

# Toolformer: Language Models Can Teach Themselves to Use Tools

Timo Schick   Jane Dwivedi-Yu   Roberto Dessì<sup>†</sup>   Roberta Raileanu  
 Maria Lomeli   Luke Zettlemoyer   Nicola Cancedda   Thomas Scialom

Meta AI Research   <sup>†</sup>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.

## 1 Introduction

Large language models achieve impressive zero- and few-shot results on a variety of natural language processing tasks (Brown et al., 2020; Chowdhery et al., 2022, i.a.) and show several emergent capabilities (Wei et al., 2022). However, all of these models have several inherent limitations that can at best be partially addressed by further scaling. These limitations include an inability to access up-to-date information on recent events (Komeili et al., 2022) and the related tendency to hallucinate facts (Maynez et al., 2020; Ji et al., 2022), difficulties in understanding low-resource languages (Lin et al., 2021), a lack of mathematical skills to perform precise calculations (Patel et al., 2021) and an unawareness of the progression of time (Dhingra et al., 2022).

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.

A simple way to overcome these limitations of today’s language models is to give them the ability to *use external tools* such as search engines, calculators, or calendars. However, existing approaches either rely on large amounts of human annotations (Komeili et al., 2022; Thoppilan et al., 2022) or limit tool use to task-specific settings only (e.g., Gao et al., 2022; Parisi et al., 2022), hindering a more widespread adoption of tool use in LMs. Therefore, we propose *Toolformer*, a model that learns to use tools in a novel way, which fulfills the following desiderata:

- The use of tools should be learned in a self-supervised way without requiring large amounts of *human annotations*. This is impor-

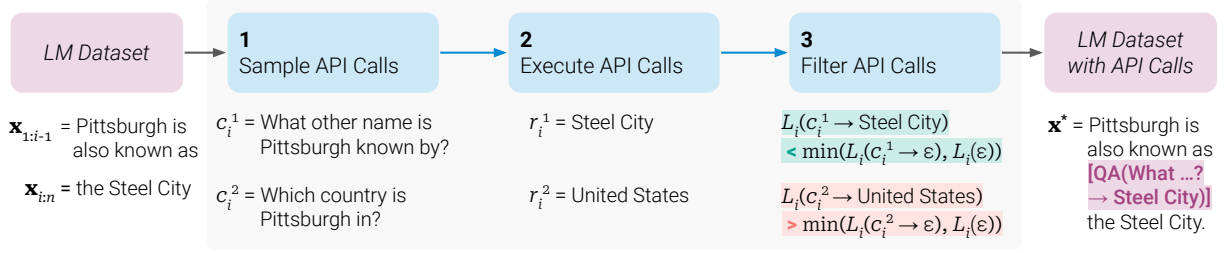


Figure 2: Key steps in our approach, illustrated for a *question answering* tool: Given an input text  $\mathbf{x}$ , we first sample a position  $i$  and corresponding API call candidates  $c_i^1, c_i^2, \dots, c_i^k$ . We then execute these API calls and filter out all calls which do not reduce the loss  $L_i$  over the next tokens. All remaining API calls are interleaved with the original text, resulting in a new text  $\mathbf{x}^*$ .

tant not only because of the costs associated with such annotations, but also because what humans find useful may be different from what a model finds useful.

- The LM should not lose any of its *generality* and should be able to decide for itself *when* and *how* to use which tool. In contrast to existing approaches, this enables a much more comprehensive use of tools that is not tied to specific tasks.

Our approach for achieving these goals is based on the recent idea of using large LMs with *in-context learning* (Brown et al., 2020) to generate entire datasets from scratch (Schick and Schütze, 2021b; Honovich et al., 2022; Wang et al., 2022): Given just a handful of human-written examples of how an API can be used, we let a LM annotate a huge language modeling dataset with potential API calls. We then use a self-supervised loss to determine which of these API calls actually help the model in predicting future tokens. Finally, we finetune the LM itself on the API calls that it considers useful. As illustrated in Figure 1, through this simple approach, LMs can learn to control a variety of tools, and to choose for themselves which tool to use when and how.

As our approach is agnostic of the dataset being used, we can apply it to the exact same dataset that was used to pretrain a model in the first place. This ensures that the model does not lose any of its generality and language modeling abilities. We conduct experiments on a variety of different downstream tasks, demonstrating that after learning to use tools, Toolformer, which is based on a pretrained GPT-J model (Wang and Komatsuzaki, 2021) with 6.7B parameters, achieves much stronger zero-shot results, clearly outperforming a much larger GPT-3 model (Brown et al., 2020) and

several other baselines on various tasks.

## 2 Approach

Our aim is to equip a language model  $M$  with the ability to use different tools by means of API calls. We require that inputs and outputs for each API can be represented as text sequences. This allows seamless insertion of API calls into any given text, using special tokens to mark the start and end of each such call.

We represent each API call as a tuple  $c = (a_c, i_c)$  where  $a_c$  is the name of the API and  $i_c$  is the corresponding input. Given an API call  $c$  with a corresponding result  $r$ , we denote the linearized sequences of the API call not including and including its result, respectively, as:

$$\begin{aligned} e(c) &= \langle \text{API} \rangle a_c (i_c) \langle / \text{API} \rangle \\ e(c, r) &= \langle \text{API} \rangle a_c (i_c) \rightarrow r \langle / \text{API} \rangle \end{aligned}$$

where “ $\langle \text{API} \rangle$ ”, “ $\langle / \text{API} \rangle$ ” and “ $\rightarrow$ ” are special tokens.<sup>1</sup> Some examples of linearized API calls inserted into text sequences are shown in Figure 1.

