# End-to-End QA Generation for Thrive@Pitt: Robust Data Preparation, Quality Assurance, and Performance Monitoring

# Contents

# 1  Thrive@Pitt

The University of Pittsburgh's Thrive@Pitt initiative (`https://www.thrive.pitt.edu`) serves as a holistic platform for supporting student well-being and success. It centralizes mental health resources, physical wellness programs, academic support, financial guidance, and other valuable services to help students flourish during their time at the university. By bringing these resources together in one accessible location, Thrive@Pitt simplifies the process of connecting with the right support, fosters a culture of proactive well-being, and addresses the diverse needs of the Pitt community.

In practice, Thrive@Pitt:

- Creates a single digital portal where students, faculty, and staff can discover campus-wide well-being and support services.

- Encourages early intervention by highlighting resources before crises escalate, thereby bolstering student retention and mental health.

- Promotes an inclusive community culture that supports holistic development—academically, emotionally, and socially.

- Evolves continuously, adding new resources and tailoring its content based on community feedback and emerging wellness challenges.

# 2  Introduction

This project addresses the information accessibility challenge by building an end-to-end automated QA generation framework that produces high-quality question-answer pairs derived from the Thrive@Pitt website content. The system incorporates retrieval-augmented generation architecture, semantic document processing with entity recognition, enhanced web crawling with anti-bot detection, and multi-dimensional quality control—all optimized specifically for university support services. These question-answer pairs will be used to finetune and train SkillBuilder.io's chatbot model.

Building on the robust well-being ecosystem of Thrive@Pitt, the University of Pittsburgh sought to further enhance access to resources through an automated conversational assistant powered by Skill-Builder.io. While Thrive@Pitt already centralizes critical information, end-users often have specific, nuanced questions—ranging from "How do I apply for emergency funding?" to "What mental health services are available during midterms?" Enabling an AI-powered chatbot to answer such queries can bring immediate, user-friendly assistance to the community. However, constructing a chatbot that both understands complex domain-specific queries and provides accurate, context-rich responses is non-trivial. Careful data preparation and a robust question-answering (QA) pipeline are crucial for the chatbot's success. The quality of training data directly impacts model performance—better question-answer pairs lead to more accurate and helpful responses from the resulting chatbot, ultimately delivering more precise information about Thrive services, programs, and resources to the university community. A robust data pipeline must:

- **Extract relevant content** from Thrive@Pitt's pages using a specialized web crawler that respects site structure and avoids detection as a bot.

- **Segment and annotate content** into meaningful chunks for more precise retrieval and generation processes.

- **Generate realistic QA pairs** that cover the breadth of Thrive@Pitt's services, ensuring the chatbot can address varied queries.

- **Minimize hallucinations** and maintain factual consistency by leveraging advanced retrieval-augmented generation (RAG) workflows.

- **Adapt to new or updated content** on the Thrive@Pitt site so that the chatbot remains accurate over time.
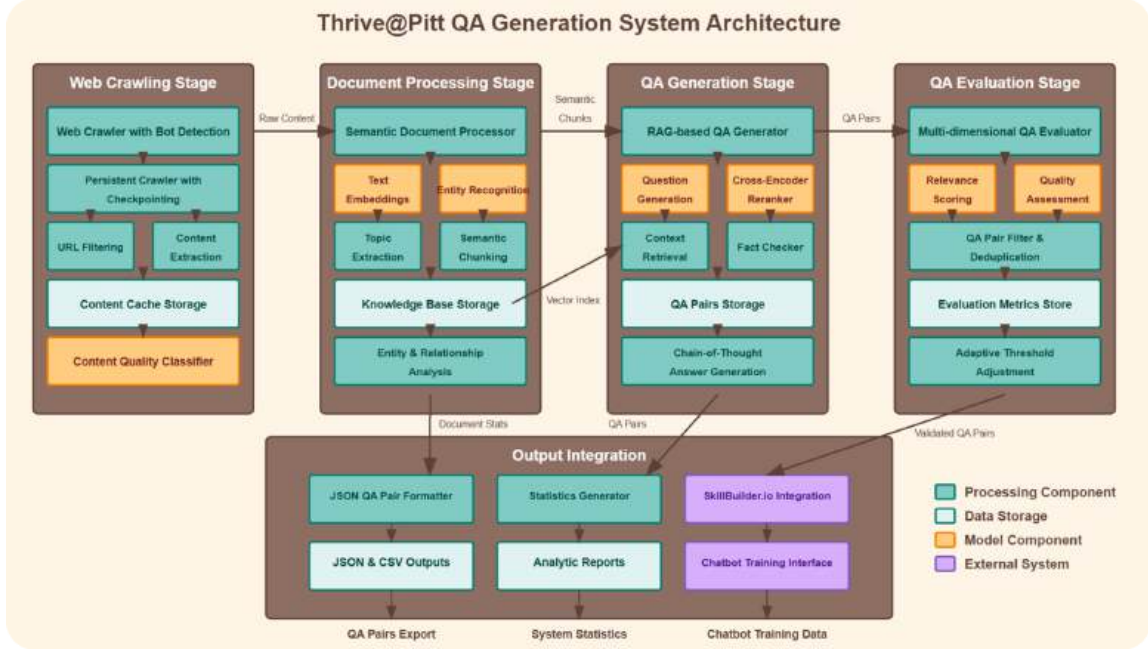
Figure 1: Thrive@Pitt QA Generation System Architecture. The diagram illustrates the complete pipeline flow across five main stages: (1) Web Crawling Stage with Bot Detection and Persistent Checkpointing; (2) Document Processing Stage with Entity Recognition and Semantic Chunking; (3) QA Generation Stage featuring RAG-based generation with Cross-Encoder Reranking and Fact Checking; (4) QA Evaluation Stage implementing Multi-dimensional Quality Assessment and Adaptive Thresholding; (5) Output Integration with formatters for JSON/CSV outputs and SkillBuilder.io integration. Color-coding distinguishes between Processing Components (teal), Data Storage (light teal), Model Components (orange), and External Systems (purple).

## 3    Project Overview

The project delivers an end-to-end automated application that employs a multi-stage approach to transform website content into curated question–answer (QA) pairs. Designed to support the University of Pittsburgh's websites, the system extracts and processes information, producing high-quality QA data optimized for training conversational agents. By integrating sophisticated AI framework with rigorous multi-dimensional evaluation metrics, the system ensures both the accuracy and practical relevance of the generated QA pairs, enhancing accessibility and engagement for users.

1. **Web Crawling:** A domain-optimized crawler systematically extracts textual content from Thrive@Pitt, ignoring duplicate or irrelevant assets and respecting robots.txt rules. To avoid bot-detection, it varies user-agent strings, adheres to realistic delays, and logs successful vs. failed attempts with checkpointing for resilience.

2. **Document Processing:** The extracted content is cleaned, normalized, and semantically chunked. Named entities and key phrases (e.g., organization names, student services) are identified to boost downstream retrieval accuracy.

3. **Retrieval-Augmented Generation (RAG):**

   - *Chunk-Level Embeddings:* Each text chunk is transformed into a vector representation, enabling semantic search and ranking.
   - *QA Generation:* An advanced language model is given a relevant subset of chunks and tasked with creating natural, domain-specific question–answer pairs.
   - *Consistency Checks:* The system employs cross-referencing techniques to confirm that generated answers match the retrieved text, reducing hallucinations and reinforcing factual correctness.

4. **Quality Evaluation and Filtering:** A multi-dimensional scoring method appraises each QA pair's relevance, completeness, factuality, and format. Lower-quality or near-duplicate pairs are discarded, ensuring the final dataset remains tightly aligned with Thrive@Pitt's official information.

5. **Output Integration:** Curated QA pairs, along with metadata (page source, timestamp, and quality metrics), are compiled into JSON or CSV formats. These can be ingested directly by chatbot frameworks (like SkillBuilder.io), enabling easy finetuning or augmentation of existing conversational models.

## 4  Key Benefits

1. **High-Fidelity Training Data.** By automatically extracting and verifying QA pairs from Thrive@Pitt content, the pipeline provides the chatbot with a rich, accurate dataset. This comprehensive coverage allows the chatbot to answer diverse user questions more reliably, without relying on manually curated or out-of-date information.

2. **Cross-Domain Applicability and Faster Deployment.** Beyond Thrive@Pitt, the QA generation system can be easily adapted to other university websites, rapidly producing curated QA data for diverse chatbot applications. This lowers the time and effort needed for manual content curation, accelerating the deployment of accurate, up-to-date support across different departmental or campus resources.

3. **Fewer Hallucinations and Errors.** The pipeline's cross-referencing and multi-step validation ensure that QA pairs match official site content. This reduces the risk of "hallucinated" or fabricated answers, improving the chatbot's trustworthiness and overall user satisfaction.

4. **Rapid Content Refresh and Adaptability.** Whenever Thrive@Pitt webpages are updated, the pipeline can re-crawl and re-process content. Newly extracted QA data is fed into the chatbot, ensuring it always reflects the latest well-being services, guidelines, and announcements.

5. **Modular Architecture for Future Enhancements.** The pipeline's design supports easy integration of newer advanced AI models or expanded retrieval techniques. This adaptability helps the chatbot remain state-of-the-art, whether adding advanced embeddings, cross-encoders, or fact-checking modules.

6. **Streamlined Data Integration and Maintenance.** Automated data cleaning, chunking, and QA generation minimize the need for manual oversight. The resulting QA pairs can be directly loaded into chatbot frameworks, reducing downtime and maintaining consistent, high-quality dialogues as new content appears.

7. **Improved Engagement and Resource Utilization.** By continually generating relevant, timely QA pairs, the pipeline equips the chatbot to direct users to the right services at the right time. This targeted guidance boosts user trust and satisfaction, encouraging more interactions with university wellness programs. In turn, it increases the overall utilization of resources by proactively connecting students, faculty, and staff with vital information before issues escalate.

## 5  System Architecture

The QA Generation System has a modular pipeline designed to automatically transform website content from any University of Pittsburgh website into a high-quality dataset of question-answer (QA) pairs. This dataset is specifically tailored for training conversational agents, to provide students with readily accessible information about university resources. Figure 1 provides a high-level overview of the system's architecture, illustrating the data flow and key components. The system comprises five primary stages: Web Crawling, Document Processing, QA Generation, QA Evaluation, and Output Integration. Each stage employs a combination of rule-based techniques and state-of-the-art neural models to ensure accuracy, relevance, and robustness.

### 5.1  Stage 1: Web Crawling

**Goal:** To collect relevant textual content from any university domain while adhering to ethical web crawling practices and mitigating the risk of bot detection.

This stage utilizes a custom-built web crawler, `PersistentCrawler`, which extends an `EnhancedWebCrawler` class. The crawler leverages libraries such as `requests` for HTTP communication and `BeautifulSoup` for HTML parsing. Key features include:

- **Bot Detection Avoidance:** The crawler employs several techniques to mimic human browsing behavior and avoid triggering anti-bot mechanisms:

  - *User-Agent Rotation:* A pool of modern user-agent strings (e.g., Chrome, Firefox, Safari) is randomly selected for each request.
  - *Request Delays:* Variable delays, following a distribution that simulates human reading times (implemented in `_get_random_delay()`), are introduced between requests. The base delay is configurable (defaulting to 2.0 seconds) and is further randomized.
  - *Referrer Headers:* Realistic `Referer` headers, mimicking navigation from previously visited pages, are included in requests.
  - *Session Management:* A `requests.Session` object with a persistent cookie jar (`cookiejar.LWPCookieJar`) is used to maintain state across requests, simulating a continuous browsing session.
  - *HEAD Requests:* Occasional HEAD requests are made before GET requests, mimicking browser behavior.
  - *Robots.txt Respect:* The crawler parses and adheres to the `robots.txt` file, respecting `Disallow` rules and `Crawl-delay` directives (implemented in `_check_robots_txt()`).

- **URL Filtering and Prioritization:** The crawler focuses on relevant content by:

  - *Domain Restriction:* Only URLs within the `thrive.pitt.edu` domain and its subdomains are crawled (implemented in `_is_valid_url()`).
  - *Content Type Filtering:* Non-HTML content (e.g., images, PDFs, documents) is ignored.
  - *Priority Paths:* A configurable list of URL path prefixes (`priority_paths`) allows prioritizing pages likely to contain valuable well-being information (e.g., "/services", "/resources").
  - *URL Normalization:* URLs are normalized to remove fragments and unnecessary query parameters, preventing redundant crawling.

- **Content Quality Assessment:** A heuristic scoring function (`_get_page_quality_score()`) evaluates the quality of each crawled page based on:

  - *Content Length:* Longer pages are generally favored, with penalties for very short pages.
  - *Presence of Meaningful Phrases:* The presence of domain-specific keywords and phrases related to well-being and university services (defined in `meaningful_phrases`) increases the score.
  - *Structural Cues:* The presence of HTML headings (detected during text extraction) is considered a positive indicator of structured content.
  - *Bot Detection Patterns:* The presence of text patterns indicative of bot detection or access denial (defined in `bot_patterns`) results in immediate rejection of the page.

  The quality score is a floating-point value between 0.0 and 1.0.

- **Checkpointing and Persistence:** The crawler's state is periodically saved to disk (implemented in `save_checkpoint()` and `load_checkpoint()`), allowing it to resume from interruptions without data loss. The checkpoint includes:

  - *Visited URLs:* A set of URLs that have been successfully crawled.
  - *Failed URLs:* A set of URLs that resulted in errors or were deemed invalid.
  - *Content Cache:* A dictionary mapping URLs to extracted and cleaned text content.
  - *Page Scores:* A dictionary mapping URLs to their calculated quality scores.

