dana Angluin, Robert Frank, Noah Amsel, Andrew Benz, and Simon Mendelsohn. 2018. Context-free transductions with neural stacks . <i>arXiv preprint</i> .	1850
1796	Francesca G.E. Happé. 1994. An advanced test of theory of mind: Understanding of story characters thoughts and feelings by able autistic, mentally handicapped, and normal children and adults . <i>Journal of Autism and Developmental Disorders</i> , 24:129–154.	1851
1797		1852
1798		1853
1799		1854
1800		1855
1801	F. Maxwell Harper and Joseph A. Konstan. 2015. The MovieLens datasets: History and context . <i>ACM Trans. Interact. Intell. Syst.</i> , 5(4).	1856
1802		1857
1803		1858
	Behnam Hedayatnia, Karthik Gopalakrishnan, Seokhwan Kim, Yang Liu, Mihail Eric, and Dilek Hakkani-Tur. 2020. Policy-driven neural response generation for knowledge-grounded dialog systems . In <i>Proceedings of the 13th International Conference on Natural Language Generation</i> , pages 412–421, Dublin, Ireland. Association for Computational Linguistics.	
	Irene Heim. 1983. On the projection problem for presuppositions. In Paul Portner and Barbara H. Partee, editors, <i>Formal Semantics - The Essential Readings</i> , pages 249–260. Blackwell, Oxford.	
	Mikael Henaff, Jason Weston, Arthur Szlam, Antoine Bordes, and Yann LeCun. 2016. Tracking the world state with recurrent entity networks . <i>arXiv preprint</i> .	
	Lisa Anne Hendricks, Kaylee Burns, Kate Saenko, Trevor Darrell, and Anna Rohrbach. 2018. Women also snowboard: Overcoming bias in captioning models . In <i>Proceedings of the European Conference on Computer Vision (ECCV)</i> , pages 771–787, Cham. Springer.	
	Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. 2021a. Measuring coding challenge competence with APPS . <i>arXiv</i> , arXiv:2105.09938.	
	Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. 2020. Aligning AI with shared human values . <i>arXiv preprint</i> .	
	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021b. Measuring massive multitask language understanding . In <i>International Conference on Learning Representations</i> .	
	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021c. Measuring mathematical problem solving with the MATH dataset . <i>arXiv preprint</i> .	
	Tom Henighan, Jared Kaplan, Mor Katz, Mark Chen, Christopher Hesse, Jacob Jackson, Heewoo Jun, Tom B. Brown, Prafulla Dhariwal, Scott Gray, Chris Hallacy, Benjamin Mann, Alec Radford, Aditya Ramesh, Nick Ryder, Daniel M. Ziegler, John Schulman, Dario Amodei, and Sam McCandlish. 2020. Scaling laws for autoregressive generative modeling . <i>arXiv preprint</i> .	
	Joseph Henrich, Steven J. Heine, and Ara Norenzayan. 2010. The weirdest people in the world? <i>Behavioral and Brain Sciences</i> , 33(2-3):61–83.	
	Karl Moritz Hermann, Tomas Kociský, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015a. Teaching machines to read and comprehend . <i>arXiv preprint arXiv:1506.03340</i> .	

1859	Karl Moritz Hermann, Tomás Kociský, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015b. Teaching machines to read and comprehend . In <i>Advances in Neural Information Processing Systems</i> , volume 28. Curran Associates, Inc.	1913
1860		1914
1861		1915
1862		1916
1863		1917
1864		1918
1865	Danny Hernandez, Jared Kaplan, Tom Henighan, and Sam McCandlish. 2021. Scaling laws for transfer . <i>arXiv preprint</i> .	1919
1866		
1867		
1868	Álvaro Hernandez, Suhan Woo, Héctor Corrales, Ignacio Parra, Euntai Kim, D. Fernández Llorca, and Miguel A. Sotelo. 2020. 3D-DEEP: 3-dimensional deep-learning based on elevation patterns for road scene interpretation . In <i>2020 IEEE Intelligent Vehicles Symposium (IV)</i> , Piscataway, NJ. Institute of Electrical and Electronics Engineers.	1920
1869		1921
1870		1922
1871		1923
1872		1924
1873		1925
1874		
1875	Jonathan Herzig, Paweł Krzysztof Nowak, Thomas Müller, Francesco Piccinno, and Julian Eisenschlos. 2020. TaPas: Weakly supervised table parsing via pre-training . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4320–4333, Online. Association for Computational Linguistics.	1926
1876		1927
1877		1928
1878		1929
1879		1930
1880		
1881		
1882	John Hewitt, Michael Hahn, Surya Ganguli, Percy Liang, and Christopher D. Manning. 2020. RNNs can generate bounded hierarchical languages with optimal memory . <i>arXiv preprint</i> .	1931
1883		1932
1884		1933
1885		1934
1886	Mireille Hildebrandt. 2018. Algorithmic regulation and the rule of law . <i>Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences</i> , 376(2128):20170355.	1935
1887		
1888		
1889		
1890	Keith J. Holyoak. 2012. Analogy and relational reasoning . In Keith J. Holyoak and Robert G. Morrison, editors, <i>The Oxford Handbook of Thinking and Reasoning</i> . Oxford University Press, Oxford.	1936
1891		1937
1892		1938
1893		1939
1894	Richard P. Honeck. 1997. <i>A Proverb in Mind: The Cognitive Science of Proverbial Wit and Wisdom</i> . Lawrence Erlbaum Associates, Mahwah, NJ.	1940
1895		
1896		
1897	Alexandra Horowitz. 2017. Smelling themselves: Dogs investigate their own odours longer when modified in an “olfactory mirror” test . <i>Behavioural Processes</i> , 143:17–24.	1941
1898		1942
1899		1943
1900		1944
1901	Mohammad Javad Hosseini, Hannaneh Hajishirzi, Oren Etzioni, and Nate Kushman. 2014. Learning to solve arithmetic word problems with verb categorization . In <i>Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 523–533, Doha, Qatar. Association for Computational Linguistics.	1945
1902		1946
1903		1947
1904		1948
1905		1949
1906		
1907		
1908	Yufang Hou. 2020. Bridging anaphora resolution as question answering . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1428–1438, Online. Association for Computational Linguistics.	1950
1909		1951
1910		1952
1911		1953
1912		1954
1913	Yufang Hou, Katja Markert, and Michael Strube. 2013. Global inference for bridging anaphora resolution . In <i>Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 907–917, Atlanta, Georgia. Association for Computational Linguistics.	1955
1914		1956
1915		1957
1916		1958
1917		
1918		
1919		
1920	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for NLP . In <i>International Conference on Machine Learning</i> , pages 2790–2799.	1959
1921		1960
1922		1961
1923		1962
1924		
1925		
1926	China Household Management Research Center, Ministry of Public Security. 2019. National name report 2018. http://news.cpd.com.cn/n18151/201901/t20190130_830962.html (Accessed 3 March 2021).	1963
1927		1964
1928		1965
1929		1966
1930		
1931	China Household Management Research Center, Ministry of Public Security. 2020. National name report 2019. https://www.mps.gov.cn/n2254314/n6409334/c6874817/content.html (Accessed 3 March 2021).	1967
1932		1968
1933		1969
1934		1970
1935		
1936	China Household Management Research Center, Ministry of Public Security. 2021. National name report 2020. https://www.mps.gov.cn/n2253534/n2253535/c7725981/content.html (Accessed 3 March 2021).	1971
1937		1972
1938		1973
1939		1974
1940		
1941	Hrisztalina Hrisztova-Gothardt and Melita Aleksa Varga. 2015. Introduction to Paremiology: A Comprehensive Guide to Proverb Studies . De Gruyter Open, Warsaw.	1975
1942		1976
1943		1977
1944		
1945	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. In <i>International Conference on Learning Representations</i> .	1978
1946		1979
1947		1980
1948		1981
1949		
1950	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models . In <i>International Conference on Learning Representations</i> .	1982
1951		1983
1952		1984
1953		1985
1954		
1955	Chengsong Huang, Qian Liu, Bill Yuchen Lin, Tianyu Pang, Chao Du, and Min Lin. 2023. Lorahub: Efficient cross-task generalization via dynamic lora composition. <i>arXiv preprint arXiv:2307.13269</i> .	1986
1956		1987
1957		1988
1958		
1959	Chengsong Huang, Qian Liu, Bill Yuchen Lin, Tianyu Pang, Chao Du, and Min Lin. 2024a. Lorahub: Efficient cross-task generalization via dynamic lora composition . <i>Preprint</i> , arXiv:2307.13269.	1989
1960		1990
1961		1991
1962		
1963	Daniel Huang, Prafulla Dhariwal, Dawn Song, and Ilya Sutskever. 2018. Gamepad: A learning environment for theorem proving . <i>arXiv preprint</i> .	1992
1964		1993
1965		1994

1966	Haoxu Huang, Fanqi Lin, Yingdong Hu, Shengjie Wang, and Yang Gao. 2024b. Copa: General robotic manipulation through spatial constraints of parts with foundation models. <i>Preprint</i> , arXiv:2403.08248.	2022
1967		2023
1968		2024
1969		2025
1970	Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. Cosmos QA: Machine reading comprehension with contextual commonsense reasoning. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2391–2401, Hong Kong, China. Association for Computational Linguistics.	2026
1971		2027
1972		2028
1973		2029
1974		2030
1975		2031
1976		2032
1977		2033
1978		2034
1979	Quzhe Huang, Zhenwei An, Nan Zhuang, Mingxu Tao, Chen Zhang, Yang Jin, Kun Xu, Kun Xu, Liwei Chen, Songfang Huang, and Yansong Feng. 2024c. Harder tasks need more experts: Dynamic routing in moe models. <i>Preprint</i> , arXiv:2403.07652.	2035
1980		2036
1981		2037
1982		2038
1983		
1984	Drew A. Hudson and Christopher D. Manning. 2019. GQA: A new dataset for real-world visual reasoning and compositional question answering. <i>arXiv preprint</i> .	2039
1985		2040
1986		2041
1987		2042
1988	Thad Hughes and Daniel Ramage. 2007. Lexical semantic relatedness with random graph walks. In <i>Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)</i> , pages 581–589, Prague, Czech Republic. Association for Computational Linguistics.	2043
1989		2044
1990		
1991		
1992		
1993		
1994		
1995	David Hume. 1739–1740. <i>A Treatise of Human Nature</i> . John Noon, London.	2045
1996		2046
1997	Dieuwke Hupkes, Verna Dankers, Mathijs Mul, and Elia Bruni. 2020. Compositionality decomposed: How do neural networks generalise? <i>Journal of Artificial Intelligence Research</i> , 67:757–795.	2047
1998		2048
1999		2049
2000		2050
2001	Annamarie W. Huttunen, Geoffrey K. Adams, and Michael L. Platt. 2017. Can self-awareness be taught? Monkeys pass the mirror test – again. <i>Proceedings of the National Academy of Sciences</i> , 114(13):3281–3283.	2051
2002		
2003		
2004		
2005		
2006	David Huynh and Stefano Mazzocchi. 2012. OpenRefine. https://openrefine.org/ .	2052
2007		2053
2008		2054
2009		2055
2010		2056
2011		2057
2012		2058
2013	Instagram Engineering. 2015. Emojineering part 1: Machine learning for emoji trends. Medium.	2059
2014		2060
2015		2061
2016		2062
2017		2063
2018		2064
2019		2065
2020		
2021	Gabriel Ilharco, Marco Túlio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt, Hanneh Hajishirzi, and Ali Farhadi. 2022. Editing models with task arithmetic. <i>arXiv preprint arXiv:2212.04089</i> .	
2022	Geoffrey Irving, Paul Christiano, and Dario Amodei. 2018. AI safety via debate. <i>arXiv preprint</i> .	
2023		
2024	Gautier Izacard and Edouard Grave. 2020. Leveraging passage retrieval with generative models for open domain question answering. <i>arXiv preprint</i> .	
2025		
2026	Kushal Jain, Adwait Deshpande, Kumar Shridhar, Felix Laumann, and Ayushman Dash. 2020. Indic-transformers: An analysis of transformer language models for Indian languages. <i>arXiv preprint</i> .	
2027		
2028		
2029		
2030		
2031	Joel Jang, Seungone Kim, Seonghyeon Ye, Doyoung Kim, Lajanugen Logeswaran, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2023. Exploring the benefits of training expert language models over instruction tuning. <i>arXiv preprint arXiv:2302.03202</i> .	
2032		
2033		
2034		
2035	Mario Jarmasz. 2012. Roget’s Thesaurus as a lexical resource for natural language processing. Master’s thesis, University of Ottawa, Ottawa.	
2036		
2037		
2038		
2039	Prashant Jayannavar, Anjali Narayan-Chen, and Julia Hockenmaier. 2020. Learning to execute instructions in a Minecraft dialogue. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 2589–2602, Online. Association for Computational Linguistics.	
2040		
2041		
2042		
2043		
2044		
2045	Paloma Jeretic, Alex Warstadt, Suvrat Bhooshan, and Adina Williams. 2020. Are natural language inference models IMPPRESSive? Learning IMPLICATURE and PRESUPPOSITION. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8690–8705, Online. Association for Computational Linguistics.	
2046		
2047		
2048		
2049		
2050		
2051		
2052	Jay J. Jiang and David W. Conrath. 1997. Semantic similarity based on corpus statistics and lexical taxonomy. In <i>Proceedings of the 10th Research on Computational Linguistics International Conference</i> , pages 19–33, Taipei, Taiwan. The Association for Computational Linguistics and Chinese Language Processing (ACLCLP).	
2053		
2054		
2055		
2056		
2057		
2058		
2059	Nanjiang Jiang and Marie-Catherine de Marneffe. 2019. Do you know that Florence is packed with visitors? Evaluating state-of-the-art models of speaker commitment. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4208–4213, Florence, Italy. Association for Computational Linguistics.	
2060		
2061		
2062		
2063		
2064		
2065		
2066	Xisen Jin, Xiang Ren, Daniel Preotiuc-Pietro, and Pengxiang Cheng. 2022. Dataless knowledge fusion by merging weights of language models. <i>arXiv preprint arXiv:2212.09849</i> .	
2067		
2068		
2069		
2070	Matt Gardner, Johannes Welbl, Nelson F. Liu. 2017. Crowdsourcing multiple choice science questions. <i>arXiv:1707.06209v1</i> .	
2071		
2072		
2073	Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C. Lawrence Zitnick, and Ross Girshick. 2016. CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning. <i>arXiv preprint</i> .	
2074		
2075		
2076		
2077		

2078	Erik Jones, Robin Jia, Aditi Raghunathan, and Percy Liang. 2020. <i>Robust encodings: A framework for combating adversarial typos</i> . <i>arXiv preprint</i> .	2134
2079		2135
2080		2136
2081	Aditya Joshi, Pushpak Bhattacharyya, and Mark J. Carman. 2017. <i>Automatic sarcasm detection: A survey</i> . <i>ACM Comput. Surv.</i> , 50(5).	2137
2082		2138
2083		2139
2084	Aditya Joshi, Vinita Sharma, and Pushpak Bhattacharyya. 2015. <i>Harnessing context incongruity for sarcasm detection</i> . In <i>Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)</i> , pages 757–762, Beijing, China. Association for Computational Linguistics.	2140
2085		2141
2086		2142
2087		2143
2088		2144
2089		
2090		
2091		
2092	Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. <i>TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension</i> . <i>arXiv preprint arXiv:1705.03551</i> , arXiv:1705.03551.	2145
2093		2146
2094		2147
2095		
2096		
2097	Armand Joulin and Tomas Mikolov. 2015. <i>Inferring algorithmic patterns with stack-augmented recurrent nets</i> . In <i>Proceedings of the 28th International Conference on Neural Information Processing Systems, NIPS’15</i> , volume 1, page 190–198, Cambridge, MA, USA. MIT Press.	2148
2098		2149
2099		2150
2100		
2101		
2102		
2103	Mihir Kale and Abhinav Rastogi. 2020. <i>Template guided text generation for task-oriented dialogue</i> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 6505–6520, Online. Association for Computational Linguistics.	2151
2104		2152
2105		2153
2106		2154
2107		
2108		
2109	Yuji Kanagawa and Tomoyuki Kaneko. 2019. <i>Rogue-Gym: A new challenge for generalization in reinforcement learning</i> . In <i>2019 IEEE Conference on Games (CoG)</i> , pages 1–8, Piscataway, NJ. Institute of Electrical and Electronics Engineers.	2155
2110		2156
2111		2157
2112		
2113		
2114	Sean Kandel, Andreas Paepcke, Joseph Hellerstein, and Jeffrey Heer. 2011. <i>Wrangler: Interactive visual specification of data transformation scripts</i> . In <i>Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, CHI ’11</i> , page 3363–3372, New York, NY, USA. Association for Computing Machinery.	2158
2115		2159
2116		2160
2117		2161
2118		2162
2119		
2120		
2121	Immanuel Kant. 1781/1787. <i>Critique of Pure Reason</i> . The Cambridge Edition of the Works of Immanuel Kant, edited by Paul Guyer and Allen W. Wood. Cambridge University Press.	2163
2122		2164
2123		2165
2124		2166
2125	Immanuel Kant. 1783. <i>Prolegomena to Any Future Metaphysics</i> , 2nd edition. Cambridge Texts in the History of Philosophy, edited by Gary Hatfield. Cambridge University Press.	2167
2126		2168
2127		2169
2128		2170
2129	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. <i>Scaling laws for neural language models</i> . <i>arXiv preprint arXiv:2001.08361</i> .	2171
2130		2172
2131		2173
2132		2174
2133		
2134	Andrej Karpathy. 2015. <i>The unreasonable effectiveness of recurrent neural networks</i> . <i>Andrej Karpathy’s blog</i> .	2175
2135		2176
2136		2177
2137	Lauri Karttunen. 2012. <i>Simple and phrasal implicatives</i> . In <i>*SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)</i> , pages 124–131, Montréal, Canada. Association for Computational Linguistics.	2178
2138		2179
2139		2180
2140		2181
2141		
2142		
2143		
2144		
2145	Nora Kassner, Benno Krojer, and Hinrich Schütze. 2020. <i>Are pretrained language models symbolic reasoners over knowledge?</i> <i>arXiv preprint</i> .	2182
2146		2183
2147		2184
2148	Nora Kassner and Hinrich Schütze. 2019. <i>Negated and misprimed probes for pretrained language models: Birds can talk, but cannot fly</i> . <i>arXiv preprint</i> .	2185
2149		2186
2150		2187
2151	Divyansh Kaushik, Eduard Hovy, and Zachary C. Lipton. 2019. <i>Learning the difference that makes a difference with counterfactually-augmented data</i> . <i>arXiv preprint</i> .	2188
2152		
2153		
2154		
2155	Zachary Kenton, Tom Everitt, Laura Weidinger, Iason Gabriel, Vladimir Mikulik, and Geoffrey Irving. 2021. <i>Alignment of language agents</i> .	2189
2156		2190
2157		2191
2158	Faisal Khan, Bilge Mutlu, and Jerry Zhu. 2011. <i>How do humans teach: On curriculum learning and teaching dimension</i> . In <i>Advances in Neural Information Processing Systems</i> , volume 24. Curran Associates, Inc.	2192
2159		2193
2160		2194
2161		2195
2162		2196
2163	Daniel Khashabi, Arman Cohan, Siamak Shakeri, Pedram Hosseini, Pouya Pezeshkpour, Maliheh Alikhani, Moin Aminnaseri, Marzieh Bitaab, Faeze Brahman, Sarik Ghazarian, Mozhdeh Gheini, Arman Kabiri, Rabeeh Karimi Mahabadi, Omid Memarrast, Ahmadreza Mosallanezhad, Erfan Noury, Shahab Raji, Mohammad Sadegh Rasooli, Sepideh Sadeghi, Erfan Sadeqi Azer, Niloofar Safi Samghabadi, Mahsa Shafaei, Saber Sheybani, Ali Tazarv, and Yadollah Yaghoobzadeh. 2020a. <i>ParsiNLU: A suite of language understanding challenges for persian</i> . <i>arXiv preprint</i> .	2197
2164		2198
2165		2199
2166		2200
2167		2201
2168		2202
2169		2203
2170		2204
2171		2205
2172		2206
2173		2207
2174		2208
2175	Daniel Khashabi, Sewon Min, Tushar Khot, Ashish Sabharwal, Oyvind Tafjord, Peter Clark, and Hannaneh Hajishirzi. 2020b. <i>UNIFIEDQA: Crossing format boundaries with a single QA system</i> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 1896–1907, Online. Association for Computational Linguistics.	2209
2176		2210
2177		2211
2178		2212
2179		2213
2180		2214
2181		2215
2182	Mikhail Khodak, Nikunj Saunshi, and Kiran Vodrahalli. 2017. <i>A large self-annotated corpus for sarcasm</i> . <i>arXiv preprint</i> .	2216
2183		2217
2184		2218
2185	Tushar Khot, Peter Clark, Michal Guerquin, Peter Jansen, and Ashish Sabharwal. 2020. <i>Qasc: A dataset for question answering via sentence composition</i> . <i>arXiv:1910.11473v2</i> .	2219
2186		2220
2187		2221
2188		2222
2189		2223

2189	Andrew Kim, Maxim Ruzmaykin, Aaron Truong, and Adam Summerville. 2019. Cooperation and code-names: Understanding natural language processing via codenames . In <i>Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment</i> , volume 15, pages 160–166, Menlo Park, CA. Association for the Advancement of Artificial Intelligence.	2244
2190		2245
2191		2246
2192		2247
2193		2248
2194		
2195		
2196		
2197	Yoon Kim, Yacine Jernite, David Sontag, and Alexander M. Rush. 2015. Character-aware neural language models . <i>arXiv preprint</i> .	2249
2198		2250
2199		2251
2200	Milton King and Paul Cook. 2020. Evaluating approaches to personalizing language models . In <i>Proceedings of the 12th Language Resources and Evaluation Conference</i> , pages 2461–2469, Marseille, France. European Language Resources Association.	2252
2201		2253
2202		2254
2203		2255
2204		2256
2205	Christo Kirov and Ryan Cotterell. 2018. Recurrent Neural Networks in Linguistic Theory: Revisiting Pinker and Prince (1988) and the Past Tense Debate . <i>Transactions of the Association for Computational Linguistics</i> , 6:651–665.	2260
2206		2261
2207		2262
2208		2263
2209		
2210	Emanuel Kitzelmann. 2010. Inductive programming: A survey of program synthesis techniques . In <i>Approaches and Applications of Inductive Programming</i> , pages 50–73, Berlin. Springer.	2264
2211		2265
2212		2266
2213		
2214	Joshua Knobe. 2003. Intentional action and side effects in ordinary language . <i>Analysis</i> , 63.	2267
2215		2268
2216	Vid Kocijan, Ana-Maria Cretu, Oana-Maria Camburu, Yordan Yordanov, and Thomas Lukasiewicz. 2019. A surprisingly robust trick for the Winograd schema challenge . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4837–4842, Florence, Italy. Association for Computational Linguistics.	2269
2217		2270
2218		2271
2219		2272
2220		2273
2221		
2222		
2223	Tomáš Kočiský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis, and Edward Grefenstette. 2018. The NarrativeQA reading comprehension challenge . <i>Transactions of the Association for Computational Linguistics</i> , 6:317–328.	2274
2224		2275
2225		2276
2226		2277
2227		2278
2228	Jan Kocoń, Piotr Miłkowski, and Kamil Kanclerz. 2021. MultiEmo: Multilingual, multilevel, multidomain sentiment analysis corpus of consumer reviews . In <i>Computational Science – ICCS 2021</i> , pages 297–312, Cham. Springer.	2279
2229		2280
2230		2281
2231		2282
2232		
2233	Alexander W. Kocurek and Ethan Jerzak. 2021. Counterlogicals as counterconventionals . <i>Journal of Philosophical Logic</i> , 50:673–704.	2283
2234		2284
2235		2285
2236	Alexander W. Kocurek, Ethan Jerzak, and Rachel Etta Rudolph. 2020. Against conventional wisdom . <i>Philosophers' Imprint</i> , 20(22):1–27.	2286
2237		
2238		
2239	Moshe Koppel and Jonathan Schler. 2004. Authorship verification as a one-class classification problem . In <i>Proceedings of the Twenty-First International Conference on Machine Learning</i> , page 62, New York, NY, USA. Association for Computing Machinery.	2287
2240		2288
2241		2289
2242		2290
2243		
537	Jarmo Korhonen. 2009. Sprichwörter und zweisprachige lexikographie: Deutsch-schwedische und deutsch-finnische wörterbücher im vergleich . In C. Földes, editor, <i>Phraseologie disziplinär und interdisziplinär</i> , pages 537–549. Gunter Narr Verlag.	2244
538		2245
539		2246
540		2247
541		2248
542	Dimitrios Kotsakos, Theodoros Lappas, Dimitrios Kotzias, Dimitrios Gunopoulos, Nattiya Kanhabua, and Kjetil Nørvåg. 2014. A burstiness-aware approach for document dating . In <i>Proceedings of the 37th International ACM SIGIR Conference on Research & Development in Information Retrieval, SIGIR '14</i> , page 1003–1006, New York, NY, USA. Association for Computing Machinery.	2249
543		2250
544		2251
545		2252
546	Samuel Kounev, Jeffrey O. Kephart, Aleksandar Milenkoski, and Xiaoyun Zhu, editors. 2017. Self-Aware Computing Systems . Springer, Cham.	2253
547		2254
548	Sarah E. Kreps, Miles McCain, and Miles Brundage. 2020. All the news that's fit to fabricate: AI-generated text as a tool of media misinformation . <i>SSRN</i> .	2255
549		2256
550	Kalpesh Krishna, Aurko Roy, and Mohit Iyyer. 2021. Hurdles to progress in long-form question answering . <i>arXiv preprint</i> .	2257
551		2258
552	Wojciech Kryscinski, Bryan McCann, Caiming Xiong, and Richard Socher. 2020. Evaluating the factual consistency of abstractive text summarization . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 9332–9346, Online. Association for Computational Linguistics.	2259
553		2260
554	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina N. Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc V. Le, and Slav Petrov. 2019. Natural questions: A benchmark for question answering research . <i>Transactions of the Association for Computational Linguistics</i> , 7:453–466.	2261
555		2262
556	Niklas Kühl, Marc Goutier, Lucas Baier, Clemens Wolff, and Dominik Martin. 2020. Human vs. supervised machine learning: Who learns patterns faster? <i>arXiv preprint</i> .	2263
557		2264
558	Heinrich Küttler, Nantas Nardelli, Alexander H. Miller, Roberta Raileanu, Marco Selvatici, Edward Grefenstette, and Tim Rocktäschel. 2020. The NetHack learning environment . <i>arXiv preprint</i> .	2265
559		2266
560	Kevin Lacker. 2020. Giving GPT-3 a Turing test . <i>Kevin Lacker's blog</i> .	2267
561		2268
562	Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. RACE: Large-scale ReADING comprehension dataset from examinations . In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 785–794, Copenhagen, Denmark. Association for Computational Linguistics.	2269
563		2270
564		2271
565		2272
566		2273
567		