Given a dataset  $\mathcal{C} = \{\mathbf{x}^1, \dots, \mathbf{x}^{|\mathcal{C}|}\}$  of plain texts, we first convert this dataset into a dataset  $\mathcal{C}^*$  augmented with API calls. This is done in three steps, illustrated in Figure 2: First, we exploit the in-context learning ability of  $M$  to sample a large number of potential API calls. We then execute these API calls and finally check whether the obtained responses are helpful for predicting future tokens; this is used as a filtering criterion. After filtering, we merge API calls for different tools, resulting in the augmented dataset  $\mathcal{C}^*$ , and finetune

<sup>1</sup>In practice, we use the token sequences “ $[$ ”, “ $]$ ” and “ $\rightarrow$ ” to represent “ $\langle \text{API} \rangle$ ”, “ $\langle / \text{API} \rangle$ ” and “ $\rightarrow$ ”, respectively. This enables our approach to work without modifying the existing LM’s vocabulary. For reasons of readability, we still refer to them as “ $\langle \text{API} \rangle$ ”, “ $\langle / \text{API} \rangle$ ” and “ $\rightarrow$ ” throughout this section.

Your task is to add calls to a Question Answering API to a piece of text. The questions should help you get information required to complete the text. You can call the API by writing "[QA(question)]" where "question" is the question you want to ask. Here are some examples of API calls:

**Input:** Joe Biden was born in Scranton, Pennsylvania.

**Output:** Joe Biden was born in [QA("Where was Joe Biden born?")] Scranton, [QA("In which state is Scranton?")] Pennsylvania.

**Input:** Coca-Cola, or Coke, is a carbonated soft drink manufactured by the Coca-Cola Company.

**Output:** Coca-Cola, or [QA("What other name is Coca-Cola known by?")] Coke, is a carbonated soft drink manufactured by [QA("Who manufactures Coca-Cola?")] the Coca-Cola Company.

**Input:**  $\mathbf{x}$

**Output:**

Figure 3: An exemplary prompt  $P(\mathbf{x})$  used to generate API calls for the question answering tool.

$M$  itself on this dataset. Each of these steps is described in more detail below.

**Sampling API Calls** For each API, we write a prompt  $P(\mathbf{x})$  that encourages the LM to annotate an example  $\mathbf{x} = x_1, \dots, x_n$  with API calls. An example of such a prompt for a question answering tool is shown in Figure 3; all prompts used are shown in Appendix A.2. Let  $p_M(z_{n+1} \mid z_1, \dots, z_n)$  be the probability that  $M$  assigns to token  $z_{n+1}$  as a continuation for the sequence  $z_1, \dots, z_n$ . We first sample up to  $k$  candidate *positions* for doing API calls by computing, for each  $i \in \{1, \dots, n\}$ , the probability

$$p_i = p_M(\langle \text{API} \rangle \mid P(\mathbf{x}), x_{1:i-1})$$

that  $M$  assigns to starting an API call at position  $i$ . Given a sampling threshold  $\tau_s$ , we keep all positions  $I = \{i \mid p_i > \tau_s\}$ ; if there are more than  $k$  such positions, we only keep the top  $k$ .

For each position  $i \in I$ , we then obtain up to  $m$  API calls  $c_i^1, \dots, c_i^m$  by sampling from  $M$  given the sequence  $[P(\mathbf{x}), x_1, \dots, x_{i-1}, \langle \text{API} \rangle]$  as a prefix and  $\langle / \text{API} \rangle$  as an end-of-sequence token.<sup>2</sup>

<sup>2</sup>We discard all examples where  $M$  does not generate the  $\langle / \text{API} \rangle$  token.

**Executing API Calls** As a next step, we execute all API calls generated by  $M$  to obtain the corresponding results. How this is done depends entirely on the API itself – for example, it can involve calling another neural network, executing a Python script or using a retrieval system to perform search over a large corpus. The response for each API call  $c_i$  needs to be a single text sequence  $r_i$ .

**Filtering API Calls** Let  $i$  be the position of the API call  $c_i$  in the sequence  $\mathbf{x} = x_1, \dots, x_n$ , and let  $r_i$  be the response from the API. Further, given a sequence  $(w_i \mid i \in \mathbb{N})$  of *weights*, let

$$L_i(\mathbf{z}) = - \sum_{j=i}^n w_{j-i} \cdot \log p_M(x_j \mid \mathbf{z}, x_{1:j-1})$$

be the weighted cross entropy loss for  $M$  over the tokens  $x_i, \dots, x_n$  if the model is prefixed with  $\mathbf{z}$ . We compare two different instantiations of this loss:

$$\begin{aligned} L_i^+ &= L_i(\mathbf{e}(c_i, r_i)) \\ L_i^- &= \min(L_i(\varepsilon), L_i(\mathbf{e}(c_i, \varepsilon))) \end{aligned}$$

where  $\varepsilon$  denotes an empty sequence. The former is the weighted loss over all tokens  $x_i, \dots, x_n$  if the API call and its result are given to  $M$  as a prefix;<sup>3</sup> the latter is the minimum of the losses obtained from (i) doing no API call at all and (ii) doing an API call, but not providing the response. Intuitively, an API call is helpful to  $M$  if providing it with both the input *and* the output of this call makes it easier for the model to predict future tokens, compared to not receiving the API call at all, or receiving only its input. Given a filtering threshold  $\tau_f$ , we thus only keep API calls for which

$$L_i^- - L_i^+ \geq \tau_f$$

holds, i.e., adding the API call and its result *reduces* the loss by at least  $\tau_f$ , compared to not doing any API call or obtaining no result from it.

