



# Supa-crawl (LLM-playground) – architecture & workflow

## Overview

The `supa-crawl` project is a Python web-crawling pipeline that wraps the [Crawl4AI](#) framework with Supabase storage. The `AdvancedWebCrawler` class configures a headless Chromium browser (via `BrowserConfig`) and creates a `CrawlerRunConfig` for each crawl <sup>1</sup>. It supports memory-adaptive and semaphore dispatchers for multi-URL crawling <sup>2</sup> and integrates an LLM extraction strategy to produce JSON summaries <sup>3</sup>. Results (URL, raw markdown, analysis) can be stored in a Supabase `pages` table using a dedicated handler <sup>4</sup>. `main.py` orchestrates tests of each dispatcher and the LLM pipeline <sup>5</sup>. Environment variables specify the Supabase URL/key and OpenAI API key <sup>6</sup>. Official Crawl4AI docs explain why LLM extraction is used (for unstructured or semantic data) <sup>7</sup> and how dispatchers offer adaptive crawling with rate limiting <sup>8</sup>. Supabase CLI docs describe how to initialize and start a local Supabase project with `supabase init` and `supabase start` <sup>9</sup> <sup>10</sup>.

## Sources used

Official documentation and code were referenced from:

- **Crawl4AI docs** – LLM extraction strategy, dispatcher behaviour and configuration <sup>7</sup> <sup>8</sup> <sup>11</sup>.
- **Supabase docs** – CLI commands (`supabase init`, `start`, `stop`) for local development <sup>9</sup> <sup>10</sup> <sup>12</sup>.
- **Repository code** – `src/config/environment.py`, `src/crawlers/async_crawler.py`, `src/storage/supabase_handler.py`, `src/models/schemas.py` and `main.py` on the `LLM-playground` branch <sup>6</sup> <sup>1</sup> <sup>4</sup> <sup>13</sup> <sup>14</sup>.

## Directory tree (relevant parts)

```
supa-crawl/
├─ src/
│   └─ config/
│       └─ environment.py          # loads env vars and creates browser/crawler
config
├─ crawlers/
│   └─ async_crawler.py          # defines AdvancedWebCrawler and crawl methods
├─ models/
│   └─ schemas.py                # Pydantic model for LLM extraction
└─ storage/
    └─ supabase_handler.py        # handles Supabase client and insertion/upsert
```

```

└─ main.py                # demo entry point testing each pipeline stage
└─ supabase/              # migrations and config for Supabase project
└─ docs/                  # documentation (not analysed)
└─ supabase-cli-commands.md # local guide for Supabase CLI commands
└─ tests/                 # pytest scripts
└─ ... (other packaging files)

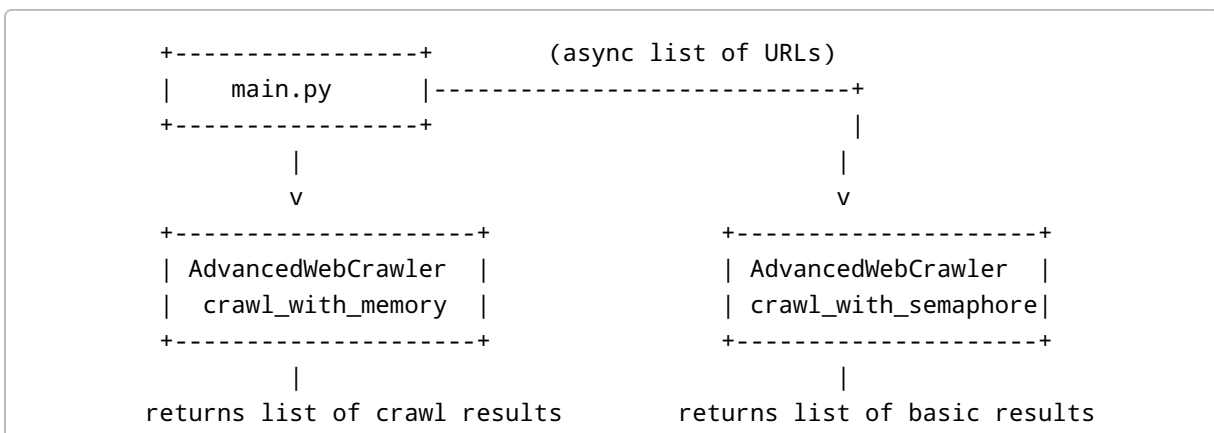
```