  Checkpoint files are stored in a configurable directory (`checkpoint_dir`, defaulting to "checkpoints") and are named using a unique identifier based on the base URL. A validation step (`validate_checkpoint()`) removes low-quality content from the cache upon loading.

The output of the Web Crawling stage is a *content cache*, a dictionary where keys are URLs and values are the extracted and cleaned text content of the corresponding pages. This cache, along with metadata such as quality scores and visit timestamps, serves as the input for the Document Processing stage.

## 5.2 Stage 2: Document Processing

**Goal:** To transform the raw, unstructured text extracted from web pages into a semantically enriched and structured representation suitable for downstream QA generation.

This stage employs a `DocumentProcessor` class, which utilizes a combination of rule-based techniques and neural models for text cleaning, semantic analysis, and chunking. The core data structures are `SemanticDocument` and `SemanticChunk`, representing entire web pages and coherent segments within them, respectively.

- **Text Normalization:** The raw text from each page is cleaned using the `_clean_text()` method. This involves:
    - *Whitespace Handling:* Excessive whitespace and newline characters are normalized.
    - *Unicode Correction:* Common Unicode character issues (e.g., incorrect apostrophes) are corrected.
    - *HTML Artifact Removal:* Remaining HTML entities and tags are removed.
    - *Section Header Handling:* Markdown-style section headers (e.g., "HEADING:", "TITLE:") are standardized.

- **Named Entity Recognition (NER):** The `analyze_entities()` method, managed by the `ResourceManager`, attempts to identify named entities (e.g., organizations, people, dates, locations) within the text. It prioritizes spaCy (`en_core_web_md` or `en_core_web_sm`) for high-quality NER. If spaCy is unavailable or fails, it falls back to NLTK's named entity chunker, and finally to a set of regular expressions for basic entity extraction. Extracted entities are stored with their text, label, and character offsets.

- **Key Phrase Extraction:** The `extract_key_phrases()` method, also managed by the `ResourceManager`, identifies important phrases and topics within the document. It primarily uses TF-IDF (Term Frequency-Inverse Document Frequency) to score n-grams (1-3 words). Stop words are removed, and a maximum number of features are considered. If TF-IDF fails, it falls back to extracting noun phrases using spaCy (if available) or, as a last resort, extracting capitalized phrases using regular expressions.

- **Semantic Chunking:** The `_create_semantic_chunks()` method divides the document into smaller, contextually coherent units (chunks). The chunking strategy adapts based on the document's structure:
    - *Section-Based Chunking:* If the document contains markdown-style section headers ("" or ""), it is split into chunks based on these headers (`_chunk_by_sections()`).
    - *Paragraph-Based Chunking:* If no section headers are found, the document is split into chunks based on paragraph breaks (""). A target chunk length (approximately 1000 characters) is maintained, with some overlap between adjacent chunks to preserve context (`_chunk_by_paragraphs()`).
    - *Sentence-Based Chunking:* If paragraph-based chunking results in very large chunks (greater than 1500 characters), a sentence-based splitting with overlap is applied (`_split_large_chunk()`). This uses a robust, fallback sentence tokenizer (`get_fallback_tokenize()`) that handles various edge cases and prioritizes NLTK, then spaCy, and finally a regex-based approach.

    Each `SemanticChunk` object stores the chunk's text, the source document URL, a unique index, and any extracted entities.

- **Embedding Generation:** Each chunk is converted into a dense vector representation (embedding) using a pre-trained sentence transformer model. The `ModelManager` class handles loading and unloading models to manage GPU memory efficiently. It prioritizes `all-mpnet-base-v2` for high-quality embeddings, falling back to smaller models (e.g., `all-MiniLM-L12-v2`, `all-MiniLM-L6-v2`) if necessary. The embedding dimension depends on the chosen model (e.g., 768 for `all-mpnet-base-v2`). Embeddings are stored as PyTorch tensors.

The output of the Document Processing stage is a `KnowledgeBase` object. This object stores:

- **documents:** A dictionary mapping URLs to `SemanticDocument` objects. Each `SemanticDocument` contains the full text, title, extracted entities, topics, and a list of `SemanticChunk` objects.

- **chunks:** A list of all `SemanticChunk` objects across all documents.

- **embeddings:** A PyTorch tensor containing the embeddings for all chunks, used for efficient similarity search.

- **chunk_text_index:** A simple text-based index mapping terms to chunk indices, used as a fallback for retrieval.

The `KnowledgeBase` provides methods for adding documents, rebuilding the search index, and performing semantic searches using cosine similarity between query embeddings and chunk embeddings. It also includes a fallback text-based search for cases where embeddings are unavailable.

### 5.3 Stage 3: QA Generation

**Goal:** To generate high-quality, contextually relevant question-answer pairs based on the processed documents and chunks in the knowledge base.

This stage utilizes a `QAGenerator` class, which implements a Retrieval-Augmented Generation (RAG) approach, combining neural language models with a retrieval mechanism to ground generated QA pairs in the source content.

- **Model Loading:** The `QAGenerator` loads several pre-trained models using the `ModelManager`:
  - *Embedding Model:* Used for retrieving relevant chunks (prioritizes `all-mpnet-base-v2`).
  - *Question Generation Model:* A T5-based model (Flan-T5-XL/XXL/Large/Base, with fallback to smaller models) used to generate questions based on context.
  - *Answer Generation Model:* A T5-based model (often the same as the question generation model) used to generate answers.
  - *Fact Checker (Optional):* A model (`vectara/hallucination_evaluation_model`) to assess the factual consistency of generated answers.
  - *Cross-Encoder Re-ranker:* A model (`cross-encoder/ms-marco-MiniLM-L-12-v2` or `cross-encoder/ms-marco-MiniLM-L-6-v2`) to re-rank retrieved chunks based on their relevance to the question.
  - *QA Evaluator (Optional):* A model (`cross-encoder/ms-marco-MiniLM-L-6-v2`) to score the relevance of QA pairs.

- **Question Generation:** The `generate_questions_from_documents()` method generates questions for a given list of URLs. It prioritizes generating questions related to key topics extracted from the documents. For each topic, it calls `generate_questions_for_topic()`, which uses a combination of neural and rule-based approaches:
  - *Neural Question Generation (`_generate_questions_neural()`):* Uses the question generation model (Flan-T5) with a variety of prompts designed to elicit diverse and natural questions. Prompts are tailored based on the identified topic type (e.g., SERVICE, PROGRAM, LOCATION, WELLNESS). Chain-of-thought prompting is also used to encourage more sophisticated question generation.
  - *Rule-Based Question Generation (`_generate_questions_rule_based()`):* Used as a fallback if neural generation fails or produces insufficient questions. This method generates questions based on sentence patterns and topic types.
  - *General Question Generation (`generate_general_questions()`):* Generates general questions about the document, not tied to specific topics.
  - *Adaptive Question Generation (`_generate_adaptive_questions()`):* Used when very few questions are generated, providing more general questions that are less dependent on specific content.

  Generated questions are deduplicated (`_deduplicate_questions()`) and cleaned (`_clean_question()`).

- **Context Retrieval:** The `_retrieve_context()` method retrieves relevant chunks from the knowledge base for a given question. It uses a hybrid approach:

- - *Vector Similarity Search:* Calculates the cosine similarity between the question embedding and the embeddings of all chunks in the knowledge base.
  - *Cross-Encoder Re-ranking:* The top-k retrieved chunks are re-ranked using a cross-encoder model, which provides a more accurate measure of relevance between the question and each chunk.
  - *Diversity Selection:* A greedy algorithm (`_select_diverse_chunks()`) selects a diverse set of chunks, prioritizing those from different source documents, to provide a broader context for answer generation.
  - *Fallback Text Search:* If vector search fails or returns too few results, a fallback text-based search (`_fallback_search()`) is used.

- **Answer Generation:** The `_generate_answer()` method generates an answer for a given question and its retrieved context. It uses the answer generation model (Flan-T5) with a detailed prompt that instructs the model to:

  - Synthesize information from multiple context chunks.
  - Maintain a confident and direct tone.
  - Avoid disclaimers ("I don't have enough information") if any relevant information is available.
  - Provide only information that is supported by the context.

  The method uses beam search for answer generation and selects the best answer from multiple candidates based on length and factual consistency. The `_merge_context_advanced()` method combines the retrieved chunks into a single context string, with clear markers indicating the source document of each chunk.

- **Fact Checking (Optional):** If a fact-checking model is available, it is used to assess the factual consistency of the generated answer with the retrieved context.

The output of the QA Generation stage is a list of `qa_pairs`, where each pair is a dictionary containing the question, answer, source URL, topic, topic type, and initial quality scores.

### 5.4 Stage 4: QA Evaluation

**Goal:** To assess the quality of the generated QA pairs and filter out low-quality or redundant pairs.

This stage evaluates each QA pair based on multiple criteria and assigns an overall quality score. The `QAGenerator` class's `_evaluate_answer()` method performs this evaluation.

- **Multi-Dimensional Scoring:** Each QA pair is scored based on the following metrics:
  - *Relevance (R):* Measures how well the answer aligns with the question. It combines a model-based relevance score (using a cross-encoder or the `qa_evaluator` model) with a keyword overlap score:
    $$R = 0.7 \times (\text{model-based relevance}) + 0.3 \times (\text{keyword match}).$$
  - *Factuality (F):* Assesses how faithfully the answer reflects the source text. It uses n-gram overlap between the answer and the retrieved context. An optional fact-checking model can override this score:
    $$F = 0.2 + 0.8 \times \frac{|\text{overlap n-grams}|}{|\text{answer n-grams}|}.$$
  - *Completeness (C):* Checks if the answer addresses all aspects of the question. Short or truncated answers receive lower scores. $C \in [0, 1]$.
  - *Formatting ($\Phi$):* Ensures the answer is free of formatting issues, disclaimers, and repetitive text. Each major formatting violation lowers the score. $\Phi \in [0, 1]$.

- **Overall Score ($O$):** A weighted sum of the individual metrics:
  $$O = 0.25R + 0.35F + 0.25C + 0.15\Phi.$$

- **Filtering:** QA pairs with an overall score below a threshold (defaulting to 0.5) are discarded. The threshold can be adjusted adaptively based on the overall quality distribution of the generated pairs.

- **Deduplication:** Near-duplicate QA pairs are identified and removed using a combination of Jaccard similarity (on word sets) and Levenshtein distance (edit distance). The `_questions_are_similar()` method implements this check. The `_deduplicate_questions()` method removes duplicate questions, and a similar approach is used for answers.

The `filter_qa_pairs()` method combines the scoring, filtering, and deduplication steps. It also includes an adaptive threshold adjustment, lowering the minimum score if very few QA pairs meet the initial threshold.

### 5.5 Stage 5: Output Integration

**Goal:** To package the filtered and evaluated QA pairs into various formats suitable for downstream use, such as training conversational AI agents.

This stage prepares the final QA dataset and provides analytics about its content and quality.

- **Data Export:** The QA pairs are exported in multiple formats:
  - *JSON:* A JSON file containing the QA pairs, along with metadata (source URL, topic, scores).
  - *CSV:* A CSV file for easy inspection and analysis in spreadsheet software.
  - *Markdown:* A Markdown file for human-readable documentation.

- **Analytics and Reporting:** The system generates summary statistics about the QA dataset, including:
  - *Total QA Pairs:* The number of QA pairs in the final dataset.
  - *Average Scores:* The average relevance, factuality, completeness, formatting, and overall scores.
  - *Score Distribution:* The number of QA pairs in different quality categories (e.g., excellent, good, average, below average).
  - *Topic Coverage:* The distribution of QA pairs across different topics.
  - *URL Coverage:* The distribution of QA pairs across different source URLs.

- **Chatbot Integration:** The generated QA pairs are designed to be directly integrated into conversational AI platforms, such as SkillBuilder.io. The JSON format is particularly suitable for this purpose. The data can also be used for fine-tuning retrieval-augmented generation (RAG) models or other language models.

The output of this stage is a set of files containing the curated QA pairs in various formats, along with a report summarizing the dataset's characteristics.

## 6 Technical Implementation

The Thrive@Pitt QA Generation System is implemented in Python, leveraging a range of state-of-the-art libraries and frameworks for natural language processing, web crawling, and machine learning. Key components include:

- **Core NLP Libraries:**
  - *Hugging Face Transformers:* Provides access to pre-trained language models (e.g., Flan-T5, BERT) and tokenizers.
  - *spaCy:* Used for named entity recognition (NER) and sentence boundary detection.
  - *NLTK:* Used for tokenization, stop word removal, and as a fallback for NER.
  - *SentenceTransformers:* Used for generating sentence and document embeddings.
  - *BeautifulSoup:* Used for parsing HTML content.

- **Web Crawling:**
  - *requests:* Used for making HTTP requests.
  - *BeautifulSoup:* Used for parsing HTML content.
  - *urllib.parse:* Used for URL parsing and manipulation.

- **QA Generation and Evaluation:**

  - *Flan-T5 (XL/XXL/Large/Base):* Used for question generation and answer generation.
  - *Cross-Encoder Re-rankers (ms-marco-MiniLM-L-12-v2, ms-marco-MiniLM-L-6-v2):* Used for re-ranking retrieved context chunks.
  - *Fact Checker (vectara/hallucination_evaluation_model):* Used for evaluating the factual consistency of answers.

- **Resource Management:**

  - *Dynamic Model Loading/Unloading:* The `ModelManager` class dynamically loads and unloads models based on memory availability, allowing the system to operate efficiently even with limited GPU resources.
  - *Quantization:* Support for 4-bit quantization (using `BitsAndBytesConfig`) is included to reduce the memory footprint of large models.
  - *Fallback Mechanisms:* The system includes fallback mechanisms for various components (e.g., using smaller models, alternative NLP libraries, rule-based methods) to ensure robustness.

