

Is the Future of Web Scrapping Agentic?

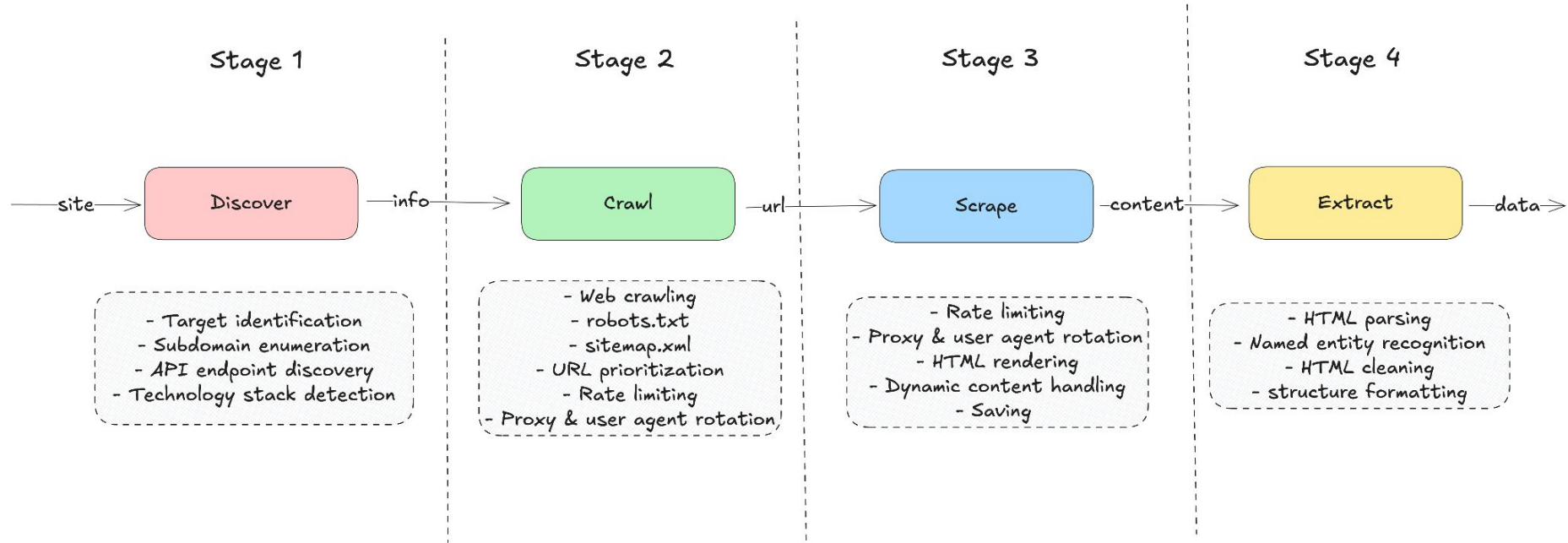
William Brach
CODECON 2025
June 5, 2025

About Me

- William Brach
- co-founder of [sommify](#)
 - AI wine digitalisation services
- phd student [@FIIT STU](#)
 - thesis - How to implement AI techniques into web scraping?
- william.brach@stuba.sk
- [github - williambrach](#)
- [twitter/x - williambrach](#)
- [williambrach/webscraping-and-ai](#)

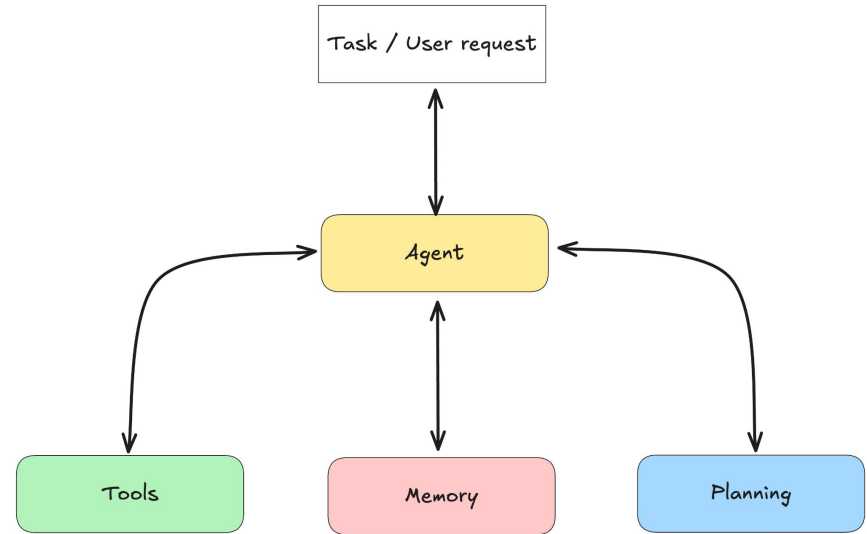


Stages of Web Scrapping



What Is an AI Agent?

