# 15 papers that shaped
# NLP & LLMs

**2017**

Attention Is All You Need

→

**2018**

BERT

→

**2020**

GPT-3: Few-Shot Learners

**2020**

RAG

←

**2020**

Scaling Laws

←

**2020**

T5

**2021**

LoRA

→

**2022**

Chain-of-Thought Prompting

→

**2022**

Self-Consistency

**2022**

Toolformer

←

**2022**

Instruction Tuning

←

**2022**

In-Context Learning & Induction Heads

**2022**

ColBERTv2

→

**2023**

LLMs as a Judge

→

**2025**

DeepSeek-R1
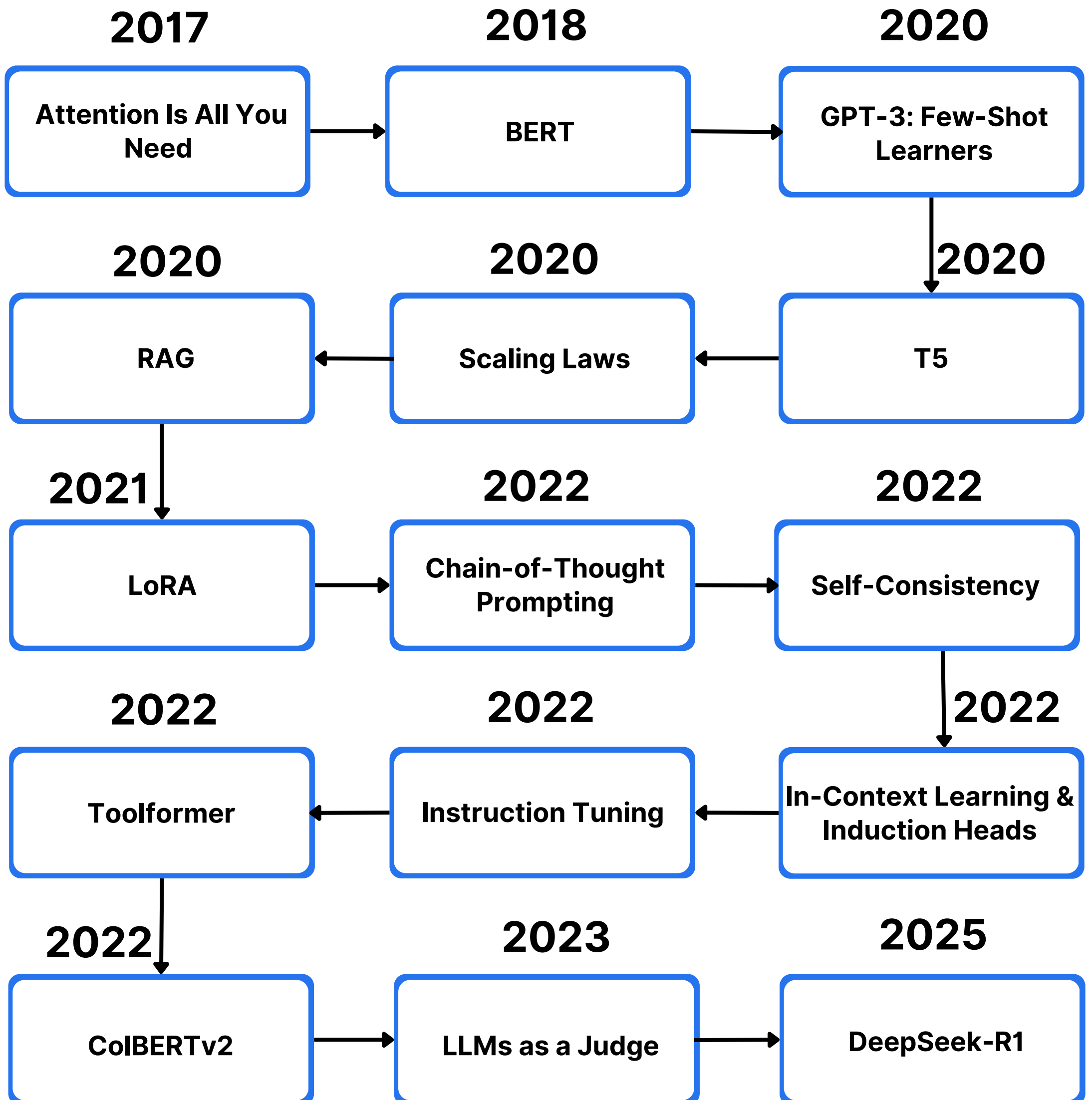
# Attention Is All You Need (2017)

Transformer backbone of all modern LLMs.