**Model Finetuning** After sampling and filtering calls for all APIs, we finally merge the remaining API calls and interleave them with the original inputs. That is, for an input text  $\mathbf{x} = x_1, \dots, x_n$  with a corresponding API call and result  $(c_i, r_i)$  at position  $i$ , we construct the new sequence  $\mathbf{x}^* =$

<sup>3</sup>We provide  $\mathbf{e}(c_i, r_i)$  as a prefix instead of inserting it at position  $i$  because  $M$  is not yet finetuned on any examples containing API calls, so inserting it in the middle of  $\mathbf{x}$  would interrupt the flow and not align with patterns in the pretraining corpus, thus hurting perplexity.

$x_{1:i-1}, e(c_i, r_i), x_{i:n}$ ; we proceed analogously for texts with multiple API calls. Doing this for all  $\mathbf{x} \in \mathcal{C}$  results in the new dataset  $\mathcal{C}^*$  augmented with API calls. We use this new dataset to finetune  $M$ , using a standard language modeling objective. Crucially, apart from inserted API calls the augmented dataset  $\mathcal{C}^*$  contains the exact same texts as  $\mathcal{C}$ , the original dataset. As a consequence, finetuning  $M$  on  $\mathcal{C}^*$  exposes it to the same content as finetuning on  $\mathcal{C}$ . Moreover, as API calls are inserted in exactly those positions and with exactly those inputs that help  $M$  predict future tokens, finetuning on  $\mathcal{C}^*$  enables the language model to decide when and how to use which tool, based purely on its own feedback.

**Inference** When generating text with  $M$  after finetuning with our approach, we perform regular decoding until  $M$  produces the “ $\rightarrow$ ” token, indicating that it next expects the response for an API call. At this point, we interrupt the decoding process, call the appropriate API to get a response, and continue the decoding process after inserting both the response and the `</API>` token.

### 3 Tools

We explore a variety of tools to address different shortcomings of regular LMs. The only constraints we impose on these tools is that (i) both their inputs and outputs can be represented as text sequences, and (ii) we can obtain a few demonstrations of their intended use. Concretely, we explore the following five tools: a question answering system, a Wikipedia search engine, a calculator, a calendar, and a machine translation system. Some examples of potential calls and return strings for the APIs associated with each of these tools are shown in Table 1. We briefly discuss all tools below; further details can be found in Appendix A.

**Question Answering** Our first tool is a question answering system based on another LM that can answer simple factoid questions. Specifically, we use *Atlas* (Izacard et al., 2022), a retrieval-augmented LM finetuned on Natural Questions (Kwiatkowski et al., 2019).

**Calculator** As a second tool, we use a calculator that can perform simple numeric calculations; we only support the four basic arithmetic operations. Results are always rounded to two decimal places.

**Wikipedia Search** Our third tool is a search engine that, given a search term, returns short text

snippets from Wikipedia. Compared to our question answering tool, this search enables a model to get more comprehensive information on a subject, but requires it to extract the relevant parts by itself. As our search engine, we use a BM25 retriever (Robertson et al., 1995; Baeza-Yates et al., 1999) that indexes the Wikipedia dump from KILT (Petroni et al., 2021).

**Machine Translation System** Our fourth tool is a machine translation system based on a LM that can translate a phrase from any language into English. More concretely, we use the 600M parameter NLLB (Costa-jussà et al., 2022) as our multilingual machine translation model that works for 200 languages (including low-resource ones). The source language is automatically detected using the *fast-Text* classifier (Joulin et al., 2016), while the target language is always set to English.

**Calendar** Our final tool is a calendar API that, when queried, returns the current date without taking any input. This provides temporal context for predictions that require some awareness of time.

## 4 Experiments

We investigate whether our approach enables a model to use tools without any further supervision and to decide for itself when and how to call which of the available tools. To test this, we select a variety of downstream tasks where we assume at least one of the considered tools to be useful, and evaluate performance in zero-shot settings (Section 4.2). Beyond that, we also ensure that our approach does not hurt the model’s core language modeling abilities; we verify this by looking at perplexity on two language modeling datasets (Section 4.3). Finally, we investigate how the ability to learn using tools is affected by model size (Section 4.4).

