

# Build an AI Agent to Automate Your Research



AMAN KHARWAL · NOVEMBER 11, 2025

AmanXai



By Aman Kharwal



## Build an **AI Agent** to Automate Your Research.

If you're a student or working in **AI**, you probably feel overwhelmed by the amount of information out there. The answers exist, but sorting through and making sense of them takes a lot of time. Imagine if you could create an AI agent to handle that first round of research for you. In this guide, we'll build an AI agent to automate research that takes your question, searches the web, reads the top results, and gives you the most relevant passages along with a short summary.

## How will our AI Agent for Research Work?

The agent we're building is a fantastic example of a simple **Retrieval-Augmented** system. The core idea isn't just to find pages with the right keywords (like a simple Ctrl+F), but to find passages with the right meaning.

We will use **vector embeddings**. Think of an embedding as a coordinate of meaning in a vast, high-dimensional space. The sentence-transformers library provides models that are experts at turning any piece of text into a list of numbers (a vector) that represents its location in that meaning space.

Our agent's entire job is to:

1. Turn your query into a vector.
2. Search the web, scrape the text, and turn all the content into more vectors.
3. Find the text vectors whose coordinates are closest to your query's coordinates.

That's it. It's just very clever math, and you can build it in about 100 lines of Python.

## Now, Let's Build an AI Agent to Automate Your Research

First, create a new Python file and install the dependencies:

```
pip install ddgs requests beautifulsoup4 sentence-transformers numpy
```

## Step 1: Import Libraries and Configure Parameters

The top of your file will define all imports and key settings:

```
1 import re
2 import urllib.parse
3 from ddgs import DDGS          # package name is 'ddgs' (d
4 import requests
5 from bs4 import BeautifulSoup
6 from sentence_transformers import SentenceTransformer
7 import numpy as np
8 import time
```

---

Next, configure constants like search results, summary size, and embedding model:

```
1 SEARCH_RESULTS = 6           # How many URLs to check
2 PASSAGES_PER_PAGE = 4         # How many passages to pull from
3 EMBEDDING_MODEL = "sentence-transformers/all-MiniLM-L6-v2"
4 TOP_PASSAGES = 5             # How many relevant passages to
5 SUMMARY_SENTENCES = 3         # How many sentences in the final
6 TIMEOUT = 8                  # How long to wait for a webpage
```

---

## Step 2: Search and Fetch Web Pages

We'll use **DuckDuckGo Search** for free, API-less search results:

```
1 def unwrap_ddg(url):
2     """If DuckDuckGo returns a redirect wrapper, extract t
```

```

3 try:
4     parsed = urllib.parse.urlparse(url)
5     if "duckduckgo.com" in parsed.netloc:
6         qs = urllib.parse.parse_qs(parsed.query)
7         uddg = qs.get("uddg")
8         if uddg:
9             return urllib.parse.unquote(uddg[0])
10    except Exception:
11        pass
12    return url
13
14 def search_web(query, max_results=SEARCH_RESULTS):
15     """Search the web and return a list of URLs."""
16     urls = []
17     with DDGS() as ddgs:
18         for r in ddgs.text(query, max_results=max_results)
19             url = r.get("href") or r.get("url")
20             if not url:
21                 continue
22             url = unwrap_ddg(url) # Clean up DDG redirect
23             urls.append(url)
24     return urls

```

Then, fetch and clean the page with requests and BeautifulSoup:

```

1 def fetch_text(url, timeout=TIMEOUT):
2     """Fetch and clean text content from a URL."""
3     headers = {"User-Agent": "Mozilla/5.0 (research-agent)"}
4     try:
5         r = requests.get(url, timeout=timeout, headers=hea
6         if r.status_code != 200:
7             return ""