2300	Brenden M. Lake and Marco Baroni. 2017. Generalization without systematicity: On the compositional skills of sequence-to-sequence recurrent networks . <i>arXiv preprint</i> .	2354
2301		2355
2302		2356
2303		2357
2304	Brenden M. Lake and Gregory L. Murphy. 2020. Word meaning in minds and machines . <i>arXiv preprint</i> .	2358
2305		2359
2306	Brenden M. Lake, Tomer D. Ullman, Joshua B. Tenenbaum, and Samuel J. Gershman. 2017. Building machines that learn and think like people . <i>Behavioral and Brain Sciences</i> , 40:e253.	2360
2307		
2308		
2309		
2310	George Lakoff and Mark Johnson. 2008. <i>Metaphors We Live By</i> . University of Chicago Press, Chicago.	2361
2311		2362
2312	Yair Lakretz, Théo Desbordes, Jean-Rémi King, Benoît Crabbé, Maxime Oquab, and Stanislas Dehaene. 2021a. Can RNNs learn recursive nested subject-verb agreements? <i>arXiv preprint</i> .	2363
2313		2364
2314		2365
2315		2366
2316	Yair Lakretz, Dieuwke Hupkes, Alessandra Vergallito, Marco Marelli, Marco Baroni, and Stanislas Dehaene. 2021b. Mechanisms for handling nested dependencies in neural-network language models and humans . <i>Cognition</i> , 213:104699. Special Issue in Honour of Jacques Mehler, Cognition’s founding editor.	2367
2317		2368
2318		2369
2319		2370
2320		
2321		
2322	Yair Lakretz, German Kruszewski, Theo Desbordes, Dieuwke Hupkes, Stanislas Dehaene, and Marco Baroni. 2019. The emergence of number and syntax units in LSTM language models . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 11–20, Minneapolis, Minnesota. Association for Computational Linguistics.	2371
2323		2372
2324		2373
2325		2374
2326		2375
2327		2376
2328		
2329		
2330		
2331	Guillaume Lample and François Charton. 2019. Deep learning for symbolic mathematics . <i>arXiv preprint</i> .	2377
2332		2378
2333	Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. ALBERT: A lite BERT for self-supervised learning of language representations . <i>arXiv preprint</i> .	2379
2334		
2335		
2336		
2337	Matthew Le, Y-Lan Boureau, and Maximilian Nickel. 2019. Revisiting the evaluation of theory of mind through question answering . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 5872–5877, Hong Kong, China. Association for Computational Linguistics.	2380
2338		2381
2339		2382
2340		2383
2341		2384
2342		2385
2343		
2344		
2345	Remi Lebret, David Grangier, and Michael Auli. 2016a. Neural text generation from structured data with application to the biography domain . <i>Preprint</i> , arXiv:1603.07771.	2386
2346		2387
2347		2388
2348		2389
2349	Rémi Lebret, David Grangier, and Michael Auli. 2016b. The wikibio corpus: A corpus of biographical texts for natural language generation . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> .	2390
2350		2391
2351		2392
2352		2393
2353		2394
		2395
	Nayeon Lee, Yejin Bang, Andrea Madotto, and Pascale Fung. 2021. Towards few-shot fact-checking via perplexity . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1971–1981, Online. Association for Computational Linguistics.	2396
		2397
		2398
		2399
		2400
		2401
		2402
		2403
	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning . <i>Preprint</i> , arXiv:2104.08691.	2404
		2405
		2406
	Iddo Lev, Bill MacCartney, Christopher Manning, and Roger Levy. 2004. Solving logic puzzles: From robust processing to precise semantics . In <i>Proceedings of the 2nd Workshop on Text Meaning and Interpretation</i> , pages 9–16, Barcelona, Spain. Association for Computational Linguistics.	2407
		2408
		2409
	Hector Levesque, Ernest Davis, and Leora Morgenstern. 2012. The winograd schema challenge . In <i>Thirteenth International Conference on the Principles of Knowledge Representation and Reasoning</i> .	2410
		2411
		2412
	Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. 2017. Zero-shot relation extraction via reading comprehension . In <i>Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)</i> , pages 333–342, Vancouver, Canada. Association for Computational Linguistics.	2413
		2414
		2415
	Ran Levy, Liat Ein-Dor, Shay Hummel, Ruty Rinott, and Noam Slonim. 2015. TR9856: A multi-word term relatedness benchmark . In <i>Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)</i> , pages 419–424, Beijing, China. Association for Computational Linguistics.	2416
		2417
		2418
		2419
	Sharon Levy, Michael Saxon, and William Yang Wang. 2021. Investigating memorization of conspiracy theories in text generation . <i>arXiv preprint</i> .	2420
		2421
		2422
	Patrick Lewis, Barlas Oguz, Ruty Rinott, Sebastian Riedel, and Holger Schwenk. 2020a. MLQA: Evaluating cross-lingual extractive question answering . In	2423
		2424

2410			Zujie Liang, Weitao Jiang, Haifeng Hu, and Jiaying Zhu. 2020b. <i>Learning to contrast the counterfactual samples for robust visual question answering</i> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 3285–3292, Online. Association for Computational Linguistics.	2466
2411				2467
2412				2468
2413				2469
2414	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020b. <i>Retrieval-augmented generation for knowledge-intensive NLP tasks</i> . <i>arXiv preprint</i> .			2470
2415				2471
2416				2472
2417				
2418				
2419				
2420	Patrick Lewis, Pontus Stenetorp, and Sebastian Riedel. 2020c. <i>Question and answer test-train overlap in open-domain question answering datasets</i> . <i>arXiv preprint</i> .			2473
2421				2474
2422				2475
2423				2476
2424	Pingzhi Li, Zhenyu Zhang, Prateek Yadav, Yi-Lin Sung, Yu Cheng, Mohit Bansal, and Tianlong Chen. 2023. Merge, then compress: Demystify efficient smoe with hints from its routing policy. <i>arXiv preprint arXiv:2310.01334</i> .			2477
2425				2478
2426				2479
2427				
2428				
2429	Tao Li, Daniel Khashabi, Tushar Khot, Ashish Sabharwal, and Vivek Srikumar. 2020a. <i>UNQOVERing stereotyping biases via underspecified questions</i> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 3475–3489, Online. Association for Computational Linguistics.			2480
2430				2481
2431				2482
2432				2483
2433				
2434				
2435	Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. 2017. <i>DailyDialog: A manually labelled multi-turn dialogue dataset</i> . In <i>Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 986–995, Taipei, Taiwan. Asian Federation of Natural Language Processing.			2484
2436				2485
2437				2486
2438				2487
2439				2488
2440				
2441				
2442	Yiwei Li, G Brian Golding, and Lucian Ilie. 2020b. <i>DELPHI: Accurate deep ensemble model for protein interaction sites prediction</i> . <i>Bioinformatics</i> , 37(7):896–904.			2489
2443				2490
2444				2491
2445				2492
2446	Yuhua Li, Zuhair A. Bandar, and David Mclean. 2003. <i>An approach for measuring semantic similarity between words using multiple information sources</i> . <i>IEEE Transactions on Knowledge and Data Engineering</i> , 15(4):871–882.			2493
2447				2494
2448				2495
2449				
2450				
2451	Chao-Chun Liang, Yu-Shiang Wong, Yi-Chung Lin, and Keh-Yih Su. 2018. <i>A meaning-based statistical English math word problem solver</i> . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 652–662, New Orleans, Louisiana. Association for Computational Linguistics.			2496
2452				2497
2453				2498
2454				2499
2455				2500
2456				2501
2457				
2458				
2459	Paul Pu Liang, Irene Mengze Li, Emily Zheng, Yao Chong Lim, Ruslan Salakhutdinov, and Louis-Philippe Morency. 2020a. <i>Towards debiasing sentence representations</i> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5502–5515, Online. Association for Computational Linguistics.			2502
2460				2503
2461				2504
2462				2505
2463				2506
2464				2507
2465				2508
2466				
2467				
2468				
2469				
2470				
2471				
2472				
2473				
2474				
2475				
2476				
2477				
2478				
2479				
2480				
2481				
2482				
2483				
2484				
2485				
2486				
2487				
2488				
2489				
2490				
2491				
2492				
2493				
2494				
2495				
2496				
2497				
2498				
2499				
2500				
2501				
2502				
2503				
2504				
2505				
2506				
2507				
2508				
2509				
2510				
2511				
2512				
2513				
2514				
2515				
2516				
2517				
2518				
2519				

2520	Tal Linzen, Emmanuel Dupoux, and Yoav Goldberg.	Shayne Longpre, Le Hou, Tu Vu, Albert Webson,	2575
2521	2016. Assessing the ability of LSTMs to learn syntax-sensitive dependencies . <i>Transactions of the Association for Computational Linguistics</i> , 4:521–535.	Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, et al. 2023. The flan collection: Designing data and methods for effective instruction tuning. <i>arXiv preprint arXiv:2301.13688</i> .	2576
2522			2577
2523			2578
2524	Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Motta, Tenghao Huang, Mohit Bansal, and Colin A Rafel. 2022. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. <i>Advances in Neural Information Processing Systems</i> , 35:1950–1965.	Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization . In <i>International Conference on Learning Representations</i> .	2579
2525			2580
2526			2581
2527			2582
2528			2583
2529		Nicholas Lourie, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Unicorn on rainbow: A universal commonsense reasoning model on a new multitask benchmark . <i>arXiv preprint</i> .	2584
2530	Jerry Liu. 2024. LlamaIndex, a data framework for your LLM applications. https://github.com/run-llama/llama_index .	Nicholas Lourie, Ronan Le Bras, and Yejin Choi. 2020. Scruples: A corpus of community ethical judgments on 32,000 real-life anecdotes . <i>arXiv preprint</i> .	2585
2531			2586
2532			2587
2533	Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2021a. What makes good in-context examples for GPT-3? <i>arXiv preprint</i> .	Kaiji Lu, Piotr Mardziel, Fangjing Wu, Preetam Amancharla, and Anupam Datta. 2020. Gender bias in neural natural language processing . In Vivek Nigam, Tajana Ban Kirigin, Carolyn Talcott, Joshua Guttmann, Stepan Kuznetsov, Boon Thau Loo, and Mitsuhiro Okada, editors, <i>Logic, Language, and Security</i> . Springer, Cham.	2588
2534			2589
2535			2590
2536			2591
2537	Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. 2020a. LogiQA: A challenge dataset for machine reading comprehension with logical reasoning . <i>arXiv preprint</i> .	Keming Lu, Hongyi Yuan, Runji Lin, Junyang Lin, Zheng Yuan, Chang Zhou, and Jingren Zhou. 2023. Routing to the expert: Efficient reward-guided ensemble of large language models. <i>arXiv preprint arXiv:2311.08692</i> .	2592
2538			2593
2539			2594
2540			2595
2541	Nelson F. Liu, Tony Lee, Robin Jia, and Percy Liang. 2021b. Can small and synthetic benchmarks drive modeling innovation? A retrospective study of question answering modeling approaches . <i>arXiv preprint</i> .		2596
2542			2597
2543			2598
2544			2599
2545	Tianyu Liu, Yizhe Zhang, Chris Brockett, Yi Mao, Zhifang Sui, Weizhu Chen, and Bill Dolan. 2021c. A token-level reference-free hallucination detection benchmark for free-form text generation . <i>arXiv preprint</i> .	Zhenyi Lu, Chenghao Fan, Wei Wei, Xiaoye Qu, Dangyang Chen, and Yu Cheng. 2024. Twin-merging: Dynamic integration of modular expertise in model merging. <i>arXiv preprint arXiv:2406.15479</i> .	2600
2546			2601
2547			2602
2548			2603
2549			2604
2550	Ye Liu, Shaika Chowdhury, Chenwei Zhang, Cornelia Caragea, and Philip S. Yu. 2020b. Interpretable multi-step reasoning with knowledge extraction on complex healthcare question answering . <i>arXiv preprint</i> .		2605
2551			2606
2552			2607
2553			2608
2554	Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020c. Multilingual denoising pre-training for neural machine translation . <i>arXiv preprint</i> .	Jelena Luketina, Nantas Nardelli, Gregory Farquhar, Jakob N. Foerster, Jacob Andreas, Edward Grefenstette, Shimon Whiteson, and Tim Rocktäschel. 2019. A survey of reinforcement learning informed by natural language . In <i>Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19</i> , page 6309–6317.	2609
2555			2610
2556			2611
2557			2612
2558			2613
2559	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 . <i>arXiv preprint</i> .	Mingyu Derek Ma, Jiao Sun, Mu Yang, Kung-Hsiang Huang, Nuan Wen, Shikhar Singh, Rujun Han, and Nanyun Peng. 2021. EventPlus: A temporal event understanding pipeline . <i>arXiv preprint</i> .	2614
2560			2615
2561			2616
2562			2617
2563	Hector Llorens, Nathanael Chambers, Naushad UzZaman, Nasrin Mostafazadeh, James Allen, and James Pustejovsky. 2015. SemEval-2015 task 5: QA TempEval - evaluating temporal information understanding with question answering . In <i>Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)</i> , pages 792–800, Denver, Colorado. Association for Computational Linguistics.	Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011a. Learning word vectors for sentiment analysis . In <i>Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies</i> .	2618
2564			2619
2565			2620
2566			2621
2567			2622
2568			2623
2569			2624
2570			2625
2571			2626
2572	Lajanugen Logeswaran, Honglak Lee, and Samy Bengio. 2018. Content preserving text generation with attribute controls . <i>arXiv preprint</i> .	Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011b. Learning word vectors for sentiment analysis . In <i>Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies</i> .	2627
2573			2628
2574			2629

2629	Andrea Madotto, Zhaojiang Lin, Genta Indra Winata, and Pascale Fung. 2021. Few-shot bot: Prompt-based learning for dialogue systems .	Philip Massey, Patrick Xia, David Bamman, and Noah A. Smith. 2015. Annotating character relationships in literary texts . <i>arXiv preprint</i> .	2683
2630			2684
2631			2685
2632	Andrea Madotto, Zihan Liu, Zhaojiang Lin, and Pascale Fung. 2020. Language models as few-shot learner for task-oriented dialogue systems . <i>arXiv preprint</i> .	Michael S Matena and Colin A Raffel. 2022. Merging models with fisher-weighted averaging. <i>Advances in Neural Information Processing Systems</i> , 35:17703–17716.	2686
2633			2687
2634			2688
2635	Eric Malmi, Sebastian Krause, Sascha Rothe, Daniil Mirylenka, and Aliaksei Severyn. 2019. Encode, tag, realize: High-precision text editing . <i>arXiv preprint</i> .	Matthew Matero, Akash Idnani, Youngseo Son, Salvatore Giorgi, Huy Vu, Mohammad Zamani, Parth Limbachiya, Sharath Chandra Guntuku, and H. Andrew Schwartz. 2019. Suicide risk assessment with multi-level dual-context language and BERT . In <i>Proceedings of the Sixth Workshop on Computational Linguistics and Clinical Psychology</i> , pages 39–44, Minneapolis, Minnesota. Association for Computational Linguistics.	2689
2636			2690
2637			2691
2638	Eric Malmi, Daniele Pighin, Sebastian Krause, and Mikhail Kozhevnikov. 2018. Automatic prediction of discourse connectives . In <i>Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)</i> , Miyazaki, Japan. European Language Resources Association.	2692	
2639			2693
2640			2694
2641			2695
2642			2696
2643			2697
2644	Jihang Mao and Wanli Liu. 2019. A BERT-based approach for automatic humor detection and scoring . In <i>Proceedings of the Iberian Languages Evaluation Forum (IberLEF 2019)</i> , pages 197–202.	Maxine. 2023. Llama-2, mo’ lora. https://crumblly.medium.com/llama-2-molora-f5f909434711 .	2698
2645			2699
2646			2700
2647			
2648	Gary Marcus. 2020. The next decade in AI: Four steps towards robust artificial intelligence . <i>arXiv preprint</i> .	Chandler May, Alex Wang, Shikha Bordia, Samuel R. Bowman, and Rachel Rudinger. 2019. On measuring social biases in sentence encoders . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 622–628, Minneapolis, Minnesota. Association for Computational Linguistics.	2701
2649			2702
2650	Gary Marcus and Ernest Davis. 2020. GPT-3, bloviator: OpenAI’s language generator has no idea what it’s talking about . <i>MIT Technology Review</i> .	2703	
2651			2704
2652			2705
2653	Mitchell Marcus, Grace Kim, Mary Ann Marcinkiewicz, Robert MacIntyre, Ann Bies, Mark Ferguson, Karen Katz, and Britta Schasberger. 1994. The Penn Treebank: Annotating predicate argument structure . In <i>Human Language Technology: Proceedings of a Workshop held at Plainsboro, New Jersey, March 8-11, 1994</i> .	2706	
2654			2707
2655			2708
2656			
2657			
2658			
2659			
2660	Katja Markert, Yufang Hou, and Michael Strube. 2012. Collective classification for fine-grained information status . In <i>Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 795–804, Jeju Island, Korea. Association for Computational Linguistics.	Andrew Mayne. 2020. OpenAI API alchemy: Emoji storytelling . <i>Andrew Mayne blog</i> .	2709
2661			2710
2662			
2663			
2664			
2665			
2666	Kim Marriott, Bongshin Lee, Matthew Butler, Ed Cutrell, Kirsten Ellis, Cagatay Goncu, Marti Hearst, Kathleen McCoy, and Danielle Albers Szafir. 2021. Inclusive data visualization for people with disabilities: A call to action . <i>Interactions</i> , 28(3):47–51.	Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization . <i>arXiv preprint</i> .	2711
2667			2712
2668			2713
2669			
2670			
2671	Rebecca Marvin and Tal Linzen. 2018. Targeted syntactic evaluation of language models . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1192–1202, Brussels, Belgium. Association for Computational Linguistics.	Eric Mays, Fred J. Damerau, and Robert L. Mercer. 1991. Context based spelling correction . <i>Information Processing & Management</i> , 27(5):517–522.	2714
2672			2715
2673			2716
2674			
2675			
2676			
2677	Javier Marín, Aritro Biswas, Ferda Ofli, Nicholas Hynes, Amaia Salvador, Yusuf Aytar, Ingmar Weber, and Antonio Torralba. 2021. Recipe1M+: A dataset for learning cross-modal embeddings for cooking recipes and food images . <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 43(1):187–203.	Momoh Karmah Mbogba, Zeeshan Haider, S.M. Chapal Hossain, Daobin Huang, Kashan Memon, Fazil Panhwar, Zeling Lei, and Gang Zhao. 2018. The application of convolution neural network based cell segmentation during cryopreservation . <i>Cryobiology</i> , 85:95–104.	2717
2678			2718
2679			2719
2680			2720
2681			2721
2682			2722
2683			
2684			
2685			
2686			
2687			
2688			
2689			
2690			
2691			
2692			
2693			
2694			
2695			
2696			
2697			
2698			
2699			
2700			
2701			
2702			
2703			
2704			
2705			
2706			
2707			
2708			
2709			
2710			
2711			
2712			
2713			
2714			
2715			
2716			
2717			
2718			
2719			
2720			
2721			
2722			
2723			
2724			
2725			
2726			
2727			
2728			
2729			
2730			
2731			
2732			
2733			
2734			
2735			
2736			
2737			

2738	R. Thomas McCoy, Robert Frank, and Tal Linzen. 2020.	2792
2739	Does Syntax Need to Grow on Trees? Sources of Hierarchical Inductive Bias in Sequence-to-Sequence Networks. <i>Transactions of the Association for Computational Linguistics</i> , 8:125–140.	2793
2740		2794
2741		2795
2742		2796
2743	R. Thomas McCoy, Ellie Pavlick, and Tal Linzen. 2019.	2797
2744	Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference. <i>arXiv preprint</i> .	2798
2745		
2746		
2747	Brendan McMahan, Eider Moore, Daniel Ramage,	2799
2748	Seth Hampson, and Blaise Aguera y Arcas. 2017.	2800
2749	Communication-efficient learning of deep networks	2801
2750	from decentralized data. In <i>Artificial intelligence and</i>	2802
2751	<i>statistics</i> .	2803
2752	Shikib Mehri and Maxine Eskenazi. 2020.	2804
2753	USR: An unsupervised and reference free evaluation metric	
2754	for dialog generation. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational</i>	
2755	<i>Linguistics</i> , pages 681–707, Online. Association for	
2756	Computational Linguistics.	
2757		
2758	Christine Palm Meister. 2007.	2805
2759	Phraseologie des schwedischen. In H. Burger et al., editor, <i>Phraseologie/Phrasology</i> , volume 2, pages 673–681. De	2806
2760	Gruyter Mouton.	2807
2761		
2762	Francisco S. Melo, Carla Guerra, and Manuel Lopes.	2808
2763	2018. Interactive optimal teaching with unknown	2809
2764	learners. In <i>Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18</i> , pages 2567–2573.	2810
2765		2811
2766		
2767	Julia Mendelsohn, Yulia Tsvetkov, and Dan Jurafsky.	2812
2768	2020. A framework for the computational linguistic	2813
2769	analysis of dehumanization. <i>Frontiers in Artificial</i>	2814
2770	<i>Intelligence</i> , 3.	2815
2771	Yuanliang Meng, Anna Rumshisky, and Alexey Romanov. 2017.	2816
2772	Temporal information extraction for question answering using syntactic dependencies in	2817
2773	an LSTM-based architecture. In <i>Proceedings of the</i>	
2774	<i>2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 887–896, Copenhagen,	
2775	Denmark. Association for Computational Linguistics.	
2776		
2777		
2778		
2779	Stephen Merity, Caiming Xiong, James Bradbury, and	2818
2780	Richard Socher. 2016.	2819
2781	Pointer sentinel mixture models. <i>arXiv preprint</i> .	2820
2782	William Merrill. 2020.	2821
2783	On the linguistic capacity of real-time counter automata. <i>arXiv preprint</i> .	2822
2784		2823
2785		
2786	Elliot Meyerson and Risto Miikkulainen. 2017.	2824
2787	Beyond shared hierarchies: Deep multitask learning through	2825
2788	soft layer ordering. <i>ArXiv</i> , abs/1711.00108.	2826
2789		2827
2790		
2791	Wolfgang Mieder. 2019. "Andere zeiten, andere lehren": Sprach-und kulturgeschichtliche betrachtungen zum sprichwort. In K. Steyer, editor, <i>Wortverbindungen - mehr oder weniger fest</i> , pages 415–438. De Gruyter, Berlin.	
2792	Rada Mihalcea and Carlo Strapparava. 2005.	2828
2793	Making computers laugh: Investigations in automatic humor recognition. In <i>Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing</i> , pages 531–538, Vancouver, British Columbia, Canada. Association for Computational Linguistics.	2829
2794		2830
2795		2831
2796		2832
2797		2833
2798		2834
2799	John Miller, Karl Krauth, Benjamin Recht, and Ludwig Schmidt. 2020.	2835
2800	The effect of natural distribution shift on question answering models. In <i>Proceedings of the 37th International Conference on Machine Learning</i> , volume 119 of <i>Proceedings of Machine Learning Research</i> , pages 6905–6916. PMLR.	2836
2801		2837
2802		2838
2803		2839
2804		2840
2805	Tristan Miller and Iryna Gurevych. 2015.	2841
2806	Automatic disambiguation of English puns. In <i>Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 719–729, Beijing, China. Association for Computational Linguistics.	2842
2807		2843
2808		2844
2809		2845
2810		2846
2811		2847
2812	Tristan Miller, Christian Hempelmann, and Iryna Gurevych. 2017.	2848
2813	SemEval-2017 task 7: Detection and interpretation of English puns. In <i>Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)</i> , pages 58–68, Vancouver, Canada. Association for Computational Linguistics.	2849
2814		2850
2815		2851
2816		2852
2817		2853
2818	David Milne and Ian H. Witten. 2008.	2854
2819	An effective, low-cost measure of semantic relatedness obtained from Wikipedia links. In <i>Proceeding of AAAI Workshop on Wikipedia and Artificial Intelligence: An Evolving Synergy</i> , page 25–30, Menlo Park. Association for the Advancement of Artificial Intelligence.	2855
2820		2856
2821		2857
2822		2858
2823		2859
2824	Republic of China Ministry of the Interior. 2018.	2860
2825	National name statistical analysis. https://www.ris.gov.tw/documents/data/5/2/107namestat.pdf	2861
2826	(Accessed 3 March 2021).	2862
2827		
2828	Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2021.	2863
2829	Cross-task generalization via natural language crowdsourcing instructions. <i>arXiv preprint arXiv:2104.08773</i> , arXiv:2104.08773.	2864
2830		2865
2831		2866
2832	Ishan Misra, Abhinav Shrivastava, Abhinav Kumar Gupta, and Martial Hebert. 2016.	2867
2833	Cross-stitch networks for multi-task learning. <i>2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 3994–4003.	2868
2834		2869
2835		2870
2836		2871
2837	Margaret Mitchell, Kees van Deemter, and Ehud Reiter. 2010.	2872
2838	Natural reference to objects in a visual domain. In <i>Proceedings of the 6th International Natural Language Generation Conference</i> . Association for Computational Linguistics.	2873
2839		2874
2840		2875
2841		2876
2842	Margaret Mitchell, Kees van Deemter, and Ehud Reiter. 2013.	2877
2843	Generating expressions that refer to visible objects. In <i>Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1174–1184, Atlanta, Georgia. Association for Computational Linguistics.	2878
2844		2879
2845		2880
2846		2881
2847		2882
2848		2883
2849		2884
2850		2885
2851		2886
2852		2887
2853		2888
2854		2889
2855		2890
2856		2891
2857		2892
2858		2893
2859		2894
2860		2895
2861		2896
2862		2897
2863		2898
2864		2899
2865		2900
2866		2901
2867		2902
2868		2903
2869		2904
2870		2905
2871		2906
2872		2907
2873		2908
2874		2909
2875		2910
2876		2911
2877		2912
2878		2913
2879		2914
2880		2915
2881		2916
2882		2917
2883		2918
2884		2919
2885		2920
2886		2921
2887		2922
2888		2923
2889		2924
2890		2925
2891		2926
2892		2927
2893		2928
2894		2929
2895		2930
2896		2931
2897		2932
2898		2933
2899		2934
2900		2935
2901		2936
2902		2937
2903		2938
2904		2939
2905		2940
2906		2941
2907		2942
2908		2943
2909		2944
2910		2945
2911		2946
2912		2947
2913		2948
2914		2949
2915		2950
2916		2951
2917		2952
2918		2953
2919		2954
2920		2955
2921		2956
2922		2957
2923		2958
2924		2959
2925		2960
2926		2961
2927		2962
2928		2963
2929		2964
2930		2965
2931		2966
2932		2967
2933		2968
2934		2969
2935		2970
2936		2971
2937		2972
2938		2973
2939		2974
2940		2975
2941		2976
2942		2977
2943		2978
2944		2979
2945		2980
2946		2981
2947		2982
2948		2983
2949		2984
2950		2985
2951		2986
2952		2987
2953		2988
2954		2989
2955		2990
2956		2991
2957		2992
2958		2993
2959		2994
2960		2995
2961		2996
2962		2997
2963		2998
2964		2999
2965		3000
2966		3001
2967		3002
2968		3003
2969		3004
2970		3005
2971		3006
2972		3007
2973		3008
2974		3009
2975		3010
2976		3011
2977		3012
2978		3013
2979		3014
2980		3015
2981		3016
2982		3017
2983		3018
2984		3019
2985		3020
2986		3021
2987		3022
2988		3023
2989		3024
2990		3025
2991		3026
2992		3027
2993		3028
2994		3029
2995		3030
2996		3031
2997		3032
2998		3033
2999		3034
3000		3035
3001		3036
3002		3037
3003		3038
3004		3039
3005		3040
3006		3041
3007		3042
3008		3043
3009		3044
3010		3045
3011		3046
3012		3047
3013		3048
3014		3049
3015		3050
3016		3051
3017		3052
3018		3053
3019		3054
3020		3055
3021		3056
3022		3057
3023		3058
3024		3059
3025		3060
3026		3061
3027		3062
3028		3063
3029		3064
3030		3065
3031		3066
3032		3067
3033		3068
3034		3069
3035		3070
3036		3071
3037		3072
3038		3073
3039		3074
3040		3075
3041		3076
3042		3077
3043		3078
3044		3079
3045		3080
3046		3081
3047		3082
3048		3083
3049		3084
3050		3085
3051		3086
3052		3087
3053		3088
3054		3089
3055		3090
3056		3091
3057		3092
3058		3093
3059		3094
3060		3095
3061		3096
3062		3097
3063		3098
3064		3099
3065		3100
3066		3101
3067		3102
3068		3103
3069		3104
3070		3105
3071		3106
3072		3107
3073		3108
3074		3109
3075		3110
3076		3111
3077		3112
3078		3113
3079		3114
3080		3115
3081		3116
3082		3117
3083		3118
3084		3119
3085		3120
3086		3121
3087		3122
3088		3123
3089		3124
3090		3125
3091		3126
3092		3127
3093		3128
3094		3129
3095		3130
3096		3131
3097		3132
3098		3133
3099		3134
3100		3135
3101		3136
3102		31

