

# Tool augmentation with LLMs

---



## Tool Augmentation with LLMs

Tool augmentation means **connecting an LLM to external tools** (APIs, functions, retrieval systems, calculators, programs) so it can perform tasks that exceed its built-in abilities.

Modern LLMs are not just text-generators — they become **reasoning agents** that know:

1. *When* to call a tool
2. *Which* tool to call
3. *How* to interpret tool results
4. *How* to combine tools + internal reasoning to solve problems

This greatly expands their capabilities.

---



## Why Do LLMs Need Tools?

LLMs have limitations:

- They cannot access the **internet** or **real-time data**.
- They cannot reliably do **math**, **code execution**, or **long-term planning**.
- They hallucinate facts instead of checking information.
- They cannot interact with the external world on their own.

**Tools solve these gaps.**

By adding tools, LLMs can:

- Fetch up-to-date information
- Execute code
- Perform exact calculations

- Access databases, vector stores, search engines
  - Control software, robots, or pipelines
- 



## How Tool Augmentation Works (High-Level)

1. **User query** arrives → LLM tries to understand what is being asked.
2. LLM decides:
  - Should I answer myself?
  - Or should I call an external tool?
3. If needed, LLM generates a **tool call** (API-like instruction).
4. External tool executes the action (search, compute, query, run code).
5. Tool returns results to the LLM.
6. LLM interprets the results and produces the final answer.

This process is often called **Tool Use, Agents, Function Calling, or LLM Tooling**.

---



## Types of Tools Commonly Used by LLMs

### ◆ 1. Search Tools / Retrieval Tools

- Web search
- Vector databases (like FAISS, Pinecone)
- Knowledge bases (Wikidata)

→ For factual accuracy and grounding.

---

### ◆ 2. Calculator / Math Engines

- Exact arithmetic
- Algebra, calculus
- Probabilistic models

→ Avoid hallucinated math.

---

### ◆ 3. Code Execution Environments

- Python sandbox
- SQL execution engines
- Shell commands

→ For data analysis, simulation, plotting.

---

### ◆ 4. Specialized APIs

- Weather APIs
- Finance APIs
- Flight/hotel search
- Mapping/geolocation
- Machine-control APIs

→ For real-world tasks.

---

### ◆ 5. Reasoning Tools (Planning, Agents, Chains)

- Multi-step reasoning frameworks (ReAct, AutoGPT, BabyAGI)
- Planning engines
- Tool calling loops

→ For complex workflows.

---

## Architectural Pattern: ReAct (Reason + Act)

A popular strategy:

1. **Reason (R):** LLM thinks step-by-step
2. **Act (A):** LLM uses a tool
3. **Observe:** Tool returns result
4. **Reason again:** Update thoughts
5. **Act again (if needed)**

This loop continues until final answer.

This enables:

- Multi-step planning
  - Complex tool use sequences
  - Self-correction
- 

## Tool Augmentation Turns LLMs into Agents

With tools, LLMs shift from “text predictors” to **autonomous agents** capable of:

- Web browsing
- Database querying
- File manipulation
- Code generation + execution
- Multi-step reasoning

Example tasks:

- “Book me a flight for tomorrow”
- “Plot the distribution of column X in the CSV”
- “Summarize the latest news”
- “Find similar research papers using dense retrieval”

Without tool augmentation, LLMs cannot perform these.

---



## Why Tool Augmentation Improves LLM Behavior

### ✓ Accuracy

Factual grounding through search.

### ✓ Reliability

Exact calculations.

### ✓ Capability

Access to specialized knowledge.

### ✓ Extendability

LLMs no longer need retraining to gain new skills — just plug in a new tool.

---

## Example: Tool-Augmented LLM Workflow

**User:** "What's the weather in Tokyo right now?"

**LLM internal reasoning:**

- I cannot know real-time weather
- I should call the weather API

**LLM → Tool:**

```
call_weather_api("Tokyo")
```

**Tool Response:**

```
{"temp": 18, "condition": "cloudy"}
```

**LLM → User:**

"It's 18°C and cloudy in Tokyo right now."

---



### In One Sentence

Tool augmentation allows LLMs to overcome their internal limitations by interacting with external tools, enabling accurate, grounded, and powerful real-world capabilities.