### 4.1 Experimental Setup

**Dataset Generation** Throughout all of our experiments, we use a subset of CCNet (Wenzek et al., 2020) as our language modeling dataset  $\mathcal{C}$  and GPT-J (Wang and Komatsuzaki, 2021) as our language model  $M$ . To reduce the computational cost of annotating  $\mathcal{C}$  with API calls, we define heuristics for some APIs to get a subset of  $\mathcal{C}$  for which API calls are more likely to be helpful than for an average text. For example, we only consider texts for the calculator tool if they contain at least three numbers. Details of the heuristics used are given in



API Name	Example Input	Example Output
Question Answering	Where was the Knights of Columbus founded?	New Haven, Connecticut
Wikipedia Search	Fishing Reel Types	Spin fishing > Spin fishing is distinguished between fly fishing and bait cast fishing by the type of rod and reel used. There are two types of reels used when spin fishing, the open faced reel and the closed faced reel.
Calculator	$27 + 4 * 2$	35
Calendar	$\varepsilon$	Today is Monday, January 30, 2023.
Machine Translation	sûreté nucléaire	nuclear safety

Table 1: Examples of inputs and outputs for all APIs used.

API	Number of Examples		
	$\tau_f = 0.5$	$\tau_f = 1.0$	$\tau_f = 2.0$
Question Answering	51,987	18,526	5,135
Wikipedia Search	207,241	60,974	13,944
Calculator	3,680	994	138
Calendar	61,811	20,587	3,007
Machine Translation	3,156	1,034	229

Table 2: Number of examples with API calls in  $\mathcal{C}^*$  for different values of our filtering threshold  $\tau_f$ .

Appendix A. For obtaining  $\mathcal{C}^*$  from  $\mathcal{C}$ , we perform all steps described in Section 2 and additionally filter out all examples for which all API calls were eliminated in the filtering step.<sup>4</sup> For the weighting function, we use

$$w_t = \frac{\tilde{w}_t}{\sum_{s \in \mathbb{N}} \tilde{w}_s} \text{ with } \tilde{w}_t = \max(0, 1 - 0.2 \cdot t)$$

to make sure that API calls happen close to where the information provided by the API is actually helpful for the model. The thresholds  $\tau_s$  and  $\tau_f$  are chosen individually for each tool to ensure a sufficiently larger number of examples; see Appendix A for details. Table 2 shows relevant statistics of our final dataset augmented with API calls.

**Model Finetuning** We finetune  $M$  on  $\mathcal{C}^*$  using a batch size of 128 and a learning rate of  $1 \cdot 10^{-5}$  with linear warmup for the first 10% of training. Details of our finetuning procedure are given in Appendix B.

**Baseline Models** Throughout the remainder of this section, we mainly compare the following models:

- **GPT-J**: A regular GPT-J model without any finetuning.
- **GPT-J + CC**: GPT-J finetuned on  $\mathcal{C}$ , our subset of CCNet *without* any API calls.
- **Toolformer**: GPT-J finetuned on  $\mathcal{C}^*$ , our subset of CCNet augmented with API calls.
- **Toolformer (disabled)**: The same model as Toolformer, but API calls are disabled during decoding.<sup>5</sup>

For most tasks, we additionally compare to OPT (66B) (Zhang et al., 2022) and GPT-3<sup>6</sup> (175B) (Brown et al., 2020), two models that are about 10 and 25 times larger than our other baseline models, respectively.

## 4.2 Downstream Tasks

We evaluate all models on a variety of downstream tasks. In all cases, we consider a prompted zero-shot setup – i.e., models are instructed to solve each task in natural language, but we do not provide any in-context examples. This is in contrast to prior work on tool use (e.g., Gao et al., 2022; Parisi et al., 2022), where models are provided with dataset-specific examples of how a tool can be used to solve a concrete task. We choose the more challenging zero-shot setup as we are interested in seeing whether Toolformer works in precisely those cases where a user does not specify in advance which tools should be used in which way for solving a specific problem.

We use standard greedy decoding, but with one modification for Toolformer: We let the model start an API call not just when  $\langle \text{API} \rangle$  is the most likely

<sup>4</sup>While this filtering alters the distribution of training examples, we assume that the remaining examples are close enough to the original distribution so that  $M$ ’s language modeling abilities remain unaffected. This assumption is empirically validated in Section 4.3.

<sup>5</sup>This is achieved by manually setting the probability of the  $\langle \text{API} \rangle$  token to 0.

<sup>6</sup>We use the original *davinci* variant that is not finetuned on any instructions.

token, but whenever it is one of the  $k$  most likely tokens. For  $k = 1$ , this corresponds to regular greedy decoding; we instead use  $k = 10$  to increase the disposition of our model to make use of the APIs that it has access to. At the same time, we only at most one API call per input to make sure the model does not get stuck in a loop where it constantly calls APIs without producing any actual output. The effect of these modifications is explored in Section 5.

#### 4.2.1 LAMA

We evaluate our models on the SQuAD, Google-RE and T-REx subsets of the LAMA benchmark (Petroni et al., 2019). For each of these subsets, the task is to complete a short statement with a missing fact (e.g., a date or a place). As LAMA was originally designed to evaluate *masked* language models (e.g., Devlin et al., 2019), we filter out examples where the mask token is not the final token, so that the remaining examples can be processed in a left-to-right fashion. To account for different tokenizations and added complexity from not informing the model that a single word is required, we use a slightly more lenient evaluation criterion than exact match and simply check whether the correct word is within the first five words predicted by the model. As LAMA is based on statements obtained directly from Wikipedia, we prevent Toolformer from using the Wikipedia Search API to avoid giving it an unfair advantage.