```

```
8     ct = r.headers.get("content-type", "")
9     if "html" not in ct.lower(): # Skip non-HTML conte
10        return ""
11
12     soup = BeautifulSoup(r.text, "html.parser")
13
14     # Remove all annoying tags
15     for tag in soup(["script", "style", "noscript", "h
16         tag.extract()
17
18     # Get all paragraph text
19     paragraphs = [p.get_text(" ", strip=True) for p in
20     text = " ".join([p for p in paragraphs if p])
21
22     if text.strip():
23         # Clean up whitespace
24         return re.sub(r"\s+", " ", text).strip()
25
26     # --- Fallback logic if <p> tags fail ---
27     meta = soup.find("meta", attrs={"name": "descripti
28     if meta and meta.get("content"):
29         return meta["content"].strip()
30     if soup.title and soup.title.string:
31         return soup.title.string.strip()
32
33 except Exception:
34     return "" # Fail silently
35 return ""
```

---

## Step 3: Chunk, Embed, and Rank Passages

We'll break long articles into smaller passages and embed them using SentenceTransformer:

```
1 def chunk_passages(text, max_words=120):
2     """Split long text into smaller passages."""
3     words = text.split()
4     if not words:
5         return []
6     chunks = []
7     i = 0
8     while i < len(words):
9         chunk = words[i : i + max_words]
10        chunks.append(" ".join(chunk))
11        i += max_words
12    return chunks
13
14 def split_sentences(text):
15     """A simple sentence splitter."""
16     parts = re.split(r'(?<=[.!?])\s+', text)
17     return [p.strip() for p in parts if p.strip()]
18
19 class ShortResearchAgent:
20     def __init__(self, embed_model=EMBEDDING_MODEL):
21         print(f"Loading embedder: {embed_model}...")
22         # This downloads the model on first run
23         self.embedder = SentenceTransformer(embed_model)
24
25     def run(self, query):
26         start = time.time()
27
28         # 1. Search
29         urls = search_web(query)
```

```

30 print(f"Found {len(urls)} urls.")
31
32 # 2. Fetch & Chunk
33 docs = []
34 for u in urls:
35     txt = fetch_text(u)
36     if not txt:
37         continue
38     chunks = chunk_passages(txt, max_words=120)
39     for c in chunks[:PAGES_PER_PAGE]:
40         docs.append({"url": u, "passage": c})
41
42 if not docs:
43     print("No documents fetched.")
44     return {"query": query, "passages": [], "summa

```

Now, initialize the embedding model and compute cosine similarity:

```

1 # 3. Embed (Turn text into numbers)
2     texts = [d["passage"] for d in docs]
3     emb_texts = self.embedder.encode(texts, convert_to_n
4     q_emb = self.embedder.encode([query], convert_to_n
5
6     # 4. Rank (Find similarity)
7     def cosine(a, b):
8         return np.dot(a, b) / (np.linalg.norm(a) * np.
9
10    sims = [cosine(e, q_emb) for e in emb_texts]
11    top_idx = np.argsort(sims)[::-1][:TOP_PASSE
12    top_passages = [{"url": docs[i]["url"], "passage":
```

## Step 4: Generate a Mini Summary

Now, extract and rank sentences within top passages to form a concise summary:

```
1      # 5. Summarize (Extractive)
2      sentences = []
3      for tp in top_passages:
4          for s in split_sentences(tp["passage"]):
5              sentences.append({"sent": s, "url": tp["ur"]
6
7      if not sentences:
8          summary = "No summary could be generated."
9      else:
10         sent_texts = [s["sent"] for s in sentences]
11         sent_embs = self.embedder.encode(sent_texts, c
12         sent_sims = [cosine(e, q_emb) for e in sent_em
13
14         top_sent_idx = np.argsort(sent_sims)[::-1][:SU
15         chosen = [sentences[idx] for idx in top_sent_i
16
17         # De-duplicate and format
18         seen = set()
19         lines = []
20         for s in chosen:
21             key = s["sent"].lower()[:80] # Check first
22             if key in seen:
23                 continue
24             seen.add(key)
25             lines.append(f"{s['sent']} (Source: {s['ur
26             summary = " ".join(lines)
27
```

```
28     elapsed = time.time() - start
29     return {"query": query, "passages": top_passages,
```

## Step 5: Run the Agent

Finally, run the AI agent with any query:

```
1 if __name__ == "__main__":
2     agent = ShortResearchAgent()
3     q = "What causes urban heat islands and how can cities
4
5     print(f"Running query: {q}\n")
6     out = agent.run(q)
7
8     print("\nTop passages:")
9     for p in out["passages"]:
10        print(f"- score {p['score']:.3f} src {p['url']}\n")
11
12    print("--- Extractive summary ---")
13    print(out["summary"])
14    print("-----")
15    print(f"\nDone in {out['time']:.1f}s")
```

Found 6 urls.

