

Leveraging AI to Expand Home Electrification and Efficiency Incentive Data Access

Dawson Chan¹, Sheel Gada¹, Kaleb Lamphiear¹, Giovanni Loia¹, Tommy Yue¹, Nick Montoni², Brian Lips²
¹Department of Chemical Engineering – University of Washington; ²NC Clean Energy Technology Center



Background

Electric utility providers and local governments offer economic incentives for renewable energy systems, energy-efficient equipment, and building improvements. The NC Clean Energy Technology Center (NCCETC) maintains the Database of State Incentives for Renewables and Efficiency (DSIRE), a comprehensive source for policies and incentives supporting the use of renewable energy. Currently, this database is manually populated and updated. Due to the volume of available incentives, it has become unfeasible to continue this method, especially for small utilities with under 30,000 customers.

The aim of our project is to use a language learning model (LLM) to create an automated web-scraping tool that finds and summarizes these incentives, returning an in-depth summary of information formatted to be entered into DSIRE.

a) Residential Energy Conservation Subsidy Exclusion (Corporate)

b) Examples of keywords chosen for scraping

Incentive Type	Derived Keywords
Rebate Program	Rebate, incentive application ...
Grant Program	Grant, funding, eligibility
Solar renewable energy credit program	Solar, efficiency, program ...

Figure 1: a) Example overview of an entry in DSIRE. b) Types of incentives as categorized within DSIRE alongside the keywords derived for scraping and analyzing webpages.