2849	Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. 2013. Playing Atari with deep reinforcement learning . <i>arXiv preprint</i> .	2906
2850		2907
2851	Mohammed Muqeeth, Haokun Liu, and Colin Raffel. 2023. Soft merging of experts with adaptive routing. <i>arXiv preprint arXiv:2306.03745</i> .	2908
2852		
2853	Yoichi Murakami and Kenji Mizuguchi. 2010. Applying the naïve Bayes classifier with kernel density estimation to the prediction of protein–protein interaction sites . <i>Bioinformatics</i> , 26(15):1841–1848.	2909
2854		2910
2855	Gregory L. Murphy. 1988. Comprehending complex concepts . <i>Cognitive Science</i> , 12(4):529–562.	2911
2856		2912
2857	Moin Nadeem, Anna Bethke, and Siva Reddy. 2020. StereoSet: Measuring stereotypical bias in pretrained language models . <i>arXiv preprint</i> .	2913
2858		2914
2859	Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, and Graham Neubig. 2018. Stress test evaluation for natural language inference . <i>arXiv preprint</i> .	2915
2860		2916
2861	Ramanujapuram Narasimhachar. 1988. <i>History of Kannada Literature: Readership Lectures</i> . Asian Educational Services, New Dehli.	2917
2862		
2863	Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don't give me the details, just the summary! Topic-aware convolutional neural networks for extreme summarization . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1797–1807, Brussels, Belgium. Association for Computational Linguistics.	2918
2864		2919
2865	Ziad S. Nasreddine, Natalie A. Phillips, Valérie Bédirian, Simon Charbonneau, Victor Whitehead, Isabelle Collin, Jeffrey L. Cummings, and Howard Chertkow. 2005. The Montreal Cognitive Assessment, MoCA: A brief screening tool for mild cognitive impairment . <i>Journal of the American Geriatrics Society</i> , 53(4):695–699.	2920
2866		2921
2867	Aida Nematzadeh, Kaylee Burns, Erin Grant, Alison Gopnik, and Tom Griffiths. 2018. Evaluating theory of mind in question answering . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2392–2400, Brussels, Belgium. Association for Computational Linguistics.	2922
2868		2923
2869	Nam Nguyen and Yunsong Guo. 2007. Comparisons of sequence labeling algorithms and extensions . In <i>Proceedings of the 24th International Conference on Machine Learning</i> , page 681–688, New York, NY, USA. Association for Computing Machinery.	2924
2870		
2871	Allen Nie, Erin Bennett, and Noah Goodman. 2019. DisSent: Learning sentence representations from explicit discourse relations . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4497–4510, Florence, Italy. Association for Computational Linguistics.	2925
2872		2926
2873	Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020. Adversarial NLI: A new benchmark for natural language understanding . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> ,	2927
2874		2928
2875		2929
2876		2930
2877		2931
2878		
2879		
2880		
2881		
2882		
2883		
2884		
2885		
2886		
2887		
2888		
2889		
2890		
2891		
2892		
2893		
2894		
2895		
2896		
2897		
2898		
2899		
2900		
2901		
2902		
2903		
2904		
2905		

- 2961 pages 4885–4901, Online. Association for Computational Linguistics. 3017
- 2962
- 2963 Marilyn Nippold, Melissa Allen, and Dixon Kirsch. 3018
- 2964 2001. [Proverb comprehension as a function of reading proficiency in preadolescents](#). *Language Speech* 3019
- 2965 and *Hearing Services in Schools*, 32:90. 3020
- 2966
- 2967 Masaaki Nishino, Sho Takase, Tsutomu Hirao, and 3021
- 2968 Masaaki Nagata. 2019. [Generating natural anagrams: Towards language generation under hard combinatorial constraints](#). 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)*, 3022
- 2969 pages 6408–6412, Hong Kong, China. Association for Computational Linguistics. 3023
- 2970
- 2971
- 2972
- 2973
- 2974
- 2975
- 2976 David Noever, Matt Ciolino, and Josh Kalin. 2020. [The chess transformer: Mastering play using generative language models](#). 3024
- 2977
- 2978
- 2979 David Nolan and Akina Mikami. 2013. ["The things that we have to do": Ethics and instrumentality in 3025](#)
- 2980 humanitarian communication. *Global Media and 3026*
- 2981 *Communication*, 9(1):53–70. 3027
- 2982
- 2983 Klaus Oberauer, Robin Hörnig, Andrea Weidenfeld, 3028
- 2984 and Oliver Wilhelm. 2005. [Effects of directionality 3029](#)
- 2985 in deductive reasoning, II. Premise integration and 3030
- 2986 conclusion evaluation. *The Quarterly Journal of Experimental Psychology Section A*, 58(7):1225–1247. 3031
- 2987
- 2988 Klaus Oberauer and Oliver Wilhelm. 2000. [Effects of 3032](#)
- 2989 directionality in deductive reasoning, I. The comprehension 3033
- 2990 of single relational premises. *Journal of Experimental Psychology: Learning, Memory, and 3034*
- 2991 Cognition
- 2992 , 26(6):1702–1712. 3035
- 2993 The Working Committee on the Revision of the 3036
- 2994 National Standard Occupational Classification. 3037
- 2995 2015. *Standard Occupational Classification of 3038*
- 2996 the People's Republic of China. China Labour 3039
- 2997 and Social Security Publishing House. http://www.jiangmen.gov.cn/bmpd/jmsrlzyhshbzj/zwfw/bmjd/jdks/content/post_2334804.html 3040
- 3000 (Accessed 4 June 2022).
- 3001 Isaac Ong, Amjad Almahairi, Vincent Wu, Wei-Lin Chi- 3041
- 3002 ang, Tianhao Wu, Joseph E. Gonzalez, M Waleed 3042
- 3003 Kadous, and Ion Stoica. 2024. [Routellm: Learning 3043](#)
- 3004 to route llms with preference data. *Preprint*, 3044
- 3005 arXiv:2406.18665.
- 3006 Silviu Oprea and Walid Magdy. 2020. [iSarcasm: A 3045](#)
- 3007 dataset of intended sarcasm. In *Proceedings of the 3046*
- 3008 58th Annual Meeting of the Association for Computational 3047
- 3009 Linguistics
- 3010 , pages 1279–1289, Online. Association for Computational Linguistics. 3048
- 3011 Peter-Michael Osera and Steve Zdancewic. 2015. [Type- 3049](#)
- 3012 and-example-directed program synthesis. In *Proceedings 3050*
- 3013 of the 36th ACM SIGPLAN Conference on Programming 3051
- 3014 Language Design and Implementation, PLDI '15
- 3015 , page 619–630, New York, NY, USA. Association for Computing Machinery. 3052
- 3016
- Oleksiy Ostapenko, Zhan Su, Edoardo Maria Ponti, Laurent Charlin, Nicolas Le Roux, Matheus Pereira, Lucas Caccia, and Alessandro Sordoni. 2024. Towards modular llms by building and reusing a library of loras. *arXiv preprint arXiv:2405.11157*. 3017
- Jahna C. Otterbacher, Dragomir R. Radev, and Airong Luo. 2002. [Revisions that improve cohesion in multi-document summaries: A preliminary study](#). In *Proceedings of the ACL-02 Workshop on Automatic Summarization*, 3018
- 3019 pages 27–44, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics. 3020
- 3021
- Artidoro Pagnoni, Vidhisha Balachandran, and Yulia Tsvetkov. 2021. [Understanding factuality in abstractive summarization with FRANK: A benchmark for factuality metrics](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 3022
- 3023 pages 4812–4829, Online. Association for Computational Linguistics. 3024
- Kartikey Pant and Tanvi Dadu. 2020. [Sarcasm detection using context separators in online discourse](#). *arXiv preprint*. 3025
- 3026
- Tae Jin Park, Naoyuki Kanda, Dimitrios Dimitriadis, Kyu J. Han, Shinji Watanabe, and Shrikanth Narayanan. 2021. [A review of speaker diarization: Recent advances with deep learning](#). *arXiv preprint*. 3027
- 3028
- Arkil Patel, Satwik Bhattacharya, and Navin Goyal. 2021. [Are NLP models really able to solve simple math word problems?](#) In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 3029
- 3030 pages 2080–2094, Online. Association for Computational Linguistics. 3031
- 3032
- 3033
- 3034
- 3035
- Anthony M. Paul. 1970. [Figurative language](#). *Philosophy & Rhetoric*, 3(4):225–248. 3036
- 3037
- Ali Payani and Faramarz Fekri. 2019. [Learning algorithms via neural logic networks](#). *arXiv preprint*. 3038
- Judea Pearl. 1988. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann, San Francisco. 3039
- 3040
- Judea Pearl. 2000. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, Cambridge. 3041
- 3042
- Devin Pelser and Hugh Murrell. 2019. [Deep and dense sarcasm detection](#). *arXiv preprint*. 3043
- Ethan Perez, Douwe Kiela, and Kyunghyun Cho. 2021. [True few-shot learning with language models](#). *arXiv preprint*. 3044
- 3045
- Carla Perez Almendros, Luis Espinosa Anke, and Steven Schockaert. 2020. [Don't patronize me! An annotated dataset with patronizing and condescending language towards vulnerable communities](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, 3046
- 3047 pages 5891–5902, Barcelona, Spain (Online). International Committee 3048
- 3049 on Computational Linguistics. 3050
- 3051

3072	Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H. Miller, and Sebastian Riedel. 2019. Language models as knowledge bases? <i>arXiv preprint</i> .	3127
3073		3128
3074		3129
3075		3130
3076	Dessislava Petrova-Antonova and Rumyana Tancheva. 2020. Data cleaning: A case study with OpenRefine and Trifacta Wrangler . In <i>Quality of Information and Communications Technology</i> , pages 32–40, Cham. Springer.	3131
3077		3132
3078		3133
3079		3134
3080		3135
3081	Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. 2021. AdapterFusion: Non-destructive task composition for transfer learning . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 487–503.	3136
3082		3137
3083		3138
3084		3139
3085		3140
3086		3141
3087	Thang M. Pham, Trung Bui, Long Mai, and Anh Nguyen. 2020. Out of order: How important is the sequential order of words in a sentence in natural language understanding tasks? <i>arXiv preprint</i> .	3142
3088		3143
3089		3144
3090		3145
3091	Steve Piantadosi. 2020. Fleet system .	3146
3092		3147
3093	Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. WiC: The word-in-context dataset for evaluating context-sensitive meaning representations . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 1267–1273, Minneapolis, Minnesota. Association for Computational Linguistics.	3148
3094		3149
3095		3150
3096		3151
3097		3152
3098		3153
3099		3154
3100		3154
3101	Tony A. Plate. 1994. <i>Distributed representations and nested compositional structure</i> . Ph.D. thesis, University of Toronto, Toronto.	3155
3102		3156
3103		3157
3104	Tony A. Plate. 2003. <i>Holographic Reduced Representations: Distributed Representation for Cognitive Structures</i> . CSLI, Stanford, CA.	3158
3105		3159
3106		3160
3107	Robert Plutchik. 1980. A general psychoevolutionary theory of emotion . In Robert Plutchik and Henry Kellerman, editors, <i>Theories of Emotion</i> , pages 3–33. Academic Press.	3161
3108		3162
3109		3163
3110		3164
3111	Nadia Polikarpova, Ivan Kuraj, and Armando Solar-Lezama. 2016. Program synthesis from polymorphic refinement types . In <i>Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI ’16</i> , page 522–538, New York, NY, USA. Association for Computing Machinery.	3165
3112		3166
3113		3167
3114		3168
3115		3169
3116		3170
3117		3171
3118	Stanislas Polu and Ilya Sutskever. 2020. Generative language modeling for automated theorem proving . <i>arXiv preprint</i> .	3172
3119		3173
3120		3174
3121	Edoardo Maria Ponti, Alessandro Sordoni, Yoshua Bengio, and Siva Reddy. 2023. Combining parameter-efficient modules for task-level generalisation. In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 687–702.	3175
3122		3176
3123		3177
3124		3178
3125		3179
3126		3180
3127	Simone Paolo Ponzetto and Michael Strube. 2007. Knowledge derived from Wikipedia for computing semantic relatedness . <i>Journal of Artificial Intelligence Research</i> , 30:181–212.	3181
3128		3181
3129		3181
3130		3181

3182	Quora, Inc. 2017. <i>Quora question pairs</i> .	3237
3183	Dragomir R. Radev, Lori Levin, and Thomas E. Payne.	3238
3184	2008. <i>The North American computational linguistics olympiad (NACLO)</i> . In <i>Proceedings of the Third Workshop on Issues in Teaching Computational Linguistics</i> , pages 87–96, Columbus, Ohio. Association for Computational Linguistics.	3239
3185		
3186		
3187		
3188		
3189	Alec Radford, Jeff Wu, Rewon Child, David Luan,	3240
3190	Dario Amodei, and Ilya Sutskever. 2019. <i>Language</i>	3241
3191	<i>models are unsupervised multitask learners</i> . <i>OpenAI blog</i> , 1(8):9.	3242
3192		3243
3193	Kira Radinsky, Eugene Agichtein, Evgeniy Gabrilovich,	3244
3194	and Shaul Markovitch. 2011. <i>A word at a time: Computing word relatedness using temporal semantic analysis</i> . In <i>WWW ’11: Proceedings of the 20th International Conference on World Wide Web</i> , page 337–346, New York, NY, USA. Association for Computing Machinery.	3245
3195		3246
3196		3247
3197		3248
3198		3249
3199		3250
3200	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. <i>Exploring the limits of transfer learning with a unified text-to-text transformer</i> . <i>Journal of Machine Learning Research</i> , 21:1–67.	3251
3201		3252
3202		3253
3203		3254
3204		
3205		
3206	Alessandro Raganato, Jose Camacho-Collados, and Roberto Navigli. 2017. <i>Word sense disambiguation: A unified evaluation framework and empirical comparison</i> . In <i>Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers</i> , pages 99–110, Valencia, Spain. Association for Computational Linguistics.	3255
3207		3256
3208		3257
3209		3258
3210		
3211		
3212		
3213		
3214	Altaf Rahman and Vincent Ng. 2012. <i>Resolving complex cases of definite pronouns: The Winograd schema challenge</i> . In <i>Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning</i> , pages 777–789, Jeju Island, Korea. Association for Computational Linguistics.	3261
3215		3262
3216		
3217		
3218		
3219		
3220		
3221	Sunny Rai and Shampa Chakraverty. 2020. <i>A survey on computational metaphor processing</i> . <i>ACM Comput. Surv.</i> , 53(2).	3263
3222		3264
3223		3265
3224	Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. <i>Know what you don’t know: Unanswerable questions for SQuAD</i> . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 784–789, Melbourne, Australia. Association for Computational Linguistics.	3266
3225		3267
3226		3268
3227		3269
3228		3270
3229		3271
3230		3272
3231	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. <i>SQuAD: 100,000+ questions for machine comprehension of text</i> . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2383–2392, Austin, Texas. Association for Computational Linguistics.	3273
3232		
3233		
3234		
3235		
3236		
3237		
3238		
3239		
3240	Prajit Ramachandran and Quoc V. Le. 2019. <i>Diversity and depth in per-example routing models</i> . In <i>International Conference on Learning Representations</i> .	3237
3241		3238
3242		3239
3243		
3244	Alexandre Ramé, Kartik Ahuja, Jianyu Zhang, Matthieu Cord, Léon Bottou, and David Lopez-Paz. 2022. Recycling diverse models for out-of-distribution generalization. <i>arXiv preprint arXiv:2212.10445</i> .	3240
3245		3241
3246		3242
3247		3243
3248	Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghab Gupta, and Pranav Khaitan. 2020. <i>Towards scalable multi-domain conversational agents: The schema-guided dialogue dataset</i> . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 34, pages 8689–8696, Menlo Park, CA. Association for the Advancement of Artificial Intelligence.	3244
3249		3245
3250		3246
3251	Ian Ravenscroft. 2019. Folk psychology as a theory. In Edward N. Zalta, editor, <i>The Stanford Encyclopedia of Philosophy</i> , Summer 2019 edition. Metaphysics Research Lab, Stanford University.	3247
3252		3248
3253		3249
3254		3250
3255	Siva Reddy, Danqi Chen, and Christopher D. Manning.	3251
3256	2019. <i>CoQA: A conversational question answering challenge</i> . <i>Transactions of the Association for Computational Linguistics</i> , 7:249–266.	3252
3257		3253
3258		3254
3259	Scott Reed and Nando de Freitas. 2015. <i>Neural programmer-interpreters</i> . <i>arXiv preprint</i> .	3255
3260		3256
3261	Marek Rei. 2017. <i>Semi-supervised multitask learning for sequence labeling</i> . <i>arXiv preprint</i> .	3257
3262		
3263	Nils Reimers and Iryna Gurevych. 2019. <i>Sentence-BERT: Sentence embeddings using siamese BERT-networks</i> . <i>arXiv preprint</i> .	3258
3264		
3265		
3266	Joseph Reisinger and Raymond J. Mooney. 2010. <i>Multi-prototype vector-space models of word meaning</i> . In <i>Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics</i> , pages 109–117, Los Angeles, California. Association for Computational Linguistics.	3259
3267		3260
3268		3261
3269		3262
3270		3263
3271		3264
3272		3265
3273	He Ren and Quan Yang. 2017. <i>Neural joke generation</i> .	3266
3274		3267
3275		3268
3276		3269
3277		3270
3278		3271
3279	Philip Resnik. 1995. <i>Using information content to evaluate semantic similarity in a taxonomy</i> . In <i>IJCAI’95: Proceedings of the 14th International Joint Conference on Artificial Intelligence, Volume 1</i> , page 448–453, San Francisco. Morgan Kaufmann.	3272
3280		3273
3281		3274
3282		3275
3283	Philip Resnik. 1999. <i>Semantic similarity in a taxonomy: An information-based measure and its application to problems of ambiguity in natural language</i> . <i>Journal of Artificial Intelligence Research</i> , 11:95–130.	3276
3284		3277
3285		3278
3286	Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S Gordon. 2011a. Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In <i>2011 AAAI Spring Symposium Series</i> .	3279
3287		3280
3288		3281
3289		3282
3290	Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S. Gordon. 2011b. <i>Choice of plausible alternatives: An evaluation of commonsense causal reasoning</i> . <i>AAAI Spring Symposium</i> .	3283
3291		3284
3292		3285
3293		3286
3294		
3295		
3296		
3297		
3298		
3299		
3300		

3291	Anna Rogers, Olga Kovaleva, and Anna Rumshisky.	Phillip Rust, Jonas Pfeiffer, Ivan Vulić, Sebastian Ruder, and Iryna Gurevych.	3344
3292	2020. Getting closer to AI complete question answering: A set of prerequisite real tasks . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 123–150, Online. Association for Computational Linguistics.		3345
3293			3346
3294			3347
3295			
3296			
3297	Alexis Ross and Ellie Pavlick. 2019. How well do NLI models capture verb veridicality? In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2230–2240, Hong Kong, China. Association for Computational Linguistics.	Amrita Saha, Rahul Aralikatte, Mitesh M. Khapra, and Karthik Sankaranarayanan. 2018. DuoRC: Towards Complex Language Understanding with Paraphrased Reading Comprehension. In <i>Meeting of the Association for Computational Linguistics (ACL)</i> .	3348
3298			3349
3299			3350
3300			3351
3301			3352
3302			
3303			
3304	Kenneth J. Rothman and Sander Greenland. 2005. Causation and causal inference in epidemiology . <i>American Journal of Public Health</i> , 95(S1):S144–S150.	Gözde Gül Şahin, Yova Kementchedjhieva, Phillip Rust, and Iryna Gurevych. 2020. Puzzling Machines: A challenge on learning from small data . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1241–1254, Online. Association for Computational Linguistics.	3353
3305			3354
3306			3355
3307	Joshua Rozner, Christopher Potts, and Kyle Mahowald. 2021. Decrypting cryptic crosswords: Semantically complex wordplay puzzles as a target for NLP . <i>arXiv preprint</i> .		3356
3308			3357
3309			3358
3310			
3311	Sebastian Ruder, Joachim Bingel, Isabelle Augenstein, and Anders Søgaard. 2017. Latent multi-task architecture learning . In <i>AAAI Conference on Artificial Intelligence</i> .	Keisuke Sakaguchi, Kevin Duh, Matt Post, and Benjamin Van Durme. 2016. Robust word recognition via semi-character recurrent neural network . <i>arXiv preprint</i> .	3359
3312			3360
3313			3361
3314			3362
3315	Rachel Rudinger, Jason Naradowsky, Brian Leonard, and Benjamin Van Durme. 2018. Gender bias in coreference resolution . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)</i> , pages 8–14, New Orleans, Louisiana. Association for Computational Linguistics.	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2020a. WinoGrande: An adversarial Winograd schema challenge at scale . In <i>Proceedings of the Thirty-Fourth AAAI Conference on Artificial Intelligence</i> , pages 8732–8740, New York, NY, USA. Association for the Advancement of Artificial Intelligence.	3363
3316			3364
3317			3365
3318			3366
3319			3367
3320			3368
3321			3369
3322			
3323	Rachel Etta Rudolph and Alexander W. Kocurek. 2020. Comparing conventions . In <i>Proceedings of Semantics and Linguistic Theory</i> , pages 294–313, Washington, D.C. Linguistic Society of America.	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2020b. WINOGRANDE: An adversarial Winograd schema challenge at scale . In <i>The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI-20</i> , pages 8732–8734, Menlo Park, CA. Association for the Advancement of Artificial Intelligence.	3370
3324			3371
3325			3372
3326			3373
3327	Rosa Rugani, Giorgio Vallortigara, Konstantinos Priftis, and Lucia Regolin. 2015. Number-space mapping in the newborn chick resembles humans' mental number line . <i>Science</i> , 347(6221):534–536.	María del Pilar Salas-Zárate, Mario Andrés Paredes-Valverde, Miguel Ángel Rodríguez-García, Rafael Valencia-García, and Giner Alor-Hernández. 2017. Automatic detection of satire in Twitter: A psycholinguistic-based approach . <i>Knowledge-Based Systems</i> , 128:20–33.	3374
3328			3375
3329			3376
3330			
3331	Joshua S. Rule. 2020. The child as hacker: Building more human-like models of learning . Ph.D. thesis, Massachusetts Institute of Technology, Cambridge, MA.	Julian Salazar, Davis Liang, Toan Q. Nguyen, and Katrin Kirchhoff. 2020. Masked language model scoring . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 2699–2712, Online. Association for Computational Linguistics.	3377
3332			3378
3333			3379
3334			3380
3335	Joshua S. Rule, Joshua B. Tenenbaum, and Steven T. Piantadosi. 2020. The child as hacker . <i>Trends in Cognitive Sciences</i> , 24(11):900–915.	Suresh Kumar Sanampudi and G. Vijaya Kumari. 2010. Temporal reasoning in natural language processing: A survey . <i>International Journal of Computer Applications</i> , 1(4):53–57.	3381
3336			3382
3337			
3338	D. E. Rumelhart, J. L. McClelland, and PDP Research Group, editors. 1986. <i>Parallel Distributed Processing. Volume 1: Foundations</i> . MIT Press, Cambridge, MA.	Evan Sandhaus. 2008. The New York Times annotated corpus LDC2008T19 . <i>Linguistic Data Consortium</i> .	3383
3339			3384
3340			3385
3341			3386
3342	Stuart J. Russell and Peter Norvig. 2002. Artificial Intelligence: A Modern Approach . Pearson, Hoboken.	Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, Manan Dey, M. Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma, Eliza Szczęchla, Taewoon	3387
3343			3388