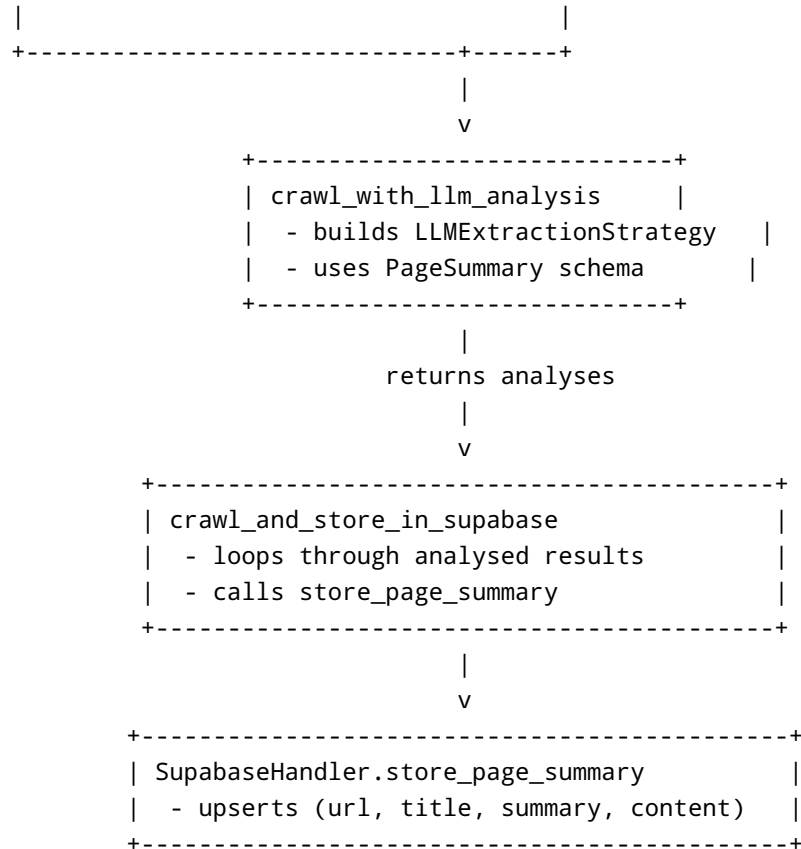
## File breakdown

Path	Role	Key functions/classes	Called by / Call
<code>src/config/environment.py</code>	Centralised environment and crawler configuration	<code>EnvironmentConfig</code> reads <code>SUPABASE_URL</code> , <code>SUPABASE_KEY</code> and <code>OPENAI_API_KEY</code> from <code>.env</code> <sup>6</sup> ; <code>CrawlerConfig.create_browser_config</code> builds a <code>BrowserConfig</code> with stealth settings <sup>15</sup> ; <code>create_crawler_run_config</code> returns a <code>CrawlerRunConfig</code> with default extraction strategy <sup>16</sup>	Imported by <code>main.py</code> for global config
<code>src/models/schemas.py</code>	Pydantic schemas for LLM extraction	<code>PageSummary</code> defines <code>title</code> and <code>summary</code> fields <sup>18</sup> ; <code>CrawlResult</code> and <code>Crawl4AIResponse</code> wrap results <sup>19</sup>	Used by <code>main.py</code> for configuring <code>LLMExtractionStrategy</code>
<code>src/crawlers/async_crawler.py</code>	Implements <code>AdvancedWebCrawler</code>	Constructor sets <code>BrowserConfig</code> and <code>CrawlerRunConfig</code> <sup>1</sup> and initialises <code>SupabaseHandler</code> <sup>21</sup> . <code>crawl_with_memory_adaptive_dispatcher</code> configures a <code>MemoryAdaptiveDispatcher</code> with rate limiting and memory thresholds <sup>2</sup> ; <code>crawl_with_semaphore_dispatcher</code> uses <code>SemaphoreDispatcher</code> for simple concurrency <sup>22</sup> . <code>crawl_with_llm_analysis</code> checks for an OpenAI key and builds an <code>LLMExtractionStrategy</code> with a multi-line prompt, schema and chunking parameters <sup>3</sup> ; it parses JSON results and returns analysis dictionaries <sup>23</sup> . <code>crawl_and_store_in_supabase</code> calls the LLM crawl and stores each summary in Supabase via <code>store_page_summary</code> <sup>24</sup> .	Called by <code>main.py</code> for demonstration of <code>SupabaseHandler</code> and <code>PageSummary</code>

Path	Role	Key functions/classes	Called by / Call
<code>src/storage/supabase_handler.py</code>	Wraps Supabase client and storage operations	<code>_initialize_client</code> creates a Supabase client when credentials exist <sup>25</sup> . <code>_extract_first_paragraph</code> trims markdown to the first meaningful paragraph for storage <sup>26</sup> . <code>store_crawl_results</code> inserts <code>url</code> and a truncated <code>content</code> into the <code>pages</code> table; optionally prepends LLM title/summary <sup>27</sup> . <code>store_page_summary</code> performs an upsert of <code>url</code> , <code>title</code> , <code>summary</code> and optional <code>content</code> using Supabase-py's <code>upsert()</code> <sup>28</sup> .	Called by <code>async_crawl_and_store</code> and tests
<code>main.py</code>	Example script demonstrating the crawler	Defines <code>urls</code> list of target pages <sup>29</sup> . Inside <code>main()</code> it initialises <code>AdvancedWebCrawler</code> , then sequentially tests memory-adaptive crawling, semaphore dispatching, LLM analysis and full pipeline storage <sup>5</sup> .	Run directly via <code>python main.py</code>
<code>supabase-cli-commands.md</code>	Local CLI reference (not official but summarises Supabase CLI usage)	Contains commands for installation, authentication, migrations, data inspection and troubleshooting	Stand-alone; used for onboarding