---

## How LLMs Are Trained for Tool Calling

Modern LLMs (GPT-4, GPT-3.5 Tools, Claude Tool-Use, Gemini with Actions) can **call external tools** such as:

- Python interpreters
- Search APIs

- Calculators
- Database query engines
- Browsers
- Retrieval systems

But an LLM doesn't magically know how to do this — it must be **trained** to understand:

1. **When** to use a tool (decision)
2. **Which** tool to call
3. **What arguments** to pass
4. **How to combine tool outputs with natural language**

This requires a combination of **supervised fine-tuning, synthetic data generation**, and sometimes **RL-based refinement**.

Let's break down the full training pipeline ↗

---

## 1 Step 1— Give the LLM “Tool APIs” as Part of Its Context

Every tool is described with:

- **Name**
- **Description of what it does**
- **Input arguments schema**
- **Output format**

Example (JSON schema):

```
{
  "name": "calculator",
  "description": "Evaluates a mathematical expression",
  "parameters": {
    "expression": "string"
}
```

```
    }  
}
```

This is included in the prompt so the LLM knows what tools exist and what inputs they expect.

## 2 Step 2 — Supervised Fine-Tuning (SFT) on Tool-Call Demonstrations

The LLM is trained using examples of **correct tool usage**.

This is standard *supervised learning*.

### Example Training Pair

**Human Input:**

| "What's the square root of 1728?"

**Correct Model Output:**

```
{  
  "tool": "calculator",  
  "arguments": {  
    "expression": "sqrt(1728)"  
  }  
}
```

Then after the tool returns a result (e.g., [41.569...](#)), the training includes:

**Model:**

| "The square root of 1728 is approximately 41.57."

The model learns:

- Detect math questions
- Choose the “calculator” tool
- Format JSON correctly

- Interpret tool results
- Produce final answers

Thousands–millions of these examples are used.

### Sources of training data:

- Human-written demonstrations
- Programmatically generated tasks
- Synthetic examples using older LLMs (self-play)
- Automatically generated API tasks (fuzzing)

---

## 3 Step 3 — Synthetic Tool Use via Self-Play (Self-Generated Examples)

Modern models like GPT-4/Claude use **self-instruct + tool-instruct**:

1. LLM generates tasks requiring tool use
2. LLM generates correct tool calls
3. A separate verifier (or execution engine) checks correctness
4. Only valid samples are added to training data

This produces:

- Many tools × many tasks combinations
- Difficult reasoning scenarios
- Rare or edge-case tool use examples

This dramatically improves tool competence.

---

## 4 Step 4 — Error Correction Training

LLMs need to handle:

- Incorrect tool calls
- Invalid JSON

- Wrong arguments
- Missing parameters
- Trying to answer without using tools

So additional training data is created:

- LLM makes incorrect calls
- A supervisor model or human edits them
- The correction is used as training data

This teaches robustness.

---

## 5 Step 5 — Reinforcement Learning (RL) for Better Tool Decisions

RLHF or RLAIF is used to teach:

- When it is better to call a tool
- When the model should **not** call a tool
- Multi-step tool planning
- Minimizing unnecessary tool calls
- Choosing between competing tools

Reward model examples:

- reward if the calculator is used for math
- reward if the search API is used for fact lookup
- reward if the answer is hallucinated without a tool
- reward for invalid tool arguments

Policy gradient methods (PPO, DPO, RRHF) refine the model.

---

## 6 Step 6 — Multi-Turn Tool Use Training

LLMs also learn sequences like:

1. Call search
2. Read results
3. Call the browser with a selected link
4. Extract data
5. Summarize

This is taught with multi-turn example conversations.

---

## 7 Step 7 — Inference-Time Tool Execution Loop

At runtime, calling a tool works like this:

1. LLM outputs JSON for tool call
2. Tool is executed by the external system
3. System returns the result as an assistant message
4. LLM continues processing using the result

This creates the illusion that the LLM “thinks” with tools — but it’s actually:

- Pattern learned from data
  - External system executing code
  - LLM integrating results into reasoning
- 



## Why This Works

Because LLMs are powerful sequence-matching systems:

- They learn patterns of when a tool is appropriate
  - They learn proper API formatting
  - They mimic demonstration patterns
  - RL helps them choose the *optimal* tool
  - Training forces them to plan multi-step actions
-

## Summary (Exam-Ready)

Large Language Models learn tool calling through a combination of supervised fine-tuning on curated tool-use examples, self-generated synthetic examples, reinforcement learning to optimize tool-selection behavior, and multi-turn training that teaches them how to use tool outputs inside reasoning. Tools are described via in-context API schemas, and the model learns to generate valid, structured tool calls and integrate tool outputs into final answers.