3400	Kim, Gunjan Chhablani, Nihal V. Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesh Sharma, Andrea Santilli, Thibault Févry, Jason Alan Fries, Ryan Teehan, Stella Biderman, Leo Gao, Tali Bers, Thomas Wolf, and Alexander M. Rush. 2022. Multitask prompted training enables zero-shot task generalization . In <i>The Tenth International Conference on Learning Representations</i> .	3458
3401		3459
3402	Megan Scudellari. 2017. Cryopreservation aims to engineer novel ways to freeze, store, and thaw organs . <i>Proceedings of the National Academy of Sciences</i> , 114(50):13060–13062.	3460
3403		3461
3404	Abigail See, Peter J. Liu, and Christopher D. Manning. 2017a. Get to the point: Summarization with pointer-generator networks . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics</i> .	3462
3405		3463
3406		3464
3407		3465
3408		3466
3409	Abigail See, Peter J. Liu, and Christopher D. Manning. 2017b. Get to the point: Summarization with pointer-generator networks . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1073–1083, Vancouver, Canada. Association for Computational Linguistics.	3467
3410		3468
3411	Thibault Sellam, Dipanjan Das, and Ankur P. Parikh. 2020. BLEURT: Learning robust metrics for text generation . <i>arXiv preprint</i> .	3469
3412		3470
3413		3471
3414		3472
3415		3473
3416	Daniel Selsam, Matthew Lamm, Benedikt Bünz, Percy Liang, Leonardo de Moura, and David L. Dill. 2018. Learning a SAT solver from single-bit supervision . <i>arXiv preprint</i> .	3474
3417		3475
3418		3476
3419	Lutfi Kerem Senel and Hinrich Schütze. 2021. Does he wink or does he nod? A challenging benchmark for evaluating word understanding of language models . <i>arXiv preprint</i> .	3477
3420		3478
3421		3479
3422		3480
3423	Zhuoqi Chen, Ming Tang, and Yuxin Chen. 2021. A simple neural network module for relational reasoning . <i>arXiv preprint</i> .	3481
3424		3482
3425	Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Jurafsky, Noah A. Smith, and Yejin Choi. 2020. Social bias frames: Reasoning about social and power implications of language . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5477–5490, Online. Association for Computational Linguistics.	3483
3426		3484
3427		3485
3428		3486
3429		3487
3430	Luzi Sennhauser and Robert C. Berwick. 2018. Evaluating the ability of LSTMs to learn context-free grammars . <i>arXiv preprint</i> .	3488
3431		3489
3432	Rico Sennrich. 2016. How grammatical is character-level neural machine translation? Assessing MT quality with contrastive translation pairs . <i>arXiv preprint</i> .	3490
3433		3491
3434		3492
3435		3493
3436	Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015. Neural machine translation of rare words with subword units . <i>arXiv preprint</i> .	3494
3437		3495
3438		3496
3439		3497
3440		3498
3441		3499
3442	Rico Sennrich and Biao Zhang. 2019. Revisiting low-resource neural machine translation: A case study . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 211–221, Florence, Italy. Association for Computational Linguistics.	3500
3443		3501
3444	Min Joon Seo, Hannaneh Hajishirzi, Ali Farhadi, and Oren Etzioni. 2014. Diagram understanding in geometry questions . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 28, Menlo Park, CA. Association for the Advancement of Artificial Intelligence.	3502
3445		3503
3446		3504
3447		3505
3448	Usman Shahid and Elena Zheleva. 2021. Counterfactual learning in networks: An empirical study of model dependence .	3506
3449		3507
3450		3508
3451	Janelle Shane. 2020. All your questions answered. AI Weirdness .	3509
3452		3510
3453		
3454		
3455		
3456		
3457		

3511	David Elliot Shaw, William R. Swartout, and C. Cordell Green. 1975. Inferring LISP programs from examples . In <i>IJCAI'75: Proceedings of the 4th International Joint Conference on Artificial Intelligence</i> , volume 1, pages 260–267. Artificial Intelligence Laboratory, Cambridge, MA.	3565
3512		3566
3513		3567
3514		3568
3515		3569
3516		3570
3517	Ekaterina Shutova. 2010. Automatic metaphor interpretation as a paraphrasing task . In <i>Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics</i> , pages 1029–1037, Los Angeles, California. Association for Computational Linguistics.	3571
3518		
3519		
3520		
3521		
3522	Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. 2016. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer . In <i>International Conference on Learning Representations</i> .	3572
3523		3573
3524		3574
3525		3575
3526		3576
3527	Ekaterina Shutova and Simone Teufel. 2010. Metaphor corpus annotated for source-target domain mappings . In <i>Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)</i> , Valletta, Malta. European Language Resources Association.	3577
3528		
3529		
3530		
3531	Damien Sileo, Tim Van De Cruys, Camille Pradel, and Philippe Muller. 2019. Mining discourse markers for unsupervised sentence representation learning . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 3477–3486, Minneapolis, Minnesota. Association for Computational Linguistics.	3578
3532		3579
3533		3580
3534		3581
3535		3582
3536		3583
3537		3584
3538		3585
3539		3586
3540	Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer . In <i>International Conference on Learning Representations</i> .	3587
3541		
3542		
3543	Shannon Zejiang Shen, Hunter Lang, Bailin Wang, Yoon Kim, and David Sontag. 2024. Learning to decode collaboratively with multiple language models. <i>arXiv preprint arXiv:2403.03870</i> .	3588
3544		
3545		
3546		
3547		
3548		
3549		
3550		
3551	Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. 2019. The woman worked as a babysitter: On biases in language generation . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3407–3412, Hong Kong, China. Association for Computational Linguistics.	3589
3552		
3553		
3554		
3555		
3556	Shaoyun Shi, Hanxiong Chen, Weizhi Ma, Jiaxin Mao, Min Zhang, and Yongfeng Zhang. 2020. Neural logic reasoning . <i>arXiv preprint</i> .	3590
3557		
3558		
3559		
3560		
3561		
3562	Han-Chin Shing, Suraj Nair, Ayah Zirikly, Meir Friedenberg, Hal Daumé III, and Philip Resnik. 2018. Expert, crowdsourced, and machine assessment of suicide risk via online postings . In <i>Proceedings of the Fifth Workshop on Computational Linguistics and Clinical Psychology: From Keyboard to Clinic</i> , pages 25–36, New Orleans, LA. Association for Computational Linguistics.	3591
3563		
3564		
3565		
3566		
3567		
3568		
3569		
3570		
3571	Damien Sileo, Tim Van de Cruys, Camille Pradel, and Philippe Muller. 2020. DiscSense: Automated semantic analysis of discourse markers . In <i>Proceedings of the 12th Language Resources and Evaluation Conference</i> , pages 991–999, Marseille, France. European Language Resources Association.	3592
3572		
3573		
3574		
3575		
3576		
3577		
3578	Damien Sileo, Wout Vossen, and Robbe Raymaekers. 2022. Zero-shot recommendation as language modeling . In <i>Advances in Information Retrieval: 44th European Conference on IR Research, ECIR 2022, Stavanger, Norway, April 10–14, 2022, Proceedings, Part II</i> , page 223–230, Cham. Springer.	3593
3579		3594
3580		3595
3581		3596
3582		3597
3583		3598
3584		
3585		
3586		
3587		
3588		
3589		
3590		
3591		
3592		
3593	Gurdeep Singh, Kaustubh Dhole, Priyadarshini P. Pai, and Sukanta Mondal. 2014. SPRINGS: Prediction of protein-protein interaction sites using artificial neural networks . <i>Journal of Proteomics & Computational Biology</i> , 1:7.	3599
3594		3600
3595		3601
3596		3602
3597		3603
3598		
3599	Rishabh Singh and Sumit Gulwani. 2015. Predicting a correct program in programming by example . In <i>Computer Aided Verification</i> , pages 398–414, Cham. Springer International Publishing.	3604
3600		3605
3601		3606
3602		3607
3603		
3604	Rishabh Singh and Sumit Gulwani. 2016. Transforming spreadsheet data types using examples . In <i>Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL '16</i> , page 343–356, New York, NY, USA. Association for Computing Machinery.	3608
3605		3609
3606		3610
3607		3611
3608		3612
3609		3613
3610		
3611		
3612		
3613		
3614	Shikhar Singh, Nuan Wen, Yu Hou, Pegah Alipoormabashi, Te-lin Wu, Xuezhe Ma, and Nanyun Peng. 2021. COM2SENSE: A commonsense reasoning benchmark with complementary sentences . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 883–898, Online. Association for Computational Linguistics.	3614
3615		3615
3616		3616
3617		3617
3618		3618
3619		3619
3620		3620

3621	Koustuv Sinha, Shagun Sodhani, Jin Dong, Joelle Pineau, and William L. Hamilton. 2019. CLUTRR: A diagnostic benchmark for inductive reasoning from text. <i>arXiv preprint</i> .	3678
3622		3679
3623		3680
3624		3681
3625	Natalia Skachkova, Thomas Trost, and Dietrich Klakow.	3682
3626	2018. Closing brackets with recurrent neural networks. In <i>Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 232–239, Brussels, Belgium. Association for Computational Linguistics.	3683
3627		3684
3628		3685
3629		3686
3630		3687
3631	Douglas R. Smith. 1984. The synthesis of LISP programs from examples: A survey. In Alan W. Biermann, Gerhard Guiho, and Yves Kodratoff, editors, <i>Automatic Program Construction Techniques</i> , pages 307–324. Macmillan, New York.	3688
3632		3689
3633		3690
3634		3691
3635		3692
3636	Paul Smolensky. 1988. On the proper treatment of connectionism. <i>Behavioral and Brain Sciences</i> , 11(1):1–23.	3693
3637		3694
3638		3695
3639	Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In <i>Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing</i> , pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.	3696
3640		3697
3641		3698
3642		3699
3643		3700
3644		3701
3645		3702
3646		3703
3647	Irene Solaiman, Miles Brundage, Jack Clark, Amanda Askell, Ariel Herbert-Voss, Jeff Wu, Alec Radford, Gretchen Krueger, Jong Wook Kim, Sarah Kreps, Miles McCain, Alex Newhouse, Jason Blazakis, Kris McGuffie, and Jasmine Wang. 2019. Release strategies and the social impacts of language models. <i>arXiv preprint</i> .	3704
3648		3705
3649		3706
3650		3707
3651		3708
3652		3709
3653		3710
3654	Mahmoud Soltani Firouz, Ali Farahmandi, and Soleiman Hosseinpour. 2021. Early detection of freeze damage in navel orange fruit using nondestructive low intensity ultrasound coupled with machine learning. <i>Food Analytical Methods</i> , 14:1140–1149.	3711
3655		3712
3656		3713
3657		3714
3658		3715
3659	Dan Sperber and Deirdre Wilson. 2002. Pragmatics, modularity and mind-reading. <i>Mind & Language</i> , 17(1-2):3–23.	3716
3660		3717
3661		3718
3662	Peter Spirtes, Clark Glymour, and Richard Scheines. 2000. <i>Causation, Prediction, and Search</i> . MIT Press, Cambridge, MA.	3719
3663		3720
3664		3721
3665	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew Dai, Andrew La, Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong,	3722
3666		3723
3667		3724
3668		3725
3669		3726
3670		3727
3671		3728
3672		3729
3673		3730
3674		3731
3675		3732
3676		3733
3677		3734
3678	Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakaş, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartłomiej Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Howald, Bryan Orinon, Cameron Diao, Cameron Dour, Catherine Stinson, Cedrick Argueta, César Ferri Ramírez, Chandan Singh, Charles Rathkopf, Chenlin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Chris Waites, Christian Voigt, Christopher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman, Dan Roth, Daniel Freeman, Daniel Khashabi, Daniel Levy, Daniel Moseguí González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilyar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ellie Pavlick, Emanuele Rodola, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engefú Manyasi, Evgenii Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca Happé, Francois Chollet, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Giambattista Parascandolo, Giorgio Mariani, Gloria Wang, Gonzalo Jaimovich-López, Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Shevlin, Hinrich Schütze, Hiromu Yakura, Hongming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geissinger, Jackson Kernion, Jacob Hilton, Jaehoon Lee, Jaime Fernández Fisac, James B. Simon, James Koppel, James Zheng, James Zou, Jan Kocoń, Jana Thompson, Janelle Wingfield, Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova, Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse Engel, Jesujoba Alabi, Jiacheng Xu, Jiaming Song, Jillian Tang, Joan Waweru, John Burden, John Miller, John U. Balis, Jonathan Batchelder, Jonathan Berant, Jörg Frohberg, Jos Rozen, Jose Hernandez-Orallo, Joseph Boudeman, Joseph Guerr, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chua, Kamil Kanclerz, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja Markert, Kaustubh D. Dhole, Kevin Gimbel, Kevin Omundi, Kory Mathewson, Kristen Chiaffullo, Ksenia Shkaruta, Kumar Shridhar, Kyle McDonell, Kyle Richardson, Laria Reynolds, Leo Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-Ochando, Louis-Philippe Morency, Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheng	3741

3742	He, Luis Oliveros Colón, Luke Metz, Lütfi Kerem	3813
3743	Şenel, Maarten Bosma, Maarten Sap, Maartje ter	3814
3744	Hoeve, Maheen Farooqi, Manaal Faruqui, Mantas	3815
3745	Mazeika, Marco Baturan, Marco Marelli, Marco	3816
3746	Maru, Maria Jose Ramírez Quintana, Marie Tolkiehn,	3817
3747	Mario Giulianelli, Martha Lewis, Martin Potthast,	3818
3748	Matthew L. Leavitt, Matthias Hagen, Mátyás Schubert,	3819
3749	Medina Orduna Baitemirova, Melody Arnaud,	3820
3750	Melvin McElrath, Michael A. Yee, Michael Cohen,	3821
3751	Michael Gu, Michael Ivanitskiy, Michael Starratt,	3822
3752	Michael Strube, Michał Śwędrowski, Michele Bevilacqua, Michihiro Yasunaga, Mihir Kale, Mike	3823
3753	Cain, Mimeo Xu, Mirac Suzgun, Mitch Walker,	3824
3754	Mo Tiwari, Mohit Bansal, Moin Aminnaseri, Mor	3825
3755	Geva, Mozhdeh Gheini, Mukund Varma T, Nanyun	3826
3756	Peng, Nathan A. Chi, Nayeon Lee, Neta Gur-Ari	3827
3757	Krakover, Nicholas Cameron, Nicholas Roberts,	3828
3758	Nick Doiron, Nicole Martinez, Nikita Nangia, Niklas	3829
3759	Deckers, Niklas Muennighoff, Nitish Shirish Keskar,	3830
3760	Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan	3831
3761	Wen, Oliver Zhang, Omar Agha, Omar Elbaghdadi,	3832
3762	Omer Levy, Owain Evans, Pablo Antonio Moreno	3833
3763	Casares, Parth Doshi, Pascale Fung, Paul Pu Liang,	3834
3764	Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao,	3835
3765	Percy Liang, Peter Chang, Peter Eckersley, Phu Mon	3836
3766	Htut, Pinyu Hwang, Piotr Miłkowski, Piyush Patil,	3837
3767	Pouya Pezeshkpour, Priti Oli, Qiaozhu Mei, Qing	3838
3768	Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta	3839
3769	Rudolph, Raefer Gabriel, Rahel Habacker, Ramon	3840
3770	Risco, Raphaël Millière, Rhythm Garg, Richard	3841
3771	Barnes, Rif A. Saurous, Riku Arakawa, Robbe	3842
3772	Raymaekers, Robert Frank, Rohan Sikand, Roman	3843
3773	Novak, Roman Sitelew, Ronan LeBras, Rosanne	3844
3774	Liu, Rowan Jacobs, Rui Zhang, Ruslan Salakhutdinov,	3845
3775	Ryan Chi, Ryan Lee, Ryan Stovall, Ryan	3846
3776	Teehan, Rylan Yang, Sahib Singh, Saif M. Mohammad,	3847
3777	Sajant Anand, Sam Dillavou, Sam Shleifer,	3848
3778	Sam Wiseman, Samuel Gruetter, Samuel R. Bowman,	3849
3779	Samuel S. Schoenholz, Sanghyun Han, Sanjeev Kwatra,	3850
3780	Sarah A. Rous, Sarik Ghazarian, Sayan	3851
3781	Ghosh, Sean Casey, Sebastian Bischoff, Sebastian	3852
3782	Gehrman, Sebastian Schuster, Sepideh Sadeghi,	3853
3783	Shadi Hamdan, Sharon Zhou, Shashank Srivastava,	3854
3784	Sherry Shi, Shikhar Singh, Shima Asaadi, Shixiang	3855
3785	Shane Gu, Shubh Pachchigar, Shubham Toshniwal,	3856
3786	Shyam Upadhyay, Shyamolima, Debnath,	3857
3787	Siamak Shakeri, Simon Thormeyer, Simone Melzi,	3858
3788	Siva Reddy, Sneha Priscilla Makini, Soo-Hwan Lee,	3859
3789	Spencer Torene, Sriharsha Hatwar, Stanislas De-	3860
3790	haene, Stefan Divic, Stefano Ermon, Stella Bider-	3861
3791	man, Stephanie Lin, Stephen Prasad, Steven T. Pi-	3862
3792	antadosi, Stuart M. Shieber, Summer Mishergi,	3863
3793	Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen,	3864
3794	Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsu Hashimoto,	3865
3795	Te-Lin Wu, Théo Desbordes, Theodore Rothschild,	3866
3796	Thomas Phan, Tianle Wang, Tiberius Nkonyili, Timo	3867
3797	Schick, Timofei Kornev, Titus Tunduny, Tobias Ger-	3868
3798	stenberg, Trenton Chang, Trishala Neeraj, Tushar	3869
3799	Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera	3870
3800	Demberg, Victoria Nyamai, Vikas Raunak, Vinay	3871
3801	Ramasesh, Vinay Uday Prabhu, Vishakh Padmumar,	3872
3802	Vivek Srikumar, William Fedus, William Saunders,	3873
3803	William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong,	3874
3804	Xinran Zhao, Xinyi Wu, Xudong Shen,	3875
3805	Yadollah Yaghoobzadeh, Yair Lakretz, Yangqiu Song,	3876
	Yasaman Bahri, Yejin Choi, Yichi Yang, Yiding	3877
	Hao, Yifu Chen, Yonatan Belinkov, Yu Hou, Yufang	3878
	Hou, Yuntao Bai, Zachary Seid, Zhuoye Zhao, Zijian	3879
	Wang, Zijie J. Wang, Zirui Wang, and Ziyi Wu. 2023.	3880
	Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. <i>Preprint</i> ,	3881
	arXiv:2206.04615.	3882
	Shashank Srivastava, Snigdha Chaturvedi, and Tom	3883
	Mitchell. 2016. Inferring interpersonal relations in narrative summaries. In <i>Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, AAAI'16</i> , page 2807–2813, Menlo Park, CA. Association for the Advancement of Artificial Intelligence.	3884
	Robert Stalnaker. 1978. Assertion. In P. Cole, editor,	3885
	<i>Pragmatics, Syntax and Semantics</i> 9, pages 315–332.	3886
	Brill, Leiden.	3887
	Trevor Standley, Amir Zamir, Dawn Chen, Leonidas	3888
	Guibas, Jitendra Malik, and Silvio Savarese. 2020.	3889
	Which tasks should be learned together in multi-task	3890
	learning? In <i>International conference on machine learning</i> , pages 9120–9132. PMLR.	3891
	Gabriel Stanovsky, Noah A. Smith, and Luke Zettlemoyer. 2019. Evaluating gender bias in machine translation. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 1679–1684, Florence, Italy. Association for Computational Linguistics.	3892
	Gerard J. Steen, Aletta G. Dorst, J. Berenike Herrmann,	3893
	Anna A. Kaal, Tina Krennmayr, and Tryntje Pasma.	3894
	2010. <i>A Method for Linguistic Metaphor Identification: From MIP to MIPVU.</i> Converging Evidence in Language and Communication Research 14. John Benjamins, Amsterdam.	3895
	Bernd Steinbach and Roman Kohut. 2002. Neural networks – a model of boolean functions. <i>5th International Workshop on Boolean Problems, Freiburg, Sept. 2002.</i>	3896
	Sebastian U. Stich. 2018. Local sgd converges	3897
	fast and communicates little. <i>arXiv preprint arXiv:1805.09767.</i>	3898
	Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M.	3899
	Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,	3900
	Dario Amodei, and Paul Christiano. 2020. Learning to summarize from human feedback. <i>arXiv preprint.</i>	3901
	Kevin Stowe, Leonardo Ribeiro, and Iryna Gurevych.	3902
	2020. Metaphoric paraphrase generation. <i>arXiv preprint.</i>	3903
	Michael Strube and Simone Paolo Ponzetto. 2006.	3904
	Wikirelate! Computing semantic relatedness using Wikipedia. In <i>AAAI'06: Proceedings of the 21st National Conference on Artificial Intelligence</i> , volume 2,	3905
	page 1419–1424. Association for the Advancement of Artificial Intelligence.	3906

3860	Saku Sugawara, Hikaru Yokono, and Akiko Aizawa.	3915
3861	2017. Prerequisite skills for reading comprehension: Multi-perspective analysis of MCTest datasets and systems. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 31, Menlo Park, CA. Association for the Advancement of Artificial Intelligence.	3916
3862		3917
3863		3918
3864		
3865		
3866		
3867	Oyvind Tafjord, Peter Clark, Matt Gardner, Wen-tau Yih, and Ashish Sabharwal. 2018. Quarel: A dataset and models for answering questions about qualitative relationships. <i>CoRR</i> , abs/1811.08048.	3919
3868		3920
3869		3921
3870		3922
3871	Oyvind Tafjord, Matt Gardner, Kevin Lin, and Peter Clark. 2019. "quartz: An open-domain dataset of qualitative relationship questions". <i>EMNLP</i> .	3923
3872		3924
3873		3925
3874	Alon Talmor, Yanai Elazar, Yoav Goldberg, and Jonathan Berant. 2019a. oLMpics – on what language model pre-training captures. <i>arXiv preprint</i> .	3926
3875		3927
3876		3928
3877	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019b. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 3929	3929
3878	4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.	3930
3879		3931
3880		3932
3881		3933
3882	Derek Tam, Mohit Bansal, and Colin Raffel. 2023. Merging by matching models in task subspaces. <i>arXiv preprint arXiv:2312.04339</i> .	3934
3883		3935
3884		3936
3885	Alex Tamkin, Miles Brundage, Jack Clark, and Deep Ganguli. 2021. Understanding the capabilities, limitations, and societal impact of large language models. <i>arXiv preprint</i> .	3937
3886		3938
3887	Jiwei Tan, Xiaojun Wan, and Jianguo Xiao. 2015. Learning to recommend quotes for writing. In <i>Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, AAAI'15</i> , page 3939	3939
3888	2453–2459, Menlo Park, CA. Association for the Advancement of Artificial Intelligence.	3940
3889		3941
3890		3942
3891	Niket Tandon, Bhavana Dalvi, Keisuke Sakaguchi, Peter Clark, and Antoine Bosselut. 2019. WIQA: A dataset for “what if...” reasoning over procedural text. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3943	3943
3892	6076–6085, Hong Kong, China. Association for Computational Linguistics.	3944
3893		3945
3894		3946
3895	Anke Tang, Li Shen, Yong Luo, Nan Yin, Lefei Zhang, and Dacheng Tao. 2024. Merging multi-task models via weight-ensembling mixture of experts. <i>Preprint</i> , 3947	3947
3896	arXiv:2402.00433.	3948
3897		3949
3898	Jan Arne Telle, José Hernández-Orallo, and Cèsar Ferri. 2019. The teaching size: Computable teachers and learners for universal languages. <i>Machine Learning</i> ,	3950
3899	108:1653–1675.	3951
3900		3952
3901		3953
3902		3954
3903		3955
3904		3956
3905	Mirac Suzgun, Sebastian Gehrmann, Yonatan Belinkov, and Stuart Shieber. 2019b. Memory-augmented recurrent neural networks can learn generalized Dyck languages. <i>arXiv preprint</i> .	3957
3906		3958
3907		3959
3908		
3909	Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, and Jason Wei. 2022. Challenging big-bench tasks and whether chain-of-thought can solve them. <i>arXiv preprint arXiv:2210.09261</i> .	
3910		
3911		
3912		
3913		
3914		

3968 Damien Teney, Ehsan Abbasnedjad, and Anton van den
3969 Hengel. 2020. [Learning what makes a difference](#)
3970 from counterfactual examples and gradient supervision.

3971 *arXiv preprint*.

3972 Paul Thagard, Keith J. Holyoak, Greg Nelson, and
3973 David Gochfeld. 1990. [Analog retrieval by constraint](#)
3974 [satisfaction](#). *Artificial Intelligence*, 46(3):259–310.

3975 Avijit Thawani, Jay Pujara, Filip Ilievski, and Pedro
3976 Szekely. 2021. [Representing numbers in NLP: A](#)
3977 [survey and a vision](#). In *Proceedings of the 2021*
3978 *Conference of the North American Chapter of the*
3979 *Association for Computational Linguistics: Human*
3980 *Language Technologies*, pages 644–656, Online. As-
3981 *sociation for Computational Linguistics*.

3982 Jesse Thomason, Shiqi Zhang, Raymond Mooney, and
3983 Peter Stone. 2015. [Learning to interpret natural](#)
3984 [language commands through human-robot dialog](#).
3985 In *Proceedings of the Twenty-Fourth International*
3986 *Joint Conference on Artificial Intelligence, IJCAI-15*,
3987 pages 1923–1929.