Results for all models can be seen in Table 3. All GPT-J models without tool use achieve similar performance. Crucially, Toolformer clearly outperforms these baseline models, improving upon the best baseline by 11.7, 5.2 and 18.6 points, respectively. It also clearly outperforms OPT (66B) and GPT-3 (175B), despite both models being much larger. This is achieved because the model independently decides to ask the question answering tool for the required information in almost all cases (98.1%); for only very few examples, it uses a different tool (0.7%) or no tool at all (1.2%).

#### 4.2.2 Math Datasets

We test mathematical reasoning abilities on ASDiv (Miao et al., 2020), SVAMP (Patel et al., 2021) and the MAWPS benchmark (Koncel-Kedziorski et al., 2016). We again account for the fact that we test all models in a zero-shot setup by using a more lenient evaluation criterion: As the required output is always a number, we simply check for the first

Model	SQuAD	Google-RE	T-REx
GPT-J	17.8	4.9	31.9
GPT-J + CC	19.2	5.6	33.2
Toolformer (disabled)	22.1	6.3	34.9
Toolformer	<b>33.8</b>	<b>11.5</b>	<b>53.5</b>
OPT (66B)	21.6	2.9	30.1
GPT-3 (175B)	26.8	7.0	39.8

Table 3: Results on subsets of LAMA. Toolformer uses the question answering tool for most examples, clearly outperforming all baselines of the same size and achieving results competitive with GPT-3 (175B).

Model	ASDiv	SVAMP	MAWPS
GPT-J	7.5	5.2	9.9
GPT-J + CC	9.6	5.0	9.3
Toolformer (disabled)	14.8	6.3	15.0
Toolformer	<b>40.4</b>	<b>29.4</b>	<b>44.0</b>
OPT (66B)	6.0	4.9	7.9
GPT-3 (175B)	14.0	10.0	19.8

Table 4: Results for various benchmarks requiring mathematical reasoning. Toolformer makes use of the calculator tool for most examples, clearly outperforming even OPT (66B) and GPT-3 (175B).

number predicted by the model.<sup>7</sup>

Table 4 shows results for all benchmarks. While GPT-J and GPT-J + CC perform about the same, Toolformer achieves stronger results even when API calls are disabled. We surmise that this is because the model is finetuned on many examples of API calls and their results, improving its own mathematical capabilities. Nonetheless, allowing the model to make API calls more than doubles performance for all tasks, and also clearly outperforms the much larger OPT and GPT-3 models. This is because across all benchmarks, for 97.9% of all examples the model decides to ask the calculator tool for help.

#### 4.2.3 Question Answering

We look at Web Questions (Berant et al., 2013), Natural Questions (Kwiatkowski et al., 2019) and TriviaQA (Joshi et al., 2017), the three question answering datasets considered by Brown et al. (2020). For evaluation, we check whether the first 20 words predicted by a model contain the correct answer instead of requiring an exact match. For Toolformer, we disable the question answering tool as

<sup>7</sup>An exception to this is if the model’s prediction contains an equation (e.g., “The correct answer is  $5+3=8$ ”), in which case we consider the first number after the “=” sign to be its prediction.

Model	WebQS	NQ	TriviaQA
GPT-J	18.5	12.8	43.9
GPT-J + CC	18.4	12.2	45.6
Toolformer (disabled)	18.9	12.6	46.7
Toolformer	<b>26.3</b>	<b>17.7</b>	<b>48.8</b>
OPT (66B)	18.6	11.4	45.7
GPT-3 (175B)	<u>29.0</u>	<u>22.6</u>	<u>65.9</u>

Table 5: Results for various question answering dataset. Using the Wikipedia search tool for most examples, Toolformer clearly outperforms baselines of the same size, but falls short of GPT-3 (175B).

this would make solving the tasks trivial, especially given that the underlying QA system was finetuned on Natural Questions.

Results are shown in Table 5. Once again, Toolformer clearly outperforms all other models based on GPT-J, this time mostly relying on the Wikipedia search API (99.3%) to find relevant information. However, Toolformer still lags behind the much larger GPT-3 (175B) model. This is likely due to both the simplicity of our search engine (in many cases, it returns results that are clearly not a good match for a given query) and the inability of Toolformer to *interact* with it, e.g., by reformulating its query if results are not helpful or by browsing through multiple of the top results. We believe that adding this functionality is an exciting direction for future work.

#### 4.2.4 Multilingual Question Answering

We evaluate Toolformer and all baseline models on MLQA (Lewis et al., 2019), a multilingual question-answering benchmark. A context paragraph for each question is provided in English, while the question can be in Arabic, German, Spanish, Hindi, Vietnamese, or Simplified Chinese. In order to solve the task, the model needs to be able to understand both the paragraph and the question, so it may benefit from translating the question into English. Our evaluation metric is the percentage of times the model’s generation, capped at 10 words, contains the correct answer.