Pipeline Overview

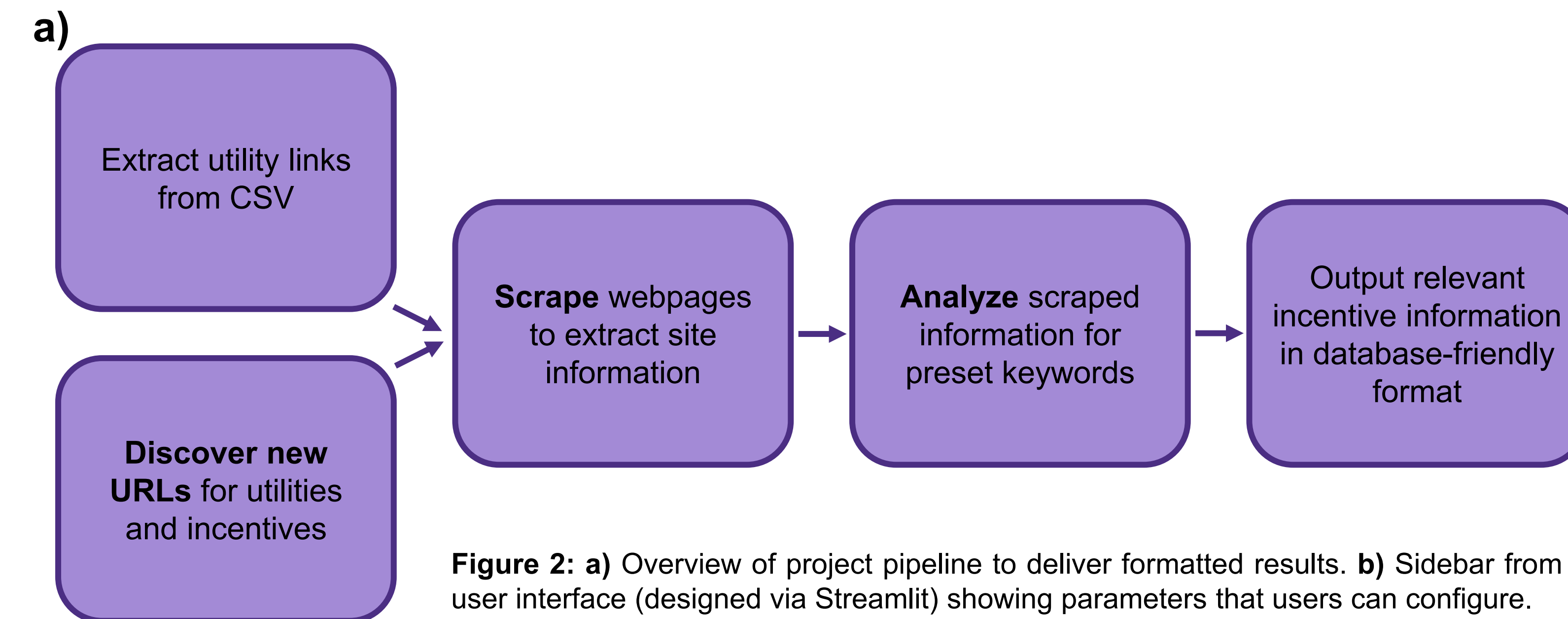
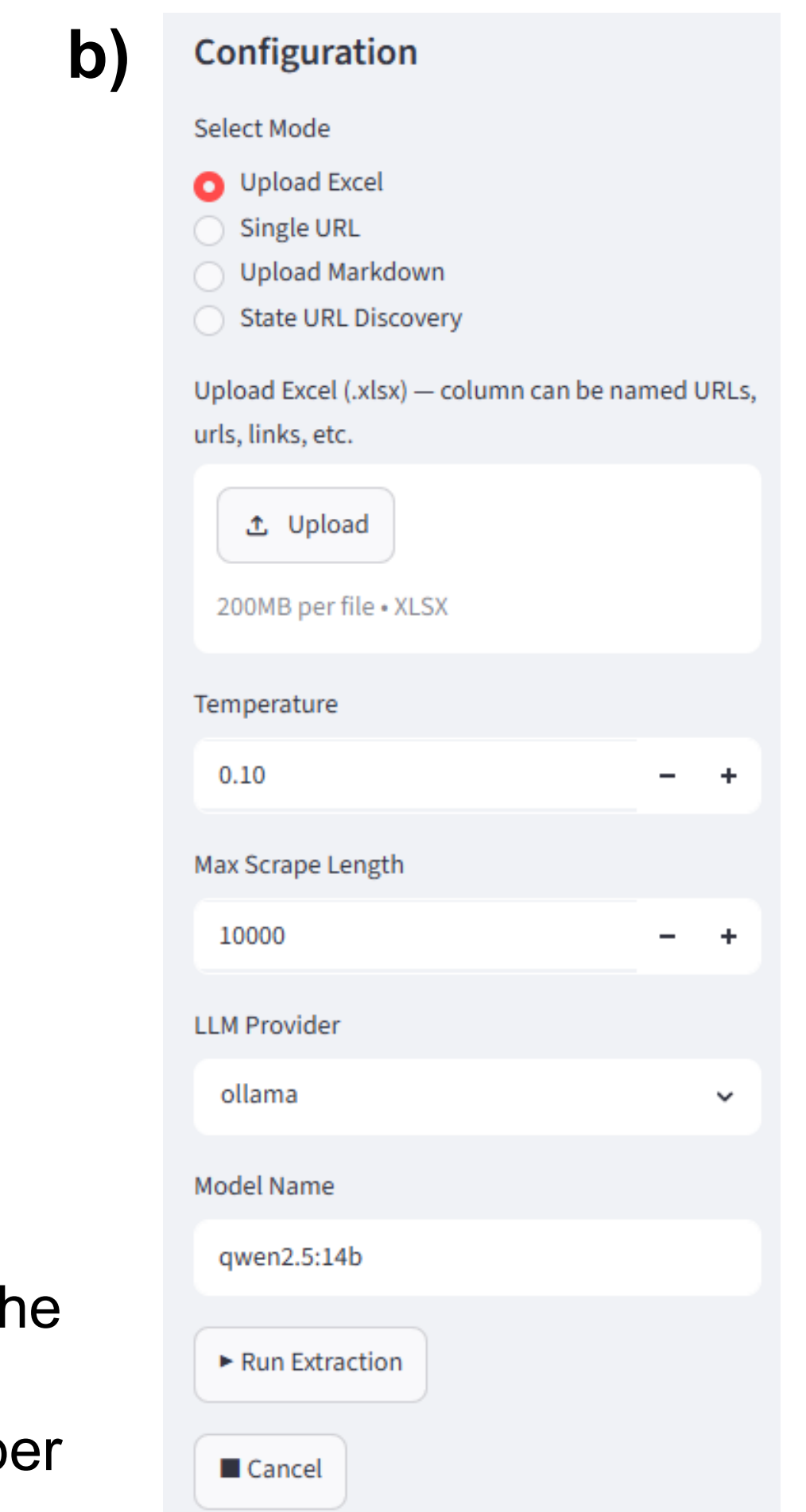


Figure 2: a) Overview of project pipeline to deliver formatted results. b) Sidebar from user interface (designed via Streamlit) showing parameters that users can configure.