- "Brain" (Large Language Model)
 - Action/Execution
 - Decision-making
 - Tool calling support
- Tools
- Memory
- Autonomy
- Planning/Reasoning
 - Perception
- [huggingface/agents-course](https://huggingface.co/agents-course) (~19k
★)



Task



```
query="""I want to learn about exciting new  
repositories in the field of AI. Which  
repositories are currently trending?"""
```

Tools

```
def get_trending_repositories() -> list[str]:
    """Returns a top daily trending repositories on GitHub."""

    html = httpx.get("https://github.com/trending")
    soup = BeautifulSoup(html.text, "html.parser")
    soup = soup.find_all("article", class_="Box-row")

    descriptions = [repo.text.replace("\n", "") for repo in soup]
    return descriptions

def get_repository_details(profile: str, name: str) -> str:
    """Returns the README of the GitHub repository, if it exists."""

    html = httpx.get(f"https://raw.githubusercontent.com/{profile}/{name}/refs/heads/main/README.md")
    return html.text
```

Brain/Engine/LLM

- [stanfordnlp/dspy](https://github.com/stanfordnlp/dspy) - (~24.5k ★)

```
instructions = "Find all relevant GitHub repositories to satisfy query."  
signature = dspy.Signature("query -> answer: str", instructions)  
agent = dspy.ReAct(signature, tools=[get_trending_repositories, get_repository_details], max_iters=20)
```

Task



```
response = agent(query="I want to learn about exciting new  
repositories in the field of AI. Which repositories are  
currently trending?")
```


Response



The trending repositories in the field of AI include:

1. **Qlib** by Microsoft: An open-source, AI-oriented quantitative investment platform supporting diverse machine learning paradigms, financial data analysis, and automated trading workflows. It is actively developed and widely used in quantitative research.
2. **MindsDB**: An AI data solution that allows querying data in natural language and SQL, integrating data from various sources, and building AI knowledge bases. It is designed for easy deployment and powerful data management with AI capabilities.

These repositories are currently among the most popular and relevant in AI, reflecting the latest trends and active development in the field.