- **Error Handling and Resilience:**

  - *Comprehensive Exception Handling:* `try-except` blocks are used throughout the pipeline to handle potential errors gracefully.
  - *Automatic Checkpointing:* The system periodically saves its state (crawled content, processed documents, generated QA pairs) to allow for recovery from interruptions.
  - *Logging:* Detailed logging is implemented using the Python `logging` module to track progress and diagnose issues.
  - *Backoff Strategies:* The `backoff` library is used to implement retry mechanisms for network requests, handling temporary network issues.

The system is designed to be modular and extensible, allowing for easy integration of new models, algorithms, and evaluation metrics.

# 7 Concerns

Below, the project identifies the primary concerns that require careful attention when deploying the chatbot platform. Each concern highlights key questions to pose to SkillBuilder.io, followed by an explanation of why it matters.

## 7.1 Hallucination Management

**Key Questions**

- How will the platform ensure that Large Language Model (LLM) responses remain grounded in curated data from university website?

- What mechanisms (e.g., citation or reference checks) exist to detect or prevent hallucinations *in real time*?

**Why It Matters:** Hallucinations occur when an LLM fabricates information not supported by its training data Ji et al. (2023); OpenAI (2023), posing a serious risk of misinformation about Thrive@Pitt services or unverified health recommendations. Ensuring the chatbot remains factually grounded preserves user trust and fosters engagement.

## 7.2 Data Updates and Exclusions

**Key Questions**

- How does the platform prevent outdated or invalid embeddings from influencing new responses after the model parameters have been fine-tuned on older information?

- Is there version control (or a similar mechanism) that tracks changes to the data over time?

- How does the platform incorporate new information—through fine-tuning or a RAG implementation—to update the knowledge base when outdated data is removed and replaced?

**Why It Matters:** University periodically updates webpages or retires outdated content. If the chatbot fails to synchronize with these changes, users could receive irrelevant or contradictory information. Proactive and transparent versioning ensures that new data supersedes the old, enabling reliable, current answers.

### 7.3 Toxicity Filtering / Adversarial questions / Guard Rails

**Key Questions**

- How does the platform handle incoming user inputs that may contain hate speech, harassment, or otherwise harmful text?

- How is training data filtered for toxicity before being used for fine-tuning the model?

- What remediation steps are triggered if content is flagged as toxic or unsafe?

**Why It Matters:** The inadvertent addition of toxic content to training data can influence the model's parameter updates, leading to undesired biases or harmful outputs. Also, chatbots sometimes receive abusive or harmful text from users. Without robust detection and filtering, the chatbot could inadvertently propagate or respond to toxic language in a way that conflicts with campus policy.

### 7.4 Performance Monitoring and Accountability

**Key Questions**

- What reporting tools are available to measure the chatbot's effectiveness (e.g., user satisfaction, resolved inquiries)?

- Can the platform log failed queries and flag them for manual review or data updates?

**Why It Matters:** Ongoing performance tracking is crucial for maintaining quality and relevance. By logging user interactions and identifying unanswered or poorly answered queries, administrators can continuously optimize the chatbot's data and responses. This iterative improvement process ensures the chatbot remains an accurate, up-to-date resource for students, staff and faculty.

## 8 Sample QA Pairs and Evaluation Results

To demonstrate the effectiveness of our QA generation system, this section presents a sample of 30 question-answer pairs extracted from the full dataset of more than 300 question-answer pairs. These examples showcase the system's ability to generate contextually relevant, factually accurate, and well-formed QA pairs across various aspects of University's resources and services.

## 8.1 Sample of the Generated QA Pairs

| Question | Answer |
|---|---|
| What is the Wellness Concierge Program? | WCP is a new initiative that offers one-on-one meetings with trained Peer Wellness Coaches (PWC). PWCs are fellow students available to support you and help you navigate your wellness jouey. The program aims to foster a community where Pitt students can grow and thrive. What to Expect from the Wellness Concierge Program? Prior to meeting the first time, you'll engage in activities and discussions with your PWC to help identify personal wellness needs, interests, and goals. |
| What is the Campus Well-Being Consortium at Pitt? | The Campus Well-Being Consortium at Pitt is a collective of faculty and staff members, student leaders and community members working together to foster an environment where students can thrive. Thrive at Pitt was formed to advance our mission: To support and maintain all aspects of student well-being by promoting and improving accessibility to a range of resources relevant to our diverse student body. |
| What is the mission of the Campus Well-Being Consortium at Pitt? | The University of Pittsburgh is committed to creating a community where all students can thrive – both as students and beyond. In 2020, a dedicated group of campus community members joined together to create the Campus Well-Being Consortium at Pitt. Thrive at Pitt was formed to advance our mission: To support and maintain all aspects of student well-being by promoting and improving accessibility to a range of resources relevant to our diverse student body. |
| What are the factors that comprise a pillar of well-being? | Well-being is multidimensional, which means it encompasses numerous aspects (or pillars) of our life. Each pillar is comprised and impacted by many individual factors (e.g. sleep, diet, mindset, exercise, relationships, etc.). When we tend to these factors through daily routines, goal setting and practice, the results can be very satisfying and rewarding. When certain parts are left out, or neglected, the outcome is likely not as good. In exploring, supporting, or enhancing your own well-being it is important to become familiar with, and give proper attention to, each pillar. |
| What are the pillars of well-being? | Well-being is multidimensional, which means it encompasses numerous aspects (or pillars) of our life. Each pillar is comprised and impacted by many individual factors (e.g. sleep, diet, mindset, exercise, relationships, etc.). When we tend to these factors through daily routines, goal setting and practice, the results can be very satisfying and rewarding. When certain parts are left out, or neglected, the outcome is likely not as good. In exploring, supporting, or enhancing your own well-being it is important to become familiar with, and give proper attention to, each pillar. |

| Question | Answer |
| --- | --- |
| How many pillars of well-being are there? | Well-being is multidimensional, which means it encompasses numerous aspects (or pillars) of our life. Each pillar is comprised and impacted by many individual factors (e.g. sleep, diet, mindset, exercise, relationships, etc.). When we tend to these factors through daily routines, goal setting and practice, the results can be very satisfying and rewarding. When certain parts are left out, or neglected, the outcome is likely not as good. In exploring, supporting, or enhancing your own well-being it is important to become familiar with, and give proper attention to, each pillar. |
| What is the Counseling Center? | The mission of the University Counseling Center (UCC) is to provide high quality services and programs that foster the academic and personal development, as well as the psychological well-being, of Pitt's diverse student body. From outreach and prevention to crisis intervention, the UCC is dedicated to supporting and empowering students to maximize their college experience and realize their fullest potential. Drop-In Services Call 412-648-7930 Group Counseling Let's Talk Sessions are drop-in format that offer informal, confidential consultation with a clinician. Consultations are typically 10-15 minutes long and are offered on a first-come, first serve basis. |
| What is the mission of the Campus Well-Being Consortium? | The Campus Well-Being Consortium is a collective of faculty and staff members, student leaders and community members working together to foster an environment where students can thrive. Thrive at Pitt was formed to advance our mission: To support and maintain all aspects of student well-being by promoting and improving accessibility to a range of resources relevant to our diverse student body. |
| What is well-being? | Highlights What is Well-Being? Well-being can be defined as an optimal and dynamic state that allows people to achieve their full potential both as an individual and a member of the greater community. Lea about well-being. Build Your Foundation Each pillar of well-bearing serves as a foundational block to support your overall well- being. Together they create a solid foundation that allows you to thrive. Explore the pillars of wellbeing. |
| What is the Thrive at Pitt website? | The University of Pittsburgh is committed to creating a community where all students can thrive – both as students and beyond. Thrive at Pitt was formed to advance our mission: To support and maintain all aspects of student well-being by promoting and improving accessibility to a range of resources relevant to our diverse student body. |

| Question | Answer |
|---|---|
| What are some of the pillars that comprise the well-being? | Well-Being Well-being is multidimensional, which means it encompasses numerous aspects (or pillars) of our life. Each pillar is comprised and impacted by many individual factors (e.g. sleep, diet, mindset, exercise, relationships, etc.). When we tend to these factors through daily routines, goal setting and practice, the results can be very satisfying and rewarding. When certain parts are left out, or neglected, the outcome is likely not as good. In exploring, supporting, or enhancing your own well-being it is important to become familiar with, and give proper attention to, each pillar. Independently, every pillar plays an essential role for supporting our well-bearing. Collectively, they provide a solid foundation that accounts for and holds up your well-well-being. |
| What is the focus of the RISE program? | The Reaching Inside Your Soul for Excellence (RISE) program helps students get acclimated to the campus and help develop necessary skills needed to succeed at the University of Pittsburgh. Students in the program receive a mentor that will work with them to take responsibility of their future, make meaningful connections, and much more during their time at Pitt. |
| Where can I find the Office of Cross Cultural and Leadership Development? | The Office of Cross Cultural and Leadership Development (CCLD) is located on the 6th floor of the William Pitt Union and is a hub of campus activity. CCLD advises some of the most active student organizations on campus, provides leadership development opportunities, and offers dynamic programming. |
| What was the Campus Well-Being Consortium created? | The Campus Well-Being Consortium is a collective of faculty and staff members, student leaders and community members working together to foster an environment where students can thrive. Thrive at Pitt was formed to advance our mission: To support and maintain all aspects of student well-being. |
| What are the main factors that impact our well-being? | Well-being influences your ability to flourish and thrive. While we all come from various backgrounds, have different experiences, and face distinct challenges, leaing to improve and maintain the various areas of your well-being can help you overcome obstacles, achieve goals, and lead a meaningful and fulfilling life. Highlights What is Well-Being? |
| What are the individual pillars? | The pillars of well-being are comprised and impacted by many individual factors (e.g. sleep, diet, mindset, exercise, relationships, etc.). When we tend to these factors through daily routines, goal setting and practice, the results can be very satisfying and rewarding. When certain parts are left out, or neglected, the outcome is likely not as good. In exploring, supporting, or enhancing your own wellbeing it is important to become familiar with, and give proper attention to, each pillar. |
| What are the physical well being services? | Physical well-being is the knowledge and practice toward positive health conditions. It involves the adoption of proactive measures such as regular physical activity, good nutrition, utilization of preventative health care, and abstaining from harmful behaviors such as alcohol consumption, tobacco and drug use. |

| Question | Answer |
|---|---|
| What is the best version of yourself? | To thrive is to continually pursue the best version of yourself. Your well-being is impacted in many ways, and giving attention to each area is essential to ensuring your basic needs are met. In addition, continual focus on improving your well-bearing can lead to numerous positive benefits. For example, individuals who maintain higher levels of well- being are more likely to: Have a more positive outlook on life and overall better mental health Have greater confidence and higher levels. |
| How does well-being measure up to our lives? | Well-being is broad and multifaceted with correlations to numerous aspects of our life. It is impacted by many factors, both psychological and physical. These factors not only impact our satisfaction with life, they also influence our desire or ability to pursue and achieve our individual and collective capabilities. Thus, when we dedicate necessary time and effort to each area of well-being, we increase our chances for success, happiness, and a life defined as thriving. |
| What is the definition of well-being? | Well-being can be defined as an optimal and dynamic state that allows people to achieve their full potential both as an individual and a member of the greater community. Lea about well-being. Build Your Foundation Each pillar of well-bearing serves as a foundational block to support your overall well- being. Together they create a solid foundation that allows you to thrive. Explore the pillars of well being. Wellness Concierge Program Personalized, peer-led support through one-on-one meetings to assist you on your well- Being jouey. Some Inspiration My mission in life is not merely to survive, but to thrive; and to do so with some passion, some humor, and some style. Maya Angelou What Does It Mean to Thrive? To thrive is to continually pursue the best version of yourself. Through continued maintenance of your well being, you can build a Foundation to help you thrive in the face of challenges and become your best self! |
| What is the foundation for well-being? | Highlights Well-being is an optimal and dynamic state that allows people to achieve their full potential both as an individual and a member of the greater community. To thrive is to continually pursue the best version of yourself. Through continued maintenance of your well-being, you can build a foundation to help you thrive in. |
| What does pursuing the best version of yourself benefit you? | To thrive is to continually pursue the best version of yourself. Through continued maintenance of your well-being, you can build a foundation to help you thrive in the face of challenges and become your best self! |
| What does environmental stewardship encompass? | Environmental stewardship involves adopting and advocating for efforts that ensure responsible and efficient resource management while increasing access to clean air and water; nutritious food; and adequate clothing, shelter, and safety in both the built and natural environments at a variety of places and scales (e.g., home, school, work, community, city, country, world). |

| Question | Answer |
|---|---|
| What is the Office of Cross Cultural and Leadership Development? | The Office of Cross Cultural and Leadership Development (CCLD) is located on the 6th floor of the William Pitt Union and is a hub of campus activity. CCLD advises some of the most active student organizations on campus, provides leadership development opportunities, and offers dynamic programming in the areas of diversity and culture, social justice, and Frateity and Sorority Life. |
| How can you improve your spiritual well-being? | Well-being influences your ability to flourish and thrive. While we all come from various backgrounds, have different experiences, and face distinct challenges, leaing to improve and maintain the various areas of your well-being can help you overcome obstacles, achieve goals, and lead a meaningful and fulfilling life. Highlights What is Well-Being? |
| What are some tips for physical well-being? | Environmental well-being is the knowledge and practice toward positive health conditions. It involves the adoption of proactive measures such as regular physical activity, good nutrition, utilization of preventative health care, and abstaining from harmful behaviors such as alcohol consumption, tobacco and drug use. 1 5 Tips for Environmental Well-Being Engage in at least 150 minutes of physical activity each week. 2, 3, 4 Set and maintain a healthy dietary patte. 5 Abstain from using alcohol, tobacco, and other drugs/substances. 6 Get 7-8 hours of sleep each night. 7 Schedule routine check-ups with a primary care physician (PCP) and other doctors. |
| What are some ways to achieve physical well-being? | Well-being is broad and multifaceted with correlations to numerous aspects of our life. It is impacted by many factors, both psychological and physical. These factors not only impact our satisfaction with life, they also influence our desire or ability to pursue and achieve our individual and collective capabilities. Thus, when we dedicate necessary time and effort to each area of well-being, we increase our chances for success, happiness, and a life defined as thriving. |
| What is physical well-being? | Physical well-being is the knowledge and practice toward positive health conditions. It involves the adoption of proactive measures such as regular physical activity, good nutrition, utilization of preventative health care, and abstaining from harmful behaviors such as alcohol consumption, tobacco and drug use. |

| Question | Answer |
|---|---|
| What is the main idea of well-being? | Highlights What is Well-Being? Well-being can be defined as an optimal and dynamic state that allows people to achieve their full potential both as an individual and a member of the greater community. Lea about well-being. Build Your Foundation Each pillar of well-bearing serves as a foundational block to support your overall well- being. Together they create a solid foundation that allows you to thrive. Explore the pillars of well being. Wellness Concierge Program Personalized, peer-led support through one-on-one meetings to assist you on your well- well being jouey. Some Inspiration My mission in life is not merely to survive, but to thrive; and to do so with some passion, some humor, and some style. Maya Angelou What Does It Mean to Thrive? To thrive is to continually pursue the best version of yourself. Through continued maintenance of your well being, you can build a Foundation to help you thrive in the face of challenges and become your best self! |

## 8.2 QA Pair Structure Explanation

To illustrate the structure of our QA pairs, let's examine an example from the dataset:

- **Question:** "What is the Wellness Concierge Program?"