- The pipeline includes different modes based on the structure of the sites and/or stored data:
- Upload Excel — provide a .xlsx with a list of URLs.
 - Single URL — attach a link or website that you wish to crawl, scrape, and analyze
 - Upload Markdown — If you have pre-scraped information from a website, you can upload the file for it to be analyzed (bypasses the scraping portion)
 - State URL Discovery — Choose which search engine, state, search topics, and results per search topic to find relevant URLs

Scraper

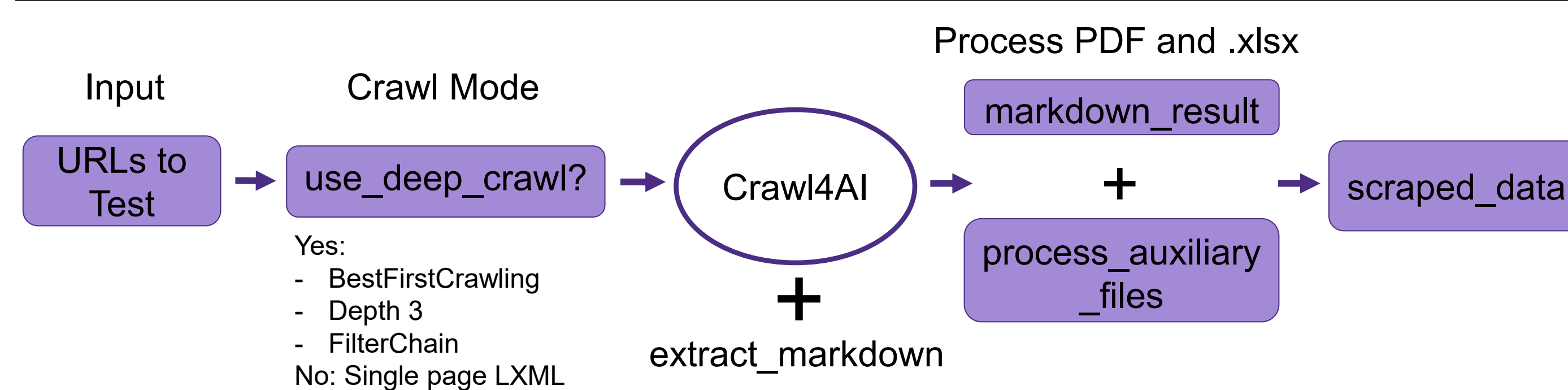


Figure 3: Scraper process flow diagram. The first portion of the pipeline was designed to gather all the relevant information found on a utility's website through a process called scraping.

- Crawling Method = BestFirstCrawlingStrategy - follows linked subpages across the site, prioritizing pages that score highest for incentive-related keywords
- max_depth - crawls up to 3 levels deep and stays within the seed domain
- FilterChain:
 - SEOFilter - helps identify pages with strong SEO characteristics (meta tags, headers, ect.)
 - DomainFilter - which domains to include or exclude

URL Discoverer

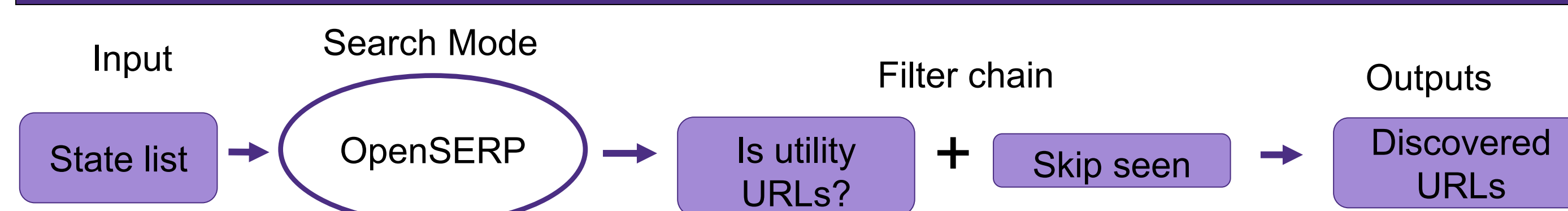
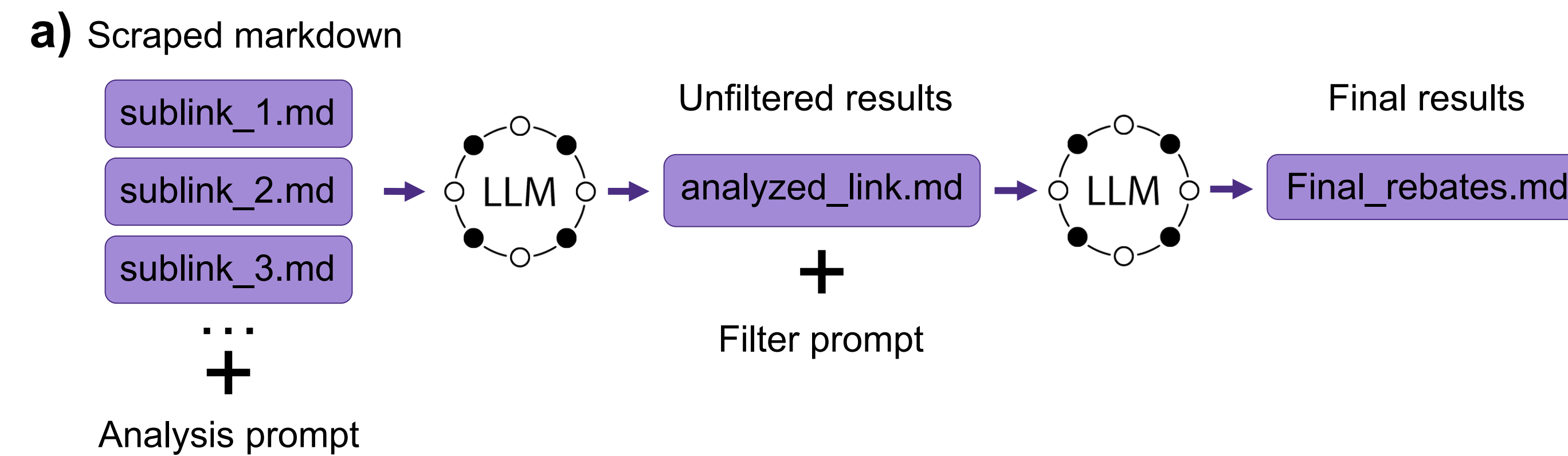


