

# **Tool Calling**

**NLP: Fall 2025**

**Anoop Sarkar**

# **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

Feb 2023

<https://arxiv.org/abs/2302.04761>

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.

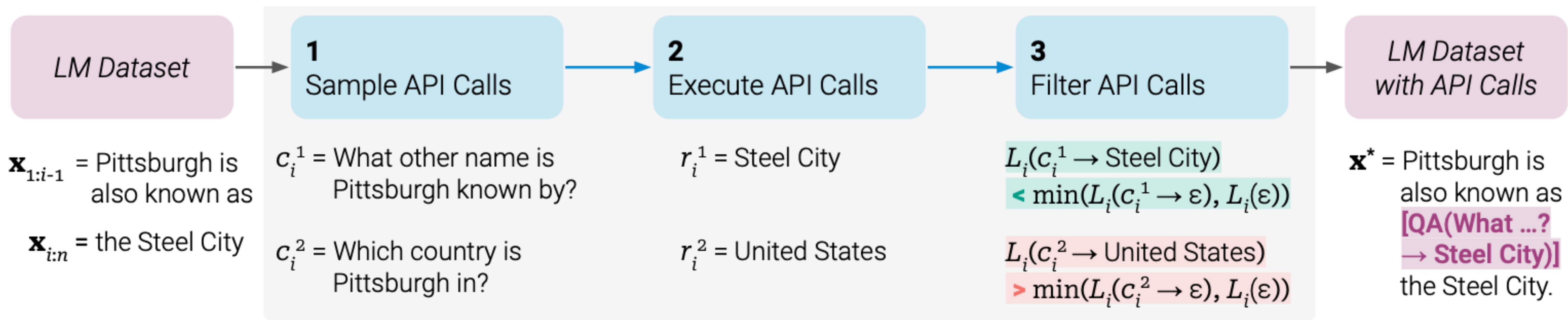


Figure 2: Key steps in our approach, illustrated for a *question answering* tool: Given an input text  $\mathbf{x}$ , we first sample a position  $i$  and corresponding API call candidates  $c_i^1, c_i^2, \dots, c_i^k$ . We then execute these API calls and filter out all calls which do not reduce the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

*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.

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	$\epsilon$	Today is Monday, January 30, 2023.
Machine Translation	sûreté nucléaire	nuclear safety

# Tools and agents

- Pipeline toolkits let you chain together multiple LLM calls
- Each output from the LLM can be used for subsequent LLM calls to break down a complex task into subgoals
- Combine multi-stage LLM calls with the following:
  - Tool calling
  - RAG
  - Automate retries
  - Fix up output syntax (e.g. JSON output)