Top passages:

- score 0.852 src <https://environment.co/what-are-urban-heat-islands-and-how-to-prevent-them/>  
heat islands is an urgent environmental goal, as climate change impacts every aspect of human life. Here's a closer look at the causes of urban heat islands and how smart city planning can address the...
- score 0.839 src <https://environment.co/what-are-urban-heat-islands-and-how-to-prevent-them/>  
We are reader-supported. When you buy through links on our site, we may earn affiliate commission. Look at many metro areas from a distance on a warm day, and you'll notice a dome-like atmospheric app...
- score 0.820 src <https://earth.org/how-cities-around-the-world-are-tackling-the-urban-heat-crisis/>  
As climate change accelerates and urbanisation intensifies, cities worldwide face a growing threat: urban heat. Rising temperatures, exacerbated by the urban heat island effect, are endangering public...
- score 0.729 src <https://www.wri.org/insights/urban-heat-effect-solutions>  
others. The reason for this comes down to the urban environment. Built infrastructure like roads, buildings and sidewalks, as well as natural infrastructure like trees and water bodies, determines how...
- score 0.708 src <https://environment.co/what-are-urban-heat-islands-and-how-to-prevent-them/>  
Fahrenheit. However, they can soar up to 15° to 20° Fahrenheit higher during the worst part of the afternoon. Excessive heat can kill. People who live in urban heat islands run higher risks of heat-re...

--- Extractive summary ---

A primary cause of urban heat islands is the built environment. (Source: <https://environment.co/what-are-urban-heat-islands-and-how-to-prevent-them/>) Other causes of urban heat islands include: According to the EPA, urban heat islands raise temperatures by an average of 1° to 7° (Source: <https://environment.co/what-are-urban-heat-islands-and-how-to-prevent-them/>) Here's a closer look at the causes of urban heat islands and how smart city planning can address them. (Source: <https://environment.co/what-are-urban-heat-islands-and-how-to-prevent-them>)

Output

## Final Words

You have just built the core logic of **RAG**, the architecture behind many of the most powerful GenAI systems today. First, you built the **Retriever** (Search, Fetch, Chunk); next, you built the **Ranker** (Embed, Cosine Similarity); and then, you built a simple **Generator** (the extractive summarizer).

The journey from a simple script to a powerful AI system is just a series of small, understandable steps. You just took the first and most important one.

I hope you liked this article on how to build an AI agent to automate your research. Feel free to ask valuable questions in the comments section below. You can follow me on [Instagram](#) for many more resources.





## Aman Kharwal

AI/ML Engineer | Published Author. My aim is to decode data science for the real world in the most simple words.



ARTICLES: 2028



### PREVIOUS POST

[The 4 Best AI Courses to Take in 2026](#)



### NEXT POST

[MLOps Project Ideas to Get You Hired](#)

## Recommended For You

By Aman Kharwal

**20 ML Projects for 2026.**

AmanXai

[20 Projects to Master ML for 2026](#)

December 17, 2025

By Aman Kharwal

**Build an AI Resume Screener with Python & Llama 3.**

AmanXai

[Build an AI Resume Screener with Python & Llama 3](#)

December 16, 2025



By Aman Kharwal

## From Student to **ML Engineer**: The 2026 Roadmap.



By Aman Kharwal

## Feature Engineering Tricks Most Courses Don't Teach.



## The Ultimate 2026 ML Engineering Roadmap

December 15, 2025

## Feature Engineering Tricks I Use as an ML Engineer

December 13, 2025 / 1 Comment

**Leave a Reply**

---