Results are shown in Table 6. Using API calls consistently improves Toolformer’s performance for all languages, suggesting that it has learned to make use of the machine translation tool. Depending on the language, this tool is used for 63.8% to 94.9% of all examples; the only exception to this is Hindi, for which the machine translation tool is used in only 7.3% of cases. However, Tool-

Model	Es	De	Hi	Vi	Zh	Ar
GPT-J	15.2	<b>16.5</b>	1.3	8.2	<b>18.2</b>	<b>8.2</b>
GPT-J + CC	15.7	14.9	0.5	8.3	13.7	4.6
Toolformer (disabled)	19.8	11.9	1.2	10.1	15.0	3.1
Toolformer	<b>20.6</b>	13.5	<b>1.4</b>	<b>10.6</b>	16.8	3.7
OPT (66B)	0.3	0.1	1.1	0.2	0.7	0.1
GPT-3 (175B)	3.4	1.1	0.1	1.7	17.7	0.1
GPT-J (All En)	24.3	27.0	23.9	23.3	23.1	23.6
GPT-3 (All En)	24.7	27.2	26.1	24.9	23.6	24.0

Table 6: Results on MLQA for Spanish (Es), German (De), Hindi (Hi), Vietnamese (Vi), Chinese (Zh) and Arabic (Ar). While using the machine translation tool to translate questions is helpful across all languages, further pretraining on CCNet deteriorates performance; consequently, Toolformer does not consistently outperform GPT-J. The final two rows correspond to models that are given contexts and questions in English.

former does not consistently outperform vanilla GPT-J. This is mainly because for some languages, finetuning on CCNet deteriorates performance; this might be due to a distribution shift compared to GPT-J’s original pretraining data.

OPT and GPT-3 perform surprisingly weak across all languages, mostly because they fail to provide an answer in English despite being instructed to do so. A potential reason for GPT-J not suffering from this problem is that it was trained on more multilingual data than both OPT and GPT-3, including the EuroParl corpus (Koehn, 2005; Gao et al., 2020). As an upper bound, we also evaluate GPT-J and GPT-3 on a variant of MLQA where both the context and the question are provided in English. In this setup, GPT-3 performs better than all other models, supporting our hypothesis that its subpar performance on MLQA is due to the multilingual aspect of the task.

#### 4.2.5 Temporal Datasets

To investigate the calendar API’s utility, we evaluate all models on TEMPLAMA (Dhingra et al., 2022) and a new dataset that we call DATESET. TEMPLAMA is a dataset built from Wikidata that contains cloze queries about facts that change with time (e.g., “Cristiano Ronaldo plays for \_\_\_”) as well as the correct answer for the years between 2010 and 2020. DATESET, described in Appendix D, is also generated through a series of templates, but populated using a combination of random dates/durations (e.g., “What day of the week was it 30 days ago?”). Critically, knowing the current date is required to answer these questions.

Model	TEMPLAMA	DATESET
GPT-J	13.7	3.9
GPT-J + CC	12.9	2.9
Toolformer (disabled)	12.7	5.9
Toolformer	<b>16.3</b>	<b>27.3</b>
OPT (66B)	14.5	1.3
GPT-3 (175B)	15.5	0.8

Table 7: Results for the temporal datasets. Toolformer outperforms all baselines, but does not make use of the calendar tool for TEMPLAMA.

For both tasks, we use the same evaluation as for the original LAMA dataset.

Results shown in Table 7 illustrate that Toolformer outperforms all baselines for both TEMPLAMA and DATESET. However, closer inspection shows that improvements on TEMPLAMA can not be attributed to the calendar tool, which is only used for 0.2% of all examples, but mostly to the Wikipedia search and question answering tools, which Toolformer calls the most. This makes sense given that named entities in TEMPLAMA are often so specific and rare that even knowing the exact date alone would be of little help. The best course of action for this dataset – first querying the calendar API to get the current date, and then querying the question answering system with this date – is not only prohibited by our restriction of using at most one API call per example, but also hard to learn for Toolformer given that all API calls in its training data are sampled independently.

For DATESET, on the other hand, the considerable improvement of Toolformer compared to other models can be fully accredited to the calendar tool, which it makes use of for 54.8% of all examples.

### 4.3 Language Modeling

In addition to verifying improved performance on various downstream tasks, we also want to ensure that language modeling performance of Toolformer does not degrade through our finetuning with API calls. To this end, we evaluate our models on two language modeling datasets: WikiText (Merity et al., 2017) and a subset of 10,000 randomly selected documents from CCNet (Wenzek et al., 2020) that were not used during training. Perplexities of various models are shown in Table 8. As one would expect, finetuning on CCNet leads to slightly improved performance on a different CCNet subset, but it slightly deteriorates performance on WikiText, presumably because the original pre-

Model	WikiText	CCNet
GPT-J	<b>9.9</b>	10.6
GPT-J + CC	10.3	<b>10.5</b>
Toolformer (disabled)	10.3	<b>10.5</b>

Table 8: Perplexities of different models on WikiText and our validation subset of CCNet. Adding API calls comes without a cost in terms of perplexity for language modeling without any API calls.

training data for GPT-J is more similar to WikiText than our randomly selected subset of CCNet. Most importantly, however, training on  $\mathcal{C}^*$  (our dataset annotated with API calls) does not lead to an increase in perplexity compared to training on  $\mathcal{C}$  when API calls are disabled at inference time.<sup>8</sup>

### 4.4 Scaling Laws

We investigate how the ability to ask external tools for help affects performance as we vary the size of our LM. To this end, we apply our approach not just to GPT-J, but also to four smaller models from the GPT-2 family (Radford et al., 2019), with 124M, 355M, 775M and 1.6B parameters, respectively. We do so using only a subset of three tools: the question answering system, the calculator, and the Wikipedia search engine. Apart from this, we follow the experimental setup described in Section 4.1.