3988 Judith Jarvis Thomson. 1976. [Killing, letting die, and](#)
3989 [the trolley problem](#). *The Monist*, 59(2):204–217.

3990 Poonam B. Thorat, Rajeshwari M. Goudar, and Sunita
3991 Barve. 2015. [Survey on collaborative filtering,](#)
3992 [content-based filtering and hybrid recommendation](#)
3993 [system](#). *International Journal of Computer Applica-*
3994 *tions*, 110:31–36.

3995 James Thorne, Andreas Vlachos, Christos
3996 Christodoulopoulos, and Arpit Mittal. 2018.
3997 [FEVER: A large-scale dataset for fact extraction](#)
3998 and [VERification](#). In *Proceedings of the 2018*
3999 *Conference of the North American Chapter of*
4000 *the Association for Computational Linguistics:*
4001 *Human Language Technologies, Volume 1 (Long*
4002 *Papers)*, pages 809–819, New Orleans, Louisiana.
4003 *Association for Computational Linguistics*.

4004 Xiaoyu Tong. 2021. [Metaphor paraphrasing and word](#)
4005 [sense disambiguation: Toward a new approach to](#)
4006 [automated metaphor](#).

4007 Xiaoyu Tong, Ekaterina Shutova, and Martha Lewis.
4008 2021. [Recent advances in neural metaphor process-](#)
4009 [ing: A linguistic, cognitive and social perspective](#).
4010 In *Proceedings of the 2021 Conference of the North*
4011 *American Chapter of the Association for Compu-*
4012 *tational Linguistics: Human Language Technologies*,
4013 pages 4673–4686, Online. *Association for Compu-*
4014 *tational Linguistics*.

4015 Shubham Toshniwal, Sam Wiseman, Karen Livescu,
4016 and Kevin Gimpel. 2021. [Learning chess blindfolded:](#)
4017 [Evaluating language models on state tracking](#). *arXiv*
4018 *preprint*.

4019 David Toubiana, Nir Sade, Lifeng Liu, Maria del
4020 Mar Rubio Wilhelmi, Yariv Brotnman, Urszula
4021 Luzarowska, John P. Vogel, and Eduardo Blumwald.
4022 2020. [Correlation-based network analysis combined](#)
4023 [with machine learning techniques highlight the role](#)

4024 of the gaba shunt in *brachypodium sylvaticum* freez-
4025 ing tolerance. *Scientific Reports*, 10:no. 4489.

4026 Andrew Trask, Felix Hill, Scott Reed, Jack Rae, Chris
4027 Dyer, and Phil Blunsom. 2018. [Neural arithmetic](#)
4028 [logic units](#). *arXiv preprint*.

4029 George Tsatsaronis, Iraklis Varlamis, and Michalis
4030 Vazirgiannis. 2010. [Text relatedness based on a word](#)
4031 [thesaurus](#). *Journal of Artificial Intelligence Research*,
4032 37(1):1–40.

4033 George Tsatsaronis, Iraklis Varlamis, Michalis Vazir-
4034 giannis, and Kjetil Nørvåg. 2009. [Omiotis: A](#)
4035 [thesaurus-based measure of text relatedness](#). In
4036 *Machine Learning and Knowledge Discovery in*
4037 *Databases*, pages 742–745, Berlin. Springer.

4038 Yulia Tsvetkov, Leonid Boytsov, Anatole Gershman,
4039 Eric Nyberg, and Chris Dyer. 2014. [Metaphor detec-](#)
4040 [tion with cross-lingual model transfer](#). In *Proceed-*
4041 *ings of the 52nd Annual Meeting of the Association*
4042 *for Computational Linguistics (Volume 1: Long Pa-*
4043 *pers)*, pages 248–258, Baltimore, Maryland. *Associa-*
4044 *tion for Computational Linguistics*.

4045 Shikhar Vashishth, Shib Sankar Dasgupta,
4046 Swayambhu Nath Ray, and Partha Talukdar.
4047 2018. [Dating documents using graph convolution](#)
4048 [networks](#). In *Proceedings of the 56th Annual*
4049 *Meeting of the Association for Computational*
4050 *Linguistics (Volume 1: Long Papers)*, pages
4051 1605–1615, Melbourne, Australia. *Association for*
4052 *Computational Linguistics*.

4053 Siddharth Vashishtha, Adam Poliak, Yash Kumar Lal,
4054 Benjamin Van Durme, and Aaron Steven White. 2020.
4055 [Temporal reasoning in natural language inference](#).
4056 In *Findings of the Association for Computational*
4057 *Linguistics: EMNLP 2020*, pages 4070–4078, Online.
4058 *Association for Computational Linguistics*.

4059 Siddharth Vashishtha, Benjamin Van Durme, and
4060 Aaron Steven White. 2019. [Fine-grained temporal](#)
4061 [relation extraction](#). *arXiv preprint*.

4062 Oleg Vasilev, Vedant Dharnidharka, and John Bohan-
4063 non. 2020. [Fill in the BLANC: Human-free quality](#)
4064 [estimation of document summaries](#). In *Proceedings*
4065 *of the First Workshop on Evaluation and Comparison*
4066 *of NLP Systems*, pages 11–20, Online. *Association*
4067 *for Computational Linguistics*.

4068 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
4069 Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
4070 Kaiser, and Illia Polosukhin. 2017. Attention is all
4071 you need. *Advances in neural information processing*
4072 *systems*, 30.

4073 Jette Viethen and Robert Dale. 2008. [The use of spa-](#)
4074 [tial relations in referring expression generation](#). In
4075 *Proceedings of the Fifth International Natural Lan-*
4076 *guage Generation Conference*, pages 59–67, Salt
4077 Fork, Ohio, USA. *Association for Computational*
4078 *Linguistics*.

4079	Oriol Vinyals, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H. Choi, Richard Powell, Timo Ewalds, Petko Georgiev, Junhyuk Oh, Dan Horgan, Manuel Kroiss, Ivo Danihelka, Aja Huang, Laurent Sifre, Trevor Cai, John P. Agapiou, Max Jaderberg, Alexander S. Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin Dalibard, David Budden, Yury Sulsky, James Molloy, Tom L. Paine, Caglar Gulcehre, Ziyu Wang, Tobias Pfaff, Yuhuai Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKinney, Oliver Smith, Tom Schaul, Timothy Lillicrap, Koray Kavukcuoglu, Demis Hassabis, Chris Apps, and David Silver. 2019. Grandmaster level in StarCraft II using multi-agent reinforcement learning. <i>Nature</i> , 575:350–354.	4136
4080		4137
4081		4138
4082		4139
4083		4140
4084		4141
4085		
4086		
4087		
4088		
4089		
4090		
4091		
4092		
4093		
4094	Ellen M. Voorhees. 2002. Overview of the trec 2002 question answering track. In <i>Proceedings of the Eleventh Text REtrieval Conference (TREC 2002)</i> .	4142
4095		4143
4096		
4097	Tu Vu, Tong Wang, Tsendsuren Munkhdalai, Alessandro Sordoni, Adam Trischler, Andrew Mattarella-Mickey, Subhransu Maji, and Mohit Iyyer. 2020. Exploring and predicting transferability across nlp tasks. <i>arXiv preprint arXiv:2005.00770</i> .	4144
4098		4145
4099		4146
4100		4147
4101		4148
4102	Henning Wachsmuth, Nona Naderi, Yufang Hou, Yonatan Bilu, Vinodkumar Prabhakaran, Tim Alberdingk Thijm, Graeme Hirst, and Benno Stein. 2017. Computational argumentation quality assessment in natural language. In <i>Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers</i> , pages 176–187, Valencia, Spain. Association for Computational Linguistics.	4149
4103		4150
4104		4151
4105		
4106		
4107		
4108		
4109		
4110		
4111	Eric Wallace, Florian Tramèr, Matthew Jagielski, and Ariel Herbert-Voss. 2020. Does GPT-2 know your phone number? <i>Berkeley Artificial Intelligence Research blog</i> .	4159
4112		4160
4113		4161
4114		4162
4115	Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020a. Asking and answering questions to evaluate the factual consistency of summaries. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5008–5020, Online. Association for Computational Linguistics.	4163
4116		4164
4117		4165
4118		4166
4119		4167
4120		4168
4121	Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019a. Superglue: A stickier benchmark for general-purpose language understanding systems. <i>Advances in neural information processing systems</i> , 32.	4169
4122		4170
4123		4171
4124		4172
4125		4173
4126		
4127	Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019b. SuperGLUE: A stickier benchmark for general-purpose language understanding systems. In <i>Advances in Neural Information Processing Systems</i> , volume 32. Curran Associates, Inc.	4174
4128		4175
4129		4176
4130		4177
4131		4178
4132		4179
4133		4180
4134	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE:	4181
4135		4182
		4183
		4184
		4185
		4186
	Zijian Wang and David Jurgens. 2018. It’s going to be okay: Measuring access to support in online communities. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 33–45, Brussels, Belgium. Association for Computational Linguistics.	4187
		4188
		4189
		4190

4191	9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3711–3719, Hong Kong, China. Association for Computational Linguistics.	4247
4192		4248
4193		4249
4194		4250
4195	Ulrike Willinger, Andreas Hergovich, Michaela Schmoeger, Matthias Deckert, Susanne Stoettner, Iris Bunda, Andrea Witting, Melanie Seidler, Reinhilde Moser, Stefanie Kacena, David Jaekle, Benjamin Loader, Christian Mueller, and Eduard Auff. 2017. Cognitive and emotional demands of black humour processing: The role of intelligence, aggressiveness and mood. <i>Cognitive Processing</i> , 18:159–167.	4251
4196		4252
4197		4253
4198		4254
4199		4255
4200		4256
4201	Ingmar Weber and Alejandro Jaimes. 2011. Who uses web search for what: And how. In <i>Proceedings of the Fourth ACM International Conference on Web Search and Data Mining, WSDM '11</i> , page 15–24, New York, NY, USA. Association for Computing Machinery.	4257
4202		4258
4203		4259
4204	David Wechsler. 2008. <i>Wechsler Adult Intelligence Scale—Fourth Edition (WAIS–IV)</i> . Pearson, San Antonio.	4260
4205		4261
4206		4262
4207	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. 2022a. Finetuned language models are zero-shot learners. In <i>International Conference on Learning Representations</i> .	4263
4208		4264
4209	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc V. Le, and Denny Zhou. 2022b. Chain of thought prompting elicits reasoning in large language models. <i>arXiv preprint arXiv:2201.11903</i> .	4265
4210		4266
4211		4267
4212		4268
4213	Ludwig Wittgenstein. 1953. <i>Philosophical investigations</i> . Basil Blackwell, Oxford.	4269
4214		4270
4215		4271
4216	Thomas Wolf. 2019. Some additional experiments extending the tech report “assessing BERT’s syntactic abilities” by Yoav Goldberg.	4272
4217		4273
4218		4274
4219	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierrette Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2019. Huggingface’s transformers: State-of-the-art natural language processing. <i>arXiv preprint</i> .	4275
4220		4276
4221		4277
4222		4278
4223	Johannes Welbl, Pontus Stenetorp, and Sebastian Riedel. 2018. Constructing datasets for multi-hop reading comprehension across documents.	4279
4224		4280
4225		4281
4226		4282
4227	Orion Weller and Kevin Seppi. 2019. Humor detection: A transformer gets the last laugh. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3621–3625, Hong Kong, China. Association for Computational Linguistics.	4283
4228		4284
4229		4285
4230	Jason Weston, Antoine Bordes, Sumit Chopra, Alexander M. Rush, Bart van Merriënboer, Armand Joulin, and Tomas Mikolov. 2015. Towards AI-complete question answering: A set of prerequisite toy tasks. <i>arXiv preprint</i> .	4286
4231		4287
4232		4288
4233		4289
4234	Aaron Steven White, Rachel Rudinger, Kyle Rawlins, and Benjamin Van Durme. 2018. Lexicosyntactic inference in neural models. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 4717–4724, Brussels, Belgium. Association for Computational Linguistics.	4290
4235		4291
4236		4292
4237	Sarah White, Elisabeth Hill, Francesca Happé, and Uta Frith. 2009. Revisiting the strange stories: Revealing mentalizing impairments in autism. <i>Child Development</i> , 80(4):1097–1117.	4293
4238		4294
4239		4295
4240		4296
4241	Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.	4297
4242		4298
4243		4299
4244		4300
4245		4301
4246	Bo Wu, Pedro Szekely, and Craig A. Knoblock. 2012. Learning data transformation rules through examples: Preliminary results. In <i>Proceedings of the Ninth International Workshop on Information Integration on the Web, IIWeb ’12</i> , New York, NY, USA. Association for Computing Machinery.	4302
4247		4303
4248		4304

4305	Chengyue Wu, Teng Wang, Yixiao Ge, Zeyu Lu, Ruisong Zhou, Ying Shan, and Ping Luo. 2023. <i>pi-</i> tuning: Transferring multimodal foundation mod- els with optimal multi-task interpolation. In <i>Inter- national Conference on Machine Learning</i> , pages 37713–37727. PMLR.	4361
4306		4362
4307		4363
4308		4364
4309		4365
4310		4366
4311	Shijie Wu, Ryan Cotterell, and Mans Hulden. 2020. <i>Ap- plying the transformer to character-level transduction.</i> <i>arXiv preprint</i> .	4367
4312		4368
4313		
4314	Xun Wu, Shaohan Huang, and Furu Wei. 2024. <i>Mix- ture of LoRA experts</i> . In <i>The Twelfth International Conference on Learning Representations</i> .	4369
4315		4370
4316		4371
4317	Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Łukasz Kaiser, Stephan Gouws, Yoshiakiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2016. <i>Google’s neural machine translation system: Bridging the gap between human and machine trans- lation</i> . <i>arXiv preprint</i> .	4372
4318		4373
4319		4374
4320		4375
4321		4376
4322		4377
4323		4378
4324		
4325		
4326		
4327		
4328		
4329	Kevin Xia, Kai-Zhan Lee, Yoshua Bengio, and Elias Bareinboim. 2021. <i>The causal-neural connection: Expressiveness, learnability, and inference</i> . <i>arXiv preprint</i> .	4379
4330		4380
4331		4381
4332		4382
4333	Yijun Xiao and William Yang Wang. 2021. <i>On hal- lucination and predictive uncertainty in conditional language generation</i> . <i>arXiv preprint</i> .	4383
4334		
4335		
4336	Jing Xu, Da Ju, Margaret Li, Y-Lan Boureau, Jason Weston, and Emily Dinan. 2020a. <i>Recipes for safety in open-domain chatbots</i> . <i>arXiv preprint</i> .	4384
4337		4385
4338		4386
4339	Jingwei Xu, Junyu Lai, and Yunpeng Huang. 2024. Me- teora: Multiple-tasks embedded lora for large lan- guage models. <i>arXiv preprint arXiv:2405.13053</i> .	4387
4340		4388
4341		4389
4342	Silei Xu, Sina Semnani, Giovanni Campagna, and Mon- ica Lam. 2020b. <i>AutoQA: From databases to QA semantic parsers with only synthetic training data</i> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 422–434, Online. Association for Computa- tional Linguistics.	4390
4343		4391
4344		4392
4345		
4346		
4347		
4348		
4349	Yang Xu, Jiawei Liu, Wei Yang, and Liusheng Huang. 2018. <i>Incorporating latent meanings of morpholog- ical compositions to enhance word embeddings</i> . In <i>Proceedings of the 56th Annual Meeting of the As- sociation for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1232–1242, Melbourne, Aus- tralia. Association for Computational Linguistics.	4393
4350		4394
4351		4395
4352		4396
4353		4397
4354		4398
4355		4399
4356	Linting Xue, Aditya Barua, Noah Constant, Rami Al- Rfou, Sharan Narang, Mihir Kale, Adam Roberts, and Colin Raffel. 2021a. <i>ByT5: Towards a token-free future with pre-trained byte-to-byte models</i> . <i>arXiv preprint</i> .	4400
4357		4401
4358		4402
4359		4403
4360		
4361	Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021b. <i>mT5: A massively multilingual pre-trained text-to-text transformer</i> . In <i>Proceedings of the 2021 Conference of the North American Chap- ter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 483–498, On- line. Association for Computational Linguistics.	4404
4362		4405
4363		4406
4364		4407
4365		4408
4366		4409
4367		
4368		
4369	Prateek Yadav, Leshem Choshen, Colin Raffel, and Mo- hit Bansal. 2023a. <i>Compeft: Compression for com- municating parameter efficient updates via sparsifica- tion and quantization</i> . <i>Preprint</i> , arXiv:2311.13171.	4370
4370		4371
4371		4372
4372		
4373	Prateek Yadav, Colin Raffel, Mohammed Muqeeth, Lucas Caccia, Haokun Liu, Tianlong Chen, Mohit Bansal, Leshem Choshen, and Alessandro Sordoni. 2024. A survey on model moerging: Recycling and routing among specialized experts for collaborative learning. <i>arXiv preprint arXiv:2408.07057</i> .	4374
4374		4375
4375		4376
4376		4377
4377		4378
4378		
4379	Prateek Yadav, Derek Tam, Leshem Choshen, Colin Raffel, and Mohit Bansal. 2023b. <i>TIES-merging: Resolving interference when merging models</i> . In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	4380
4380		4381
4381		4382
4382		4383
4383		
4384	Xinru Yan and Ted Pedersen. 2017. <i>Who’s to say what’s funny? A computer using language models and deep learning, that’s who!</i> <i>arXiv preprint</i> .	4385
4385		4386
4386		
4387	Diyi Yang, Alon Lavie, Chris Dyer, and Eduard Hovy. 2015. <i>Humor recognition and humor anchor extrac- tion</i> . In <i>Proceedings of the 2015 Conference on Em- pirical Methods in Natural Language Processing</i> , pages 2367–2376, Lisbon, Portugal. Association for Computational Linguistics.	4388
4388		4389
4389		4390
4390		4391
4391		4392
4392		
4393	Enneng Yang, Zhenyi Wang, Li Shen, Shiwei Liu, Guib- ing Guo, Xingwei Wang, and Dacheng Tao. 2023. Adamerging: Adaptive model merging for multi-task learning. <i>arXiv preprint arXiv:2310.02575</i> .	4394
4394		4395
4395		4396
4396		
4397	Kaiyu Yang and Jia Deng. 2019. <i>Learning to prove theorems via interacting with proof assistants</i> . <i>arXiv preprint</i> .	4398
4398		4399
4399		
4400	Scott Cheng-Hsin Yang and Patrick Shafto. 2017. <i>Ex- plainable artificial intelligence via Bayesian teaching</i> . <i>Workshop on Teaching Machines, Robots, and Hu- mans, NIPS 2017</i> .	4401
4401		4402
4402		4403
4403		
4404	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Ben- gio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. <i>Hotpotqa: A dataset for diverse, explainable multi-hop question answer- ing</i> . In <i>Proceedings of the 2018 Conference on Em- pirical Methods in Natural Language Processing</i> .	4405
4405		4406
4406		4407
4407		4408
4408		4409
4409		
4410	Qinyuan Ye, Juan Zha, and Xiang Ren. 2022. Elic- iting and understanding cross-task skills with task-level mixture-of-experts. <i>arXiv preprint arXiv:2205.12701</i> .	4411
4411		4412
4412		4413
4413		

4414	Eric Yeh, Daniel Ramage, Christopher D. Manning, Eneko Agirre, and Aitor Soroa. 2009. WikiWalk: Random walks on Wikipedia for semantic relatedness . In <i>Proceedings of the 2009 Workshop on Graph-based Methods for Natural Language Processing (TextGraphs-4)</i> , pages 41–49, Suntec, Singapore. Association for Computational Linguistics.	4471
4415		4472
4416		4473
4417		
4418		
4419		
4420		
4421	Yelp, Inc. 2018. Yelp open dataset .	
4422		
4423	Yang Yi, Yih Wen-tau, and Christopher Meek. 2015. WikiQA: A Challenge Dataset for Open-Domain Question Answering . Association for Computational Linguistics, page 2013–2018.	
4424		
4425		
4426	Wenpeng Yin and Yulong Pei. 2015. Optimizing sentence modeling and selection for document summarization . In <i>Proceedings of the 24th International Conference on Artificial Intelligence, IJCAI-15</i> , page 1383–1389.	
4427		
4428		
4429		
4430		
4431	Tao Yu, Rui Zhang, Heyang Er, Suyi Li, Eric Xue, Bo Pang, Xi Victoria Lin, Yi Chern Tan, Tianze Shi, Zihan Li, Youxuan Jiang, Michihiro Yasunaga, Sungrok Shim, Tao Chen, Alexander Fabbri, Zifan Li, Luyao Chen, Yuwen Zhang, Shreya Dixit, Vincent Zhang, Caiming Xiong, Richard Socher, Walter Lasecki, and Dragomir Radev. 2019a. CoSQL: A conversational text-to-SQL challenge towards cross-domain natural language interfaces to databases . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 1962–1979, Hong Kong, China. Association for Computational Linguistics.	
4432		
4433		
4434		
4435		
4436		
4437		
4438		
4439		
4440		
4441		
4442		
4443		
4444		
4445		
4446	Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. 2018. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-SQL task . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3911–3921, Brussels, Belgium. Association for Computational Linguistics.	
4447		
4448		
4449		
4450		
4451		
4452		
4453		
4454		
4455	Tao Yu, Rui Zhang, Michihiro Yasunaga, Yi Chern Tan, Xi Victoria Lin, Suyi Li, Heyang Er, Irene Li, Bo Pang, Tao Chen, Emily Ji, Shreya Dixit, David Proctor, Sungrok Shim, Jonathan Kraft, Vincent Zhang, Caiming Xiong, Richard Socher, and Dragomir Radev. 2019b. SParC: Cross-domain semantic parsing in context . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4511–4523, Florence, Italy. Association for Computational Linguistics.	
4456		
4457		
4458		
4459		
4460		
4461		
4462		
4463		
4464		
4465	Weihao Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. 2020. ReClor: A reading comprehension dataset requiring logical reasoning . <i>CoRR</i> , arXiv:2002.04326.	
4466		
4467		
4468	Xiang Yu, Ngoc Thang Vu, and Jonas Kuhn. 2019c. Learning the Dyck language with attention-based Seq2Seq models . In <i>Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 138–146, Florence, Italy. Association for Computational Linguistics.	
4469		
4470		
4471		
4472		
4473		
4474	Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yuet-ting Zhuang, and Dacheng Tao. 2019d. ActivityNet-QA: A dataset for understanding complex web videos via question answering . <i>arXiv preprint</i> .	
4475		
4476		
4477		
4478	Eliezer Yudkowsky. 2008. Artificial intelligence as a positive and negative factor in global risk . In Nick Bostrom and Milan M. Ćirković, editors, <i>Global Catastrophic Risks</i> , pages 308–345. Oxford University Press, Oxford.	
4479		
4480		
4481		
4482		
4483	Ted Zadouri, Ahmet Üstün, Arash Ahmadian, Beyza Ermiş, Aycı̄ Locatelli, and Sara Hooker. 2023. Pushing mixture of experts to the limit: Extremely parameter efficient moe for instruction tuning. <i>arXiv preprint arXiv:2309.05444</i> .	
4484		
4485		
4486		
4487		
4488	Amir Zamir, Alexander Sax, Bokui (William) Shen, Leonidas J. Guibas, Jitendra Malik, and Silvio Savarese. 2018. Taskonomy: Disentangling task transfer learning . In <i>2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 3712–3722.	
4489		
4490		
4491		
4492		
4493		
4494	Wojciech Zaremba and Ilya Sutskever. 2014. Learning to execute . <i>arXiv preprint</i> .	
4495		
4496	Poorya Zaremoddi, Wray L. Buntine, and Gholamreza Haffari. 2018. Adaptive knowledge sharing in multi-task learning: Improving low-resource neural machine translation . In <i>Annual Meeting of the Association for Computational Linguistics</i> .	
4497		
4498		
4499		
4500		
4501	Omnia Zayed, John Philip McCrae, and Paul Buitelaar. 2020. Figure me out: A gold standard dataset for metaphor interpretation . In <i>Proceedings of the 12th Language Resources and Evaluation Conference</i> , pages 5810–5819, Marseille, France. European Language Resources Association.	
4502		
4503		
4504		
4505		
4506		
4507	Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2018. From recognition to cognition: Visual commonsense reasoning . <i>arXiv preprint</i> .	
4508		
4509		
4510	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019a. HellaSwag: Can a machine really finish your sentence? <i>arXiv preprint arXiv:1905.07830</i> .	
4511		
4512		
4513		
4514	Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. 2019b. Defending against neural fake news . <i>arXiv preprint</i> .	
4515		
4516		
4517		
4518	Zihao Zeng, Yibo Miao, Hongcheng Gao, Hao Zhang, and Zhijie Deng. 2024. Adamoe: Token-adaptive routing with null experts for mixture-of-experts language models . <i>Preprint</i> , arXiv:2406.13233.	
4519		
4520		
4521		
4522	Dongxiang Zhang, Lei Wang, Luming Zhang, Bing Tian Dai, and Heng Tao Shen. 2020a. The gap of semantic parsing: A survey on automatic math word problem	
4523		
4524		