⚙️ No recurrence
🚀 Massive parallelism
🧠 Foundation for GPT, BERT

---

# Attention Is All You Need

**Ashish Vaswani***
Google Brain
avaswani@google.com

**Noam Shazeer***
Google Brain
noam@google.com

**Niki Parmar***
Google Research
nikip@google.com

**Jakob Uszkoreit***
Google Research
usz@google.com

**Llion Jones***
Google Research
llion@google.com

**Aidan N. Gomez*** [†]
University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser***
Google Brain
lukaszkaiser@google.com

**Illia Polosukhin*** [‡]
illia.polosukhin@gmail.com

## Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after

# BERT (2018)

Brought transfer learning to NLP.

- 📦 Pretrain + fine-tune
- 🔄 Bidirectional context
- 🧰 Spawned countless variants

# BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Jacob Devlin    Ming-Wei Chang    Kenton Lee    Kristina Toutanova

Google AI Language

{jacobdevlin,mingweichang,kentonl,kristout}@google.com

## Abstract

We introduce a new language representation model called **BERT**, which stands for **B**idirectional **E**ncoder **R**epresentations from **T**ransformers. Unlike recent language representation models (Peters et al., 2018a; Radford et al., 2018), BERT is designed to pre-train deep bidirectional representations from unlabeled text by jointly conditioning on both left and right context in all layers. As a result, the pre-trained BERT model can be fine-tuned with just one additional output layer to create state-of-the-art models for a wide range of tasks, such as question answering and language inference, without substantial task-specific architecture modifications.

BERT is conceptually simple and empirically powerful. It obtains new state-of-the-art results on eleven natural language processing tasks, including pushing the GLUE score to

There are two existing strategies for applying pre-trained language representations to downstream tasks: *feature-based* and *fine-tuning*. The feature-based approach, such as ELMo (Peters et al., 2018a), uses task-specific architectures that include the pre-trained representations as additional features. The fine-tuning approach, such as the Generative Pre-trained Transformer (OpenAI GPT) (Radford et al., 2018), introduces minimal task-specific parameters, and is trained on the downstream tasks by simply fine-tuning *all* pre-trained parameters. The two approaches share the same objective function during pre-training, where they use unidirectional language models to learn general language representations.

We argue that current techniques restrict the power of the pre-trained representations, especially for the fine-tuning approaches. The ma-

# GPT-3: Few-Shot Learners (2020)

Scale unlocks capabilities.

📈 175B parameters
🧠 Emergent reasoning
🪄 Prompt-based interface

---

# Language Models are Few-Shot Learners

Tom B. Brown*        Benjamin Mann*        Nick Ryder*        Melanie Subbiah*

Jared Kaplan[†]        Prafulla Dhariwal        Arvind Neelakantan        Pranav Shyam        Girish Sastry

Amanda Askell        Sandhini Agarwal        Ariel Herbert-Voss        Gretchen Krueger        Tom Henighan

Rewon Child        Aditya Ramesh        Daniel M. Ziegler        Jeffrey Wu        Clemens Winter

Christopher Hesse        Mark Chen        Eric Sigler        Mateusz Litwin        Scott Gray

Benjamin Chess        Jack Clark        Christopher Berner

Sam McCandlish        Alec Radford        Ilya Sutskever        Dario Amodei

OpenAI

## Abstract

Recent work has demonstrated substantial gains on many NLP tasks and benchmarks by pre-training on a large corpus of text followed by fine-tuning on a specific task. While typically task-agnostic in architecture, this method still requires task-specific fine-tuning datasets of thousands or tens of thousands of examples. By contrast, humans can generally perform a new language task from only a few examples or from simple instructions – something which current NLP systems still largely struggle to do. Here we show that scaling up language models greatly improves task-agnostic, few-shot performance, sometimes even reaching competitiveness with prior state-of-the-art fine-tuning approaches. Specifically, we train GPT-3, an autoregressive language model with 175 billion parameters, 10x more than any previous non-sparse language model, and test its performance in the few-shot setting. For all tasks, GPT-3 is applied without any gradient updates or fine-tuning, with tasks and few-shot demonstrations specified purely via text interaction with the model. GPT-3 achieves strong performance on many NLP datasets, including translation, question-answer

# T5 (2020)

Unified all NLP as text-to-text.

