

Closing note

Task 1: RAG works perfectly in both scripts and added duckduckgo as options based sources for RAG retrieval Engine configurations for that in both scripts:

```
rag_article_retriever_engine="Duckduckgo",
```

Task 3: Reworked image handling in both scripts to fetch images from various sources. Configurations for that in both scripts:

```
image_source = "stock",
stock_primary_source = "pixabay",
secondary_source_image= True,
image_api=True,
huggingface_model="stabilityai/stable-diffusion-xl-base-1.0",
```

image_source = "stock" for stock sources

image_source = "imageai" for ai images

This feature toggles between stock sources and ai source

stock_primary_source - gives the user the chance to pick from any stock option

secondary_source_image - determines if secondary sources is needed if no image is found in primary source

Huggingface_model - change ai image generation model here

All image sources here

```
# Image sources api keys
unsplash_api_key=os.getenv('UNSPLASH_API_KEY',''),
pexels_api_key=os.getenv('PEXELS_API_KEY',''),
pixabay_api_key=os.getenv('PIXABAY_API_KEY',''),
giphy_api_key=os.getenv('GIPHY_API_KEY',''),
huggingface_api_key=os.getenv('HUGGINGFACE_API_KEY',''),
```

Change image captioning here

```
# Image captioning instance
image_caption_instance="openai/clip-vit-base-patch32",
```

Task 5: Improve text quality , structure of articles for SEO purpose and search engines (both scripts) now generation of table is now dynamic based on context

```

522
523 PARAGRAPH_WITH_HEADING_PROMPT = """You are an expert SEO content writer with over 10 years of experience creating high-ranking, reader-focused content for Fortune 500 companies and
524
525 Your task is to write paragraph {current_paragraph} of {paragraphs_per_section} for the section titled "{heading}" (Section {section_number} of {total_sections}). Target {articleaud
526
527 Requirements:
528 1. Write ONE cohesive paragraph (3-5 sentences, 80-120 words) focusing on a single main point.
529 2. Use {articlanguage} language, {voicetone} tone, and {pointofview} point of view.
530 3. Include specific details or actionable examples to engage {articleaudience}.
531 4. Integrate the primary keyword ({primary_keyword}) and LSI keywords ({lsi_keywords}) naturally, targeting 1-2% keyword density.
532 - Example LSI keywords for "SEO": "search engine ranking," "keyword research," "on-page optimization."
533 5. Use HTML tags for scannability and emphasis (no Markdown), selecting the most effective structure based on the content and trigger words (or their synonyms) in {current_points} c
534 - '<p>': Use for narrative or descriptive content. Trigger words: explain or describe or discuss or explore or introduce or define or elaborate or clarify or highlight or summar
535 - Example:
536   <h3>SEO Basics</h3>
537   <p>Exploring <strong>keyword research</strong> helps small businesses improve their <strong>search engine ranking</strong>. By identifying high-traffic terms, businesses can
538   <strong>: Emphasize 1-2 key phrases/keywords per paragraph.
539   <em>: Use 0-1 times for subtle emphasis.
540 - Example:
541   <h3>SEO Basics</h3>
542   <p>Keyword research is vital for <em>targeted content creation</em>, ensuring your site ranks for relevant terms. It involves analyzing search trends to optimize content effe
543   <ul><li>: Include 2-4 item unordered list for tips, benefits, or non-sequential items. Trigger words: tips or benefits or features or advantages or reasons or examples or ide
544   - Example:
545   <h3>SEO Basics</h3>
546   <p>Effective <strong>on-page optimization</strong> offers key benefits for small businesses: <ul><li>Optimized title tags improve click-through rates.</li><li>Clear meta desc
547   <ol><li>: Include 2-4 item ordered list for steps, processes, or sequences. Trigger words: steps or process or procedure or sequence or instructions or stages or phases or or
548   - Example:
549   <h3>SEO Basics</h3>
550   <p>To enhance <strong>SEO</strong>, follow this process for <strong>keyword research</strong>: <ol><li>Use tools like Google Keyword Planner to find relevant terms.</li><li>A
551   <table> with 'thead', 'tr', 'th', 'td': Use for comparisons, summaries, or structured data. Trigger words: compare or comparison or contrast or summarize or summary c
552   - Example:
553   <h3>SEO Basics</h3>
554   <p>A breakdown of <strong>SEO</strong> strategies clarifies their roles: <table><thead><tr><th>Strategy</th><th>Focus Area</th></tr></thead><tbody><tr><td>On-page SEO</td><td>
555 6. Start with a clear topic sentence and end with a conclusion or transition.
556 7. Ensure cohesion and alignment with {flow_instruction}.
557 8. Use '<h3>' for sections 1-2 or '<h4>' for deeper sections in the heading.
558 9. Maintain a professional, human-like tone, avoiding exaggeration.
559 10. Select the HTML structure based on trigger words or their synonyms in {current_points} or {flow_instruction}. If no trigger words or synonyms are present, default to '<p>'.
560
561 Section points to cover across {paragraphs_per_section} paragraphs:
562 {all_points}
563
564 Focus for this paragraph ({current_paragraph} of {paragraphs_per_section}):
565 {current_points}
566
567 Structure content to fit the section's flow: {flow_instruction}.
568

```

Approach to doing that was through prompt and lifting the ban on generation of tables and list as paragraph. Now the context of the discussion in the article determines the presentation pattern. Did that also in all the prompts including PAA.