- **Answer:** "WCP is a new initiative that offers one-on-one meetings with trained Peer Wellness Coaches (PWC). PWCs are fellow students available to support you and help you navigate your wellness jouey. The program aims to foster a community where Pitt students can grow and thrive. What to Expect from the Wellness Concierge Program? Prior to meeting the first time, you'll engage in activities and discussions with your PWC to help identify personal wellness needs, interests, and goals"

- **Source URL:** https://thrive.pitt.edu/wellness-concierge-program

- **Topic:** "wellness"

- **Topic Type:** "WELLNESS"

- **Scores:**
    - *Relevance:* 0.9997 - Indicates nearly perfect relevance of the answer to the question
    - *Factuality:* 0.9787 - Indicates high factual accuracy with respect to source content
    - *Completeness:* 0.8000 - Indicates good but not exhaustive coverage of all aspects
    - *Formatting:* 1.0000 - Indicates perfect formatting with no artifacts or issues
    - *Overall:* 0.9425 - The weighted combination of all metrics, showing excellent quality

This QA pair exemplifies the high standards achieved by our generation system, with strong relevance, factuality, and appropriate formatting. The answer directly addresses the question and provides valuable information about the Wellness Concierge Program's purpose and process.

## 8.3 Evaluation Metrics Analysis

The evaluation metrics output provide insights into the overall quality and characteristics of the generated QA pairs:

| Metric | Value |
|---|---|
| Total QA Pairs | 225 |

Table 2: Total QA Pairs Generated

| Quality Metric | Average Score |
|----------------|---------------|
| Relevance | 0.547 |
| Factuality | 0.957 |
| Completeness | 0.665 |
| Formatting | 1.000 |
| Overall | 0.788 |

Table 3: Average Quality Scores

| Quality Tier | Count |
|--------------|-------|
| Excellent (Overall $\geq$ 0.8) | 125 |
| Good ($0.7 \leq$ Overall $< 0.8$) | 40 |
| Average ($0.6 \leq$ Overall $< 0.7$) | 58 |
| Below Average (Overall $< 0.6$) | 2 |

Table 4: Quality Distribution

- **Overall Quality:** This is a weighted combination of all metrics, reflecting the comprehensive quality of each QA pair. A higher overall score indicates a better-quality QA pair. We categorize pairs as "Excellent" (0.8+), "Good" (0.7-0.8), "Average" (0.6-0.7), or "Below Average" (<0.6).

- **Relevance:** This metric evaluates how directly the answer addresses the question posed. Higher relevance scores indicate answers that provide information specifically pertinent to the question, without extraneous or tangential content.

- **Factuality:** This measures how accurately the answer reflects information contained in the source material. High factuality scores indicate minimal hallucination or fabrication, with content faithfully representing the source material.

- **Completeness:** This assesses whether the answer covers all aspects of the question comprehensively. Higher completeness scores indicate answers that address multiple dimensions of the question without significant omissions.

- **Formatting:** This evaluates the structural quality of the answer, including paragraph breaks, lists, and other formatting elements. Perfect formatting scores indicate answers that are well-structured and easy to read.

- **Topic Coverage:** We analyze the distribution of topics among the QA pairs to ensure comprehensive coverage across key domains, including well-being concepts, university initiatives, and available resources.

- **URL Distribution:** This examines the spread of source URLs used to generate answers, ensuring content is drawn from diverse sections of the website.

These metrics collectively help us assess the system's strengths and identify areas for improvement.