🔤 One framework for many tasks
🧠 Simple and powerful

## Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer

Colin Raffel*                                           CRAFFEL@GMAIL.COM
Noam Shazeer*                                           NOAM@GOOGLE.COM
Adam Roberts*                                           ADAROB@GOOGLE.COM
Katherine Lee*                                          KATHERINELEE@GOOGLE.COM
Sharan Narang                                           SHARANNARANG@GOOGLE.COM
Michael Matena                                          MMATENA@GOOGLE.COM
Yanqi Zhou                                              YANQIZ@GOOGLE.COM
Wei Li                                                  MWEILI@GOOGLE.COM
Peter J. Liu                                            PETERJLIU@GOOGLE.COM
*Google, Mountain View, CA 94043, USA*

## Abstract

Transfer learning, where a model is first pre-trained on a data-rich task before being fine-tuned on a downstream task, has emerged as a powerful technique in natural language processing (NLP). The effectiveness of transfer learning has given rise to a diversity of approaches, methodology, and practice. In this paper, we explore the landscape of transfer learning techniques for NLP by introducing a unified framework that converts all text-based language problems into a text-to-text format. Our systematic study compares pre-training objectives, architectures, unlabeled data sets, transfer approaches, and other factors on dozens of language understanding tasks. By combining the insights from our exploration with scale and our new "Colossal Clean Crawled Corpus", we achieve state-of-the-art results on many benchmarks covering summarization, question answering, text classification, and

# Scaling Laws (2020)

Showed how performance scales.

📊 Predictable gains
🧠 Blueprint for large models

---

# Scaling Laws for Neural Language Models

**Jared Kaplan** *

Johns Hopkins University, OpenAI

jaredk@jhu.edu

**Sam McCandlish***

OpenAI

sam@openai.com

**Tom Henighan**

OpenAI

henighan@openai.com

**Tom B. Brown**

OpenAI

tom@openai.com

**Benjamin Chess**

OpenAI

bchess@openai.com

**Rewon Child**

OpenAI

rewon@openai.com

**Scott Gray**

OpenAI

scott@openai.com

**Alec Radford**

OpenAI

alec@openai.com

**Jeffrey Wu**

OpenAI

jeffwu@openai.com

**Dario Amodei**

OpenAI

damodei@openai.com

## Abstract

We study empirical scaling laws for language model performance on the cross-entropy loss. The loss scales as a power-law with model size, dataset size, and the amount of compute used for training, with some trends spanning more than seven orders of magnitude. Other architectural details such as network width or depth have minimal effects within a wide range. Simple equations govern the dependence of overfitting on model/dataset size and the dependence of training speed on model size. These relationships allow us to determine the optimal allocation of a fixed compute budget. Larger models are significantly more sample-efficient, such that optimally compute-efficient training involves training very large models on a relatively modest amount of data and stopping significantly before convergence.

# RAG (2020)

Combined retrieval and generation.

🔍 External knowledge
📚 Factual grounding

---

# Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks

Patrick Lewis[†‡], Ethan Perez[*],

Aleksandra Piktus[†], Fabio Petroni[†], Vladimir Karpukhin[†], Naman Goyal[†], Heinrich Küttler[†],

Mike Lewis[†], Wen-tau Yih[†], Tim Rocktäschel[†‡], Sebastian Riedel[†‡], Douwe Kiela[†]

[†]Facebook AI Research; [‡]University College London; [*]New York University;
plewis@fb.com

## Abstract

Large pre-trained language models have been shown to store factual knowledge in their parameters, and achieve state-of-the-art results when fine-tuned on downstream NLP tasks. However, their ability to access and precisely manipulate knowledge is still limited, and hence on knowledge-intensive tasks, their performance lags behind task-specific architectures. Additionally, providing provenance for their decisions and updating their world knowledge remain open research problems. Pre-trained models with a differentiable access mechanism to explicit non-parametric memory have so far been only investigated for extractive downstream tasks. We explore a general-purpose fine-tuning recipe for retrieval-augmented generation (RAG) — models which combine pre-trained parametric and non-parametric memory for language generation. We introduce RAG models where the parametric memory is a pre-trained seq2seq model and the non-parametric memory is a dense vector index of Wikipedia, accessed with a pre-trained neural retriever. We compare two RAG formulations, one which conditions on the same retrieved passages across the whole generated sequence, and another which can use different passages per token. We fine-tune and evaluate our models on a wide range of knowledge-intensive NLP tasks and set the state of the art on three open domain QA tasks, outperforming parametric seq2seq models and task-specific retrieve-and-extract architectures. For language generation tasks, we find that RAG models generate more specific, diverse and factual language than a state-of-the-art parametric-only seq2seq baseline.