4525		solvers. <i>IEEE Transactions on Pattern Analysis & Machine Intelligence</i> , 42(09):2287–2305.	4580
4526			4581
4527		Haoran Zhang, Amy X. Lu, Mohamed Abdalla, Matthew McDermott, and Marzyeh Ghassemi. 2020b. Hurtful words: Quantifying biases in clinical contextual word embeddings . In <i>Proceedings of the ACM Conference on Health, Inference, and Learning, CHIL '20</i> , page 110–120, New York, NY, USA. Association for Computing Machinery.	4582
4528			4583
4529			
4530			
4531			
4532			
4533			
4534		Hongming Zhang, Xinran Zhao, and Yangqiu Song. 2020c. WinoWhy: A deep diagnosis of essential commonsense knowledge for answering Winograd schema challenge . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5736–5745, Online. Association for Computational Linguistics.	4588
4535			4589
4536			4590
4537			4591
4538			4592
4539			
4540			
4541		Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. 2019a. Multi-agent reinforcement learning: A selective overview of theories and algorithms . <i>arXiv preprint</i> .	4593
4542			4594
4543			4595
4544			4596
4545			4597
4546		Li Zhang, Qing Lyu, and Chris Callison-Burch. 2020d. Reasoning about goals, steps, and temporal ordering with WikiHow . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 4630–4639, Online. Association for Computational Linguistics.	4598
4547			4599
4548			4600
4549			4601
4550		Meishan Zhang, Yue Zhang, and Guohong Fu. 2016. Tweet sarcasm detection using deep neural network . In <i>Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers</i> , pages 2449–2460, Osaka, Japan. The COLING 2016 Organizing Committee.	4602
4551			4603
4552			4604
4553			4605
4554			4606
4555			
4556		Sheng Zhang, Xiaodong Liu, Jingjing Liu, Jianfeng Gao, Kevin Duh, and Benjamin Van Durme. 2018a. Record: Bridging the gap between human and machine commonsense reading comprehension . <i>arXiv preprint arXiv:1810.12885</i> .	4607
4557			4608
4558			4609
4559			4610
4560			4611
4561		Shiwei Zhang, Xiuzhen Zhang, Jeffrey Chan, and Paolo Rosso. 2019b. Irony detection via sentiment-based transfer learning . <i>Information Processing & Management</i> , 56(5):1633–1644.	4612
4562			4613
4563			4614
4564			4615
4565		Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification . In <i>Advances in Neural Information Processing Systems</i> .	4616
4566			4617
4567			4618
4568			
4569		Yan Zhang, Jonathon Hare, and Adam Prügel-Bennett. 2018b. Learning to count objects in natural images for visual question answering . <i>arXiv preprint</i> .	4619
4570			4620
4571			4621
4572			4622
4573			4623
4574		Yuan Zhang, Jason Baldridge, and Luheng He. 2019c. PAWS: Paraphrase Adversaries from Word Scrambling . In <i>Proc. of NAACL</i> .	4624
4575			
4576			
4577			
4578		Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. Gender bias in coreference resolution: Evaluation and debiasing methods . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for</i>	
4579			
		<i>Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)</i> , pages 15–20, New Orleans, Louisiana. Association for Computational Linguistics.	
		Tony Z. Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models . <i>arXiv preprint</i> .	4584
			4585
			4586
			4587
		Xinyu Zhao, Guoheng Sun, Ruisi Cai, Yukun Zhou, Pingzhi Li, Peihao Wang, Bowen Tan, Yexiao He, Li Chen, Yi Liang, et al. 2024a. Model-glue: Democratized llm scaling for a large model zoo in the wild . <i>arXiv preprint arXiv:2410.05357</i> .	4588
			4589
		Ziyu Zhao, Leilei Gan, Guoyin Wang, Wangchunshu Zhou, Hongxia Yang, Kun Kuang, and Fei Wu. 2024b. Loraretriever: Input-aware lora retrieval and composition for mixed tasks in the wild . <i>Preprint</i> , arXiv:2402.09997.	4590
			4591
			4592
		Ben Zhou, Daniel Khashabi, Qiang Ning, and Dan Roth. 2019. "Going on a vacation" takes longer than "going for a walk": A study of temporal commonsense understanding . <i>arXiv preprint</i> .	4593
			4594
			4595
			4596
			4597
		Chunting Zhou, Graham Neubig, Jiatao Gu, Mona Diab, Paco Guzman, Luke Zettlemoyer, and Marjan Ghazvininejad. 2020. Detecting hallucinated content in conditional neural sequence generation . <i>arXiv preprint</i> .	4602
			4603
			4604
			4605
			4606
		Jing Zhou, Zongyu Lin, Yanan Zheng, Jian Li, and Zhilin Yang. 2022. Not all tasks are born equal: Understanding zero-shot generalization . In <i>The Eleventh International Conference on Learning Representations</i> .	4607
			4608
			4609
			4610
			4611
		Haichao Zhu, Li Dong, Furu Wei, Wenhui Wang, Bing Qin, and Ting Liu. 2019. Learning to ask unanswerable questions for machine reading comprehension . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4238–4248, Florence, Italy. Association for Computational Linguistics.	4612
			4613
			4614
			4615
			4616
			4617
			4618
		Xiaojin Zhu. 2015. Machine teaching: An inverse problem to machine learning and an approach toward optimal education . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 29, Menlo Park, CA. Association for the Advancement of Artificial Intelligence.	4619
			4620
			4621
			4622
			4623
			4624
		Zhemin Zhu, Delphine Bernhard, and Iryna Gurevych. 2010. Sentence simplification by monolingual machine translation . In <i>Proceedings of the 5th International Conference on Natural Language Processing</i> .	4625
			4626
			4627
			4628
		Alan Zucconi. 6 Jan. 2016. The secrets of colour interpolation .	4629
			4630

4631 Appendix

4632 A Extended Related Work

4633 **Model Merging.** Model merging (Yadav et al.,
4634 2023b; Choshen et al., 2022; Wortsman et al., 2022;
4635 Ramé et al., 2022; Matena and Raffel, 2022; Ilharco
4636 et al., 2022; Tam et al., 2023; Jin et al., 2022; Yang
4637 et al., 2023; Zhao et al., 2024a) consolidates mul-
4638 tiple independently trained models with identical
4639 architectures into a unified model that preserves
4640 multi-model capabilities. While simple parameter
4641 averaging suffices for models within a linearly con-
4642 nected low-loss parameter space (McMahan et al.,
4643 2017; Stich, 2018; Frankle et al., 2020; Wortsman
4644 et al., 2021; Li et al., 2023), more sophisticated
4645 techniques are necessary for complex scenarios.
4646 For instance, task vectors facilitate merging expert
4647 models trained on diverse domains (Ilharco et al.,
4648 2022). Additionally, methods like weighted merg-
4649 ing using Fisher Importance Matrices (Matena and
4650 Raffel, 2022; Tam et al., 2023) and TIES-Merging,
4651 which addresses sign disagreements and redun-
4652 dancy (Yadav et al., 2023b) offers improved per-
4653 formance. As a non-adaptive expert aggregation
4654 method, merging serves as a fundamental baseline
4655 for numerous Model Editing with Regularization
4656 (MoErging) techniques.

4657 **Multitask Learning (MTL)** research offers val-
4658 uable insights for decentralized development. Notably,
4659 investigations into task-relatedness (Standley
4660 et al., 2020; Bingel and Søgaard, 2017; Achille
4661 et al., 2019; Vu et al., 2020; Zamir et al., 2018;
4662 Mou et al., 2016) provide guidance for design-
4663 ing routing mechanisms, while MTL architectures
4664 addressing the balance between shared and task-
4665 specific knowledge (Misra et al., 2016; Ruder et al.,
4666 2017; Meyerson and Miikkulainen, 2017; Zare-
4667 moodi et al., 2018; Sun et al., 2019b) offer strate-
4668 gies for combining expert contributions in a decen-
4669 tralized manner.

4670 **MoE for Multitask Learning.** Recent research
4671 has extensively investigated mixture-of-experts
4672 (MoE) models for multitask learning, achieving
4673 promising results in unseen task generalization.
4674 These approaches generally fall into two categories:
4675 (1) Example Routing: Studies like Muqeeth et al.
4676 (2023); Zadouri et al. (2023); Wang et al. (2022a)
4677 train routers to dynamically select experts for each
4678 input, while Caccia et al. (2023) demonstrate the
4679 efficacy of routing at a finer granularity by splitting
4680 expert parameters into blocks. (2) Task Routing:

4681 Ponti et al. (2023) employs a trainable skill ma-
4682 trix to assign tasks to specific parameter-efficient
4683 modules, while Gupta et al. (2022) leverages task-
4684 specific routers selected based on domain knowl-
4685 edge. Ye et al. (2022) proposes a layer-wise expert
4686 selection mechanism informed by task represen-
4687 tations derived from input embeddings. Such ap-
4688 proaches leverage task-specific representation to
4689 allow the router to effectively select the most suit-
4690 able experts for unseen tasks. While these studies
4691 differ from our setting by assuming simultaneous
4692 data access, they offer valuable insights applicable
4693 to our exploration of creating routing mechanisms
4694 over expert models.

4695 B LLM for Task Instruction Generation.

4696 B.1 Prompt Template

4697 We use the following prompt with 3 randomly se-
4698 lected samples for each task to generate its descrip-
4699 tion. The prompt is then fed into the gpt-4-turbo
4700 OpenAI API to get the generated task descriptions.

*The following are three pairs of input-output
examples from one task. Generate the task
instruction in one sentence that is most
possibly used to command a language
model to produce them. In the instruction,
remember to point out the skill or knowledge
required for the task to guide the language
model.*

- Input:

- Output:

- Input:

- Output:

- Input:

- Output:

4701 B.2 Examples of the Generated Instructions

4702 We provide several examples of LLM-generated
4703 instructions in this section.

4704 **WikiBio** (Lebret et al., 2016a) (T0 Held-In):

- 4705 • Create a short biography using the provided
4706 facts, demonstrating knowledge in historical
4707 and biographical writing.
- 4708 • Write a short biography based on the given
4709 factual bullet points, demonstrating profi-

4711	<p>ciency in summarizing and transforming structured data into coherent narrative text.</p>	4755
4712		4756
4713	<p>CommonGen (Lin et al., 2020b) (T0 Held-In):</p> <ul style="list-style-type: none">• Generate a coherent sentence using all the given abstract concepts, requiring the skill of concept integration to form a meaningful sentence.• Generate a coherent sentence by creatively combining a given set of abstract concepts.	4757
4714		4758
4715		4759
4716		4760
4717		4761
4718		4762
4719		4763
4720	<p>COPA (Huang et al., 2024b) (T0 Held-Out):</p> <ul style="list-style-type: none">• Identify the most logically consistent sentence from two given options based on the provided context, demonstrating reasoning and causal relationship skills.• Generate the most likely outcome for a given scenario by choosing between two provided options based on contextual clues and causal reasoning.	4764
4721		4765
4722		4766
4723		4767
4724		4768
4725		4769
4726		4770
4727		4771
4728		4772
4729	<p>Date Understanding (Srivastava et al., 2023) (BigBench-Hard):</p> <ul style="list-style-type: none">• Calculate the date based on the given information and present it in MM/DD/YYYY format, ensuring that you accurately account for day, month, and year changes.	4773
4730		4774
4731		4775
4732		4776
4733		4777
4734		4778
4735	<p>Hindu Mythology Trivia (Srivastava et al., 2023) (BigBench-Lite):</p> <ul style="list-style-type: none">• Generate the correct answer by making use of your knowledge in Hindu mythology and culture.	4779
4736		4780
4737		4781
4738		4782
4739		4783
4740	<h2>C Demonstrating Compositional Generation</h2>	4784
4741		4785
4742	<p>In addition to significant improvements on held-in tasks, GLIDER demonstrates strong performance on held-out tasks, showcasing its generalization capability. To further examine this ability to handle unseen tasks by composing experts, we provide specific task examples illustrating the association between selected experts and the evaluated task. As Figure 2 shows, GLIDER primarily selects two experts for the COPA (T0 held-out) task, corresponding to CosmosQA and QuaRel. The following three examples from these tasks demonstrate their close semantic relationship:</p>	4786
4743		4787
4744		4788
4745		4789
4746		4790
4747		4791
4748		4792
4749		4793
4750		4794
4751		4795
4752		4796
4753		4797
4754	<ul style="list-style-type: none">• COPA:	
4755	<ul style="list-style-type: none">– <u>Question</u>: Everyone in the class turned to stare at the student. Select the most plausible cause: - The student’s phone rang. - The student took notes.– <u>Answer</u>: The student’s phone rang.	
4756		
4757		
4758		
4759		
4760	<ul style="list-style-type: none">• CosmosQA:	
4761	<ul style="list-style-type: none">– <u>Question</u>: That idea still weirds me out . I made a blanket for the baby ’s older sister before she was born but I completely spaced that this one was on the way , caught up in my own dramas and whatnot . Luckily , I had started a few rows in white just to learn a stitch ages ago , and continuing that stitch will make an acceptable woobie , I think . According to the above context, choose the best option to answer the following question. Question: What did I make for the baby . Options: A. I made a carseat . B. None of the above choices . C. I made a crib . D. I finished a pair of booties .– <u>Answer</u>: D.	
4762		
4763		
4764		
4765		
4766		
4767		
4768		
4769		
4770		
4771		
4772		
4773		
4774		
4775		
4776		
4777	<ul style="list-style-type: none">• QuaRel:	
4778	<ul style="list-style-type: none">– <u>Question</u>: Here’s a short story: A piece of thread is much thinner than a tree so it is (A) less strong (B) more strong. What is the most sensical answer between “Thread” and “Tree”? – <u>Answer</u>: Thread.	
4779		
4780		
4781		
4782		
4783		
4784	<h2>D Datasets and Metric</h2>	
4785	<p>The specific details of all the datasets we use in this work are provided in this section.</p>	
4786		
4787	<h3>D.1 T0 Held-In Datasets</h3>	
4788	<ul style="list-style-type: none">• CommonsenseQA (Talmor et al., 2019b) under MIT License, evaluated by accuracy.	
4789		
4790	<ul style="list-style-type: none">• DREAM (Sun et al., 2019a) under MIT License, evaluated by accuracy.	
4791		
4792	<ul style="list-style-type: none">• QUAIL (Rogers et al., 2020) under CC BY-SA 4.0, evaluated by accuracy.	
4793		
4794	<ul style="list-style-type: none">• QuaRTz (Tafjord et al., 2019) under Apache 2.0 License, evaluated by accuracy.	
4795		
4796	<ul style="list-style-type: none">• Social IQA (Sap et al., 2019) under MIT License, evaluated by accuracy.	
4797		

4798	• WiQA (Tandon et al., 2019) under <i>Apache 2.0 License</i> , evaluated by accuracy.	4838
4799		4839
4800	• Cosmos QA (Huang et al., 2019) under <i>MIT License</i> , evaluated by accuracy.	4840
4801		
4802	• QASC (Khot et al., 2020) under <i>Apache 2.0 License</i> , evaluated by accuracy.	
4803		
4804	• Quarel (Tafjord et al., 2018) under <i>Apache 2.0 License</i> , evaluated by accuracy.	
4805		
4806	• SciQ (Johannes Welbl, 2017) under <i>MIT License</i> , evaluated by accuracy.	
4807		
4808	• Wiki Hop (Welbl et al., 2018) under <i>CC BY-SA 3.0</i> , evaluated by accuracy.	
4809		
4810	• Adversarial QA (Bartolo et al., 2020) under <i>Apache 2.0 License</i> , evaluated by F1 score.	
4811		
4812	• Quoref (Dasigi et al., 2019) under <i>Apache 2.0 License</i> , evaluated by F1 score.	
4813		
4814	• DuoRC (Saha et al., 2018) under <i>MIT License</i> , evaluated by F1 score.	
4815		
4816	• ROPEs (Lin et al., 2019b) under <i>Apache 2.0 License</i> , evaluated by F1 score.	
4817		
4818	• Hotpot QA (Yang et al., 2018) under <i>CC BY-SA 4.0</i> , evaluated by exact match and F1 score.	
4819		
4820	• Wiki QA (Yi et al., 2015) under <i>MIT License</i> , evaluated by mean average precision (MAP) and mean reciprocal rank (MRR).	
4821		
4822	• Common Gen (Lin et al., 2020c) under <i>MIT License</i> , evaluated by BLEU and ROUGE scores.	
4823		
4824	• Wiki Bio (Lebret et al., 2016b) under <i>CC BY-SA 3.0</i> , evaluated by BLEU score.	
4825		
4826	• Amazon (Blitzer et al., 2007) under <i>Proprietary License</i> , evaluated by accuracy.	
4827		
4828	• App Reviews (Maas et al., 2011a) under <i>Proprietary License</i> , evaluated by accuracy.	
4829		
4830	• IMDB (Maas et al., 2011b) under <i>Proprietary License</i> , evaluated by accuracy.	
4831		
4832	• Rotten Tomatoes (Zhu et al., 2010) under <i>Proprietary License</i> , evaluated by accuracy.	
4833		
4834	• Yelp (Yelp, Inc., 2018) under <i>Apache 2.0 License</i> , evaluated by accuracy.	
4835		
4836		
4837		
	• CNN Daily Mail (Hermann et al., 2015a) under <i>Apache 2.0 License</i> , evaluated by ROUGE score.	4838
		4839
	• Gigaword (Graff et al., 2003) under <i>LDC License</i> , evaluated by ROUGE score.	4840
	• MultiNews (Fabbri et al., 2019) under <i>MIT License</i> , evaluated by ROUGE score.	4843
		4844
	• SamSum (Gliwa et al., 2019) under <i>CC BY-SA 4.0</i> , evaluated by ROUGE score.	4845
		4846
	• XSum (See et al., 2017a) under <i>Apache 2.0 License</i> , evaluated by ROUGE score.	4847
		4848
	• AG News (Zhang et al., 2015) under <i>CC BY-SA 3.0</i> , evaluated by accuracy.	4849
		4850
	• DBPedia (Auer et al., 2007) under <i>CC BY-SA 3.0</i> , evaluated by accuracy.	4851
		4852
	• TREC (Voorhees, 2002) under <i>NIST License</i> , evaluated by accuracy.	4853
		4854
	• MRPC (Dolan and Brockett, 2005) under <i>Apache 2.0 License</i> , evaluated by accuracy and F1 score.	4855
		4856
	• PAWS (Zhang et al., 2019c) under <i>Apache 2.0 License</i> , evaluated by accuracy and F1 score.	4857
	• QQP (Quora, Inc., 2017) under <i>Quora Terms of Service</i> , evaluated by accuracy and F1 score.	4860
		4861
		4862
	D.2 T0 Held-Out Datasets	4863
	Held-out Tasks	4864
	• ANLI (Nie et al., 2020) under <i>MIT License</i> , evaluated by accuracy.	4865
		4866
	• CB (de Marneffe et al., 2017) under <i>CC-BY-SA License</i> , evaluated by accuracy.	4867
		4868
	• RTE (Dagan et al., 2005; Bar-Haim et al., 2006; Giampiccolo et al., 2007; Bentivogli et al., 2009) under <i>Apache 2.0 License</i> , evaluated by accuracy.	4869
		4870
	• WSC (Levesque et al., 2012) under <i>Creative Commons License</i> , evaluated by accuracy.	4871
		4872
	• Winogrande (Sakaguchi et al., 2020a) under <i>Apache License 2.0</i> , evaluated by accuracy.	4873
		4874

- 4877 • **WiC** ([Pilehvar and Camacho-Collados, 2019](#)) 4920
 4878 under *CC BY-SA 4.0 License*, evaluated by 4921
 4879 accuracy.
- 4880 • **COPA** ([Roemmele et al., 2011a](#)) under *BSD- 4922
 4881 2-Clause License*, evaluated by accuracy. 4923
 4882 • **HellaSwag** ([Zellers et al., 2019a](#)) under *MIT 4924
 4883 License*, evaluated by accuracy.
- 4884 • **Story Cloze** ([Mostafazadeh et al., 2016b](#)) 4925
 4885 under *CC-BY 4.0 License*, evaluated by accuracy. 4926
 4886
- ### D.3 BigBench-Hard Datasets
- 4887 • **abstract_narrative_understanding** ([Ghosh 4928
 4888 and Srivastava, 2021; Holyoak, 2012; Nippold 4929
 4889 et al., 2001; Tan et al., 2015; Wang et al., 4930
 4890 2020b; Mostafazadeh et al., 2016a](#)) under 4931
 4891 *Apache 2.0 License*, evaluated by accuracy 4932
 4892
- 4893 • **abstraction_and_reasoning_corpus** ([Chollet, 4933
 4894 2019; Brown et al., 2020; Chollet, 2020](#)) 4934
 4895 under *Apache 2.0 License*, evaluated by accuracy 4935
- 4896 • **anachronisms** ([Otterbacher et al., 2002; 4936
 4897 Popescu and Strapparava, 2015; Llorens et al., 4937
 4898 2015; Meng et al., 2017; Geva et al., 2021](#)) 4938
 4899 under *Apache 2.0 License*, evaluated by accuracy 4939
 4900
- 4901 • **analogical_similarity** ([Plate, 2003, 1994; 4940
 4902 Thagard et al., 1990; Gentner et al., 1993](#)) 4941
 4903 under *Apache 2.0 License*, evaluated by accuracy 4942
 4904
- 4905 • **analytic_entailment** ([Hume, 1739–1740; 4943
 4906 Kant, 1781/1787; Wittgenstein, 1953; Quine, 4944
 4907 1951; Grice and Strawson, 1956; Bolukbasi 4945
 4908 et al., 2016; Kocurek et al., 2020; Rudolph and 4946
 4909 Kocurek, 2020; Kocurek and Jerzak, 2021](#)) 4947
 4910 under *Apache 2.0 License*, evaluated by accuracy 4948
- 4911 • **arithmetic** ([Brown et al., 2020; Saxton et al., 4949
 4912 2019](#)) under *Apache 2.0 License*, evaluated by 4950
 4913 accuracy.
- 4914 • **ascii_word_recognition** ([Child et al., 2019; 4951
 4915 Chen et al., 2020](#)) under *Apache 2.0 License*, 4952
 4916 evaluated by accuracy.
- 4917 • **authorship_verification** ([Bischoff et al., 2020; 4953
 4918 Koppel and Schler, 2004](#)) under *Apache 2.0 4954
 4919 License*, evaluated by accuracy.
- **auto_categorization** under *Apache 2.0 License*, 4955
 evaluated by accuracy
- **bbq_lite** ([Crawford, 2017; Khashabi et al., 4956
 2020b; Li et al., 2020a](#)) under *Apache 2.0 License*, 4957
 evaluated by accuracy
- **bias_from_probabilities** ([Bender et al., 2021; 4958
 Abid et al., 2021](#)) under *Apache 2.0 License*, 4959
 evaluated by accuracy
- **boolean_expressions** ([Habernal et al., 2018; 4960
 Yu et al., 2020; Dua et al., 2019; Liu et al., 4961
 2020a; Sinha et al., 2019; Wang et al., 2019b; 4962
 Steinbach and Kohut, 2002; Saxton et al., 4963
 2019; Payani and Fekri, 2019; Trask et al., 4964
 2018; Selsam et al., 2018; Allamanis et al., 4965
 2016; Evans et al., 2018; Shi et al., 2020](#)) under 4966
Apache 2.0 License, evaluated by accuracy
- **bridging_anaphora_resolution_barqa** ([Hou, 4967
 2020; Hou et al., 2013; Markert et al., 2012; 4968
 Rajpurkar et al., 2016](#)) under *Apache 2.0 4969
 License*, evaluated by accuracy
- **causal_judgment** ([Gordon, 2010; Bosselut 4970
 et al., 2019; Halpern, 2016; Knobe, 2003](#)) under 4971
Apache 2.0 License, evaluated by accuracy
- **cause_and_effect** ([Gordon, 2010](#)) under 4972
Apache 2.0 License, evaluated by accuracy
- **checkmate_in_one** ([Alexander, 2020; Am- 4973
 manabrolu et al., 2019; Dambekodi et al., 4974
 2020; Ammanabrolu et al., 2020](#)) under 4975
Apache 2.0 License, evaluated by accuracy
- **chess_state_tracking** ([Weston et al., 2015; 4976
 Côté et al., 2018; Toshniwal et al., 2021; 4977
 Alexander, 2020; Chen, 2020; Noever et al., 4978
 2020; Swingle et al., 2021](#)) under *Apache 2.0 4979
 License*, evaluated by accuracy
- **chinese_remainder_theorem** under *Apache 4980
 2.0 License*, evaluated by accuracy
- **cifar10_classification** under *Apache 2.0 Li- 4981
 cense*, evaluated by accuracy
- **codenames** ([Kim et al., 2019](#)) under *Apache 4982
 2.0 License*, evaluated by accuracy
- **color** ([Gibson et al., 2017; Zucconi, 6 Jan. 4983
 2016](#)) under *Apache 2.0 License*, evaluated by 4984
 accuracy