## Understanding the Slide: Training an LLM for Tool/API Calling

This slide shows **how a language model is trained to learn when and how to call tools (APIs)** by comparing three settings:

1. No API
2. Input-Only API
3. Full API (executed correctly)

The training objective uses a **log-likelihood loss**, and “helpfulness” is computed by comparing the model’s losses across these settings.

## The Core Idea

The model is shown three versions of the same example:

- One where **no API call** is made.
- One where an **API call is written but not executed**.
- One where an **API call is written AND executed**, and the true result is inserted.

The model’s loss on each version tells us **whether the tool call was helpful**.

This is used to **train the LLM to prefer correct API usage**.

## Helpfulness Score — What It Measures

The goal is to measure:

Does adding API/tool calls actually make the model more helpful?

We compare **three losses**:

### 1. No-API loss

The model completes the task *without* any API call traces.

### 2. Input-Only API loss

The model *sees* the API call in the prompt,

**but the API result is hidden** (represented as "?").

### 3. Full-API loss

The model *sees both*:

- the API call
- the API result (e.g., "Pennsylvania")

If the model performs *better* with the API results,  
its loss should be **lower**.

---

## 📌 Three Modes Explained

### 1. ★ = No API

Text only.

Pittsburgh is also known as

The model is trained normally; no API is involved.

Loss:

$$\mathcal{L}(\text{No API}) = 2.5$$

Higher loss means: **model struggled more to predict the next tokens.**

---

## 2. ★ = Input-Only API

The model is shown the API call as *text*, but the API is **not executed**.

[QA("In which state is Pittsburgh?") → ?]

Pittsburgh is also known as

This teaches the model **the structure/syntax of tool calls**,

but not yet the usefulness of the output.

Loss:

$$\mathcal{L}(\text{Input Only}) = 2.4$$

Loss goes slightly down → the model benefits a little from the tool call format.

## 3. ★ = Full API

The model sees a tool call **and the correct tool output**.

[QA("In which state is Pittsburgh?") → Pennsylvania]

Pittsburgh is also known as

Now the model gets **real information** from the tool.

Loss:

$$\mathcal{L}(\text{Full API}) = 2.6$$

Loss is *higher*, meaning the API output *changed the distribution* of what the model expected.

This indicates the tool **significantly shifted the context**, which is good for learning.



## Computing “Helpfulness”

They calculate “helpfulness” of using the tool:

$$\text{Helpfulness} = \min(2.5, 2.4) - 2.6 = -0.2$$

Explanation:

- The best *baseline* the model could do without executing the API is:  
 $\min(\text{No API, Input Only}) = \min(2.5, 2.4) = 2.4$
- But the model with the **correct API output** has loss:  
 $\mathcal{L}(\text{Full API}) = 2.6$
- So the helpfulness = baseline – full API = **negative value**.

Meaning:

A negative helpfulness score means the tool provided **non-trivial new information** the model did not already know or expect.

This is exactly what we want:

→ **Tool use should meaningfully change the model's predictions.**

Over-training, the model learns:

- **When** a tool API call is needed
- **How** to format it
- **How** to use the tool's result to improve its answer

---

## The Workflow on the Left (Pink + Blue Boxes)

This is the pipeline for generating the training dataset:

1. **Sample API calls**
  - Ask the base model to produce potential tool calls.
2. **Execute API calls**
  - Actually run these calls (e.g., calculator, search, QA).
3. **Filter API calls**
  - Keep only sensible and successful ones.
  - Remove hallucinated or broken calls.
4. **Create LM dataset with API calls**

- Combine:
    - No-API versions
    - Input-Only versions
    - Full-API versions
  - Used for training the final tool-augmented model.
- 

## Summary in One Sentence

This slide illustrates how an LLM is trained to use tools by comparing its predictive loss with and without API calls; executing tools provides new information that changes the model's internal distribution, and this "helpfulness" signal is used to teach the model when tool usage improves its output.