# LoRA (2021)

Fine-tune big models cheaply.

💸 Low-cost adaptation
🧰 Enterprise-ready

# LoRA: Low-Rank Adaptation of Large Language Models

**Edward Hu**[*]  **Yelong Shen**[*]  **Phillip Wallis**  **Zeyuan Allen-Zhu**
**Yuanzhi Li**  **Shean Wang**  **Lu Wang**  **Weizhu Chen**
Microsoft Corporation
{edwardhu, yeshe, phwallis, zeyuana,
yuanzhil, swang, luw, wzchen}@microsoft.com
yuanzhil@andrew.cmu.edu
(Version 2)

## ABSTRACT

An important paradigm of natural language processing consists of large-scale pre-training on general domain data and adaptation to particular tasks or domains. As we pre-train larger models, full fine-tuning, which retrains all model parameters, becomes less feasible. Using GPT-3 175B as an example – deploying independent instances of fine-tuned models, each with 175B parameters, is prohibitively expensive. We propose **Low-R**ank **A**daptation, or LoRA, which freezes the pre-trained model weights and injects trainable rank decomposition matrices into each layer of the Transformer architecture, greatly reducing the number of trainable parameters for downstream tasks. Compared to GPT-3 175B fine-tuned with Adam, LoRA can reduce the number of trainable parameters by 10,000 times and the GPU memory requirement by 3 times. LoRA performs on-par or better than fine-tuning in model quality on RoBERTa, DeBERTa, GPT-2, and GPT-3, despite having fewer trainable parameters, a higher training throughput, and, unlike adapters, *no additional inference latency*. We also provide an empirical investigation into rank-deficiency in language model adaptation, which sheds light on the efficacy of LoRA. We release a package that facilitates the integration of LoRA with PyTorch models and provide our implementations and model checkpoints for RoBERTa, DeBERTa, and GPT-2 at `https://github.com/microsoft/LoRA`.

# 1  INTRODUCTION

# CoT Prompting (2022)

Prompting unlocks reasoning.

🧠 Multi-step thinking
✏️ Better complex-task performance

---

## Chain-of-Thought Prompting Elicits Reasoning in Large Language Models

Jason Wei        Xuezhi Wang        Dale Schuurmans        Maarten Bosma

Brian Ichter        Fei Xia        Ed H. Chi        Quoc V. Le        Denny Zhou

Google Research, Brain Team
{jasonwei,dennyzhou}@google.com

### Abstract

We explore how generating a *chain of thought*—a series of intermediate reasoning steps—significantly improves the ability of large language models to perform complex reasoning. In particular, we show how such reasoning abilities emerge naturally in sufficiently large language models via a simple method called *chain-of-thought prompting*, where a few chain of thought demonstrations are provided as exemplars in prompting.

Experiments on three large language models show that chain-of-thought prompting improves performance on a range of arithmetic, commonsense, and symbolic reasoning tasks. The empirical gains can be striking. For instance, prompting a PaLM 540B with just eight chain-of-thought exemplars achieves state-of-the-art accuracy on the GSM8K benchmark of math word problems, surpassing even finetuned GPT-3 with a verifier.

| Standard Prompting | Chain-of-Thought Prompting |
|---|---|
| **Model Input**<br><br>Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?<br><br>A: The answer is 11.<br><br>Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples | **Model Input**<br><br>Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?<br><br>A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.<br><br>Q: The cafeteria had 23 apples. If they used 20 to |

# Self-Consistency (2022)

Majority vote boosts reasoning.