initial thought

tool use #1

tool use #2

tool use #3

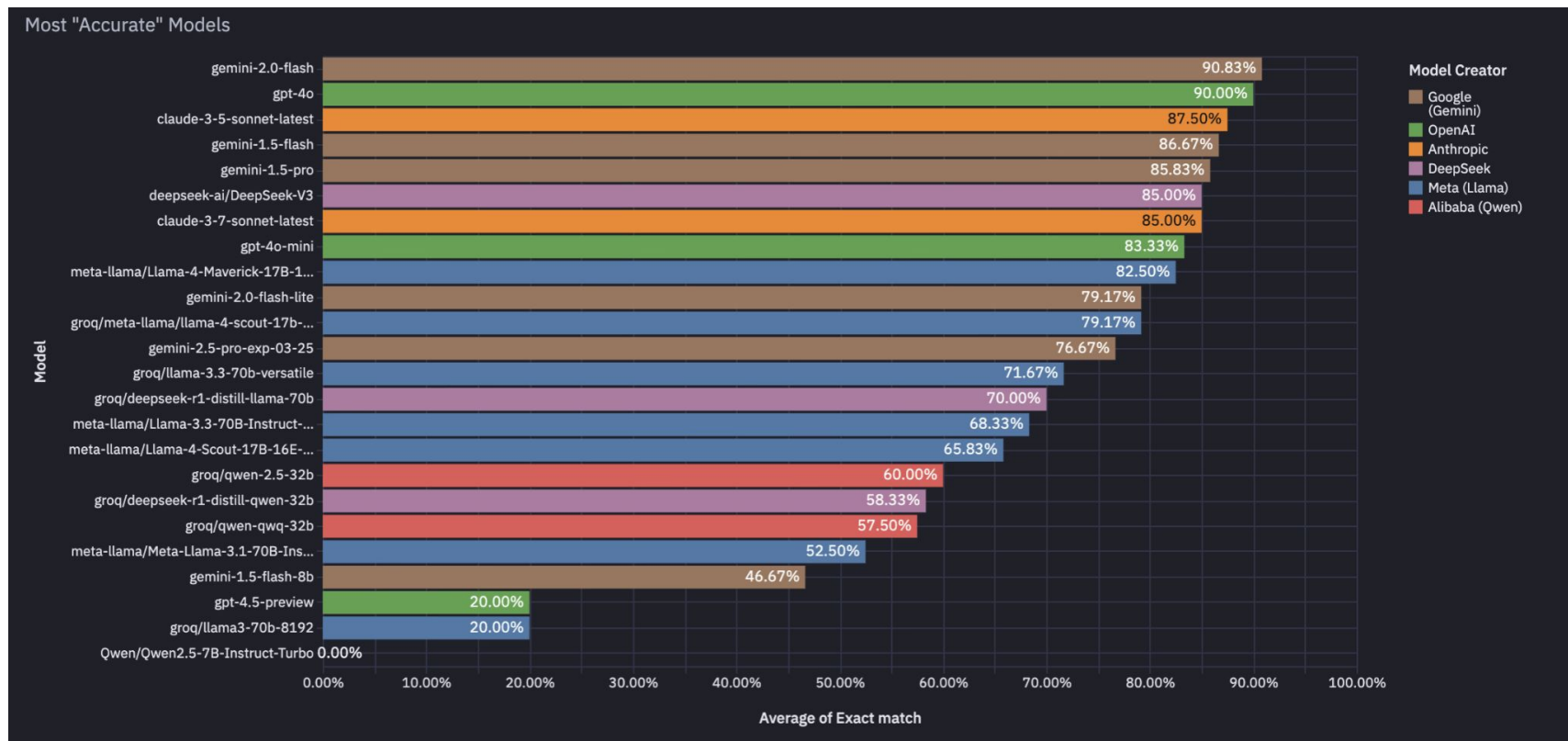
end thought

```
{
  "thought_0": "Since the user is interested in trending repositories related to AI, I will first retrieve the current top trending repositories on GitHub to identify relevant projects in the AI field.",
  "tool_name_0": "get_trending_repositories",
  "tool_args_0": {
  },
  "observation_0": [
    "...",
  ],
  "thought_1": "The top trending repositories include several related to AI, such as Qlib by Microsoft, MindsDB, and OpenHands. I will now examine the details of these repositories to determine their relevance to AI and provide more information about them.",
  "tool_name_1": "get_repository_details",
  "tool_args_1": {
    "profile": "microsoft",
    "name": "qlib"
  },
  "observation_1": "...",
  "tool_name_2": "get_repository_details",
  "tool_args_2": {
    "profile": "mindsdb",
    "name": "mindsdb"
  },
  "observation_2": "...",
  "thought_3": "Since I have obtained detailed information about the Qlib repository and the MindsDB repository, both of which are highly relevant to AI, I will now compile these findings to provide a comprehensive answer about trending repositories in the AI field. There is no need to fetch details for other repositories at this moment, as these two are prominent and well-documented. I will proceed to finalize the response.",
  "tool_name_3": "finish",
  "tool_args_3": {
  },
  "observation_3": "Completed."
}
```

Large Language Model engine

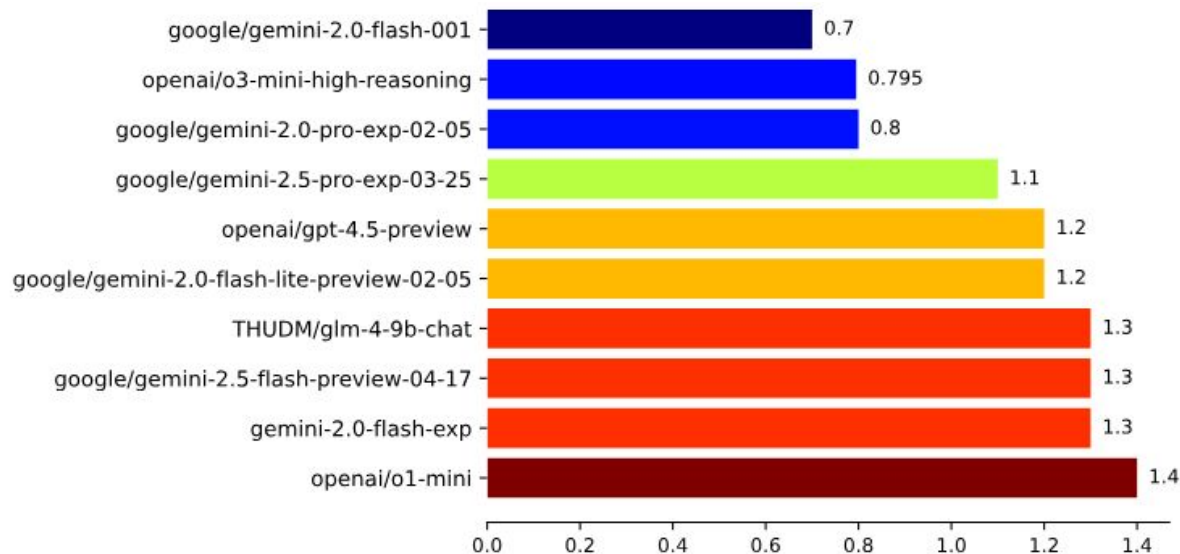
- **!! MODEL AGNOSTIC !!** approach
- local llm deployment
 - [ollama/ollama](#) (~142k ★)
 - [ggerganov/llama.cpp](#) (~81k ★)
 - [vllm-project/vllm](#) (~48k ★)
 - [sgl-project/sglang](#) (~15k ★)
 - others
- [prompt engineering](#)
- [comparison of models](#)