---

## What does $-\log p(\text{query} \mid \text{PREFIX}^*)$ mean?

### 1. What is $p(\text{query} \mid \text{PREFIX}^*)$ ?

This is the **probability assigned by the language model** to producing the correct continuation ("query") after being given the input ("PREFIX★").

- **PREFIX★** = everything the model has already seen (prompt + tool call steps)
- **query** = the next token(s) we want the model to generate

So the model is asked:

"Given this PREFIX★, what is the probability that the correct next text is the 'query'?"

This is just the usual LM objective: predicting next tokens.

---

### 2. Why the negative log?

Language models are trained using the **negative log-likelihood (NLL)**:

$$L = -\log p(\text{query} \mid \text{PREFIX}^*)$$

The negative log does 3 important things:

✓ Converts **probability** into a **loss value**

High probability → low loss

Low probability → high loss

✓ Turns a product of probabilities into a sum

Since generating text involves many tokens, logs avoid numerical underflow.

✓ Makes the objective convex (well-behaved for optimization)

---

## Example

Let:

- prefix = "Pittsburgh is also known as"
- model predicts:
  - probability 0.01 that next token is "the"
  - probability 0.30 that next token is "Steel"
  - probability 0.40 that next token is "City" (after "The")

Then:

$$p(\text{The Steel City} \mid \text{prefix}) = 0.3 \times 0.4 = 0.12$$

Loss is:

$$L = -\log(0.12) \approx 2.12$$

---

### 3. Why this is used in API/tool-trained LLMs

When training LLMs that do tool calling, different versions of the prefix are possible:

★ No API

PREFIX★ contains no tool call

★ Input-only

PREFIX★ contains the *tool call string* but *not* the tool result

★ Full API

PREFIX★ contains tool call + *executed tool response*

Each gives the model a different amount of information, so the LM assigns different probabilities to the correct continuation.

The loss tells us how good the model is at predicting the right answer **with** and **without** tools.

---

## Perplexity (PPL) in Language Models

### 1 What is Perplexity?

Perplexity measures **how well a language model predicts a sequence of tokens**.

It reflects the model's **uncertainty** or "**surprise**" when generating text.

- **Low perplexity → the model is confident → better model**
- **High perplexity → the model is confused → worse model**

**Intuition:**

Perplexity is the average number of choices the model thinks it has at each step when predicting the next token.

A model with:

- $\text{PPL} = 10 \rightarrow$  "I am choosing among ~10 likely words on average."
  - $\text{PPL} = 100 \rightarrow$  "I am very unsure — many possible next words."
- 

### 2 Mathematical Definition

For a sequence of tokens ( $x_1, x_2, \dots, x_N$ ):

$$\text{Perplexity} = \exp \left( -\frac{1}{N} \sum_{i=1}^N \log p(x_i | x_{<i}) \right)$$

This is:

$$\text{PPL} = e^{\text{Cross-Entropy}}$$

Perplexity is the **exponentiated average negative log-likelihood**.

---

### 3 Why “Perplexity”?

Because it expresses:

- **How surprised** the model is by the test data
- In the scale of **effective vocabulary size**

If perplexity = 20,

It means the model behaves as if it is choosing among **20 equally likely tokens on average**.

---

### 4 Example

Suppose a model is predicting the next word in a sentence:

“The capital of France is \_\_”

A good model assigns:

- high probability of **Paris**
- low probability to **Tokyo, banana, carrot**, etc.

If the model assigns:

- ( $p(\text{Paris}) = 0.7$ )
  - the negative log probability is small
  - perplexity is low
  - good performance.

If the model assigns:

- ( $p(\text{Paris}) = 0.05$ )
  - the model is “perplexed”
  - negative log probability is large

- high perplexity
  - poor performance.
- 

## 5 When is Perplexity Used?

Perplexity is widely used to evaluate:

- **Language models (GPT, LLaMA, BERT masked LM loss)**
- **Model checkpoints during pretraining**
- **Comparing quality across models or datasets**

It should be computed on:

- Held-out validation data
  - Not on training data (to avoid overfitting)
- 

## 6 Limitations of Perplexity

Although useful for training:

- Perplexity does **not** measure **truthfulness, reasoning, instruction following, or helpfulness**.
- LLMs today (GPT-4, Claude, etc.) are judged using **human preference, RLHF, and task-specific evaluations** rather than perplexity alone.
- Smaller models can get lower PPL than larger ones yet still perform worse on downstream tasks.