- 4963 • **com2sense** ([Singh et al., 2021](#)) under *Apache*
4964 *2.0 License*, evaluated by accuracy
4965 • **common_morpheme** ([Devlin et al., 2018](#);
4966 [Wu et al., 2016](#); [Won et al., 2021](#); [Xu et al.,
4967 2018](#); [Edmiston and Stratos, 2018](#); [El-Kishky
4968 et al., 2019](#)) under *Apache 2.0 License*, eval-
4969 uated by accuracy
4970 • **context_definition_alignment** ([Senel and
4971 Schütze, 2021](#); [Reimers and Gurevych, 2019](#))
4972 under *Apache 2.0 License*, evaluated by accu-
4973 racy
4974 • **convinceme** ([Lin et al., 2021b](#); [Levy et al.,
4975 2021](#); [Clark et al., 2018](#); [Maynez et al., 2020](#);
4976 [Wang et al., 2020a](#); [Kenton et al., 2021](#); [Xu
4977 et al., 2020a](#); [Tamkin et al., 2021](#)) under
4978 *Apache 2.0 License*, evaluated by accuracy
4979 • **coqa_conversational_question_answering** ([Reddy
4980 et al., 2019](#); [Radford et al., 2019](#); [Brown et al.,
4981 2020](#)) under *Apache 2.0 License*, evaluated by
4982 accuracy
4983 • **crash_blossom** under *Apache 2.0 License*,
4984 evaluated by accuracy
4985 • **crass_ai** ([Schölkopf et al., 2021](#); [Teney et al.,
4986 2020](#); [Liang et al., 2020b](#); [Pearl, 2000](#); [Shahid
4987 and Zheleva, 2021](#); [Xia et al., 2021](#); [Priol et al.,
4988 2020](#)) under *Apache 2.0 License*, evaluated by
4989 accuracy
4990 • **cryobiology_spanish** ([Scudellari, 2017](#);
4991 [Soltani Firouz et al., 2021](#); [Toubiana et al.,
4992 2020](#); [Mbogba et al., 2018](#)) under *Apache 2.0*
4993 *License*, evaluated by accuracy
4994 • **cryptonite** ([Efrat et al., 2021](#); [Raganato et al.,
4995 2017](#); [Sakaguchi et al., 2020b](#); [Miller and
4996 Gurevych, 2015](#); [Miller et al., 2017](#); [Joshi
4997 et al., 2017](#); [Oprea and Magdy, 2020](#); [Friedlan-
4998 der and Fine, 2018](#); [Lewis et al., 2020c](#)) under
4999 *Apache 2.0 License*, evaluated by accuracy
5000 • **cs_algorithms** under *Apache 2.0 License*,
5001 evaluated by accuracy
5002 • **cycled_letters** ([Brown et al., 2020](#)) under
5003 *Apache 2.0 License*, evaluated by accuracy
5004 • **dark_humor_detection** ([Weller and Seppi,
5005 2019](#); [Fan et al., 2020](#); [Willinger et al., 2017](#);
5006 [Yang et al., 2015](#); [Mihalcea and Strapparava,
5007 2005](#)) under *Apache 2.0 License*, evaluated by
5008 accuracy
5009 • **date_understanding** ([Vashishth et al., 2018](#);
5010 [Chambers, 2012](#); [Kotsakos et al., 2014](#);
5011 [Vashishtha et al., 2020, 2019](#)) under *Apache*
5012 *2.0 License*, evaluated by accuracy
5013 • **disambiguation_qa** ([Zhao et al., 2018](#);
5014 [Rudinger et al., 2018](#)) under *Apache 2.0 Li-
5015 cense*, evaluated by accuracy
5016 • **discourse_marker_prediction** ([Malmi et al.,
5017 2018](#); [Nie et al., 2019](#); [Sileo et al., 2019, 2020](#))
5018 under *Apache 2.0 License*, evaluated by accu-
5019 racy
5020 • **disfl_qa** ([Gupta et al., 2021](#); [Rajpurkar et al.,
5021 2018](#)) under *Apache 2.0 License*, evaluated by
5022 accuracy
5023 • **diverse_social_bias** ([Sheng et al., 2019](#);
5024 [Nadeem et al., 2020](#); [Hendrycks et al., 2020](#);
5025 [Sap et al., 2020](#); [Gehman et al., 2020](#); [Boluk-
5026 basi et al., 2016](#); [Caliskan et al., 2017](#); [May
5027 et al., 2019](#); [Liang et al., 2020a](#); [Barocas and
5028 Selbst, 2016](#); [Cho et al., 2019](#); [Blodgett et al.,
5029 2020](#); [Merity et al., 2016](#); [Socher et al., 2013](#);
5030 [Poria et al., 2019](#)) under *Apache 2.0 License*,
5031 evaluated by accuracy
5032 • **dyck_languages** ([Chomsky and Schützen-
5033 berger, 1959](#); [Suzgun et al., 2019b](#); [Hao et al.,
5034 2018](#); [Hewitt et al., 2020](#); [Hahn, 2020](#); [Suzgun
5035 et al., 2019a](#); [Sennhauser and Berwick, 2018](#);
5036 [Skachkova et al., 2018](#); [Bhattamishra et al.,
5037 2020a](#); [Yu et al., 2019c](#); [Ebrahimi et al., 2020](#);
5038 [Ackerman and Cybenko, 2020](#); [Bhattamishra
5039 et al., 2020b](#)) under *Apache 2.0 License*, eval-
5040 uated by accuracy
5041 • **dynamic_counting** ([Suzgun et al., 2019a](#);
5042 [Skachkova et al., 2018](#); [Bhattamishra et al.,
5043 2020a](#); [Suzgun et al., 2019b](#); [Yu et al., 2019c](#);
5044 [Ebrahimi et al., 2020](#); [Ackerman and Cy-
5045 benko, 2020](#); [Bhattamishra et al., 2020b](#);
5046 [Sennhauser and Berwick, 2018](#); [Merrill, 2020](#);
5047 [Karpathy, 2015](#)) under *Apache 2.0 License*,
5048 evaluated by accuracy
5049 • **elementary_math_qa** ([Amini et al., 2019](#);
5050 [Ling et al., 2017](#); [Hendrycks et al., 2021c](#); [Pa-
5051 tel et al., 2021](#); [Zhang et al., 2020a](#); [Hendrycks
5052 et al., 2021b](#)) under *Apache 2.0 License*, eval-
5053 uated by accuracy

- 5054 • **emojis_emotion_prediction** ([Shoeb and de Melo, 2020](#); [Plutchik, 1980](#)) under *Apache 2.0 License*, evaluated by accuracy 5099
- 5055 5056 5057 5058 5059 5060 5061 5062 5063 5064 5065 5066 5067 5068 5069 5070 5071 5072 5073 5074 5075 5076 5077 5078 5079 5080 5081 5082 5083 5084 5085 5086 5087 5088 5089 5090 5091 5092 5093 5094 5095 5096 5097 5098 5099 5100 5101 5102 5103 5104 5105 5106 5107 5108 5109 5110 5111 5112 5113 5114 5115 5116 5117 5118 5119 5120 5121 5122 5123 5124 5125 5126 5127 5128 5129 5130 5131 5132 5133 5134 5135 5136 5137 5138 5139 5140 5141 5142
- **emojis_emotion_prediction** ([Shoeb and de Melo, 2020](#); [Plutchik, 1980](#)) under *Apache 2.0 License*, evaluated by accuracy
 - **empirical_judgments** ([Kant, 1781/1787, 1783](#); [Spirtes et al., 2000](#); [Pearl, 1988](#); [Goldberger, 1972](#); [Rothman and Greenland, 2005](#); [Roemmel et al., 2011b](#); [Wang et al., 2019b](#); [Evans et al., 2019](#)) under *Apache 2.0 License*, evaluated by accuracy
 - **english_proverbs** ([Gyasi Obeng, 1996](#); [Honeck, 1997](#); [Hrisztova-Gotthardt and Aleksa Varga, 2015](#)) under *Apache 2.0 License*, evaluated by accuracy
 - **english_russian_proverbs** ([Bodrova, 2007](#); [Gvarjalaze and Mchedlishvili, 1971](#); [Wik](#)) under *Apache 2.0 License*, evaluated by accuracy
 - **entailed_polarity** ([Karttunen, 2012](#)) under *Apache 2.0 License*, evaluated by accuracy
 - **entailed_polarity_hindi** ([Karttunen, 2012](#)) under *Apache 2.0 License*, evaluated by accuracy
 - **epistemic_reasoning** ([Ravenscroft, 2019](#); [Call and Tomasello, 2008](#); [Bugnyar et al., 2016](#); [Stalnaker, 1978](#); [Sperber and Wilson, 2002](#); [Nematzadeh et al., 2018](#); [Le et al., 2019](#); [Jiang and de Marneffe, 2019](#); [Ross and Pavlick, 2019](#); [Bowman et al., 2015](#); [de Marneffe et al., 2012](#); [Jeretic et al., 2020](#)) under *Apache 2.0 License*, evaluated by accuracy
 - **evaluating_information_essentiality** ([Rajpurkar et al., 2018](#); [Hosseini et al., 2014](#); [Levy et al., 2017](#); [Yin and Pei, 2015](#); [de Marneffe et al., 2008](#); [Zhu et al., 2019](#)) under *Apache 2.0 License*, evaluated by accuracy
 - **fact_checker** ([Thorne et al., 2018](#); [Lee et al., 2021](#)) under *Apache 2.0 License*, evaluated by accuracy
 - **factuality_of_summary** ([Eyal et al., 2019](#); [Wang et al., 2020a](#); [Durmus et al., 2020](#); [Vasiliyev et al., 2020](#); [See et al., 2017b](#); [Hermann et al., 2015b](#); [Narayan et al., 2018](#); [Pagnoni et al., 2021](#); [Kryscinski et al., 2020](#)) under *Apache 2.0 License*, evaluated by accuracy
 - **fantasy_reasoning** ([Wang et al., 2019b, 2018](#); [McCann et al., 2018](#); [Bhagavatula et al., 2019](#);
 - **few_shot_nlg** ([Rastogi et al., 2020](#); [Kale and Rastogi, 2020](#)) under *Apache 2.0 License*, evaluated by accuracy
 - **figure_of_speech_detection** ([Potamias et al., 2020](#)) under *Apache 2.0 License*, evaluated by accuracy
 - **forecasting_subquestions** under *Apache 2.0 License*, evaluated by accuracy
 - **gem** ([Gehrman et al., 2021](#)) under *Apache 2.0 License*, evaluated by accuracy
 - **gender_inclusive_sentences_german** under *Apache 2.0 License*, evaluated by accuracy
 - **gender_sensitivity_chinese** ([on the Revision of the National Standard Occupational Classification, 2015](#); [Household Management Research Center, 2021, 2020, 2019](#); [Qimington, 2016](#); [Ministry of the Interior, 2018](#)) under *Apache 2.0 License*, evaluated by accuracy
 - **gender_sensitivity_english** ([Bordia and Bowman, 2019](#); [Marcus et al., 1994](#); [Caliskan et al., 2017](#); [Bolukbasi et al., 2016](#); [Ruderger et al., 2018](#); [Lu et al., 2020](#); [Gonen and Goldberg, 2019](#); [Hall Maudslay et al., 2019](#); [Fellbaum, 1998](#)) under *Apache 2.0 License*, evaluated by accuracy
 - **general_knowledge** ([Shane, 2020](#); [Dhingra et al., 2017](#); [Rajpurkar et al., 2016, 2018](#); [Lacker, 2020](#)) under *Apache 2.0 License*, evaluated by accuracy
 - **geometric_shapes** ([Bostock et al., 2011](#); [Marriott et al., 2021](#); [Boillot, 2019](#)) under *Apache 2.0 License*, evaluated by accuracy
 - **goal_step_wikihow** ([Zhang et al., 2020d](#)) under *Apache 2.0 License*, evaluated by accuracy
 - **gre_reading_comprehension** ([Lai et al., 2017](#)) under *Apache 2.0 License*, evaluated by accuracy
 - **hh_alignment** under *Apache 2.0 License*, evaluated by accuracy

- 5143 • **high_low_game** under *Apache 2.0 License*,
5144 evaluated by accuracy
5145 • **hindi_question_answering** ([Brown et al., 2020](#); [Radford et al., 2019](#); [Jain et al., 2020](#);
5146 [Lewis et al., 2020a](#); [Artetxe et al., 2020](#); [Rajpurkar et al., 2016](#)) under *Apache 2.0 License*,
5147 evaluated by accuracy
5148
5149
- 5150 • **hinglish_toxicity** under *Apache 2.0 License*,
5151 evaluated by accuracy
5152 • **human_organs_senses** under *Apache 2.0 Li-*
5153 *cense*, evaluated by accuracy
5154 • **hyperbaton** ([Forsyth, 2014](#)) under *Apache*
5155 *2.0 License*, evaluated by accuracy
5156 • **identify_math_theorems** ([Gao et al., 2021](#);
5157 [Black et al., 2022](#); [Brown et al., 2020](#); [Radford](#)
5158 [et al., 2019](#); [Wang and Komatsuzaki, 2021](#)) un-
5159 der *Apache 2.0 License*, evaluated by accuracy
5160 • **identify_odd_metaphor** ([Lakoff and John-
5161 son, 2008](#); [Gao et al., 2018](#)) under *Apache 2.0*
5162 *License*, evaluated by accuracy
5163 • **implicatures** ([Davis, 2019](#); [George and](#)
5164 [Mamidi, 2020](#)) under *Apache 2.0 License*,
5165 evaluated by accuracy
5166 • **implicit_relations** ([Cain and Oakhill, 1999](#);
5167 [Bayat and Çetinkaya, 2020](#); [Srivastava et al.,](#)
5168 [2016](#); [Lin et al., 2019a](#); [Massey et al., 2015](#))
5169 under *Apache 2.0 License*, evaluated by accu-
5170 racy
5171 • **intent_recognition** ([Brown et al., 2020](#);
5172 [Winata et al., 2021](#); [Madotto et al., 2020](#);
5173 [Coucke et al., 2018](#); [Madotto et al., 2021](#)) un-
5174 der *Apache 2.0 License*, evaluated by accuracy
5175 • **international_phonetic_alphabet_nli** ([Williams](#)
5176 [et al., 2018](#)) under *Apache 2.0 License*, eval-
5177 uated by accuracy
5178 • **international_phonetic_alphabet_transliterate** ([Brown](#)
5179 [et al., 2020](#); [Liu et al., 2020c](#); [Williams et al.,](#)
5180 [2018](#)) under *Apache 2.0 License*, evaluated by
5181 accuracy
5182 • **intersect_geometry** ([Weston et al., 2015](#);
5183 [Agrawal et al., 2015](#); [Trask et al., 2018](#); [Seo](#)
5184 [et al., 2014](#); [Hosseini et al., 2014](#); [Polu and](#)
5185 [Sutskever, 2020](#); [Yang and Deng, 2019](#)) under
5186 *Apache 2.0 License*, evaluated by accuracy
5187 • **irony_identification** ([Zhang et al., 2019b](#);
5188 [Ghanem et al., 2020](#); [Salas-Zárate et al., 2017](#))
5189 under *Apache 2.0 License*, evaluated by accu-
5190 racy
5191 • **kanji_ascii** under *Apache 2.0 License*, eval-
5192 uated by accuracy
5193 • **kannada** ([Prentice and Fathman, 1975](#);
5194 [Narasimhachar, 1988](#); [Liu et al., 2021b](#); [Lev](#)
5195 [et al., 2004](#); [Lin et al., 2021a](#)) under *Apache*
5196 *2.0 License*, evaluated by accuracy
5197 • **key_value_maps** under *Apache 2.0 License*,
5198 evaluated by accuracy
5199 • **language_games** under *Apache 2.0 License*,
5200 evaluated by accuracy
5201 • **linguistic_mappings** ([McCoy et al., 2018](#),
5202 [2020](#); [Mulligan et al., 2021](#); [Rumelhart et al.,](#)
5203 [1986](#); [Kirov and Cotterell, 2018](#); [Berko, 1958](#);
5204 [Baayen et al., 1995](#)) under *Apache 2.0 License*,
5205 evaluated by accuracy
5206 • **list_functions** ([Rule et al., 2020](#); [Rule, 2020](#);
5207 [Green et al., 1974](#); [Shaw et al., 1975](#); [Bier-](#)
5208 [mann, 1978](#); [Green, 1981](#); [Smith, 1984](#); [Feser](#)
5209 [et al., 2015](#); [Osera and Zdancewic, 2015](#); [Po-](#)
5210 [likarpova et al., 2016](#); [Cropper et al., 2020](#);
5211 [Graves et al., 2014](#); [Reed and de Freitas, 2015](#);
5212 [Joulin and Mikolov, 2015](#); [Balog et al., 2016](#);
5213 [Bošnjak et al., 2017](#); [Gaunt et al., 2016](#); [Chen](#)
5214 [et al., 2019b](#); [Kitzelmann, 2010](#); [Flener and](#)
5215 [Schmid, 2008](#); [Gulwani et al., 2017a](#); [Devlin](#)
5216 [et al., 2017](#); [Ellis et al., 2020](#); [Cropper and](#)
5217 [Muggleton, 2016](#); [Piantadosi, 2020](#)) under
5218 *Apache 2.0 License*, evaluated by accuracy
5219 • **logical_args** under *Apache 2.0 License*, eval-
5220 uated by accuracy
5221 • **logical_fallacy_detection** ([Brown et al.,](#)
5222 [2020](#); [Hendrycks et al., 2021b](#); [Bender et al.,](#)
5223 [2021](#); [Wachsmuth et al., 2017](#); [Yu et al., 2020](#);
5224 [Obipi et al., 2018](#); [Oberauer et al., 2005](#); [Ober-](#)
5225 [auer and Wilhelm, 2000](#)) under *Apache 2.0*
5226 *License*, evaluated by accuracy
5227 • **logical_sequence** ([Saxton et al., 2019](#); [Lin](#)
5228 [et al., 2020a](#); [Bowman and Dahl, 2021](#)) under
5229 *Apache 2.0 License*, evaluated by accuracy
5230 • **long_context_integration** under *Apache 2.0*
5231 *License*, evaluated by accuracy

- 5232 • **mathematical_induction** (Hendrycks et al.,
5233 2021c; Patel et al., 2021) under *Apache 2.0*
5234 License, evaluated by accuracy 5277
5235 • **matrixshapes** under *Apache 2.0 License*, eval-
5236 uated by accuracy 5278
5237 • **metaphor_boolean** (Lakoff and Johnson,
5238 2008; Bizzoni and Lappin, 2018) under
5239 *Apache 2.0 License*, evaluated by accuracy 5279
5240 • **metaphor_understanding** (Paul, 1970; Tong
5241 et al., 2021; Radford et al., 2019; Rai and
5242 Chakraverty, 2020; Shutova, 2010; Stowe
5243 et al., 2020; Mohler et al., 2016; Shutova and
5244 Teufel, 2010; Birke and Sarkar, 2006; Zayed
5245 et al., 2020; Steen et al., 2010; Tong, 2021;
5246 Bizzoni and Lappin, 2018; Tsvetkov et al.,
5247 2014) under *Apache 2.0 License*, evaluated by
5248 accuracy 5280
5249 • **minute_mysteries_qa** (Sugawara et al., 2017;
5250 Dunietz et al., 2020; Kočiský et al., 2018;
5251 Mostafazadeh et al., 2016a; Frermann et al.,
5252 2018) under *Apache 2.0 License*, evaluated by
5253 accuracy 5281
5254 • **misconceptions** (Irving et al., 2018;
5255 Atanasova et al., 2020; Boller and George,
5256 1989) under *Apache 2.0 License*, evaluated by
5257 accuracy 5282
5258 • **mnist_ascii** under *Apache 2.0 License*, eval-
5259 uated by accuracy 5283
5260 • **modified_arithmetic** (Brown et al., 2020) un-
5261 der *Apache 2.0 License*, evaluated by accuracy 5284
5262 • **moral_permissibility** (Hendrycks et al.,
5263 2020; Lourie et al., 2020; Thomson, 1976)
5264 under *Apache 2.0 License*, evaluated by accu-
5265 racy 5285
5266 • **movie_dialog_same_or_different** (Park
5267 et al., 2021) under *Apache 2.0 License*,
5268 evaluated by accuracy 5286
5269 • **movie_recommendation** (Sileo et al., 2022;
5270 Thorat et al., 2015; Barkan and Koenig-
5271 stein, 2016; Harper and Konstan, 2015) under
5272 *Apache 2.0 License*, evaluated by accuracy 5287
5273 • **mult_data_wrangling** (Bender et al., 2021;
5274 Tamkin et al., 2021; Singh and Gulwani, 2015;
5275 Cropper et al., 2016; Wu et al., 2012; Gul-
5276 wani et al., 2015; Contreras-Ochando et al.,
5277 2018, 2020; Petrova-Antonova and Tancheva,
5278 2020; Huynh and Mazzocchi, 2012; Kandel
5279 et al., 2011; Bhupatiraju et al., 2017; Ellis
5280 and Gulwani, 2017; Gulwani et al., 2012; Gul-
5281 wani, 2011; Singh and Gulwani, 2016) under
5282 *Apache 2.0 License*, evaluated by accuracy 5283
5283 • **multitemo** (Kocoń et al., 2021) under *Apache*
5284 *2.0 License*, evaluated by accuracy 5284
5285 • **multistep_arithmetic** (Flanagan and Dixon,
5286 2014) under *Apache 2.0 License*, evaluated by
5287 accuracy 5285
5288 • **muslim_violence_bias** (Abid et al., 2021;
5289 Bender et al., 2021) under *Apache 2.0 License*,
5290 evaluated by accuracy 5288
5291 • **natural_instructions** (Mishra et al., 2021) un-
5292 der *Apache 2.0 License*, evaluated by accuracy 5291
5293 • **navigate** (Graves et al., 2016; Henaff et al.,
5294 2016; Geva et al., 2020b; Chen et al., 2019a;
5295 Kryscinski et al., 2020; Côté et al., 2018;
5296 Luketina et al., 2019; Thawani et al., 2021;
5297 Lake and Baroni, 2017) under *Apache 2.0 Li-*
5298 *cense*, evaluated by accuracy 5298
5299 • **nonsense_words_grammar** under *Apache*
5300 *2.0 License*, evaluated by accuracy 5299
5301 • **object_counting** (Rugani et al., 2015; Wang
5302 et al., 2019c; Zhang et al., 2018b; Brown et al.,
5303 2020) under *Apache 2.0 License*, evaluated by
5304 accuracy 5301
5305 • **odd_one_out** (Resnik, 1995, 1999; Jiang and
5306 Conrath, 1997; Li et al., 2003; Banerjee and
5307 Pedersen, 2003; Jarmasz, 2012; Hughes and
5308 Ramage, 2007; Tsatsaronis et al., 2010, 2009;
5309 Morris and Hirst, 1991; Strube and Ponzetto,
5310 2006; Ponzetto and Strube, 2007; Gabrilovich
5311 and Markovitch, 2007; Milne and Witten,
5312 2008; Yeh et al., 2009; Radinsky et al., 2011;
5313 Cilibrasi and Vitanyi, 2007; Deerwester et al.,
5314 1990; Reisinger and Mooney, 2010; El-Yaniv
5315 and Yanay, 2013) under *Apache 2.0 License*,
5316 evaluated by accuracy 5305
5317 • **paragraph_segmentation** under *Apache 2.0*
5318 *License*, evaluated by accuracy 5317
5319 • **parsinlu_qa** (Khashabi et al., 2020a) under
5320 *Apache 2.0 License*, evaluated by accuracy 5319

- 5321 • **penguins_in_a_table** (Herzig et al., 2020) un- 5364
 5322 der Apache 2.0 License, evaluated by accuracy 5365
 5323 • **periodic_elements** under Apache 2.0 License, 5366
 5324 evaluated by accuracy 5367
 5325 • **persian_idioms** under Apache 2.0 License, 5368
 5326 evaluated by accuracy 5369
 5327 • **phrase_relatedness** (Asaadi et al., 2019; 5370
 5328 Levy et al., 2015; Ein Dor et al., 2018) under 5371
 5329 Apache 2.0 License, evaluated by accuracy 5372
 5330 • **physical_intuition** under Apache 2.0 License, 5373
 5331 evaluated by accuracy 5374
 5332 • **physics** under Apache 2.0 License, evaluated 5375
 5333 by accuracy 5376
 5334 • **physics_questions** (Ling et al., 2017; Amini 5377
 5335 et al., 2019) under Apache 2.0 License, evaluated 5378
 5336 by accuracy 5379
 5337 • **polish_sequence_labeling** (Nguyen and Guo, 5380
 5338 2007; Rei, 2017; Gu et al., 2018) under 5381
 5339 Apache 2.0 License, evaluated by accuracy 5382
 5340 • **presuppositions_as_nli** (Heim, 1983; 5383
 5341 de Marneffe et al., 2019; White et al., 2018; 5384
 5342 Jeretic et al., 2020) under Apache 2.0 License, 5385
 5343 evaluated by accuracy 5386
 5344 • **program_synthesis** (Gulwani et al., 2017b) 5387
 5345 under Apache 2.0 License, evaluated by accuracy 5388
 5346
 5347 • **protein_interacting_sites** (Dhole et al., 2014; 5389
 5348 Singh et al., 2014; Li et al., 2020b; Murakami 5390
 5349 and Mizuguchi, 2010) under Apache 2.0 Li- 5391
 5350 cense, evaluated by accuracy 5392
 5351 • **python_programming_challenge** (Allama- 5393
 5352 nis et al., 2018; Alon et al., 2020; Hendrycks 5394
 5353 et al., 2021a) under Apache 2.0 License, eval- 5395
 5354 uated by accuracy
 5355 • **qa_wikidata** (Radford et al., 2019; 5396
 5356 Kwiatkowski et al., 2019; Weber and 5397
 5357 Jaimes, 2011) under Apache 2.0 License, 5398
 5358 evaluated by accuracy 5399
 5359 • **question_answer_creation** under Apache 2.0 5400
 5360 License, evaluated by accuracy 5401
 5361 • **question_selection** (Rajpurkar et al., 2016) 5402
 5362 under Apache 2.0 License, evaluated by accu- 5403
 5363 racy 5404
 5364 • **real_or_fake_text** (Dugan et al., 2020; Ip- 5405
 5365 polito et al., 2020; Solaiman et al., 2019; 5406
 5366 Zellers et al., 2019b; Brown et al., 2020; 5407
 5367 Bakhtin et al., 2019; Sandhaus, 2008; Fan 5408
 5368 et al., 2018; Marín et al., 2021) under Apache 5409
 5369 2.0 License, evaluated by accuracy 5410
 5370 • **reasoning_about_colored_objects** (Hen- 5411
 5371 dricks et al., 2018; Hosseini et al., 2014; 5412
 5372 Winograd, 1972; Wang et al., 2016; Jayan- 5413
 5373 navar et al., 2020; Suhr et al., 2019; Thomason 5414
 5374 et al., 2015; Mitchell et al., 2010; Viethen and 5415
 5375 Dale, 2008; Gatt et al., 2009; Mitchell et al., 5416
 5376 2013; Liang et al., 2018) under Apache 2.0 5417
 5377 License, evaluated by accuracy 5418
 5378 • **rephrase** under Apache 2.0 License, evaluated 5419
 5379 by accuracy 5420
 5380 • **riddle_sense** (Lin et al., 2021a; Talmor et al., 5421
 5381 2019b) under Apache 2.0 License, evaluated 5422
 5382 by accuracy 5423
 5383 • **roots_optimization_and_games** (Lample 5424
 5384 and Charton, 2019; Polu and Sutskever, 2020; 5425
 5385 Amos and Kolter, 2017; Agrawal et al., 2020) 5426
 5386 under Apache 2.0 License, evaluated by 5427
 5387 accuracy 5428
 5388 • **ruin_names** (Attardo, 2017; Ren and Yang, 5429
 5389 2017; Amin and Burghardt, 2020; Annamoradnejad and Zoghi, 2020; Blinov et al., 5430
 5390 2019; Yan and Pedersen, 2017; Frolovs, 2019) 5431
 5391 under Apache 2.0 License, evaluated by accu- 5432
 5392 racy 5433
 5393 • **salient_translation_error_detection** under 5434
 5394 Apache 2.0 License, evaluated by accuracy 5435
 5395
 5396 • **scientific_press_release** under Apache 2.0 Li- 5436
 5397 cense, evaluated by accuracy 5438
 5398 • **self_awareness** (Yudkowsky, 2008; Chella 5439
 5399 et al., 2020; Schick et al., 2021; Kounev et al., 5440
 5400 2017; Huttunen et al., 2017; Wallace et al., 5441
 5401 2020; Clark and Jackson, 1994; Horowitz, 5442
 5402 2017; Branwen, 2020; Chu et al., 2017) under 5443
 5403 Apache 2.0 License, evaluated by accuracy 5444
 5404 • **self_evaluation_courtroom** (Hildebrandt, 5445
 5405 2018; Daley, 2021; King and Cook, 2020) 5446
 5406 under Apache 2.0 License, evaluated by 5447
 5407 accuracy 5448