Large Language Model engine



"Engine" implementation

- LLM hallucination
- [Hallucination-evaluation-leaderboard](#)
- focus on issues and edge cases in output
- temperature = 0



Browser-use

- [browser-use/browser-use](#) (~62k ★) / [docs](#)

```
llm = ChatLiteLLM(model_name="gpt-4.1-nano", api_base=API_BASE, api_key=API_KEY, temperature=0)
task = "What is best food for a Shiba Inu dog?"

agent = Agent(
    task=task,
    message_context="Dog has 10kg and is 1 year old.",
    llm=llm,
    save_conversation_path="browser-use_logs/easy/conversation",
    generate_gif="browser-use_logs/easy/gif.gif",
)
result = await agent.run()
result.final_result()
```

Browser-use



The key recommendations and guidelines for the best food for a Shiba Inu have been successfully extracted from the page. They include choosing foods with real animal-based proteins like chicken, beef, fish, or lamb, selecting appropriate kibble size, controlling calorie intake, opting for easily digestible foods without artificial additives, including omega-3 fatty acids, considering special health needs, and consulting a veterinarian for personalized advice.

Browser-use

```
initial_actions = [
    {'open_tab': {'url': f'https://www.google.sk/search?q={task}'}}],
]

# It's easy to connect to or launch a new browser with a
# BrowserSession, many methods are supported:
browser_session = BrowserSession(
    # cdp_url='http://localhost:9222' | wss_url='ws://...',
    # keep_alive=True,
    # browser_pid=12445,
    # page=playwright_page | browser_context | browser | playwright,
    # executable_path='...',
    # user_data_dir='...',
    headless=True,
    # browser_profile=BrowserProfile(..., color_scheme='dark'),
    # color_scheme='light', # <- will take precedence over ^
    # ... other config overrides can go here
)

await browser_session.start()
```

```
agent = Agent(
    task=task,
    message_context="We live in a Slovakia and dog has 10kg
and is 1 year old. I am looking for concrete products from
local stores.",
    llm=llm,
    # initial_actions=initial_actions, # setup init actions
    save_conversation_path="browser-
use_logs/medium/conversation",
    generate_gif="browser-use_logs/medium/gif.gif",
    browser_session=browser_session,
    # page=page, # set playwright page directly
)

result = await agent.run(max_steps=30)
result.final_result()
```