Figure 4 shows that the ability to leverage the provided tools only emerges at around 775M parameters: smaller models achieve similar performance both with and without tools. An exception to this is the Wikipedia search engine used mostly for QA benchmarks; we hypothesize that this is because the API is comparably easy to use. While models become better at solving tasks *without* API calls as they grow in size, their ability to make good use of the provided API improves at the same time. As a consequence, there remains a large gap between predictions with and without API calls even for our biggest model.

## 5 Analysis

**Decoding Strategy** We investigate the effect of our modified decoding strategy introduced in Section 4.2, where instead of always generating the

<sup>8</sup>We do not evaluate the perplexity of Toolformer with API calls enabled as computing the probability  $p_M(x_t \mid x_1, \dots, x_{t-1})$  of token  $x_t$  given  $x_1, \dots, x_{t-1}$  would require marginalizing over all potential API calls that the model could make at position  $t$ , which is intractable.



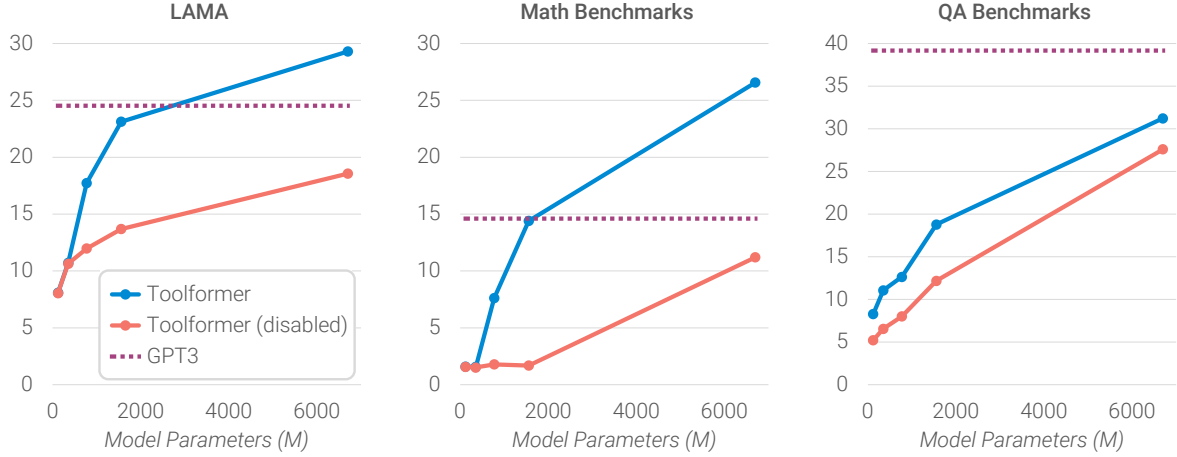


Figure 4: Average performance on LAMA, our math benchmarks and our QA benchmarks for GPT-2 models of different sizes and GPT-J finetuned with our approach, both with and without API calls. While API calls are not helpful to the smallest models, larger models learn how to make good use of them. Even for bigger models, the gap between model predictions with and without API calls remains high.

most likely token, we generate the `<API>` token if it is one of the  $k$  most likely tokens. Table 9 shows performance on the T-REx subset of LAMA and on WebQS for different values of  $k$ . As expected, increasing  $k$  leads to the model doing API calls for more examples – from 40.3% and 8.5% with  $k = 1$  (i.e., regular greedy decoding) to 98.1% and 100% for  $k = 10$ . While for T-REx, there is already a clear improvement in performance with greedy decoding, on WebQS our model only starts to make a substantial number of API calls as we slightly increase  $k$ . Interestingly, for  $k = 1$  the model is calibrated to some extent: It decides to call APIs for examples that it would perform particularly badly on without making API calls. This can be seen from the fact that performance on examples where it decides *not* to make an API call (44.3 and 19.9) is higher than average performance if no API calls are made at all (34.9 and 18.9). However, this calibration is lost for higher values of  $k$ .

**Data Quality** We qualitatively analyze some API calls generated with our approach for different APIs. Table 10 shows some examples of texts from CCNet augmented with API calls, as well as the corresponding score  $L_i^- - L_i^+$  that is used as a filtering criterion, and whether the API calls made by the model are intuitively useful in the given context. As can be seen, high values of  $L_i^- - L_i^+$  typically correspond to useful API calls, whereas low values correspond to API calls that do not provide any information that is useful for predicting future tokens. There are some exceptions, e.g., an API call for

$k$	T-REx				WebQS			
	All	AC	NC	%	All	AC	NC	%
0	34.9	–	34.9	0.0	18.9	–	18.9	0.0
1	47.8	53.0	44.3	40.3	19.3	17.1	19.9	8.5
3	52.9	58.0	29.0	82.8	<b>26.3</b>	26.5	6.6	99.3
10	<b>53.5</b>	54.0	22.5	98.1	<b>26.3</b>	26.4	–	100.0

Table 9: Toolformer results on the T-REx subset of LAMA and on WebQS for different values of  $k$  used during decoding. Numbers shown are overall performance (All), performance on the subset where the model decides to make an API call (AC) and all remaining examples (NC), as well as the percentage of examples for which the model decides to call an API (%).

“Fast train success” in the fourth example that does not give any relevant information but still reduces perplexity. However, some amount of noise in the API calls that are not filtered can actually be useful as it forces the model finetuned on  $C^*$  to not always blindly follow the results of each call it makes.