# 9 Code

```
1  import requests
2  from bs4 import BeautifulSoup
3  import json
4  from typing import List, Dict, Tuple, Set, Any, Optional, Union, Generator, Callable
5  from transformers import (
6      AutoTokenizer, AutoModelForSeq2SeqLM, AutoModelForQuestionAnswering,
7      AutoModelForSequenceClassification, T5ForConditionalGeneration,
8      pipeline, BitsAndBytesConfig
9  )
10 import torch
11 from torch.utils.data import Dataset, DataLoader
12 import nltk
13 from nltk.tokenize import sent_tokenize
14 from nltk.corpus import stopwords
```

```python
15  import spacy
16  import logging
17  import os
18  import sys
19  import re
20  import gc
21  import time
22  from tqdm import tqdm
23  from datetime import datetime
24  import backoff
25  from concurrent.futures import ThreadPoolExecutor, as_completed
26  import hashlib
27  import traceback
28  import signal
29  from urllib.parse import urljoin, urlparse
30  import random
31  from http import cookiejar
32  from sklearn.feature_extraction.text import TfidfVectorizer
33  from sklearn.metrics.pairwise import cosine_similarity
34  from sentence_transformers import SentenceTransformer, CrossEncoder, util
35  import argparse
36  import copy
37  import numpy as np
38  from collections import Counter
39  from itertools import combinations
40  import math
41  import uuid
42
43  # Configure GPU usage
44  device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
45  print(f"Using device: {device}")
46
47  # Set up logging
48  logging.basicConfig(
49      level=logging.INFO,
50      format='%(asctime)s - %(levelname)s - %(message)s',
51      handlers=[logging.FileHandler('qa_generator.log'), logging.StreamHandler()]
52  )
53  logger = logging.getLogger(__name__)
54
55  #----------------------------------------------------------------------
56  # Resource Setup and Initialization
57  #----------------------------------------------------------------------
58
59  class ResourceManager:
60
61      def __init__(self):
62          self.nltk_available = False
63          self.spacy_available = False
64          self.spacy_model = None
65          self.stopwords = set()
66          self.nltk_data_dir = os.path.join(os.getcwd(), "nltk_data")
67
68      def setup_nltk(self):
69          try:
70              # Create NLTK data directory if it doesn't exist
71              if not os.path.exists(self.nltk_data_dir):
72                  os.makedirs(self.nltk_data_dir)
73
74              # Add to NLTK's search path
75              nltk.data.path.append(self.nltk_data_dir)
76
77              # List of resources to download
78              nltk_resources = [
79                  'punkt_tab',
80                  'stopwords',
81                  'wordnet',
82                  'averaged_perceptron_tagger',
83                  'maxent_ne_chunker',
84                  'words'
85              ]
86
87              # Download resources that aren't already available
```

```
88             for resource in nltk_resources:
89                 try:
90                     # Check if resource exists before downloading
91                     try:
92                         nltk.data.find(f'{resource}')
93                         logger.info(f"NLTK resource already available: {resource}")
94                     except LookupError:
95                         # Download if not found
96                         nltk.download(resource, download_dir=self.nltk_data_dir, quiet=
    True)
97                         logger.info(f"NLTK resource downloaded: {resource}")
98                 except Exception as e:
99                     logger.warning(f"Error downloading NLTK resource {resource}: {str(e)
    }")
100
101             # Verify punkt is available for sentence tokenization
102             try:
103                 nltk.data.find('tokenizers/punkt')
104                 self.nltk_tokenizer = nltk.data.load('tokenizers/punkt/english.pickle')
105                 logger.info("NLTK punkt tokenizer loaded successfully")
106             except LookupError:
107                 logger.warning("Could not load punkt tokenizer, will use fallback
    methods")
108                 self.nltk_tokenizer = None
109
110             # Initialize stopwords
111             try:
112                 self.stopwords = set(stopwords.words('english'))
113                 logger.info("NLTK stopwords loaded successfully")
114             except:
115                 # Fallback stopwords if NLTK fails
116                 self.stopwords = {"the", "a", "an", "in", "on", "at", "is", "are", "and"
    , "to", "of", "for", "with"}
117                 logger.info("Using fallback stopwords")
118
119             self.nltk_available = True
120             logger.info("NLTK setup complete and working properly")
121
122         except Exception as e:
123             logger.error(f"Error setting up NLTK: {str(e)}")
124             # Set up basic fallback stopwords
125             self.stopwords = {"the", "a", "an", "in", "on", "at", "is", "are", "and", "
    to", "of", "for", "with"}
126             logger.info("Using basic fallback stopwords due to NLTK setup error")
127             self.nltk_available = False
128
129
130     def setup_spacy(self):
131         """Set up spaCy with fallbacks for different model availability."""
132         try:
133             # Try to load the English model - starting with larger one
134             try:
135                 import spacy
136                 self.spacy_model = spacy.load("en_core_web_md")
137                 logger.info("Loaded spaCy medium model (en_core_web_md)")
138             except:
139                 # Try smaller model
140                 try:
141                     self.spacy_model = spacy.load("en_core_web_sm")
142                     logger.info("Loaded spaCy small model (en_core_web_sm)")
143                 except:
144                     # If not installed, try to download
145                     logger.info("Downloading spaCy model")
146                     os.system("python -m spacy download en_core_web_sm")
147                     self.spacy_model = spacy.load("en_core_web_sm")
148                     logger.info("Downloaded and loaded spaCy small model")
149
150             # Test the model
151             test_text = "The University of Pittsburgh offers student services."
152             doc = self.spacy_model(test_text)
153             entities = [ent.text for ent in doc.ents]
154             if len(entities) > 0 or len(doc) > 0:
155                 self.spacy_available = True
```

```python
                logger.info("spaCy setup complete and working properly")
            else:
                logger.warning("spaCy loaded but not functioning properly")

        except Exception as e:
            logger.error(f"Error setting up spaCy: {str(e)}")
            logger.info("NER functionality will be limited")

    def setup_huggingface_access(self):
        """Set up HuggingFace access with token."""
        try:
            from huggingface_hub import login

            # Try environment variable first
            token = os.environ.get("HUGGINGFACE_TOKEN")

            # If not available, use a default (should be replaced with user's token)
            if not token:
                token = "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"

            if token:
                login(token=token)
                logger.info("Authenticated with HuggingFace Hub")
            else:
                logger.warning("No HuggingFace token found, some models may not be
accessible")

        except Exception as e:
            logger.error(f"Error setting up HuggingFace access: {str(e)}")

    def get_fallback_tokenize(self):
        """Get a robust sentence tokenization function that works with or without NLTK.
"""
        def tokenize_text(text):
            if not text:
                return []

            # Method 1: Try using NLTK's punkt tokenizer directly if available
            if self.nltk_available and hasattr(self, 'nltk_tokenizer') and self.
nltk_tokenizer:
                try:
                    return self.nltk_tokenizer.tokenize(text)
                except Exception as e:
                    logger.warning(f"NLTK tokenizer failed: {str(e)}")

            # Method 2: Try nltk.sent_tokenize which might work even if we couldn't load
 punkt directly
            if self.nltk_available:
                try:
                    from nltk.tokenize import sent_tokenize
                    return sent_tokenize(text)
                except Exception as e:
                    logger.warning(f"NLTK sent_tokenize failed: {str(e)}")

            # Method 3: Try spaCy if available
            if self.spacy_available and self.spacy_model:
                try:
                    # Use spaCy's sentence boundary detection
                    doc = self.spacy_model(text[:10000])  # Limit text length for
performance
                    return [sent.text for sent in doc.sents]
                except Exception as e:
                    logger.warning(f"spaCy sentence tokenization failed: {str(e)}")

            # Method 4: Strong regex-based fallback approach
            logger.info("Using regex-based sentence tokenization as fallback")

            try:
                # Handle common abbreviations to prevent false splits
                text = re.sub(r'\b(Mr|Mrs|Ms|Dr|Prof|Inc|Ltd|Co|Sr|Jr|Ph\.D|M\.D|B\.A|M
\.A|i\.e|e\.g)\.',
                              lambda m: m.group(0).replace('.', '<PERIOD>'), text)
```

```python
                    # Handle decimal numbers and URLs to prevent false splits
                    text = re.sub(r'(\d+)\.(\d+)', r'\1<PERIOD>\2', text)  # Decimal numbers
                    text = re.sub(r'(www\.)|(http\.)', r'\1<PERIOD>', text)  # URLs

                    sentences = []

                    # First split by punctuation + space + capital letter
                    temp_sentences = re.split(r'(?<=[.!?])\s+(?=[A-Z])', text)

                    # Then handle end-of-text punctuation in each segment
                    for segment in temp_sentences:
                        # Split segments that might end with punctuation
                        end_splits = re.split(r'(?<=[.!?])$', segment)
                        sentences.extend([s for s in end_splits if s])

                    # Clean up sentences
                    clean_sentences = []
                    for s in sentences:
                        if not s.strip():
                            continue

                        # Restore periods
                        s = s.replace('<PERIOD>', '.')

                        # Add ending punctuation if missing
                        if not re.search(r'[.!?]$', s):
                            s = s + '.'

                        clean_sentences.append(s)

                    # If we still have no sentences, try a simpler approach - split on
        paragraph breaks
                    if not clean_sentences and '\n\n' in text:
                        paragraphs = text.split('\n\n')
                        for p in paragraphs:
                            if p.strip():
                                clean_sentences.append(p.strip())

                    # If all else fails, treat the whole text as one sentence
                    if not clean_sentences and text.strip():
                        clean_sentences = [text.strip()]

                    return clean_sentences

            except Exception as e:
                logger.warning(f"Regex tokenization failed: {str(e)}")

                # Absolute last resort: simple period splitting
                try:
                    return [s.strip() + '.' for s in text.split('. ') if s.strip()]
                except:
                    # If everything fails, return original text as a single sentence
                    return [text] if text else []

        return tokenize_text

    def analyze_entities(self, text):
        """Extract named entities from text using available NER tools."""
        entities = []

        # Try spaCy first (best quality)
        if self.spacy_available and self.spacy_model:
            try:
                doc = self.spacy_model(text[:5000])  # Limit length for performance
                for ent in doc.ents:
                    entities.append({
                        "text": ent.text,
                        "label": ent.label_,
                        "start": ent.start_char,
                        "end": ent.end_char
                    })
            except Exception as e:
                logger.warning(f"spaCy entity extraction failed: {str(e)}")
```

```python
        # If no entities found or spaCy not available, try NLTK
        if not entities and self.nltk_available:
            try:
                from nltk import word_tokenize, pos_tag, ne_chunk
                from nltk.chunk import tree2conlltags

                # Process with NLTK NER
                tokens = word_tokenize(text[:3000])  # Limit length
                pos_tags = pos_tag(tokens)
                ne_tree = ne_chunk(pos_tags)

                # Extract named entities
                iob_tags = tree2conlltags(ne_tree)
                current_entity = {"text": "", "label": "", "start": 0}
                char_index = 0

                for word, pos, tag in iob_tags:
                    if tag != "O":  # Part of a named entity
                        entity_label = tag.split("-")[1]

                        if tag.startswith("B-"):  # Beginning of entity
                            # Save previous entity if exists
                            if current_entity["text"]:
                                entities.append(current_entity.copy())

                            # Start new entity
                            current_entity = {
                                "text": word,
                                "label": entity_label,
                                "start": char_index
                            }
                        elif tag.startswith("I-"):  # Continuation of entity
                            current_entity["text"] += " " + word
                    else:
                        # End of entity
                        if current_entity["text"]:
                            entities.append(current_entity.copy())
                            current_entity = {"text": "", "label": "", "start": 0}

                    # Update character index (approximate)
                    char_index += len(word) + 1

                # Add final entity if exists
                if current_entity["text"]:
                    entities.append(current_entity)

            except Exception as e:
                logger.warning(f"NLTK entity extraction failed: {str(e)}")

        # If still no entities, use regex patterns for basic extraction
        if not entities:
            # Simple patterns for common entity types
            patterns = [
                (r'\b[A-Z][a-z]+ (University|College|School)\b', 'ORG'),  # Educational
    institutions
                (r'\b[A-Z][a-z]+ (Center|Service|Office|Department)\b', 'ORG'),  #
    University services
                (r'\b[A-Z][a-z]+ (Hall|Building|Library|Center)\b', 'FAC'),  # Campus
    facilities
                (r'\b(January|February|March|April|May|June|July|August|September|
    October|November|December) \d{1,2}(st|nd|rd|th)?,? \d{4}\b', 'DATE'),  # Dates
                (r'\b\d{1,2}:\d{2} (AM|PM|am|pm)\b', 'TIME'),  # Times
                (r'\b[A-Z][a-z]+ [A-Z][a-z]+\b', 'PERSON')  # Potential names
            ]

            for pattern, label in patterns:
                for match in re.finditer(pattern, text):
                    entities.append({
                        "text": match.group(0),
                        "label": label,
                        "start": match.start(),
                        "end": match.end()
```

```
364                                })
365
366                return entities
367
368        def extract_key_phrases(self, text, top_n=10):
369            """Extract key phrases that represent important concepts in the text."""
370            if not text or len(text) < 20:
371                return []
372
373            # Clean text
374            text = re.sub(r'\s+', ' ', text).strip()
375
376            # Split into sentences
377            tokenize = self.get_fallback_tokenize()
378            sentences = tokenize(text)
379
380            if not sentences:
381                return []
382
383            try:
384                # Use TF-IDF to find important n-grams
385                vectorizer = TfidfVectorizer(
386                    ngram_range=(1, 3),
387                    stop_words=list(self.stopwords) if self.stopwords else 'english',
388                    max_features=100
389                )
390
391                # Get matrix
392                tfidf_matrix = vectorizer.fit_transform(sentences)
393
394                # Get feature names
395                feature_names = vectorizer.get_feature_names_out()
396
397                # Calculate scores for each feature across all sentences
398                feature_scores = tfidf_matrix.sum(axis=0).A1
399
400                # Sort features by score
401                sorted_features = sorted(zip(feature_names, feature_scores), key=lambda x: x
        [1], reverse=True)
402
403                # Filter out single stop words and very short phrases
404                key_phrases = []
405                for phrase, score in sorted_features:
406                    if len(phrase) > 3 and not (len(phrase.split()) == 1 and phrase in self.
        stopwords):
407                        key_phrases.append({
408                            "text": phrase,
409                            "score": float(score),
410                            "type": "PHRASE"
411                        })
412
413                        if len(key_phrases) >= top_n:
414                            break
415
416                return key_phrases
417
418            except Exception as e:
419                logger.error(f"Error extracting key phrases: {str(e)}")
420
421                # Fallback: extract noun phrases if available
422                if self.spacy_available and self.spacy_model:
423                    try:
424                        doc = self.spacy_model(text[:5000])
425                        noun_phrases = [{"text": chunk.text, "score": 0.5, "type": "
        NOUN_PHRASE"}
426                                        for chunk in doc.noun_chunks]
427                        return noun_phrases[:top_n]
428                    except:
429                        pass
430
431                # Last resort: just return capitalized phrases
432                cap_phrases = re.findall(r'\b[A-Z][a-zA-Z]*(?:\s+[A-Z][a-zA-Z]*)+\b', text)
433                return [{"text": phrase, "score": 0.5, "type": "CAP_PHRASE"}
```

```
434                     for phrase in set(cap_phrases)][:top_n]