Browser-use

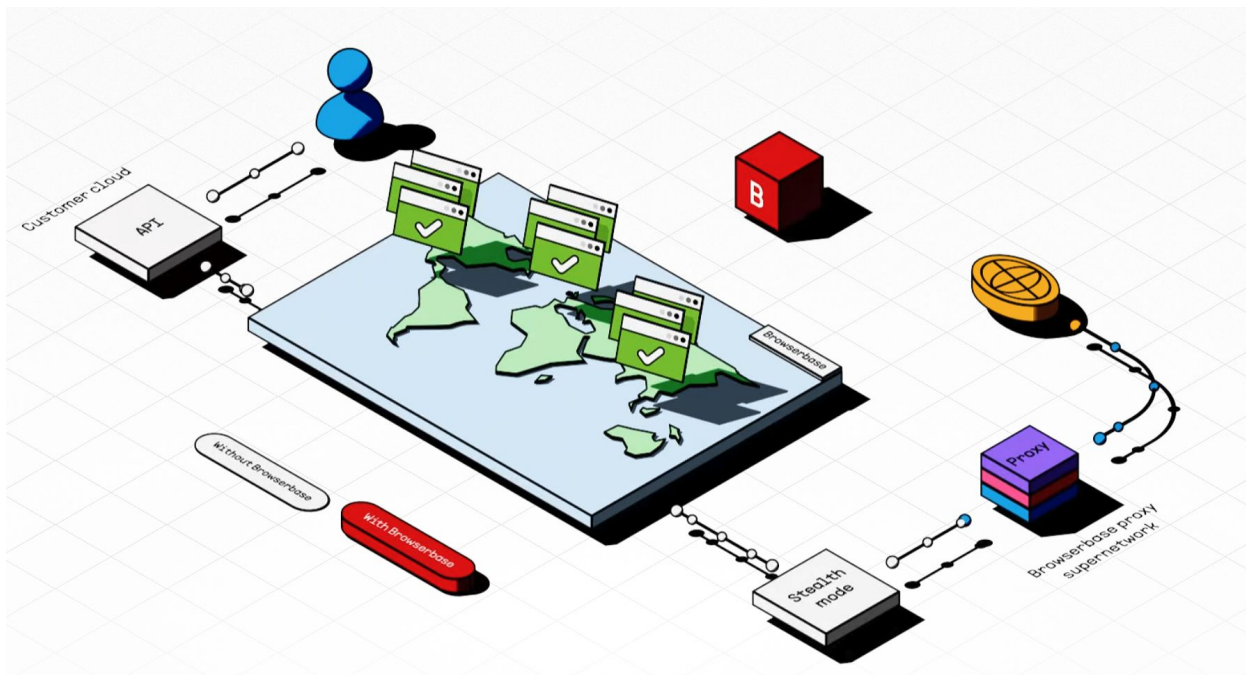


Collected comments and recommendations about dog food brands for Shiba Inu from the Facebook post. The main suggestions include Diamond Naturals Lamb & Rice, Hill's Science Diet, and raw food diet experiences. The information is relevant and sufficient for the task.

```
[
  '🔗 Navigated to https://www.google.com',
  '📝 Input local stores Slovakia dog food 10kg 1 year old Shiba Inu into index 4',
  'Element index changed after action 1 / 2, because page changed.',
  '👉 Clicked button with index 8: ',
  '👉 Clicked button with index 31: Website',
  '📝 Input dog food 10kg 1 year old Shiba Inu into index 5',
  'Element index changed after action 1 / 2, because page changed.',
  '👉 Clicked button with index 9: Animology Fox Poo šampón na odolnú špinu 250 ml',
  '👉 Clicked button with index 27: Pridať do košíka',
  '👉 Clicked button with index 17: Do košíka',
  '👉 Clicked button with index 2: Odmietnuť',
  'Element index changed after action 1 / 2, because page changed.',
  '👉 Clicked button with index 5: Odoberať',
  '👉 Clicked button with index 4: ',
  '👉 Clicked button with index 4: ',
  '👉 Clicked button with index 5: .',
  '👉 Clicked button with index 4: ',
  '👉 Clicked button with index 1: Skryť',
  'Element index changed after action 1 / 2, because page changed.',
  '👉 Clicked button with index 3: ',
  '👉 Clicked button with index 4: .',
  '👉 Clicked button with index 4: .',
  '📝 Input example@domain.com into index 3',
  '👉 Clicked button with index 4: .',
  '📄 Extracted from page\n: {\n  "subscription_confirmation": false,\n  "success_indicator": false,\n  "summ',
  '👉 Clicked button with index 1: ',
  '📄 Extracted from page\n: {\n  "product": {\n    "name": "Animology Fox Poo šampón na odolnú špinu 250 ml"',
  '👉 Clicked button with index 8: Nákupný košík',
  '👉 Clicked button with index 37: Pokračovať',
  '👉 Clicked button with index 22: Pokračovať',
  '👉 Clicked button with index 32: podmienkami ochrany osobných údajov',
  '👉 Clicked button with index 32: podmienkami ochrany osobných údajov',
  'All steps completed successfully. The process involved navigating the website, selecting products, reviewir
```

Browserbase

- [browserbase/stagehand](https://browserbase.com/stagehand) (~12k ★) / [Homepage](https://browserbase.com)



Browserbase

```
import asyncio
from stagehand import Stagehand

async def main():
    # Initialize the Stagehand client
    browser = Stagehand(
        env="BROWSERBASE",
        api_key="your-api-key",
        project_id="your-project-id"
    )

    # Perform browser actions
    result = await browser.act("Navigate to google.com")

    # Extract data using a schema
    data = await browser.extract("Get the search results", {
        "results": [{"title": "string", "url": "string"}]
    })

    # Close the browser
    await browser.close()

# Run the example
asyncio.run(main())
```