Task 6: Improved image positioning across the entire article giving the user to determine is positioning based of images these are the options.

```

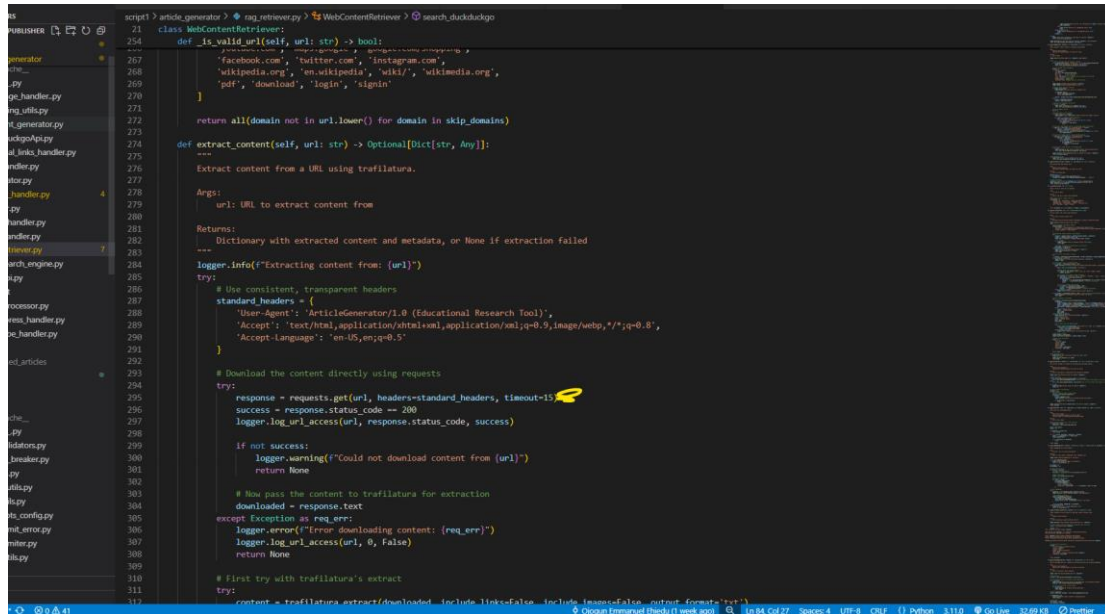
# positioning
youtube_position = "random", #random / first_heading / end
keytakeaways_position = "random", #random / before_conclusion / middle
image_position = "end", #random / under_first_heading /middle / end
paa_image_position = "random", # random / first_heading / end

```

Task 4: To improve generation time implement parallel body generation

1. Implement content generation in parallel: Script currently processes sections sequentially, split them into independent async tasks Use `asyncio.gather()` or multithreading where safe. Like blocknotes, conclusion should be generated at once.

2. Refactor I/O Layers: to Async with `aiohttp` Replace blocking HTTP requests with `aiohttp` for non-blocking concurrency. This reduces idle time while waiting on APIs or external services. Example of blocking requests is this

A screenshot of a code editor with a dark theme. The left sidebar shows a file explorer with various Python files. The main editor area displays a Python script for a class named 'WebContentRetriever'. The script includes methods for validating URLs, extracting content, and handling errors. A yellow arrow points to a 'timeout=15' parameter in a 'requests.get' call. The status bar at the bottom indicates the file is 'script1.py' and the editor is running Python 3.11.0.

```
21 class WebContentRetriever:
22     def __init__(self, url: str) -> bool:
23         """
24         Returns:
25         bool: True if the URL is valid, False otherwise
26         """
27         skip_domains = [
28             'facebook.com', 'twitter.com', 'instagram.com',
29             'wikipedia.org', 'en.wikipedia.org', 'wikimedia.org',
30             'pdf', 'download', 'login', 'signin'
31         ]
32         return all(domain not in url.lower() for domain in skip_domains)
33
34     def extract_content(self, url: str) -> Optional[Dict[str, Any]]:
35         """
36         Extract content from a URL using trafilatura.
37
38         Args:
39             url: URL to extract content from
40
41         Returns:
42             Dictionary with extracted content and metadata, or None if extraction failed
43         """
44         logger.info(f"Extracting content from: {url}")
45
46         # Use consistent, transparent headers
47         standard_headers = {
48             'User-Agent': 'ArticleGenerator/1.0 (Educational Research Tool)',
49             'Accept': 'text/html,application/xhtml+xml,application/xml;q=0.9,image/webp,*/*;q=0.8',
50             'Accept-Language': 'en-US,en;q=0.5'
51         }
52
53         # Download the content directly using requests
54         try:
55             response = requests.get(url, headers=standard_headers, timeout=15)
56             success = response.status_code == 200
57             logger.log_url_access(url, response.status_code, success)
58
59             if not success:
60                 logger.warning(f"Could not download content from {url}")
61                 return None
62
63             # Now pass the content to trafilatura for extraction
64             downloaded = response.text
65             except Exception as req_err:
66                 logger.error(f"Error downloading content: {req_err}")
67                 logger.log_url_access(url, 0, False)
68                 return None
69
70             # First try with trafilatura's extract
71             try:
72                 content = trafilatura.extract(downloaded, include_links=False, include_images=False, output_format='text')
73             except Exception as trafilatura_err:
74                 logger.error(f"Error extracting content with trafilatura: {trafilatura_err}")
75                 return None
76
77             return {'url': url, 'content': content}
```

They were used in several places across both scripts that will help too.

3. **Cache embedding:** this will help to automatically return similar embeddings that has been worked on before.
4. **Persistent Caching Between Runs:** Store previously generated sections or embeddings to disk (e.g., SQLite, Redis, or simple JSON files). On next run, load them quickly without recomputing.

Note: Major improvement comes from implementing 1 and 2 suggestion as with that every thing works at once and better handling of requests with this improvement there should be improvement in entire generation.