435

436 #---------------------------------------------------------------------
437 # Enhanced Model Manager
438 #---------------------------------------------------------------------
439
440 class ModelManager:
441     """Model manager with dynamic loading/unloading and quantization support."""
442
443     def __init__(self, use_gpu: bool = True, memory_threshold: float = 0.95):
444         """Initialize model manager with memory management."""
445         self.use_gpu = use_gpu
446         self.memory_threshold = memory_threshold
447         self.device = torch.device("cuda" if use_gpu and torch.cuda.is_available() else
    "cpu")
448         self.loaded_models = {}
449
450         # Track how many models we have memory for
451         if self.use_gpu and torch.cuda.is_available():
452             self.total_gpu_memory = torch.cuda.get_device_properties(0).total_memory
453             logger.info(f"GPU memory: {self.total_gpu_memory / 1e9:.2f} GB")
454
455             available_memory = self.total_gpu_memory * self.memory_threshold
456             self.max_models = max(2, min(10, int(available_memory / 1.5e9)))
457             logger.info(f"Estimated capacity: {self.max_models} models can be loaded
    simultaneously")
458         else:
459             self.max_models = 2  # Conservative default for CPU
460
461     def _check_memory(self):
462         """Check if we have enough GPU memory available."""
463         if not self.use_gpu or not torch.cuda.is_available():
464             return True
465
466         # Get current memory usage
467         allocated = torch.cuda.memory_allocated()
468         reserved = torch.cuda.memory_reserved()
469         total = self.total_gpu_memory
470
471         # Calculate percentage used
472         used_fraction = (allocated + reserved) / total
473
474         # Log warning if close to threshold
475         if used_fraction > self.memory_threshold * 0.8:
476             logger.warning(f"GPU memory usage high: {used_fraction:.1%}")
477
478         # Return True if we have enough memory
479         return used_fraction < self.memory_threshold
480
481     def load_model(self, model_name: str, model_class, model_path: str, **kwargs):
482         """Load a model with intelligent memory management and quantization."""
483         # If model already loaded, return it
484         if model_name in self.loaded_models:
485             logger.info(f"Model {model_name} already loaded")
486             # Update last used timestamp
487             self.loaded_models[model_name]['last_used'] = time.time()
488             return self.loaded_models[model_name]['model']
489
490         # Check if we need to unload models to free memory
491         if len(self.loaded_models) >= self.max_models:
492             logger.info(f"Maximum models loaded ({len(self.loaded_models)}), unloading
    least recently used")
493             self._unload_least_used()
494
495         # Check if we have enough memory
496         if not self._check_memory():
497             logger.warning("Memory usage high, forcing garbage collection")
498             gc.collect()
499             if self.use_gpu and torch.cuda.is_available():
500                 torch.cuda.empty_cache()
501
502             # Check again after cleaning
503             if not self._check_memory():
```

```
504                # Unload more aggressively
505                if self.loaded_models:
506                    logger.warning("Still low on memory, unloading all models")
507                    self.cleanup()
508                else:
509                    logger.error("Not enough memory to load model even after cleanup")
510                    raise MemoryError("Not enough GPU memory available")
511
512        # Apply quantization if requested
513        quantization_config = None
514        if kwargs.pop('quantize', False) and self.use_gpu:
515            quantization_config = BitsAndBytesConfig(
516                load_in_4bit=True,
517                bnb_4bit_use_double_quant=True,
518                bnb_4bit_quant_type="nf4",
519                bnb_4bit_compute_dtype=torch.float16
520            )
521            logger.info(f"Using 4-bit quantization for {model_name}")
522
523        # Try to load the model with retries for network issues
524        @backoff.on_exception(backoff.expo,
525                              (requests.RequestException, OSError),
526                              max_tries=3, max_time=60)
527        def load_with_retry():
528            if model_class == SentenceTransformer:
529                # Special handling for sentence transformers
530                model = SentenceTransformer(model_path, device=self.device)
531            else:
532                # Default loading with quantization if specified
533                if quantization_config:
534                    model = model_class.from_pretrained(
535                        model_path,
536                        quantization_config=quantization_config,
537                        device_map="auto" if self.use_gpu else None,
538                        **kwargs
539                    )
540                else:
541                    # Handle device placement
542                    if 'device_map' not in kwargs and self.use_gpu:
543                        if hasattr(model_class, 'from_pretrained') and 'auto' in dir(
    model_class):
544                            # This model supports auto device mapping
545                            kwargs['device_map'] = "auto"
546                            model = model_class.from_pretrained(model_path, **kwargs)
547                        else:
548                            # Manual device placement
549                            model = model_class.from_pretrained(model_path, **kwargs).to
    (self.device)
550                    else:
551                        model = model_class.from_pretrained(model_path, **kwargs)
552
553                        # Explicitly move to device if not using auto mapping
554                        if 'device_map' not in kwargs and self.use_gpu:
555                            try:
556                                model = model.to(self.device)
557                            except Exception as e:
558                                logger.warning(f"Could not move model to {self.device}:
    {str(e)}")
559
560            return model
561
562        # Try to load the model
563        try:
564            logger.info(f"Loading model {model_name} from {model_path}")
565            model = load_with_retry()
566            self.loaded_models[model_name] = {
567                'model': model,
568                'last_used': time.time(),
569                'size': self._estimate_model_size(model)
570            }
571            logger.info(f"Successfully loaded {model_name}")
572            return model
573        except Exception as e:
```

```
574            logger.error(f"Error loading model {model_name}: {str(e)}")
575
576            # Try smaller fallback models if primary ones fail
577            if model_path == "google/flan-t5-xl":
578                logger.info("Trying smaller T5 model instead")
579                return self.load_model(model_name, model_class, "google/flan-t5-large",
    **kwargs)
580            elif model_path == "google/flan-t5-large":
581                logger.info("Trying even smaller T5 model")
582                return self.load_model(model_name, model_class, "google/flan-t5-base",
    **kwargs)
583            elif "v3-large" in model_path or "-large-" in model_path:
584                logger.info("Trying base model instead of large")
585                smaller_path = model_path.replace("large", "base")
586                return self.load_model(model_name, model_class, smaller_path, **kwargs)
587
588            raise
589
590    def _estimate_model_size(self, model):
591        """Estimate the size of a model in bytes (rough approximation)."""
592        try:
593            # Get model parameters
594            params = sum(p.numel() for p in model.parameters())
595
596            # Estimate bytes per parameter
597            bytes_per_param = 4  # Conservative estimate
598            if hasattr(model, 'dtype'):
599                if model.dtype in [torch.float16, torch.bfloat16]:
600                    bytes_per_param = 2
601                elif model.dtype == torch.int8:
602                    bytes_per_param = 1
603
604            # Return estimated size
605            return params * bytes_per_param
606        except:
607            # If we can't estimate, use a default value
608            return 5e8
609
610    def _unload_least_used(self):
611        """Unload the least recently used model."""
612        if not self.loaded_models:
613            return
614
615        # Find least recently used
616        lru_model = min(self.loaded_models.items(), key=lambda x: x[1]['last_used'])
617        model_name = lru_model[0]
618
619        # Unload it
620        self.unload_model(model_name)
621
622    def unload_model(self, model_name: str):
623        """Unload a specific model from memory."""
624        if model_name not in self.loaded_models:
625            return
626
627        logger.info(f"Unloading model {model_name}")
628
629        try:
630            # Get model
631            model = self.loaded_models[model_name]['model']
632
633            # Move to CPU first to free GPU memory
634            if self.use_gpu:
635                try:
636                    model = model.to('cpu')
637                except:
638                    pass
639
640            # Delete model
641            del model
642
643            # Remove from loaded models
644            del self.loaded_models[model_name]
```

```python
            # Force garbage collection
            gc.collect()
            if self.use_gpu and torch.cuda.is_available():
                torch.cuda.empty_cache()

            logger.info(f"Successfully unloaded {model_name}")

        except Exception as e:
            logger.error(f"Error unloading model {model_name}: {str(e)}")

    def cleanup(self):
        """Unload all models and free memory."""
        if not self.loaded_models:
            return

        logger.info(f"Cleaning up, unloading {len(self.loaded_models)} models")

        # Get all model names first
        model_names = list(self.loaded_models.keys())

        # Unload each model
        for model_name in model_names:
            self.unload_model(model_name)

        # Final garbage collection
        gc.collect()
        if self.use_gpu and torch.cuda.is_available():
            torch.cuda.empty_cache()

        logger.info("All models unloaded")

#----------------------------------------------------------------------
# Enhanced Crawler with Anti-Bot Detection and Content Quality Filters
#----------------------------------------------------------------------

class EnhancedWebCrawler:
    """
    Web crawler with anti-bot detection avoidance, content quality
    assessment, and intelligent page prioritization.
    """

    def __init__(self, base_url: str, delay: float = 2.0, checkpoint_dir: str = "
    checkpoints"):
        """Initialize the crawler with robust settings."""
        self.base_url = base_url
        self.base_domain = urlparse(base_url).netloc
        self.visited_urls = set()
        self.failed_urls = set()
        self.delay = delay
        self.content_cache = {}
        self.checkpoint_dir = checkpoint_dir

        # Create checkpoint directory if it doesn't exist
        os.makedirs(checkpoint_dir, exist_ok=True)

        # Create a unique ID for this crawler instance based on the base URL
        self.checkpoint_id = hashlib.md5(base_url.encode('utf-8')).hexdigest()

        # For tracking page quality
        self.page_scores = {}

        # Create persistent session with cookies
        self.session = requests.Session()
        self.session.cookies = cookiejar.LWPCookieJar()

        # Rotate user agents to avoid detection - using more modern user agents
        self.user_agents = [
            'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like
    Gecko) Chrome/120.0.0.0 Safari/537.36',
            'Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_7) AppleWebKit/605.1.15 (KHTML
    , like Gecko) Version/16.5 Safari/605.1.15',
```

```python
            'Mozilla/5.0 (Windows NT 10.0; Win64; x64; rv:120.0) Gecko/20100101 Firefox
    /120.0',
            'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like
    Gecko) Chrome/120.0.0.0 Safari/537.36 Edg/120.0.0.0',
            'Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_7) AppleWebKit/537.36 (KHTML,
    like Gecko) Chrome/120.0.0.0 Safari/537.36',
            'Mozilla/5.0 (X11; Linux x86_64) AppleWebKit/537.36 (KHTML, like Gecko)
    Chrome/120.0.0.0 Safari/537.36',
            'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like
    Gecko) Chrome/120.0.0.0 Safari/537.36 OPR/106.0.0.0',
            'Mozilla/5.0 (iPad; CPU OS 16_5 like Mac OS X) AppleWebKit/605.1.15 (KHTML,
    like Gecko) Version/16.5 Mobile/15E148 Safari/604.1'
        ]

        # Headers to mimic real browsers
        self.browser_headers = {
            'Accept': 'text/html,application/xhtml+xml,application/xml;q=0.9,image/avif,
    image/webp,image/apng,*/*;q=0.8,application/signed-exchange;v=b3;q=0.7',
            'Accept-Language': 'en-US,en;q=0.9',
            'Accept-Encoding': 'gzip, deflate, br',
            'Connection': 'keep-alive',
            'Upgrade-Insecure-Requests': '1',
            'Sec-Fetch-Dest': 'document',
            'Sec-Fetch-Mode': 'navigate',
            'Sec-Fetch-Site': 'none',
            'Sec-Fetch-User': '?1',
            'sec-ch-ua': '"Not_A Brand";v="99", "Google Chrome";v="109", "Chromium";v
    ="109"',
            'sec-ch-ua-mobile': '?0',
            'sec-ch-ua-platform': '"Windows"'
        }

        # Actual bot detection patterns - much more focused to avoid false positives
        self.bot_patterns = [
            r'captcha\s+verification',
            r'security\s+check\s+failed',
            r'automated\s+access\s+detected',
            r'access\s+denied.*automated',
            r'blocked.*bot',
            r'cloudflare.*ray\s+id',
            r'your\s+IP\s+has\s+been\s+blocked'
        ]
        self.bot_regex = re.compile('|'.join(self.bot_patterns), re.IGNORECASE)

        # Content quality indicators
        self.meaningful_phrases = [
            # University/education terms
            'university', 'campus', 'college', 'academic', 'student', 'faculty', 'staff'
    ,
            'course', 'class', 'program', 'degree', 'major', 'minor', 'education',
            'research', 'study', 'learning',

            # Thrive/wellness specific terms
            'wellness', 'health', 'counseling', 'support', 'resource', 'service',
            'assistance', 'help', 'aid', 'benefit', 'initiative', 'thrive', 'success',
            'wellbeing', 'mental health', 'physical health', 'emotional health',

            # Academic support terms
            'advising', 'tutoring', 'mentoring', 'center', 'office', 'department',
            'financial aid', 'scholarship', 'grant', 'funding', 'career', 'job',
            'internship', 'opportunity'
        ]

        # Paths to prioritize - will be populated by caller
        self.priority_paths = []

        # Cache for robots.txt rules
        self.robots_rules = None
        self.crawl_delay = 0

        # Keep track of visit times to ensure natural browsing patterns
        self.last_visit_time = {}
        self.global_last_visit = time.time() - self.delay
```

```python
778
779    def _get_checkpoint_path(self, file_type: str) -> str:
780        """Get the path for a specific checkpoint file."""
781        return os.path.join(self.checkpoint_dir, f"{file_type}_{self.checkpoint_id}.json
       ")
782
783    def save_checkpoint(self) -> None:
784        """Save crawler state to checkpoint files."""
785        try:
786            # Save content cache
787            cache_path = self._get_checkpoint_path("content")
788            with open(cache_path, 'w', encoding='utf-8') as f:
789                json.dump(self.content_cache, f)
790
791            # Save visited URLs
792            visited_path = self._get_checkpoint_path("visited")
793            with open(visited_path, 'w', encoding='utf-8') as f:
794                json.dump(list(self.visited_urls), f)
795
796            # Save failed URLs
797            failed_path = self._get_checkpoint_path("failed")
798            with open(failed_path, 'w', encoding='utf-8') as f:
799                json.dump(list(self.failed_urls), f)
800
801            # Save page scores
802            scores_path = self._get_checkpoint_path("scores")
803            with open(scores_path, 'w', encoding='utf-8') as f:
804                json.dump(self.page_scores, f)
805
806            logger.info(f"Checkpoint saved: {len(self.content_cache)} pages")
807
808        except Exception as e:
809            logger.error(f"Error saving checkpoint: {str(e)}")
810
811    def load_checkpoint(self) -> bool:
812        """Load crawler state from checkpoint files."""
813        try:
814            # Check if checkpoint files exist
815            cache_path = self._get_checkpoint_path("content")
816            visited_path = self._get_checkpoint_path("visited")
817            failed_path = self._get_checkpoint_path("failed")
818            scores_path = self._get_checkpoint_path("scores")
819
820            if not all(os.path.exists(path) for path in [cache_path, visited_path,
       failed_path]):
821                logger.info("Incomplete checkpoint files, starting fresh")
822                return False
823
824            # Load content cache
825            with open(cache_path, 'r', encoding='utf-8') as f:
826                self.content_cache = json.load(f)
827
828            # Load visited URLs
829            with open(visited_path, 'r', encoding='utf-8') as f:
830                self.visited_urls = set(json.load(f))
831
832            # Load failed URLs
833            with open(failed_path, 'r', encoding='utf-8') as f:
834                self.failed_urls = set(json.load(f))
835
836            # Load page scores if available
837            if os.path.exists(scores_path):
838                with open(scores_path, 'r', encoding='utf-8') as f:
839                    loaded_scores = json.load(f)
840                    # Convert string keys back to float values
841                    self.page_scores = {k: float(v) if isinstance(v, str) else v
842                                        for k, v in loaded_scores.items()}
843
844            logger.info(f"Checkpoint loaded: {len(self.content_cache)} pages, {len(self.
       visited_urls)} visited URLs")
845            return True
846
847        except Exception as e:
```

```python
848            logger.error(f"Error loading checkpoint: {str(e)}")