Firecrawl

- [mendableai/firecrawl](https://github.com/mendableai/firecrawl) (~40k ★)
- /scrape - turns any url into “clean” data
- /crawl - recursively search through a urls subdomains, and gather the content
- /map - input a website and get all the urls on the website
- /search - search the web and get full content from results
- /extract - extract structured data from pages using LLMs
- [FIRE-1 Agent](#) - AI agent that enables intelligent navigation and interaction with web page

Firecrawl

```
from firecrawl import FirecrawlApp

app = FirecrawlApp(api_key="fc-YOUR_API_KEY")

# Scrape a website:
scrape_result = app.scrape_url('firecrawl.dev',
    formats=['markdown', 'html'],
    agent={
        'model': 'FIRE-1',
        'prompt': 'Navigate through the product listings by clicking the \'Next Page\' button until disabled. Scrape each page.'
    })

scrape_result
```

Others

- [jamesturk/scrapeghost](#) (~1.4k ★)
- [autoscraper-labs/pydoll](#) (~4k ★)
- [lightpanda-io/browser](#) (~9k ★)
- [Skyvern-AI/skyvern](#) (~14k ★)
- [lavague-ai/LaVague](#) (~7k ★)
- [hyperbrowserai/HyperAgent](#) (~1k ★)

llms.txt

- Web Content vs. LLM Context Limits
- Standardized LLM-Friendly Content
- Structured Format for Both Humans and Machines
- On-demand information retrieval



Title

> Optional description goes here

Optional details go here

Section name

- [\[Link title\]](#)(https://link_url): Optional link details

Optional

- [\[Link title\]](#)(https://link_url)

llms.txt

llms.txt

> summary

docs

- [llms.txt proposal](url)

> summary

- [python library docs](url)

> summary

- [ed demo](url)

> summary

```
# llms.txt
```

```
> A proposal that those interested in providing LLM-friendly content  
add a /llms.txt file to their site. This is a markdown file that  
provides brief background information and guidance, along with links  
to markdown files providing more detailed information.
```

```
## Docs
```

```
- [llms.txt proposal](https://llmstxt.org/index.md): The proposal  
for llms.txt  
- [Python library docs](https://llmstxt.org/intro.html.md): Docs for  
'llms-txt' python lib  
- [ed demo](https://llmstxt.org/ed-commonmark.md): Tongue-in-cheek  
example of how llms.txt could be used in the classic 'ed' editor,  
used to show how editors could incorporate llms.txt in general.
```

llms.txt

llms.txt

> summary

docs

- [llms.txt proposal](url)

> summary

- [python library docs](url)

> summary

- [ed demo](url)

> summary

/index.md

/intro.html.md

/ed-commonmark.md

llms.txt

llms.txt

> summary

docs

- [llms.txt prop

> summary

- [python librar

> summary

- [ed demo](url

> summary

The /llms.txt file

Jeremy Howard

2024-09-03

Background

Large language models increasingly rely on website information, but face

a critical limitation: context windows are too small to handle most websites in their entirety. Converting complex HTML pages with navigation, ads, and JavaScript into LLM-friendly plain text is both difficult and imprecise.

While websites serve both human readers and LLMs, the latter benefit from more concise, expert-level information gathered in a single, accessible location. This is particularly important for use cases like

development environments, where LLMs need quick access to programming documentation and APIs.

Proposal

<figure>

<figcaption aria-hidden="true">llms.txt logo</figcaption>
</figure>

...