🔄 More reliable outputs
📈 Better logic tasks

# SELF-CONSISTENCY IMPROVES CHAIN OF THOUGHT REASONING IN LANGUAGE MODELS

**Xuezhi Wang**[†‡]   **Jason Wei**[†]   **Dale Schuurmans**[†]   **Quoc Le**[†]   **Ed H. Chi**[†]
**Sharan Narang**[†]   **Aakanksha Chowdhery**[†]   **Denny Zhou**[†§]
[†]Google Research, Brain Team
[‡]xuezhiw@google.com, [§]dennyzhou@google.com

## ABSTRACT

Chain-of-thought prompting combined with pre-trained large language models has achieved encouraging results on complex reasoning tasks. In this paper, we propose a new decoding strategy, *self-consistency*, to replace the naive greedy decoding used in chain-of-thought prompting. It first samples a diverse set of reasoning paths instead of only taking the greedy one, and then selects the most consistent answer by marginalizing out the sampled reasoning paths. Self-consistency leverages the intuition that a complex reasoning problem typically admits multiple different ways of thinking leading to its unique correct answer. Our extensive empirical evaluation shows that self-consistency boosts the performance of chain-of-thought prompting with a striking margin on a range of popular arithmetic and commonsense reasoning benchmarks, including GSM8K (+17.9%), SVAMP (+11.0%), AQuA (+12.2%), StrategyQA (+6.4%) and ARC-challenge (+3.9%).

## 1 INTRODUCTION

Although language models have demonstrated remarkable success across a range of NLP tasks, their ability to demonstrate reasoning is often seen as a limitation, which cannot be overcome solely by increasing model scale (Rae et al., 2021; BIG-bench collaboration, 2021, *inter alia*). In an effort to address this shortcoming, Wei et al. (2022) have proposed *chain-of-thought prompting*, where a language model is prompted to generate a series of short sentences that mimic the reasoning process a person might employ in solving a task. For example, given the question *"If there are 3 cars in the parking lot and 2 more cars arrive, how many cars are in the parking lot?"*, instead of directly responding with *"5"*, a language model would be prompted to respond with the entire chain-of-thought: *"There are 3 cars in the parking lot already. 2 more arrive. Now there are 3 + 2 = 5 cars. The answer is 5."*. It has been observed that chain-of-thought prompting significantly improves model performance across a variety of multi-step reasoning tasks (Wei et al., 2022).

# In-Context Learning & Induction Heads (2022)

Explained how LLMs learn from context.

🧩 Mechanistic insights
🔍 Key to interpretability

---

## Language Models (Mostly) Know What They Know

Saurav Kadavath,* Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, Scott Johnston, Sheer El-Showk, Andy Jones, Nelson Elhage, Tristan Hume, Anna Chen, Yuntao Bai, Sam Bowman, Stanislav Fort, Deep Ganguli, Danny Hernandez, Josh Jacobson, Jackson Kernion, Shauna Kravec, Liane Lovitt, Kamal Ndousse, Catherine Olsson, Sam Ringer, Dario Amodei, Tom Brown, Jack Clark, Nicholas Joseph, Ben Mann, Sam McCandlish, Chris Olah, Jared Kaplan*

Anthropic

### Abstract

We study whether language models can evaluate the validity of their own claims and predict which questions they will be able to answer correctly. We first show that larger models are well-calibrated on diverse multiple choice and true/false questions when they are provided in the right format. Thus we can approach self-evaluation on open-ended sampling tasks by asking models to first propose answers, and then to evaluate the probability "P(True)" that their answers are correct. We find encouraging performance, calibration, and scaling for P(True) on a diverse array of tasks. Performance at self-evaluation further improves when we allow models to consider many of their own samples before predicting the validity of one specific possibility. Next, we investigate whether models can be trained to predict "P(IK)", the probability that "I know" the answer to a question, without reference to any particular proposed answer. Models perform well at predicting P(IK) and partially generalize across tasks, though they struggle with calibration of P(IK) on new tasks. The predicted P(IK) probabilities also increase appropriately in the presence of relevant source

# Instruction Tuning (2022)

Made LLMs follow natural language.