- 5408 • **self_evaluation_tutoring** (Zhang et al., 5454
 5409 2019a; Irving et al., 2018) under *Apache 2.0 5455
 5410 License*, evaluated by accuracy
- 5411 • **semantic_parsing_in_context_sparc** (Yu 5456
 5412 et al., 2019b, 2018, 2019a) under *Apache 2.0 5457
 5413 License*, evaluated by accuracy
- 5414 • **semantic_parsing_spider** (Yu et al., 2018, 5458
 5415 2019b,a) under *Apache 2.0 License*, evaluated 5459
 5416 by accuracy
- 5417 • **sentence_ambiguity** under *Apache 2.0 Li- 5460
 5418 cense*, evaluated by accuracy
- 5419 • **similarities_abstraction** (Nasreddine et al., 5461
 5420 2005; Wechsler, 2008) under *Apache 2.0 Li- 5462
 5421 cense*, evaluated by accuracy
- 5422 • **simp_turing_concept** (Böhm, 1964; Sun 5463
 5423 et al., 2020; Devlin et al., 2018; Radford et al., 5464
 5424 2019; Brown et al., 2020; Vaswani et al., 2017; 5465
 5425 Hendrycks et al., 2021b; Xu et al., 2020b; 5466
 5426 Izacard and Grave, 2020; Zhu, 2015; Cak- 5467
 5427 mak and Thomaz, 2014; Goodman and Frank, 5468
 5428 2016; Degen et al., 2020; Khan et al., 2011; 5469
 5429 Basu and Christensen, 2013; Yang and Shafto, 5470
 5430 2017; Melo et al., 2018; Telle et al., 2019; 5471
 5431 Chater and Vitányi, 2003; Hupkes et al., 2020; 5472
 5432 Lakretz et al., 2019; Toshniwal et al., 2021; 5473
 5433 Bender and Koller, 2020; Kühl et al., 2020; 5474
 5434 Marcus and Davis, 2020; Sinha et al., 2019; 5475
 5435 McClelland et al., 2019) under *Apache 2.0 5476
 5436 License*, evaluated by accuracy
- 5437 • **simple_ethical_questions** (Hendrycks et al., 5477
 5438 2020; Lourie et al., 2020) under *Apache 2.0 5478
 5439 License*, evaluated by accuracy
- 5440 • **simple_text_editing** (Branwen, 2020; Malmi 5479
 5441 et al., 2019; Faltings et al., 2020) under *Apache 5480
 5442 2.0 License*, evaluated by accuracy
- 5443 • **snarks** (Brown et al., 2020; Devlin et al., 5481
 5444 2018; Lan et al., 2019; Liu et al., 2019; Rad- 5482
 5445 ford et al., 2019; Annamoradnejad and Zoghi, 5483
 5446 2020; Chen and Soo, 2018; Mao and Liu, 5484
 5447 2019; Weller and Seppi, 2019; Khodak et al., 5485
 5448 2017; Ghosh et al., 2020; González-Ibáñez 5486
 5449 et al., 2011; Joshi et al., 2015; McCoy et al., 5487
 5450 2019; Kaushik et al., 2019; Gardner et al., 5488
 5451 2020; Sennrich, 2016; Burlot et al., 2018; 5489
 5452 Naik et al., 2018; Zhang et al., 2016; Felbo 5490
 5453 et al., 2017; Pant and Dadu, 2020; Pelser 5491
 5454 and Murrell, 2019) under *Apache 2.0 License*, 5492
 5455 evaluated by accuracy
- 5456 • **social_support** (Wang and Jurgens, 2018) un- 5493
 5457 der *Apache 2.0 License*, evaluated by accuracy
- 5458 • **social_iqa** (Sap et al., 2019; Bisk et al., 2020; 5494
 5459 Talmor et al., 2019b; Zellers et al., 2018) un- 5495
 5460 der *Apache 2.0 License*, evaluated by accuracy
- 5461 • **spelling_bee** (Ginsberg, 2014) under *Apache 5496
 5462 2.0 License*, evaluated by accuracy
- 5463 • **sports_understanding** under *Apache 2.0 Li- 5497
 5464 cense*, evaluated by accuracy
- 5465 • **squad_shifts** (Miller et al., 2020; Brown et al., 5498
 5466 2020; Rajpurkar et al., 2016; Baumgartner 5499
 5467 et al., 2020; McAuley et al., 2015) under 5500
 5468 *Apache 2.0 License*, evaluated by accuracy
- 5469 • **subject_verb_agreement** (Lakretz et al., 5501
 5470 2021b, 2019; Linzen et al., 2016; Gulordava 5502
 5471 et al., 2018; Marvin and Linzen, 2018; Gold- 5503
 5472 berg, 2019; Lakretz et al., 2021a; Wolf, 2019) 5504
 5473 under *Apache 2.0 License*, evaluated by accu- 5505
 5474 racy
- 5475 • **sudoku** (Wang et al., 2019d; Russell and 5506
 5476 Norvig, 2002; Garcez and Lamb, 2020; Huang 5507
 5477 et al., 2018; Hendrycks et al., 2021c) under 5508
 5478 *Apache 2.0 License*, evaluated by accuracy
- 5479 • **sufficient_information** under *Apache 2.0 Li- 5509
 5480 cense*, evaluated by accuracy
- 5481 • **suicide_risk** (Gaur et al., 2019; Mohammadi 5510
 5482 et al., 2019; Matero et al., 2019; Shing et al., 5511
 5483 2018) under *Apache 2.0 License*, evaluated by 5512
 5484 accuracy
- 5485 • **swahili_english_proverbs** under *Apache 2.0 5513
 5486 License*, evaluated by accuracy
- 5487 • **swedish_to_german_proverbs** (Hanzén, 5514
 5488 2007; Korhonen, 2009; Meister, 2007; 5515
 5489 Mieder, 2019) under *Apache 2.0 License*, 5516
 5490 evaluated by accuracy
- 5491 • **taboo** (Joshi et al., 2017) under *Apache 2.0 5517
 5492 License*, evaluated by accuracy
- 5493 • **talkdown** (Wang and Potts, 2019; Mendel- 5518
 5494 sohn et al., 2020; Fiske, 1993; Nolan 5519
 5495 and Mikami, 2013; Breitfeller et al., 2019; 5520
 5496 Perez Almendros et al., 2020) under *Apache 5521
 5497 2.0 License*, evaluated by accuracy

- 5498 • **temporal_sequences** (Elazar et al., 2021; 5544
5499 Pustejovsky et al., 2004; Sanampudi and Ku- 5545
5500 mari, 2010; Han et al., 2020; Ma et al., 2021; 5546
5501 Brown et al., 2020; Petroni et al., 2019) under 5502
5502 *Apache 2.0 License*, evaluated by accuracy 5503
5503 • **tense** (Logeswaran et al., 2018) under *Apache 5504
5504 2.0 License*, evaluated by accuracy 5505
5505 • **text_navigation_game** (Vinyals et al., 2019; 5506
5506 Küttler et al., 2020; Kanagawa and Kaneko, 5507
5507 2019; Noever et al., 2020) under *Apache 2.0 5508
5508 License*, evaluated by accuracy 5509
5509 • **timedial** (Qin et al., 2021; Li et al., 2017; 5510
5510 Zhou et al., 2019) under *Apache 2.0 License*, 5511
5511 evaluated by accuracy 5512
5512 • **topical_chat** (Gopalakrishnan et al., 2019; 5513
5513 Mehri and Eskenazi, 2020; Gopalakrishnan 5514
5514 et al., 2020; Hedayatnia et al., 2020) under 5515
5515 *Apache 2.0 License*, evaluated by accuracy 5516
5516 • **tracking_shuffled_objects** (Liu et al., 2020b; 5517
5517 Dong et al., 2020) under *Apache 2.0 License*, 5518
5518 evaluated by accuracy 5519
5519 • **training_on_test_set** under *Apache 2.0 5520
5520 License*, evaluated by accuracy 5521
5521 • **truthful_qa** (Brown et al., 2020; Sellam et al., 5522
5522 2020; Amodei et al., 2016; Leike et al., 2018; 5523
5523 Kenton et al., 2021; Clark et al., 2018; Bhak- 5524
5524 thavatsalam et al., 2021; Hendrycks et al., 5525
5525 2021b; Khashabi et al., 2020b; Kreps et al., 5526
5526 2020; Maynez et al., 2020; Gabriel et al., 5527
5527 2020; Wang et al., 2020a; Stiennon et al., 5528
5528 2020; Lewis et al., 2020b; Krishna et al., 2021; 5529
5529 Gehrmann et al., 2021; Xu et al., 2020a; Di- 5530
5530 nan et al., 2019; Tamkin et al., 2021; Bowman 5531
5531 and Dahl, 2021) under *Apache 2.0 License*, 5532
5532 evaluated by accuracy 5533
5533 • **twenty_questions** (Rajpurkar et al., 2016; 5534
5534 Choi et al., 2018; Reddy et al., 2019; Ali- 5535
5535 annejadi et al., 2019; Clark et al., 2019; Zhang 5536
5536 et al., 2019a) under *Apache 2.0 License*, 5537
5537 evaluated by accuracy 5538
5538 • **understanding_fables** (Reimers and 5539
5539 Gurevych, 2019; Salazar et al., 2020; Wolf 5540
5540 et al., 2019) under *Apache 2.0 License*, 5541
5541 evaluated by accuracy 5542
5542 • **undo_permutation** (Pham et al., 2020) under 5543
5543 *Apache 2.0 License*, evaluated by accuracy 5544
5544 • **unit_conversion** (Hendrycks et al., 2021c; 5545
5545 Geva et al., 2020a) under *Apache 2.0 License*, 5546
5546 evaluated by accuracy 5547
5547 • **unit_interpretation** under *Apache 2.0 Li- 5548
5548 cense*, evaluated by accuracy 5549
5549 • **unnatural_in_context_learning** (Brown 5550
5550 et al., 2020; Kaplan et al., 2020; Henighan 5551
5551 et al., 2020; Hernandez et al., 2021; Bahri 5552
5552 et al., 2021; Wang et al., 2019b; Hernandez 5553
5553 et al., 2020; Hendrycks et al., 2021c,a; Liu 5554
5554 et al., 2021a; Zhao et al., 2021; Perez et al., 5555
5555 2021) under *Apache 2.0 License*, evaluated by 5556
5556 accuracy 5557
5557 • **unqover** (Li et al., 2020a; Caliskan et al., 5558
5558 2017; Ruderger et al., 2018; Zhao et al., 2018; 5559
5559 Dev et al., 2020; Stanovsky et al., 2019; 5560
5560 Nadeem et al., 2020; Sheng et al., 2019; 5561
5561 Zhang et al., 2020b) under *Apache 2.0 Li- 5562
5562 cense*, evaluated by accuracy 5563
5563 • **web_of_lies** under *Apache 2.0 License*, 5564
5564 evaluated by accuracy 5565
5565 • **what_is_the_tao** under *Apache 2.0 License*, 5566
5566 evaluated by accuracy 5567
5567 • **which_wiki_edit** under *Apache 2.0 License*, 5568
5568 evaluated by accuracy 5569
5569 • **word_problems_on_sets_and_graphs** (Bency 5570
5570 et al., 2019; Mnih et al., 2013; Russell and 5571
5571 Norvig, 2002; Besold et al., 2017; Clark et al., 5572
5572 2020; Wang et al., 2018; Lacker, 2020) under 5573
5573 *Apache 2.0 License*, evaluated by accuracy 5574
5574 • **word_sorting** (Grover et al., 2019) under 5575
5575 *Apache 2.0 License*, evaluated by accuracy 5576
5576 • **word_unscrambling** (Nishino et al., 2019; 5577
5577 Rozner et al., 2021; Jones et al., 2020; Mays 5578
5578 et al., 1991; Edizel et al., 2019; Sakaguchi 5579
5579 et al., 2016; Kim et al., 2015; Xue et al., 5580
5580 2021a; Wu et al., 2020; Rust et al., 2020) un- 5581
5581 der *Apache 2.0 License*, evaluated by accuracy 5582
5582 • **yes_no_black_white** under *Apache 2.0 Li- 5583
5583 cense*, evaluated by accuracy 5584
5584 **D.4 BigBench-Lite Datasets** 5585
5585 • **auto_debugging** (Zaremba and Sutskever, 5586
5586 2014) under *Apache 2.0 License*, evaluated 5587
5587 by accuracy

5588	• bbq_lite_json (Crawford, 2017; Khashabi et al., 2020b; Li et al., 2020a) under <i>Apache 2.0 License</i> , evaluated by accuracy	5631
5589		5632
5590		
5591	• code_line_description (Alon et al., 2018) under <i>Apache 2.0 License</i> , evaluated by accuracy	5633
5592		5634
5593	• conceptual_combinations (Fodor, 1975; Fodor and Pylyshyn, 1988; Smolensky, 1988; Lake et al., 2017; Lake and Murphy, 2020; Marcus, 2020; Henrich et al., 2010; Murphy, 1988) under <i>Apache 2.0 License</i> , evaluated by accuracy	5635
5594		5636
5595		
5596	• conlang_translation (Canfield, 2010; Şahin et al., 2020; Sennrich and Zhang, 2019) under <i>Apache 2.0 License</i> , evaluated by accuracy	5637
5597		5638
5598		5639
5599	• emoji_movie (Cruse, 2015; Instagram Engineering, 2015; Chandra Guntuku et al., 2019; Eisner et al., 2016; Mayne, 2020; Boillot, 2019) under <i>Apache 2.0 License</i> , evaluated by accuracy	5640
5600		
5601	• formal_fallacies_syllogisms_negation (Kassner and Schütze, 2019; Talmor et al., 2019a; Betz et al., 2020) under <i>Apache 2.0 License</i> , evaluated by accuracy	5641
5602		5642
5603		5643
5604	• hindu_knowledge under <i>Apache 2.0 License</i> , evaluated by accuracy	5644
5605		5645
5606		
5607	• known_unknowns (Liu et al., 2021c; Xiao and Wang, 2021; Shuster et al., 2021; Zhou et al., 2020; Dziri et al., 2021) under <i>Apache 2.0 License</i> , evaluated by accuracy	5646
5608		5647
5609		5648
5610	• language_identification (Brown, 2014) under <i>Apache 2.0 License</i> , evaluated by accuracy	5649
5611		5650
5612	• linguistics_puzzles (Bozhanov and Derzhanski, 2013; Radev et al., 2008; Sennrich and Zhang, 2019; Clark et al., 2018; Şahin et al., 2020) under <i>Apache 2.0 License</i> , evaluated by accuracy	5651
5613		5652
5614		5653
5615		5654
5616	• logic_grid_puzzle under <i>Apache 2.0 License</i> , evaluated by accuracy	5655
5617		
5618	• logical_deduction under <i>Apache 2.0 License</i> , evaluated by accuracy	5656
5619		5657
5620	• misconceptions_russian (Thorne et al., 2018; Lee et al., 2020) under <i>Apache 2.0 License</i> , evaluated by accuracy	5658
5621		
5622		
5623		
5624		
5625		
5626		
5627		
5628		
5629		
5630		
5631	• novel_concepts (Santoro et al., 2021) under <i>Apache 2.0 License</i> , evaluated by accuracy	5665
5632		5666
5633	• operators (Brown et al., 2020; Kassner et al., 2020; Hendrycks et al., 2021c; Saxton et al., 2019) under <i>Apache 2.0 License</i> , evaluated by accuracy	5667
5634		5668
5635		5669
5636		5670
5637	• parsinlu_reading_comprehension (Khashabi et al., 2020a; Xue et al., 2021b; Rajpurkar et al., 2016) under <i>Apache 2.0 License</i> , evaluated by accuracy	5671
5638		5672
5639		5673
5640		5674
5641	• play_dialog_same_or_different (Park et al., 2021) under <i>Apache 2.0 License</i> , evaluated by accuracy	
5642		
5643		
5644	• repeat_copy_logic (Graves et al., 2014) under <i>Apache 2.0 License</i> , evaluated by accuracy	
5645		
5646	• strange_stories (Happé, 1994; White et al., 2009) under <i>Apache 2.0 License</i> , evaluated by accuracy	
5647		
5648		
5649	• strategyqa (Geva et al., 2021) under <i>Apache 2.0 License</i> , evaluated by accuracy	
5650		
5651	• symbol_interpretation (Brown et al., 2020; Johnson et al., 2016; Santoro et al., 2017; Hudson and Manning, 2019; Sennrich et al., 2015) under <i>Apache 2.0 License</i> , evaluated by accuracy	
5652		
5653		
5654		
5655		
5656	• vitaminc_fact_verification (Schuster et al., 2021) under <i>Apache 2.0 License</i> , evaluated by accuracy	
5657		
5658		
5659	• winowhy (Zhang et al., 2020c; Devlin et al., 2018; Liu et al., 2019; Kocijan et al., 2019; Rahman and Ng, 2012; Sakaguchi et al., 2020b) under <i>Apache 2.0 License</i> , evaluated by accuracy	
5660		
5661		
5662		
5663		
5664	E Efficiency Analysis	
5665	GLIDER introduces minimal computational overhead by requiring only two lightweight operations per LoRA layer: a single cosine similarity calculation between the query’s task embedding and global routing vectors and a simple vector addition to combine this with the local routing score. With typical values for N (experts) and d_g (embedding dimension) in the hundreds, this amounts to just $(N \times d_g + d_g)$ FLOPs per layer, which is negligible compared to the base model’s computation.	
5666		
5667		
5668		
5669		
5670		
5671		
5672		
5673		
5674		

5675 **F Detailed Performance**

5676 We list detailed performance of all baselines for
5677 each task in Table 4, 6, 5, 7, and 7.

Method	Avg	RTE	H-Swag	COPA	WIC	Winogrande	CB	StoryCloze	ANLI-R1	ANLI-R2	ANLI-R3	WSC
Multi-Task Fine-Tuning	51.6	60.1	26.9	74.8	51.3	50.9	52.7	85.1	34.7	33	33.5	64.9
Oracle Expert	57.2	66.9	36.8	89.6	52.4	57.6	59.9	96.9	34.5	34.8	36.7	63.6
Merged Experts	45.4	55.7	25.7	61.8	50.3	53.3	45.6	63.8	33.1	33.4	33.4	43.5
LoRA Hub	46.9	52.1	27.2	75.1	50.4	50.6	33.5	84.3	33.2	34.2	33.6	41.1
Glider	57.3	64.8	29.9	91.6	50.6	60.2	69.6	96.2	36.0	34.3	38.1	58.5

Table 4: Complete results on T0 Held-Out datasets.

Expert	Avg	Boolean Expression	Causal Judgment	Date Understanding	Disambiguator QA	Formal Fallacies	Geometric Shapes	Hyperbaton
Multi-Task Fine-Tuning	34.9	49.6	55.3	35.2	55.4	51.5	10.6	50
Oracle Expert	42.2	64.4	59.5	41.7	65.1	51.7	30.1	52.7
Merged Experts	35.3	60.8	58.4	38.5	45.3	50.1	10	50
LoRA Hub	32.0	59.2	49.5	36.6	30.2	50.9	24.2	50.0
Glider	34.9	52.8	59.5	39.6	57.0	50.1	10.3	50.0
Expert	Logical Detection	Movie Recommendation	Multistep Arithmetic	Navigate Object Counting	Penguins in a Table	Reasoning about Colored Objects	Ruin Names	
Multi-Task Fine-Tuning	47.9	34.8	0	50	22.5	32.9	42	19.6
Oracle Expert	45.8	49.2	1.6	50	28.1	36.9	53.2	49.6
Merged Experts	44.3	23	0.4	50	25.4	35.6	47.5	26.6
LoRA Hub	27.5	32.0	1.2	50.0	0.0	30.9	19.3	21.2
Glider	40.8	31.2	1.6	50.0	19.7	34.2	48.4	24.3
Expert	Salient Translation Error Detection	Snarks Understanding	Sports Sequences	Temporal Sequences	Track Shuffled Objects	Web of Lies	Word Sorting	
Multi-Task Fine-Tuning	27.8	46.4	50.2	16.1	17.4	51.6	0.5	
Oracle Expert	27	61.3	50.9	28.2	20.1	59.2	2.8	
Merged Experts	24.9	44.8	51.7	19.5	17	51.6	0.9	
LoRA Hub	24.9	43.6	50.3	12.8	19.7	55.6	0.0	
Glider	25.6	48.1	51.4	12.8	16.1	52.8	0.0	

Table 5: BIG-bench Hard (BBH) results of different methods in T0 Held-In setting

Expert	Avg	BBQ Lite	Conceptual Combinations	Conlong Translation	Formal Fallacies	Hindu Knowledge
Multi-Task Fine-Tuning	36.6	40.8	44.7	26	51.5	40.6
Oracle Expert	43.5	55.3	62.1	29.8	51.6	46.3
Merged Experts	36	42.5	33	28.9	50.1	40
LoRA Hub	32.4	40.1	39.8	0.2	50.1	29.1
Glider	37.5	48.4	44.7	17.1	50.1	48.6

Expert	Known Unknowns	Linguistic Puzzles	Logic Grid Puzzle	Logical Detection	Novel Concepts	Operators
Multi-Task Fine-Tuning	47.8	0	35.9	48.1	40.6	1
Oracle Expert	65.2	0	41.7	45.4	43.8	8.6
Merged Experts	45.7	0	39.6	44.3	28.1	7.1
LoRA Hub	45.7	0.0	32.8	20.1	28.1	5.7
Glider	52.2	0.0	39.6	40.8	43.8	2.9

Expert	Play Dialog Same or Different	Repeat Copy Logic	Strange Stories	Strategy QA	Vitamin C Fact Verification	Winowhy
Multi-Task Fine-Tuning	45.8	0	47.7	52.5	54.2	44.3
Oracle Expert	63.3	0	68.4	56.1	51.1	50.5
Merged Experts	36.9	0	56.3	54.3	61.3	44.3
LoRA Hub	47.5	0.0	43.7	52.9	49.9	44.3
Glider	36.9	0.0	63.8	53.6	49.8	44.5

Table 6: BIG-bench Lite (BBL) results of different methods in T0 Held-In setting

Expert	Avg	Boolean Expression	Causal Judgment	Date Understanding	Disambiguator QA	Formal Fallacies	Geometric Shapes	Hyperbaton
Multi-Task Fine-Tuning	38.9	50	61.1	36.6	65.9	52.2	9.7	51.1
Oracle Expert	45.5	66	59.5	42.3	65.1	52.9	30.1	69.3
Merged Experts	34.6	53.6	56.8	36.9	45.7	50	12	52.2
LoRA Hub	32.8	55.2	54.2	26.8	30.2	50.0	19.5	50.0
Glider	35.3	50.8	58.4	35.2	50.4	50.5	10.3	48.2
Expert	Logical Detection	Movie Recommendation	Multistep Arithmetic	Navigate	Object Counting	Penguins in a Table	Reasoning about Colored Objects	Ruin Names
Multi-Task Fine-Tuning	49.6	32.8	0	50	35.7	39.6	56.6	19
Oracle Expert	48.8	49.2	1.6	54.6	45.7	37.6	53.5	49.6
Merged Experts	42.9	23.2	0.8	50	24.1	34.2	44.5	28.3
LoRA Hub	31.4	33.8	0.0	50.0	0.0	29.5	25.1	24.6
Glider	40.1	28.2	0.4	50.0	41.1	32.2	46.9	33.0
Expert	Salient Translation Error Detection	Snarks	Sports Understanding	Temporal Sequences	Track Shuffled Objects	Web of Lies	Word Sorting	
Multi-Task Fine-Tuning	39.2	59.7	51.2	27.2	15.7	53.6	0	0
Oracle Expert	31.5	61.3	52.5	45.3	20.3	61.6	2.9	2.9
Merged Experts	26.4	40.9	49.9	22.4	16.3	49.6	1.2	1.2
LoRA Hub	12.4	48.1	50.3	100.0	19.1	48.8	0.0	0.0
Glider	23.8	41.4	49.6	8.8	17.2	52.0	0.1	0.1

Table 7: BIG-bench Hard (BBH) results of different methods in the FLAN setting.

Expert	Avg	BBQ Lite	Conceptual Combinations	Conlong Translation	Formal Fallacies	Hindu Knowledge
Multi-Task Fine-Tuning	45.4	66.9	72.8	27.9	52.2	40
Oracle Expert	46.5	61.3	62.1	36	52.9	46.3
Merged Experts	34	40.2	27.2	29.1	50	36.6
LoRA Hub	31.8	36.0	29.1	2.6	50.0	25.7
Glider	35.5	49.4	35.9	10.1	50.5	47.4

Expert	Known Unknowns	Linguistic Puzzles	Logic Grid Puzzle	Logical Detection	Novel Concepts	Operators
Multi-Task Fine-Tuning	58.7	0	42.8	49.9	37.5	13.3
Oracle Expert	65.2	0	42.5	48.6	50	12.4
Merged Experts	43.5	0	37.9	42.9	34.4	7.6
LoRA Hub	52.2	0.0	29.3	26.2	34.4	0.0
Glider	41.3	0.0	34.4	40.1	25.0	3.3

Expert	Play Dialog Same or Different	Repeat Copy Logic	Strange Stories	Strategy QA	Vitamin C Fact Verification	Winowhy
Multi-Task Fine-Tuning	44.4	0	75.9	65.7	78.5	45.3
Oracle Expert	63.3	0	74.1	56.1	66.7	53.8
Merged Experts	36.9	0	46.6	52.1	47.9	44.3
LoRA Hub	63.1	0.0	28.2	46.7	51.4	44.3
Glider	37.6	0.0	61.5	52.6	66.5	48.3

Table 8: BIG-bench Lite (BBL) results of different methods in the FLAN setting.