## Workflow diagram





## Configuration & environment variables

- `.env` – must define `SUPABASE_URL`, `SUPABASE_KEY` and `OPENAI_API_KEY`. `EnvironmentConfig` loads these and warns if missing <sup>30</sup>.
- **Crawler config** – `BrowserConfig` parameters (headless, browser type, user agent) and `CrawlerRunConfig` settings (cache mode, word count, extraction strategy) follow official patterns <sup>1</sup>.
- **Dispatchers** – choose between `MemoryAdaptiveDispatcher` (pauses based on RAM usage; includes rate limiter and monitor) and `SemaphoreDispatcher` (fixed concurrency). Crawl4AI docs explain that dispatchers provide adaptive crawling, rate limiting and real-time monitoring <sup>8</sup>.
- **LLM extraction** – define a `PageSummary` schema and specify `LLMExtractionStrategy` parameters like provider, API token, instruction, chunking and extraction type <sup>3</sup>. Crawl4AI docs note that LLM extraction supports various providers and uses chunking to manage token limits <sup>7</sup> <sup>31</sup>.
- **Supabase client** – initialised via `create_client(SUPABASE_URL, SUPABASE_KEY)` <sup>25</sup>. To run Supabase locally, official CLI commands are: `supabase init` to set up a project and `supabase start` to launch Postgres, Auth, Storage and Studio <sup>9</sup> <sup>10</sup>; `supabase stop` stops the services <sup>12</sup>.

## Limitations & next actions

- **Manual URL list** – `main.py` hard-codes test URLs <sup>29</sup>. For production use, implement URL ingestion (e.g., read from a file or database).
- **Schema rigidity** – `PageSummary` contains only `title` and `summary` <sup>18</sup>. To extract specific fields (e.g., prices or code snippets), define a new Pydantic schema and update the `LLMExtractionStrategy` prompt accordingly.
- **Error handling** – exceptions in `async_crawler.py` are printed but not retried beyond the dispatcher's rate limiter. Adding retry logic or logging would improve robustness.
- **Supabase table schema** – the `pages` table is assumed to have `url`, `title`, `summary` and `content` columns. Migrations aren't included; the `supabase` directory likely contains SQL files to create this schema but was not analysed.
- **Client dependencies** – the code expects the `supabase` Python client to be installed. `supabase_handler.py` exits if it cannot import `create_client` <sup>32</sup>.

Next steps:

1. **Automate ingestion** – integrate a queue or scheduler to supply URLs for crawling, rather than a static list.
  2. **Custom extraction** – extend `schemas.py` and modify the LLM prompt to pull domain-specific data (e.g., product prices, code examples).
  3. **Database migrations** – ensure migrations in `supabase` folder align with the fields used in `store_page_summary`; document how to run them (`supabase db push` and `supabase migration create`).
  4. **Monitoring & logging** – integrate logging libraries and consider using the Crawl4AI monitor API for real-time dashboards.
  5. **Cost optimisation** – since LLM extraction can be slow and costly <sup>7</sup>, implement a fallback to schema-based extraction when structured data is available.
-

1 2 3 20 21 22 23 24 **async\_crawler.py**

[https://github.com/Gaya56/supa-crawl/blob/LLM-playground/src/crawlers/async\\_crawler.py](https://github.com/Gaya56/supa-crawl/blob/LLM-playground/src/crawlers/async_crawler.py)

4 25 26 27 28 32 **supabase\_handler.py**

[https://github.com/Gaya56/supa-crawl/blob/LLM-playground/src/storage/supabase\\_handler.py](https://github.com/Gaya56/supa-crawl/blob/LLM-playground/src/storage/supabase_handler.py)

5 14 29 **main.py**

<https://github.com/Gaya56/supa-crawl/blob/LLM-playground/main.py>

6 15 16 17 30 **environment.py**

<https://github.com/Gaya56/supa-crawl/blob/LLM-playground/src/config/environment.py>

7 31 **LLM Strategies - Crawl4AI Documentation (v0.7.x)**

<https://docs.crawl4ai.com/extraction/llm-strategies/>

8 **Multi-URL Crawling - Crawl4AI Documentation (v0.7.x)**

<https://docs.crawl4ai.com/advanced/multi-url-crawling/>

9 10 12 **Supabase CLI | Supabase Docs**

<https://supabase.com/docs/guides/local-development/cli/getting-started>

11 **Browser, Crawler & LLM Config - Crawl4AI Documentation (v0.7.x)**

<https://docs.crawl4ai.com/core/browser-crawler-config/>

13 18 19 **schemas.py**

<https://github.com/Gaya56/supa-crawl/blob/LLM-playground/src/models/schemas.py>