Figure 5: URL Discoverer process flow diagram. As an alternative to inputting known utility URLs, the URL discoverer was created to locate electric utilities that have not been captured in DSIRE by finding their official websites.

- Search topics - Utility ownership types (cooperative, municipal, investor-owned)
- Engine choice - Bing / DuckDuckGo tolerate automation; Google gives best results but triggers CAPTCHA
- Three-pass filter - block aggregators, keep utility domain signals, brand-name fallback on title keywords
- Merge - combines newly discovered URLs into the database with domain-level deduplication, skipping sites already in the database and backing up the old file.

Analyzer



b)

Model	Context window	Parameter count	Performance
Llama_3.1:8b	128k	8 billion	Good lightweight baseline, but struggles to filter
Gemma_4:31b	256k	31 billion	Conservative analysis
Gpt_5.4-pro	~1.05M	Undisclosed	Best overall model but \$\$\$

Figure 4: a) Analyzer process flow diagram. The scraped website information is fed into a series of LLMs, with the initial analyzer designed to find the program name, URL, concrete rebates, eligibility, and any additional utility information. These results are then fed into another LLM designed to filter and discard any sections deemed irrelevant. b) Three different models were tested. A model complexity increased, performance improved, but cost to run increased. Testing more models is recommended.

Recommendations & Limitations

- Parameters that should not be modified by user:
- **Search Depth** - Set at 3. Depth = 2 misses information; Depth = 4 is too slow
 - **AI Prompts** - Both analyzer prompts have been tuned to minimize errors
 - **Temperature** - Range of 0 (most rigid) to 1 (most creative). Set at 0.1 for both functions to keep responses grounded in source material
- Parameters that user can modify:
- **LLM** - Large models are more reliable, but may require advanced hardware
 - **LLM Time** - Maximum time model can spend per link
 - **Truncation Length** - Must change to fit each model's capacity
- Further investigation:
- **Paid API or Dedicated PC** - Provide the compute power externally
 - **Multi-user deployment** - Deploy via Streamlit to act as accessible website

Results

Figure 6: Example incentive output. After each scrape, final incentives are returned in both Excel (a) and markdown format (b). Exact categories included may vary by incentive, but general fields include Program Name, Program Type, Financial Details, Eligibility, Application Process, and Sector.

To quantify our pipeline's effectiveness, we used a subset of scraped markdown files from the city of Anaheim, CA (selected for seepage complexity and organization. 61 of the 65 definitive rebates were accurately pulled for a **success rate of 94%**. 17 loosely associated rebates which did not fit the specifications of the prompt were missed.

Acknowledgments

This work was supported by Dr. Dave Beck and Chino Ruttanasupagid. Special thanks to our stakeholders at the NC Clean Energy Technology Center

References

1. Klesel, M., & Wittmann, H. F. (2025). Retrieval-Augmented Generation (RAG). Business & Information Systems Engineering. <https://doi.org/10.1007/s12599-025-00945-3>
2. Wang, S., Khrantsova, E., Zhuang, S., & Zuccon, G. (2024). FeB4RAG: Evaluating Federated Search in the Context of Retrieval Augmented Generation. Proceedings of SIGIR 2024. doi:10.1145/3626772.3657853
3. Blázquez-Ochando, M., Prieto-Gutiérrez, J. J., & Ovalle-Perandones, M. A. (2025). Prompt engineering for bibliographic web-scraping. Scientometrics, 130(7), 3433–3453. <https://doi.org/10.1007/s11192-025-05372-5>