## 6 Related Work

**Language Model Pretraining** There are various approaches that augment language models with some form of additional textual information during pretraining, including various forms of metadata (Keskar et al., 2019), HTML tags (Aghajanyan et al., 2021), Wikipedia markup (Schick et al., 2022), or related texts obtained from an information retrieval system (Gua et al., 2020; Borgeaud et al., 2021; Izacard et al., 2022). For all of these

Example	$L_i^- - L_i^+$	Useful
The Flodden Window (a war memorial dedicated to The Middleton Archers), in the Grade I-listed Church of St Leonard in Middleton is said to be the oldest war memorial in the United Kingdom. <b>&lt;API&gt; WikiSearch(War memorial Flodden) → Battle of Flodden &gt; Commemoration &gt; The stained-glass Flodden Window in Middleton Parish Church [...] was constructed by Sir Richard Assheton in memory of the Battle of Flodden and the archers from Middleton who fought in it. &lt;/API&gt;</b> Sir Richard Assheton of Middleton (who built St Leonard) was granted knighthood [...]	5.49	✓
Note: The WL will be open on Friday, <b>&lt;API&gt; Calendar() → Today is Thursday, March 9, 2017. &lt;/API&gt;</b> March 10, and Sunday, March 19 for regular hours.	2.11	✓
The Nile has an approximate length of <b>&lt;API&gt; QA(What is the approximate length of the Nile?) → 6,853 km &lt;/API&gt;</b> 6,853 kilometers, the White Nile being its main source.	2.08	✓
If Venus had an atmosphere similar to Earth’s then you would expect Venus’ mean temperature to be 499 K (1.74 x 287) rather than 735 K which is <b>&lt;API&gt; Calculator(735 / 499) → 1.47 &lt;/API&gt;</b> 1.47 (735 / 499) times hotter than it should be.	1.59	✓
You are here: Home / Featured / Catch this fast train to success! <b>&lt;API&gt; WikiSearch(Fast train success) → Fast Train &gt; It also peaked at #23 on the Canadian CHUM singles chart, on June 26, 1971. The success of this single established Myles Goodwyn as the band’s main songwriter, and made it possible for April Wine to record a second album. &lt;/API&gt;</b> Don’t wait weeks and [...]	0.92	✗
Os Melhores Escolas em Jersey 2020 <b>&lt;API&gt; MT(Os Melhores Escolas em Jersey) → The Best Schools in Jersey &lt;/API&gt;</b> On this page you can search for Universities, Colleges and Business schools in Jersey	0.70	✓
Enjoy these pictures from the <b>&lt;API&gt; Calendar() → Today is Friday, April 19, 2013. &lt;/API&gt;</b> Easter Egg Hunt.	0.33	✓
85 patients (23%) were hospitalised alive and admitted to a hospital ward. Of them, <b>&lt;API&gt; Calculator(85 / 23) → 3.70 &lt;/API&gt;</b> 65% had a cardiac aetiology [...]	-0.02	✗
But hey, after the <b>&lt;API&gt; Calendar() → Today is Saturday, June 25, 2011. &lt;/API&gt;</b> Disneyland fiasco with the fire drill, I think it’s safe to say Chewey won’t let anyone die in a fire.	-0.41	✗
The last time I was with <b>&lt;API&gt; QA(Who was last time I was with?) → The Last Time &lt;/API&gt;</b> him I asked what he likes about me and he said he would tell me one day.	-1.23	✗

Table 10: Examples of API calls for different tools, sorted by the value of  $L_i^- - L_i^+$  that is used as a filtering criterion. High values typically correspond to API calls that are intuitively useful for predicting future tokens.

approaches, additional information is *always* provided, regardless of whether it is helpful or not. In contrast, Toolformer learns for itself to explicitly asks for the right information.

**Tool Use** Several approaches aim to equip LMs with the ability to use external tools such as search engines (Komeili et al., 2022; Thoppilan et al., 2022; Lazaridou et al., 2022; Shuster et al., 2022; Yao et al., 2022), web browsers (Nakano et al., 2021), calculators (Cobbe et al., 2021; Thoppilan et al., 2022), translation systems (Thoppilan et al., 2022) and Python interpreters (Gao et al., 2022). The way these models learn to use tools can roughly be divided into two approaches: Either they rely on large amounts of human supervision (Komeili et al., 2022; Nakano et al., 2021; Thoppilan et al., 2022) or they work by prompting the language model in a few-shot setup tailored towards a specific task where it is known a priori which tools needs to be

used (Gao et al., 2022; Lazaridou et al., 2022; Yao et al., 2022). In contrast, the self-supervised nature of Toolformer enables it to learn how and when to use tools without requiring a specific prompt that shows task-specific examples of how a tool could be used. Perhaps most closely related to our work is TALM (Parisi et al., 2022), an approach that uses a similar self-supervised objective for teaching a model to use a calculator and a search engine, but explores this only in settings where a model is finetuned for downstream tasks.

**Bootstrapping** The idea of using self-training and bootstrapping techniques to improve models has been investigated in various contexts, ranging from word sense disambiguation (Yarowsky, 1995), relation extraction (Brin, 1999; Agichtein and Gravano, 2000), parsing (McClosky et al., 2006; Reichart and Rappoport, 2007), sequence generation (He et al., 2020), few-shot text classi-