**Important:**

| Perplexity measures predictive likelihood, not alignment or usefulness.

---

## 7 Summary Table

| Concept    | Meaning                                   |
|------------|---|
| Perplexity | A metric measuring how confused the LM is |

| Concept       | Meaning   |
|---------------|---|
| Formula       | $\text{Perplexity} = \exp\left(-\frac{1}{N} \sum_{i=1}^N \log p(x_i   x_{<i})\right)$ |
| Low PPL       | Better prediction, more confident model   |
| High PPL      | More uncertainty, poorer model  |
| Related to    | Cross-entropy, negative log-likelihood  |
| Not measuring | Reasoning, alignment, truthfulness  |

## 🔧 Limitations of Toolformer

Toolformer was influential, but it has several important weaknesses that limit its performance and scalability.

### 1 Self-supervised tool-use depends on the base LM's intelligence

Toolformer generates its own training data by sampling possible tool calls and filtering them.

This works **only if the base model already has a partial understanding** of:

- When a tool should be used
- How tools behave
- What arguments do tools expect

If the base LLM is weak, the sampled tool-calls become:

- irrelevant
- syntactically incorrect
- semantically incoherent

**Result:** Low-quality training data → poor tool-use ability.

### 2 Single-step tool usage (no multi-step reasoning)

Toolformer assumes that:

- The LLM inserts **single API calls** in the middle of its output
- Every tool call is **one-shot**
- The LM never reasons across *multiple tool calls*

This means:

- No multi-hop retrieval
- No use of tools in a chain
- No iterative planning with tools

Thus, it fails on tasks requiring:

- multi-step computation
  - sequential API interactions
  - aggregation or reasoning over tool outputs
- 

### 3 Tools are shallow and pre-defined

Toolformer only supports a small set of **simple APIs**, such as:

- calculator
- web search snippet fetcher
- calendar lookup

It cannot learn:

- new APIs
- dynamic tools
- tools requiring schemas
- tools that return structured or multi-modal output

The system is **not extensible** without re-training.

---

### 4 No grounding or accuracy verification

Toolformer never checks:

- whether the API result is correct
- whether the use of the tool improves factuality

It simply uses a **perplexity-drop heuristic** to decide whether the API call is "useful".

This can include **wrong or misleading API calls** if they reduce perplexity.

---

## 5 Poor handling of complex inputs and instructions

Toolformer only learns API call patterns in contexts that look **similar to its training examples**.

It struggles with:

- ambiguous instructions
- multi-constraint tool usage
- composing different tools

The model doesn't *generalize* well outside the patterns it saw during self-supervised data generation.

---

## 6 Memoryless: No tool state or history

Tool calls do not maintain:

- state
- context
- persistent variables

Meaning:

- Cannot store retrieved values for later use
- Cannot update tool outputs across turns
- Cannot simulate program-like sequences

This makes it unsuitable for:

- planning

- database-style queries
  - agent-like behavior
- 

## 7 Limited scalability of the sampling-and-filtering process

Toolformer creates training data through three expensive steps:

1. Sample many potential tool calls
2. Execute each tool call
3. Filter them based on likelihood changes

For large toolsets or huge corpora, this becomes **computationally expensive**.

---

## 8 No explicit alignment or human feedback

Toolformer trains purely via LM likelihood.

There is:

- no reward model
- no preference optimization
- no safety filtering

So it cannot guarantee:

- safe tool use
  - ethical constraints
  - avoiding harmful or unauthorized API invocation
- 

## ★ Summary in One Line

Toolformer is limited because it relies on a moderately strong base model, supports only shallow one-shot tool calls, lacks multi-step reasoning and grounding, uses rigid APIs, and scales poorly due to heuristic self-supervision.

---

## Limitations of Toolformer — 5-Line Answer

Toolformer depends on a **base LLM's own ability** to propose useful tool calls, so weak models cannot self-generate good training data.

It struggles with **multi-step reasoning** because tool calls are inserted independently without global planning.

The approach is sensitive to **noisy demonstrations**, producing incorrect or irrelevant tool calls if the synthetic data is poor.

Toolformer cannot **adapt dynamically** to new or unseen tools without regenerating training data.

Its tool-use remains **shallow and brittle**, lacking robust decision-making compared to modern agent-style LLMs.

---