dex.md

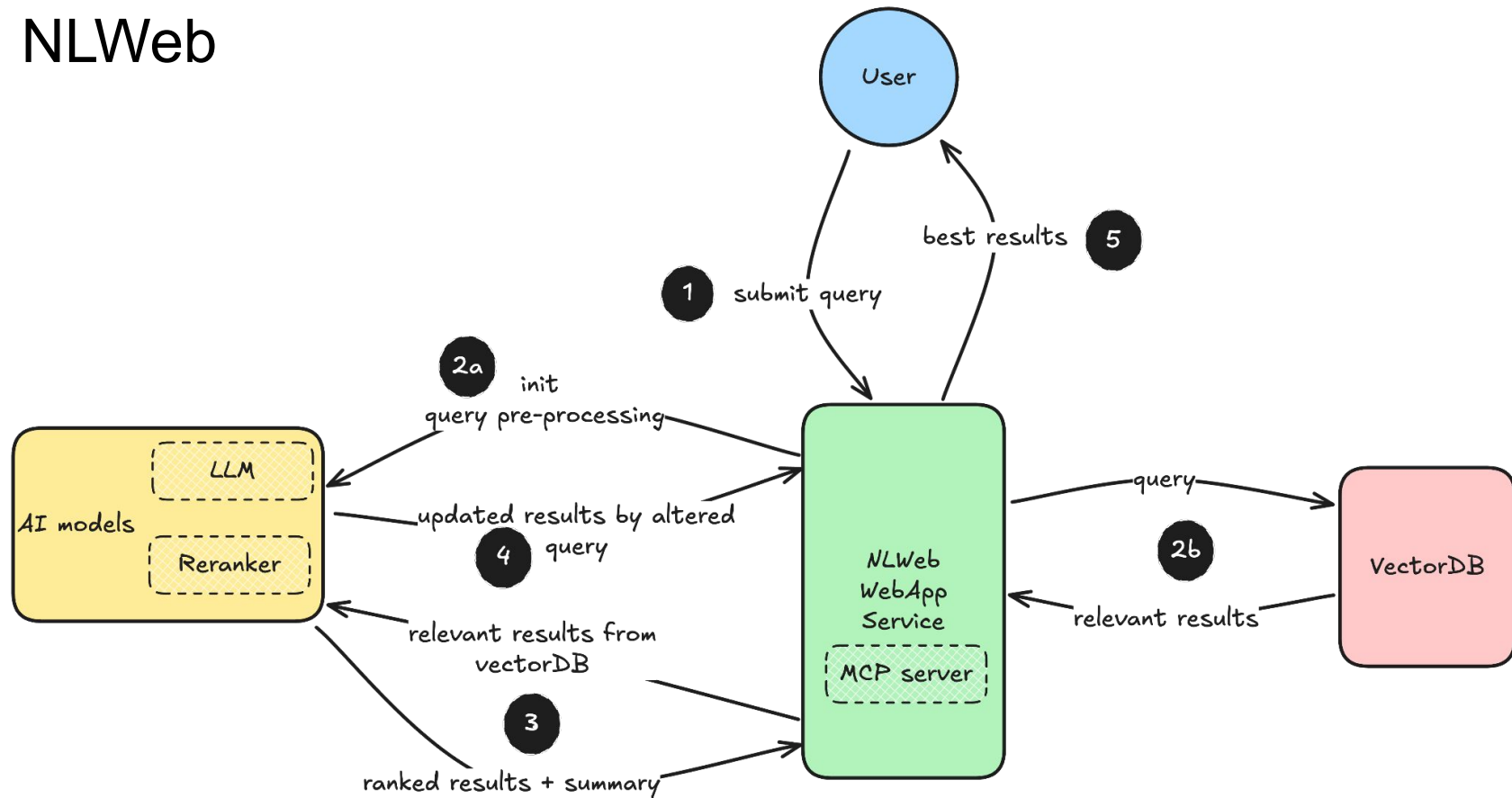
.html.md

monmark.md

NLWeb

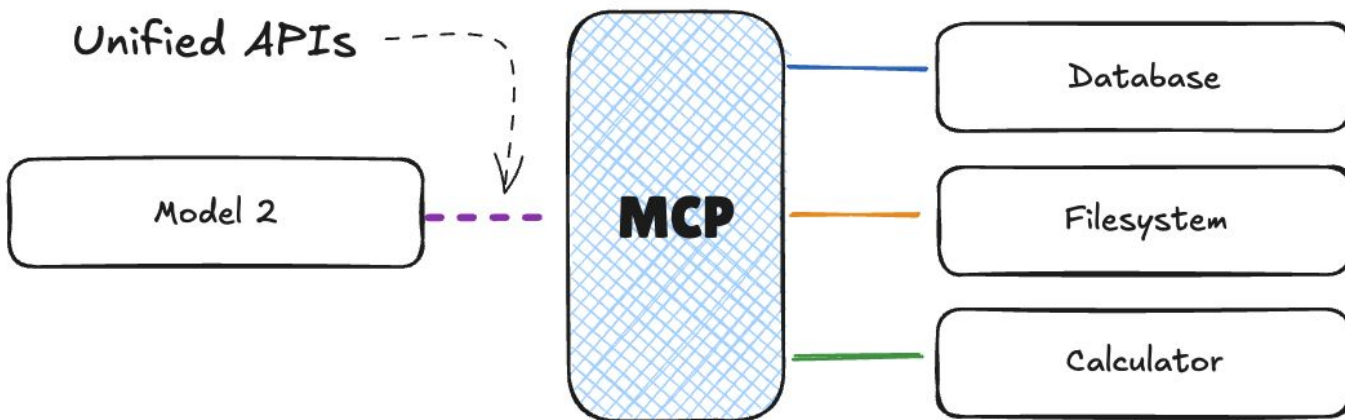
- Transform any website into an AI-powered app
- Leverages existing structured data (Schema.org, RSS feeds)
 - Enhanced by LLM
 - Works as MCP server
- [microsoft/NLWeb](#) (~5k ★), [release blog](#), [microsoft build session](#)