👨‍🏫 Conversational skills
🧰 No retraining needed

---

## Training language models to follow instructions with human feedback

**Long Ouyang***    **Jeff Wu***    **Xu Jiang***    **Diogo Almeida***    **Carroll L. Wainwright***

**Pamela Mishkin***    **Chong Zhang**    **Sandhini Agarwal**    **Katarina Slama**    **Alex Ray**

**John Schulman**    **Jacob Hilton**    **Fraser Kelton**    **Luke Miller**    **Maddie Simens**

**Amanda Askell[†]**    **Peter Welinder**    **Paul Christiano*[†]**

**Jan Leike***    **Ryan Lowe***

OpenAI

### Abstract

Making language models bigger does not inherently make them better at following a user's intent. For example, large language models can generate outputs that are untruthful, toxic, or simply not helpful to the user. In other words, these models are not *aligned* with their users. In this paper, we show an avenue for aligning language models with user intent on a wide range of tasks by fine-tuning with human feedback. Starting with a set of labeler-written prompts and prompts submitted through the OpenAI API, we collect a dataset of labeler demonstrations of the desired model behavior, which we use to fine-tune GPT-3 using supervised learning. We then collect a dataset of rankings of model outputs, which we use to further fine-tune this supervised model using reinforcement learning from human feedback. We call the resulting models *InstructGPT*. In human evaluations on our prompt distribution, outputs from the 1.3B parameter InstructGPT model are preferred to outputs from the 175B GPT-3, despite having 100x fewer parameters. Moreover, InstructGPT models show improvements in truthfulness and reductions

# Toolformer (2023)

Models teach themselves to use tools.

🛠️ API calling
🤖 Planning abilities

# Toolformer: Language Models Can Teach Themselves to Use Tools

**Timo Schick**    **Jane Dwivedi-Yu**    **Roberto Dessì**[†]    **Roberta Raileanu**

**Maria Lomeli**    **Luke Zettlemoyer**    **Nicola Cancedda**    **Thomas Scialom**

Meta AI Research    [†]Universitat Pompeu Fabra

## Abstract

Language models (LMs) exhibit remarkable abilities to solve new tasks from just a few examples or textual instructions, especially at scale. They also, paradoxically, struggle with basic functionality, such as arithmetic or factual lookup, where much simpler and smaller models excel. In this paper, we show that LMs can teach themselves to *use external tools* via simple APIs and achieve the best of both worlds. We introduce *Toolformer*, a model trained to decide which APIs to call, when to call them, what arguments to pass, and how to best incorporate the results into future token prediction. This is done in a self-supervised way, requiring nothing more than a handful of demonstrations for each API. We incorporate a range of tools, including a calculator, a Q&A system, a search engine, a translation system, and a calendar. Toolformer achieves substantially improved zero-shot performance across a variety of downstream tasks, often competitive with much larger models, without sacrificing its core language modeling abilities.

The New England Journal of Medicine is a registered trademark of [QA("Who is the publisher of The New England Journal of Medicine?") → Massachusetts Medical Society] the MMS.

Out of 1400 participants, 400 (or [Calculator(400 / 1400) → 0.29] 29%) passed the test.

The name derives from "la tortuga", the Spanish word for [MT("tortuga") → turtle] turtle.

The Brown Act is California's law [WikiSearch("Brown Act") → The Ralph M. Brown Act is an act of the California State Legislature that guarantees the public's right to attend and participate in meetings of local legislative bodies.] that requires legislative bodies, like city councils, to hold their meetings open to the public.

Figure 1: Exemplary predictions of Toolformer. The model autonomously decides to call different APIs (from top to bottom: a question answering system, a calculator, a machine translation system, and a Wikipedia search engine) to obtain information that is useful for completing a piece of text.

# ColBERTv2 (2022)

Late interaction for retrieval.

🔍 Efficient + accurate search
⚙ Scales to billions

# ColBERTv2:
# Effective and Efficient Retrieval via Lightweight Late Interaction

**Keshav Santhanam**[*]
Stanford University

**Omar Khattab**[*]
Stanford University

**Jon Saad-Falcon**
Georgia Institute of Technology

**Christopher Potts**
Stanford University

**Matei Zaharia**
Stanford University

## Abstract

Neural information retrieval (IR) has greatly advanced search and other knowledge-intensive language tasks. While many neural IR methods encode queries and documents into single-vector representations, late interaction models produce multi-vector representations at the granularity of each token and decompose relevance modeling into scalable token-level computations. This decomposition has been shown to make late interaction more effective, but it inflates the space footprint of these models by an order of magnitude. In this work, we introduce ColBERTv2, a retriever that couples an aggressive residual compression mechanism with a denoised supervision strategy to simultaneously improve the quality and space footprint of late interaction. We evaluate ColBERTv2 across a wide range of benchmarks, establishing state-of-the-art quality within and outside the training domain while reducing the space footprint of late interaction models by 6–10×.