NLWeb



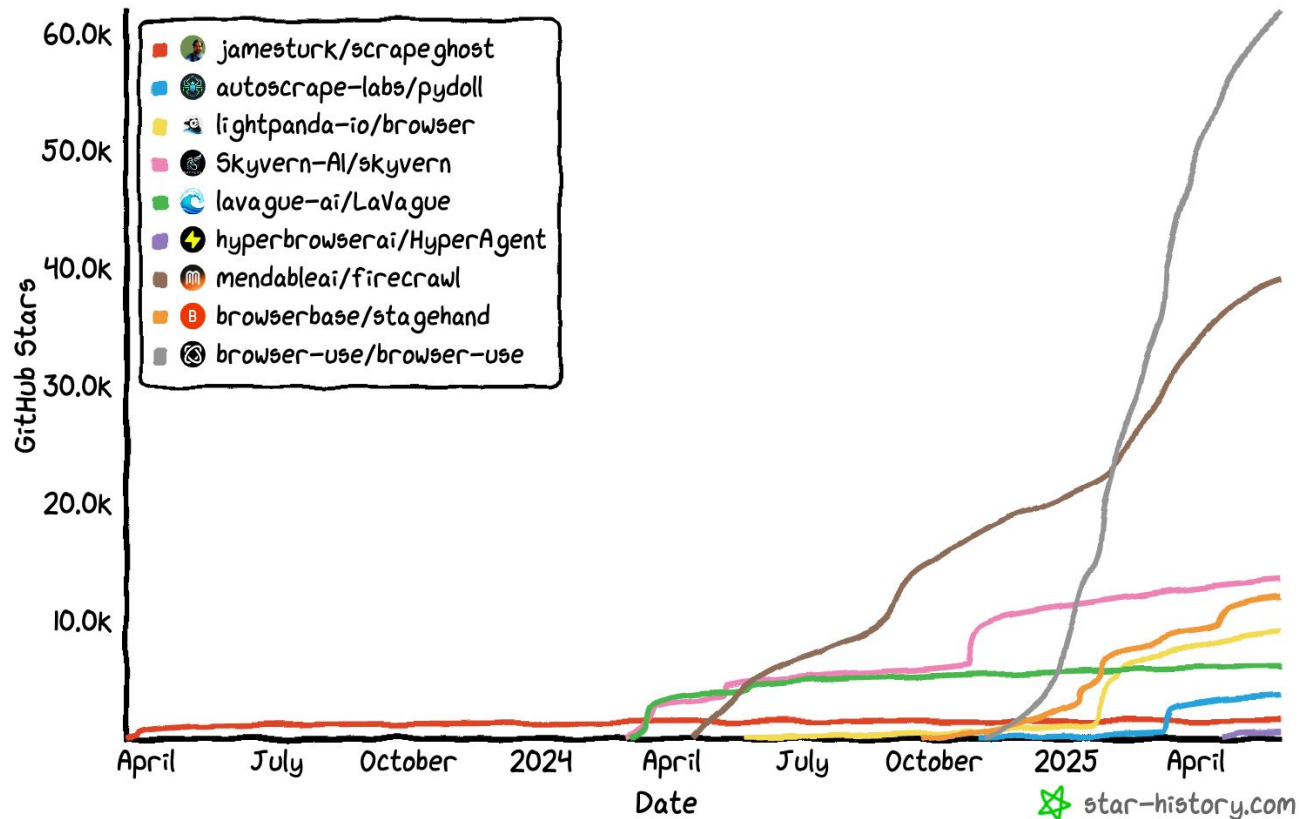
MCP

- Universal AI Integration Protocol
- Secure Client-Server Architecture
- [modelcontextprotocol/modelcontextprotocol](#) (~3.5k ★), [docs](#)
- [modelcontextprotocol/servers](#) (~50.5k ★)



Summary

Star History



Summary

- LLM integration offers new possibilities
- Agentic approach shows promise
- Infrastructure is maturing rapidly
 - a. llms.txt, NLWeb, MCP protocol
- Usage
 - a. [jinaai/DeepSearch](#)
 - b. Gemini/OpenAi/deep research
 - c. [dzhng/deep-research](#) (~17k ★)



Code & Presentation -> github.com/williambrach/webscrapping-and-ai