## 1 Introduction

Neural information retrieval (IR) has quickly dominated the search landscape over the past 2–3 years, dramatically advancing not only passage and document search (Nogueira and Cho, 2019) but also

relevance is estimated using rich yet scalable interactions between these two sets of vectors. ColBERT produces an embedding for every token in the query (and document) and models relevance as the sum of maximum similarities between each query vector and all vectors in the document.

By decomposing relevance modeling into token-level computations, late interaction aims to reduce the burden on the encoder: whereas single-vector models must capture complex query–document relationships within one dot product, late interaction encodes meaning at the level of tokens and delegates query–document matching to the interaction mechanism. This added expressivity comes at a cost: existing late interaction systems impose an order-of-magnitude larger *space footprint* than single-vector models, as they must store billions of small vectors for Web-scale collections. Considering this challenge, it might seem more fruitful to focus instead on addressing the fragility of single-vector models (Menon et al., 2022) by introducing new supervision paradigms for negative mining (Xiong et al., 2020), pretraining (Gao and Callan, 2021), and distillation (Qu et al., 2021). Indeed, recent single-vector models with highly-tuned supervision strategies (Ren et al., 2021b; Formal et al., 2021a) sometimes perform on-par or

# LLMs as a Judge (2023)

LLMs evaluate other LLMs.

🧑‍⚖️ 85% human-level agreement
📊 Automates eval pipelines

---

# Judging LLM-as-a-Judge
# with MT-Bench and Chatbot Arena

**Lianmin Zheng**[1*]  **Wei-Lin Chiang**[1*]  **Ying Sheng**[4*]  **Siyuan Zhuang**[1]

**Zhanghao Wu**[1]  **Yonghao Zhuang**[3]  **Zi Lin**[2]  **Zhuohan Li**[1]  **Dacheng Li**[13]

**Eric P. Xing**[35]  **Hao Zhang**[12]  **Joseph E. Gonzalez**[1]  **Ion Stoica**[1]

[1] UC Berkeley  [2] UC San Diego  [3] Carnegie Mellon University  [4] Stanford  [5] MBZUAI

## Abstract

Evaluating large language model (LLM) based chat assistants is challenging due to their broad capabilities and the inadequacy of existing benchmarks in measuring human preferences. To address this, we explore using strong LLMs as judges to evaluate these models on more open-ended questions. We examine the usage and limitations of LLM-as-a-judge, including position, verbosity, and self-enhancement biases, as well as limited reasoning ability, and propose solutions to mitigate some of them. We then verify the agreement between LLM judges and human preferences by introducing two benchmarks: MT-bench, a multi-turn question set; and Chatbot Arena, a crowdsourced battle platform. Our results reveal that strong LLM judges like GPT-4 can match both controlled and crowdsourced human preferences well, achieving over 80% agreement, the same level of agreement between humans. Hence, LLM-as-a-judge is a scalable and explainable way to approximate human preferences, which are otherwise very expensive to obtain. Additionally, we show our benchmark and traditional benchmarks complement each other by evaluating several variants of LLaMA and Vicuna. The MT-bench questions, 3K expert votes, and 30K conversations with human preferences are publicly available at `https://github.com/lm-sys/FastChat/tree/main/fastchat/llm_judge`.

## 1  Introduction

# DeepSeek-R1 (2025)

RL trains structured reasoning.

🧠 Step-by-step thinking
🔭 Glimpse into LLM 2.0

deepseek

## DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning

DeepSeek-AI

research@deepseek.com

## Abstract

We introduce our first-generation reasoning models, DeepSeek-R1-Zero and DeepSeek-R1. DeepSeek-R1-Zero, a model trained via large-scale reinforcement learning (RL) without supervised fine-tuning (SFT) as a preliminary step, demonstrates remarkable reasoning capabilities. Through RL, DeepSeek-R1-Zero naturally emerges with numerous powerful and intriguing reasoning behaviors. However, it encounters challenges such as poor readability, and language mixing. To address these issues and further enhance reasoning performance, we introduce DeepSeek-R1, which incorporates multi-stage training and cold-start data before RL. DeepSeek-R1 achieves performance comparable to OpenAI-o1-1217 on reasoning tasks. To support the research community, we open-source DeepSeek-R1-Zero, DeepSeek-R1, and six dense models (1.5B, 7B, 8B, 14B, 32B, 70B) distilled from DeepSeek-R1 based on Qwen and Llama.