849            return False
850
851    def validate_checkpoint(self) -> None:
852        """Validate and clean checkpoint data."""
853        if not self.content_cache:
854            return
855
856        # Check cached content for quality
857        invalid_urls = []
858
859        for url, content in tqdm(self.content_cache.items(), desc="Validating checkpoint
    "):
860            if not content or not self._is_valid_content(content, url):
861                invalid_urls.append(url)
862
863        # Remove invalid content
864        for url in invalid_urls:
865            logger.warning(f"Removing invalid content for {url}")
866            del self.content_cache[url]
867            if url in self.page_scores:
868                del self.page_scores[url]
869
870        if invalid_urls:
871            logger.info(f"Removed {len(invalid_urls)} invalid pages from checkpoint")
872
873            # Save updated checkpoint
874            self.save_checkpoint()
875
876    def _is_valid_content(self, content: str, url: str) -> bool:
877        """Check if content is valid (not an error page, has meaningful content)."""
878        # Check for empty content
879        if not content or len(content) < 50:  # Reduced minimum length
880            logger.debug(f"Content too short for {url}")
881            return False
882
883        # Check only for definitive bot detection patterns
884        if self.bot_regex.search(content):
885            logger.warning(f"Bot detection pattern found for {url}")
886            return False
887
888        # Check if content has some basic text
889        word_count = len(re.findall(r'\b\w+\b', content))
890        if word_count < 10:
891            logger.debug(f"Too few words ({word_count}) in {url}")
892            return False
893
894        return True
895
896    def _check_robots_txt(self):
897        """Parse robots.txt to respect crawling guidelines."""
898        if self.robots_rules is not None:
899            return
900
901        try:
902            # Get the robots.txt file
903            robots_url = urljoin(self.base_url, "/robots.txt")
904            logger.info(f"Checking robots.txt at {robots_url}")
905
906            # Use a simple GET request with a short timeout
907            user_agent = random.choice(self.user_agents)
908            robots_response = requests.get(
909                robots_url,
910                timeout=5,
911                headers={
912                    'User-Agent': user_agent,
913                    'Accept': 'text/plain,text/html;q=0.9,*/*;q=0.8'
914                }
915            )
916
917            if robots_response.status_code == 200:
918                # Parse robots.txt content
919                content = robots_response.text
```

```python
                self.robots_rules = []

                # Simple parsing of disallow rules and crawl-delay
                current_agent = None
                for line in content.split('\n'):
                    line = line.strip()

                    # Skip empty lines and comments
                    if not line or line.startswith('#'):
                        continue

                    # Check for User-agent directive
                    if line.lower().startswith('user-agent:'):
                        agent = line.split(':', 1)[1].strip().lower()
                        current_agent = agent

                    # Check for Disallow directive for relevant user agents
                    elif line.lower().startswith('disallow:') and (current_agent in ['*'
, 'bot', '']):
                        path = line.split(':', 1)[1].strip()
                        if path:
                            self.robots_rules.append(path)

                    # Check for Crawl-delay directive for relevant user agents
                    elif line.lower().startswith('crawl-delay:') and (current_agent in [
'*', 'bot', '']):
                        try:
                            delay = float(line.split(':', 1)[1].strip())
                            # Use the largest crawl-delay
                            self.crawl_delay = max(self.crawl_delay, delay)
                        except ValueError:
                            pass

                logger.info(f"Found {len(self.robots_rules)} disallow rules in robots.
txt")
                if self.crawl_delay > 0:
                    logger.info(f"Found crawl-delay: {self.crawl_delay} seconds")
                    # Update delay
                    self.delay = max(self.delay, self.crawl_delay)
            else:
                logger.info(f"No robots.txt found or accessible, status: {
robots_response.status_code}")
                self.robots_rules = []

        except Exception as e:
            logger.warning(f"Error parsing robots.txt: {str(e)}")
            self.robots_rules = []  # Continue without rules

    def _is_valid_url(self, url: str) -> bool:
        """Check if URL is valid and should be crawled."""
        # Parse URL
        try:
            parsed = urlparse(url)
        except:
            return False

        # Basic URL validation
        if not all([parsed.scheme, parsed.netloc]) or parsed.scheme not in ['http', '
https']:
            return False

        # Check domain
        url_domain = parsed.netloc.lower()
        if not (url_domain == self.base_domain or
                url_domain.endswith('.' + self.base_domain) or
                self.base_domain.endswith('.' + url_domain)):
            return False

        # Check file extension - skip media, documents, etc.
        path = parsed.path.lower()
        if re.search(r'\.(jpg|jpeg|png|gif|pdf|doc|docx|ppt|pptx|zip|rar|exe|css|js|xml|
json)$', path):
            return False
```

```
987
988          # Skip common utility paths
989          if re.search(r'/(login|logout|signin|signout|register|profile|cart|checkout|
      search\?)', path):
990              return False
991
992          # Check robots.txt rules
993          if self.robots_rules:
994              for rule in self.robots_rules:
995                  if path.startswith(rule):
996                      return False
997
998          return True
999
1000     def _get_random_delay(self):
1001         """Get a random delay to mimic human browsing patterns."""
1002
1003         # Base delay factors - more variance
1004         delay_options = [
1005             self.delay * random.uniform(0.7, 1.0),
1006             self.delay * random.uniform(1.0, 1.5),
1007             self.delay * random.uniform(1.5, 3.0)
1008         ]
1009
1010         # Choose a delay with weighting toward the middle
1011         weights = [0.3, 0.4, 0.3]
1012         return random.choices(delay_options, weights=weights)[0]
1013
1014     @backoff.on_exception(
1015         backoff.expo,
1016         (requests.RequestException, requests.exceptions.Timeout, requests.exceptions.
      ConnectionError),
1017         max_tries=3,
1018         giveup=lambda e: isinstance(e, requests.exceptions.HTTPError) and e.response.
      status_code in [403, 404, 405, 410],
1019     )
1020     def _get_page_content(self, url: str) -> Tuple[str, List[str]]:
1021         """Get content and links from a page with anti-bot detection measures."""
1022
1023         current_time = time.time()
1024
1025         # Global delay
1026         global_elapsed = current_time - self.global_last_visit
1027         if global_elapsed < self.delay:
1028             time.sleep(self.delay - global_elapsed)
1029
1030         # Per-URL delay
1031         if url in self.last_visit_time:
1032             url_elapsed = current_time - self.last_visit_time[url]
1033             url_min_delay = self.delay * 3
1034             if url_elapsed < url_min_delay:
1035                 time.sleep(url_min_delay - url_elapsed)
1036
1037         # Update timing records
1038         self.last_visit_time[url] = time.time()
1039         self.global_last_visit = time.time()
1040
1041         # Prepare headers with a random user agent
1042         headers = self.browser_headers.copy()
1043         headers['User-Agent'] = random.choice(self.user_agents)
1044
1045         # Add referer for more realistic browsing
1046         if self.visited_urls:
1047             # Use a previously visited URL as referer
1048             potential_referers = list(self.visited_urls)
1049             if len(potential_referers) > 5:
1050                 potential_referers = potential_referers[-5:]
1051             referer = random.choice(potential_referers)
1052             headers['Referer'] = referer
1053
1054             # Adjust Sec-Fetch-Site based on domain relationship
1055             referer_domain = urlparse(referer).netloc
1056             current_domain = urlparse(url).netloc
```

```python
            if referer_domain == current_domain:
                headers['Sec-Fetch-Site'] = 'same-origin'
            else:
                headers['Sec-Fetch-Site'] = 'cross-site'

        # Create debug directory if enabled and doesn't exist
        debug_dir = "debug_html"
        if not os.path.exists(debug_dir) and len(self.visited_urls) < 10:
            os.makedirs(debug_dir, exist_ok=True)

        try:
            # Add slight variation in the request to appear more natural
            if random.random() < 0.2:
                try:
                    self.session.head(
                        url,
                        headers=headers,
                        timeout=random.uniform(1.0, 3.0),
                        allow_redirects=True
                    )
                    # Small pause between HEAD and GET
                    time.sleep(random.uniform(0.1, 0.5))
                except:

                    pass

            # Request the page
            response = self.session.get(
                url,
                headers=headers,
                timeout=random.uniform(10.0, 20.0),
                allow_redirects=True
            )

            # Check HTTP status
            if response.status_code != 200:
                logger.warning(f"HTTP error {response.status_code} for {url}")
                self.failed_urls.add(url)
                return "", []

            # Check content type
            content_type = response.headers.get('content-type', '').lower()
            if 'text/html' not in content_type and 'application/xhtml+xml' not in
    content_type:
                logger.warning(f"Skipping non-HTML content ({content_type}) at {url}")
                self.failed_urls.add(url)
                return "", []

            # Get HTML content
            html_content = response.text

            # Save raw HTML for debugging
            if len(self.visited_urls) < 5 and debug_dir:
                filename = os.path.join(debug_dir, f"page_{len(self.visited_urls)}.html"
    )
                try:
                    with open(filename, 'w', encoding='utf-8') as f:
                        f.write(html_content)
                    logger.debug(f"Saved debug HTML to {filename}")
                except Exception as debug_e:
                    logger.debug(f"Could not save debug HTML: {str(debug_e)}")

            # Extract main content
            clean_text = self._extract_clean_text(html_content, url)

            # Validate content - much less strict validation
            if not self._is_valid_content(clean_text, url):
                logger.warning(f"Invalid content for {url}")
                self.failed_urls.add(url)
                return "", []

            # Extract links
            links = []
```

```python
            try:
                soup = BeautifulSoup(html_content, 'html.parser')

                for a_tag in soup.find_all('a', href=True):
                    href = a_tag['href']

                    # Skip empty or javascript links
                    if not href or href.startswith(('javascript:', '#', 'mailto:', 'tel:
    ')):
                        continue

                    # Convert to absolute URL
                    absolute_url = urljoin(url, href)

                    # Validate URL
                    if self._is_valid_url(absolute_url):
                        # Normalize URL - remove fragments and some query params
                        parsed = urlparse(absolute_url)

                        # Keep only the base URL for uniqueness
                        clean_url = f"{parsed.scheme}://{parsed.netloc}{parsed.path}"

                        # Only add unique URLs
                        if clean_url not in links:
                            links.append(clean_url)
            except Exception as e:
                logger.error(f"Error extracting links from {url}: {str(e)}")

            # Save content to cache
            self.content_cache[url] = clean_text

            # Calculate page quality
            quality = self._get_page_quality_score(clean_text, url)
            self.page_scores[url] = quality

            # Log success
            logger.info(f"Successfully crawled {url} - {len(clean_text)} chars, {len(
    links)} links, quality: {quality:.2f}")

            return clean_text, links

        except Exception as e:
            logger.error(f"Error fetching {url}: {str(e)}")
            self.failed_urls.add(url)
            return "", []

    def _extract_clean_text(self, html: str, url: str) -> str:
        """Extract clean, main content text from HTML with semantic structure."""
        try:
            # Parse HTML
            soup = BeautifulSoup(html, 'html.parser')

            # Remove non-content elements
            for element in soup.find_all(['script', 'style', 'noscript', 'svg', 'iframe'
    , 'form', 'nav', 'footer']):
                element.decompose()

            # Try to find main content
            main_content = None

            # Check for common content containers - prioritize semantic elements
            for selector in [
                'main', 'article', '#content', '.content', '#main-content', '.main-
    content',
                '.page-content', '.container', 'div.entry-content', '.article', 'section
    ',
                'body'
            ]:
                elements = soup.select(selector)
                if elements:
                    main_content = elements[0]
                    break
```

```python
                if not main_content:
                    main_content = soup.find('body')

                if not main_content:
                    logger.warning(f"Could not extract main content from {url}")
                    return ""

                # Extract title
                title_text = ""
                title = soup.find('title')
                if title and title.text:
                    title_text = f"TITLE: {title.text.strip()}\n\n"

                # Get all text with proper spacing
                body_text = main_content.get_text(separator='\n\n', strip=True)

                # Combine title and body
                full_text = title_text + body_text

                # Clean up text
                full_text = re.sub(r'\n{3,}', '\n\n', full_text)  # Remove excessive
    newlines

                return full_text.strip()

        except Exception as e:
            logger.error(f"Error extracting content from {url}: {str(e)}")
            return ""

    def _get_page_quality_score(self, content: str, url: str) -> float:
        """Calculate a quality score for the page (0-1)."""
        if not content:
            return 0.0

        # Base score
        score = 0.5

        # Adjust based on content length (longer is usually better)
        words = len(content.split())
        if words < 100:
            score -= 0.2
        elif words < 300:
            score -= 0.1
        elif words > 500:
            score += 0.1
        elif words > 1000:
            score += 0.2

        # Adjust based on URL priority
        parsed = urlparse(url)
        path = parsed.path.lower()

        # Priority paths get bonus
        if self.priority_paths and any(path.startswith(priority) for priority in self.
    priority_paths):
            score += 0.15

        # Home page or about page get bonus
        if path == '/' or path == '/index.html' or path.startswith('/about'):
            score += 0.1

        # Check for meaningful phrases
        matches = sum(1 for phrase in self.meaningful_phrases if phrase in content.lower
    ())
        phrase_score = min(0.2, matches * 0.02)  # Cap at 0.2
        score += phrase_score

        # Presence of structured content (headings) is usually good
        if "HEADING:" in content:
            score += 0.05

        # Bonus for having a proper title
        if "TITLE:" in content:
```

```python
                score += 0.05

        # Ensure score is in range [0, 1]
        return max(0.0, min(1.0, score))

    def _prioritize_urls(self, urls: List[str]) -> List[str]:
        """Prioritize URLs for crawling based on content value heuristics."""
        if not urls:
            return []

        # Score each URL for priority
        scored_urls = []

        for url in urls:
            score = 0.5  # Base score

            # Parse URL
            parsed = urlparse(url)
            path = parsed.path.lower()

            # Prioritize URLs based on path patterns

            # High priority paths
            if any(priority in path for priority in ['/about', '/services', '/resources'
, '/programs']):
                score += 0.4

            # Medium priority paths
            elif any(priority in path for priority in ['/wellness', '/health', '/student
', '/academic']):
                score += 0.3

            # Lower priority but still valuable
            elif any(priority in path for priority in ['/news', '/events', '/contact', '
/faq']):
                score += 0.2

            # Deprioritize pagination and archive pages
            if re.search(r'/(page|p)/\d+', path) or re.search(r'/\d{4}/(0\d|1[0-2])',
path):
                score -= 0.2

            # Prioritize shorter paths
            path_depth = path.count('/')
            if path_depth <= 1:
                score += 0.1
            elif path_depth >= 4:
                score -= 0.1

            # Check if URL has query parameters
            if parsed.query:
                score -= 0.1

            # Add to scored list
            scored_urls.append((url, score))

        # Sort by score and return URLs only
        sorted_urls = [url for url, score in sorted(scored_urls, key=lambda x: x[1],
reverse=True)]

        return sorted_urls

    def crawl(self, max_pages: int = 30, max_workers: int = 1):
        """Crawl website with prioritization, quality filters, and improved anti-bot
measures."""
        if max_pages <= 0:
            return {}

        # Start with base URL
        urls_to_visit = [self.base_url]

        # Add important paths to initial crawl list
        for path in self.priority_paths:
```

```
1333              if path:
1334                  urls_to_visit.append(urljoin(self.base_url, path))
1335
1336          # Add common paths that might exist
1337          common_paths = ["/index.html", "/home", "/about", "/contact", "/resources", "/
     services"]
1338          for path in common_paths:
1339              urls_to_visit.append(urljoin(self.base_url, path))
1340
1341          # Deduplicate
1342          urls_to_visit = list(dict.fromkeys(urls_to_visit))
1343
1344          # Initialize results
1345          results = {}
1346          if self.content_cache:
1347              results = self.content_cache.copy()
1348              logger.info(f"Starting with {len(results)} cached pages")
1349
1350          # Check robots.txt first
1351          self._check_robots_txt()
1352
1353          # Try alternate domains if base URL has www. or not
1354          alternate_urls = []
1355          parsed_base = urlparse(self.base_url)
1356          base_domain = parsed_base.netloc
1357
1358          if base_domain.startswith('www.'):
1359              # Try non-www version
1360              alt_domain = base_domain[4:]
1361              alt_url = f"{parsed_base.scheme}://{alt_domain}{parsed_base.path}"
1362              alternate_urls.append(alt_url)
1363          else:
1364              # Try www version
1365              alt_domain = f"www.{base_domain}"
1366              alt_url = f"{parsed_base.scheme}://{alt_domain}{parsed_base.path}"
1367              alternate_urls.append(alt_url)
1368
1369          # Also try https if the original URL is http
1370          if parsed_base.scheme == 'http':
1371              https_url = f"https://{base_domain}{parsed_base.path}"
1372              alternate_urls.append(https_url)
1373
1374          # Add alternate URLs to visit list
1375          for alt_url in alternate_urls:
1376              if alt_url not in urls_to_visit:
1377                  urls_to_visit.append(alt_url)
1378
1379          # Crawl sequentially
1380          with tqdm(total=max_pages, desc="Crawling pages", initial=len(results)) as pbar:
1381              attempts = 0
1382              max_attempts = max_pages * 3
1383
1384              while urls_to_visit and len(results) < max_pages and attempts < max_attempts
     :
1385                  # First prioritize URLs
1386                  urls_to_visit = self._prioritize_urls(urls_to_visit)
1387
1388                  # Get next URL to visit
1389                  url = urls_to_visit.pop(0)
1390
1391                  # Skip if already visited
1392                  if url in self.visited_urls:
1393                      continue
1394
1395                  # Get page content
1396                  content, links = self._get_page_content(url)
1397                  self.visited_urls.add(url)
1398                  attempts += 1
1399
1400                  if content:
1401                      # Store content
1402                      results[url] = content
1403                      pbar.update(1)
```

```
1404
1405                    # Add new links to queue
1406                    for link in links:
1407                        if (link not in self.visited_urls and
1408                            link not in self.failed_urls and
1409                            link not in urls_to_visit):
1410                            urls_to_visit.append(link)
1411                else:
1412                    self.failed_urls.add(url)
1413
1414                # Save checkpoint periodically
1415                if len(results) % 5 == 0:
1416                    self.content_cache = results
1417                    self.save_checkpoint()
1418
1419                # Add a randomized delay between requests
1420                time.sleep(self._get_random_delay())
1421
1422        # Update and save final state
1423        self.content_cache = results
1424        self.save_checkpoint()
1425
1426        # Log results
1427        logger.info(f"Crawling complete: {len(results)} pages, {len(self.visited_urls)}
      visited, {len(self.failed_urls)} failed")
1428
1429        # Return results sorted by quality
1430        sorted_results = {}
1431        for url, content in sorted(results.items(),
1432                                   key=lambda x: self.page_scores.get(x[0], 0),
1433                                   reverse=True):
1434            sorted_results[url] = content
1435
1436        return sorted_results
1437
1438 #----------------------------------------------------------------------
1439 # Persistent Crawler with Checkpointing
1440 #----------------------------------------------------------------------
1441
1442 class PersistentCrawler(EnhancedWebCrawler):
1443     """Enhanced crawler with checkpoint capabilities."""
1444
1445     def __init__(self, base_url: str, delay: float = 1.0, checkpoint_dir: str = "
      checkpoints"):
1446         """Initialize crawler with checkpoint support."""
1447         super().__init__(base_url, delay)
1448         self.checkpoint_dir = checkpoint_dir
1449
1450         # Create checkpoint directory if it doesn't exist
1451         os.makedirs(checkpoint_dir, exist_ok=True)
1452
1453         # Create a unique ID for this crawler instance based on the base URL
1454         self.checkpoint_id = hashlib.md5(base_url.encode('utf-8')).hexdigest()
1455
1456     def _get_checkpoint_path(self, file_type: str) -> str:
1457         """Get the path for a specific checkpoint file."""
1458         return os.path.join(self.checkpoint_dir, f"{file_type}_{self.checkpoint_id}.json
      ")
1459
1460     def save_checkpoint(self) -> None:
1461         """Save crawler state to checkpoint files."""
1462         try:
1463             # Save content cache
1464             cache_path = self._get_checkpoint_path("content")
1465             with open(cache_path, 'w', encoding='utf-8') as f:
1466                 json.dump(self.content_cache, f)
1467
1468             # Save visited URLs
1469             visited_path = self._get_checkpoint_path("visited")
1470             with open(visited_path, 'w', encoding='utf-8') as f:
1471                 json.dump(list(self.visited_urls), f)
1472
1473             # Save failed URLs
```

```python
            failed_path = self._get_checkpoint_path("failed")
            with open(failed_path, 'w', encoding='utf-8') as f:
                json.dump(list(self.failed_urls), f)

            # Save page scores
            scores_path = self._get_checkpoint_path("scores")
            with open(scores_path, 'w', encoding='utf-8') as f:
                json.dump(self.page_scores, f)

            logger.info(f"Checkpoint saved: {len(self.content_cache)} pages")

        except Exception as e:
            logger.error(f"Error saving checkpoint: {str(e)}")

    def load_checkpoint(self) -> bool:
        """Load crawler state from checkpoint files."""
        try:
            # Check if checkpoint files exist
            cache_path = self._get_checkpoint_path("content")
            visited_path = self._get_checkpoint_path("visited")
            failed_path = self._get_checkpoint_path("failed")
            scores_path = self._get_checkpoint_path("scores")

            if not all(os.path.exists(path) for path in [cache_path, visited_path,
    failed_path]):
                logger.info("Incomplete checkpoint files, starting fresh")
                return False

            # Load content cache
            with open(cache_path, 'r', encoding='utf-8') as f:
                self.content_cache = json.load(f)

            # Load visited URLs
            with open(visited_path, 'r', encoding='utf-8') as f:
                self.visited_urls = set(json.load(f))

            # Load failed URLs
            with open(failed_path, 'r', encoding='utf-8') as f:
                self.failed_urls = set(json.load(f))

            # Load page scores if available
            if os.path.exists(scores_path):
                with open(scores_path, 'r', encoding='utf-8') as f:
                    self.page_scores = json.load(f)

            logger.info(f"Checkpoint loaded: {len(self.content_cache)} pages, {len(self.
    visited_urls)} visited URLs")
            return True

        except Exception as e:
            logger.error(f"Error loading checkpoint: {str(e)}")
            return False

    def validate_checkpoint(self) -> None:
        """Validate and clean checkpoint data."""
        if not self.content_cache:
            return

        # Check cached content for quality
        invalid_urls = []

        for url, content in tqdm(self.content_cache.items(), desc="Validating checkpoint
    "):
            if not content or not self._is_valid_content(content, url):
                invalid_urls.append(url)

        # Remove invalid content
        for url in invalid_urls:
            logger.warning(f"Removing invalid content for {url}")
            del self.content_cache[url]
            if url in self.page_scores:
                del self.page_scores[url]
```

```python
1544            if invalid_urls:
1545                logger.info(f"Removed {len(invalid_urls)} invalid pages from checkpoint")
1546
1547                # Save updated checkpoint
1548                self.save_checkpoint()
1549
1550 #----------------------------------------------------------------------
1551 # Semantic Document Processor
1552 #----------------------------------------------------------------------
1553
1554 class SemanticDocument:
1555     """A document with semantic understanding for better processing."""
1556
1557     def __init__(self, text: str, url: str = "", title: str = ""):
1558         """Initialize a semantic document."""
1559         self.text = text
1560         self.url = url
1561         self.title = title
1562         self.chunks = []
1563         self.entities = []
1564         self.topics = []
1565         self.embedding = None
1566
1567     def __str__(self):
1568         """String representation."""
1569         return f"Document({self.title or self.url}, {len(self.text)} chars, {len(self.
    chunks)} chunks)"
1570
1571     def add_chunk(self, chunk):
1572         """Add a semantic chunk to the document."""
1573         self.chunks.append(chunk)
1574
1575     def add_entity(self, entity):
1576         """Add an entity to the document."""
1577         self.entities.append(entity)
1578
1579     def add_topic(self, topic):
1580         """Add a topic to the document."""
1581         self.topics.append(topic)
1582
1583     def set_embedding(self, embedding):
1584         """Set the document's embedding vector."""
1585         self.embedding = embedding
1586
1587     def get_summary(self):
1588         """Get a summary of key document statistics."""
1589         return {
1590             "url": self.url,
1591             "title": self.title,
1592             "length": len(self.text),
1593             "chunks": len(self.chunks),
1594             "entities": len(self.entities),
1595             "topics": [t["text"] for t in self.topics[:5]] if self.topics else []
1596         }
1597
1598 class SemanticChunk:
1599     """A semantic chunk of content optimized for QA processing."""
1600
1601     def __init__(self, text: str, doc_url: str = "", index: int = 0):
1602         """Initialize a semantic chunk."""
1603         self.text = text
1604         self.doc_url = doc_url
1605         self.index = index
1606         self.entities = []
1607         self.embedding = None
1608
1609     def __str__(self):
1610         """String representation."""
1611         return f"Chunk({self.index}, {len(self.text)} chars, {len(self.entities)}
    entities)"
1612
1613     def add_entity(self, entity):
1614         """Add an entity to the chunk."""
```

```python
1615            self.entities.append(entity)
1616
1617        def set_embedding(self, embedding):
1618            """Set the chunk's embedding vector."""
1619            self.embedding = embedding
1620
1621        def to_dict(self):
1622            """Convert to dictionary for serialization."""
1623            return {
1624                "text": self.text,
1625                "doc_url": self.doc_url,
1626                "index": self.index,
1627                "entities": self.entities,
1628                "embedding": self.embedding.tolist() if isinstance(self.embedding, torch.
        Tensor) else self.embedding
1629            }
1630
1631    class DocumentProcessor:
1632        """Process documents with semantic understanding and chunking."""
1633
1634        def __init__(self, resource_manager, model_manager: ModelManager = None, use_gpu:
        bool = True):
1635            """Initialize document processor with resources."""
1636            self.resource_manager = resource_manager
1637            self.model_manager = model_manager or ModelManager(use_gpu=use_gpu)
1638            self.use_gpu = use_gpu
1639            self.device = torch.device("cuda" if use_gpu and torch.cuda.is_available() else
        "cpu")
1640
1641            # Load embedding model
1642            self.embedding_model = None
1643            self.embedding_dimension = 0
1644
1645            # For tokenization
1646            self.tokenize_fn = resource_manager.get_fallback_tokenize()
1647
1648        def load_models(self):
1649            """Load necessary models for document processing."""
1650            try:
1651                # Try to load a more powerful embedding model first
1652                embedding_model_options = [
1653                    "sentence-transformers/all-mpnet-base-v2",
1654                    "sentence-transformers/all-MiniLM-L12-v2",
1655                    "sentence-transformers/all-MiniLM-L6-v2"
1656                ]
1657
1658                for model_name in embedding_model_options:
1659                    try:
1660                        logger.info(f"Loading embedding model: {model_name}")
1661                        self.embedding_model = self.model_manager.load_model(
1662                            "document_embeddings",
1663                            SentenceTransformer,
1664                            model_name
1665                        )
1666
1667                        # Test the model
1668                        test_embedding = self.embedding_model.encode("test",
        convert_to_tensor=True)
1669                        self.embedding_dimension = test_embedding.shape[0]
1670                        logger.info(f"Embedding model loaded successfully: {model_name},
        dimension: {self.embedding_dimension}")
1671                        break
1672                    except Exception as e:
1673                        logger.warning(f"Failed to load embedding model {model_name}: {str(e
        )}")
1674                        continue
1675
1676                if self.embedding_model is None:
1677                    logger.error("Could not load any embedding model")
1678
1679            except Exception as e:
1680                logger.error(f"Error loading document processing models: {str(e)}")
1681
```

```python
1682    def process_document(self, text: str, url: str = "", title: str = "") ->
        SemanticDocument:
1683        """Process a document into semantic chunks with entity recognition."""
1684        if not text:
1685            return None
1686
1687        # Clean the text first
1688        text = self._clean_text(text)
1689
1690        # Initialize document
1691        doc = SemanticDocument(text, url, title)
1692
1693        # Extract title if not provided
1694        if not title:
1695            title_match = re.search(r'TITLE: (.*?)(\n|$)', text)
1696            if title_match:
1697                doc.title = title_match.group(1).strip()
1698
1699        # Extract entities
1700        entities = self.resource_manager.analyze_entities(text)
1701        for entity in entities:
1702            doc.add_entity(entity)
1703
1704        # Extract topics/key concepts
1705        topics = self.resource_manager.extract_key_phrases(text, top_n=10)
1706        for topic in topics:
1707            doc.add_topic(topic)
1708
1709        # Create semantic chunks
1710        chunks = self._create_semantic_chunks(text, url)
1711        for chunk in chunks:
1712            doc.add_chunk(chunk)
1713
1714        # Create document embedding if model available
1715        if self.embedding_model:
1716            try:
1717                # Use title + first part of text for document-level embedding
1718                summary_text = (doc.title + ". " if doc.title else "") + text[:1000]
1719                doc.set_embedding(self.embedding_model.encode(summary_text,
        convert_to_tensor=True))
1720            except Exception as e:
1721                logger.error(f"Error creating document embedding: {str(e)}")
1722
1723        return doc
1724
1725    def _clean_text(self, text: str) -> str:
1726        """Clean text for better processing."""
1727        # Remove excessive whitespace
1728        text = re.sub(r'\s+', ' ', text)
1729
1730        # Fix Unicode characters
1731        text = text.replace('        ', "'").replace('        ', '"').replace('      ', '"
        ')
1732        text = text.replace(' ', ' ').replace('&amp;', '&').replace('&quot;', '"')
1733
1734        # Fix spacing after periods
1735        text = re.sub(r'\.([A-Z])', r'. \1', text)
1736
1737        # Handle marked section headers
1738        text = re.sub(r'HEADING: ', '\n## ', text)
1739        text = re.sub(r'TITLE: ', '# ', text)
1740        text = re.sub(r'CONTENT: ', '\n', text)
1741
1742        return text.strip()
1743
1744    def _create_semantic_chunks(self, text: str, url: str = "") -> List[SemanticChunk]:
1745        """Create semantic chunks from text preserving context."""
1746        # This is a key function for improving QA quality through better chunking
1747
1748        chunks = []
1749
1750        # Different chunking strategies based on document structure
1751        if '##' in text or '#' in text:
```

```python
                # Document has section markers, use them for chunking
                chunks = self._chunk_by_sections(text, url)
            else:
                # Try to identify sections by headings and paragraphs
                chunks = self._chunk_by_paragraphs(text, url)

        # If we get very large chunks, split them further
        max_chunk_size = 1500  # About 300-400 words typically
        new_chunks = []

        for i, chunk in enumerate(chunks):
            if len(chunk.text) > max_chunk_size:
                # Split large chunks with overlap
                sub_chunks = self._split_large_chunk(chunk.text, url, chunk.index)
                new_chunks.extend(sub_chunks)
            else:
                new_chunks.append(chunk)

        # Analyze entities for each chunk
        if new_chunks:
            self._analyze_chunk_entities(new_chunks)

            # Create embeddings for each chunk if model available
            if self.embedding_model:
                try:
                    # Prepare all chunks for batch encoding
                    texts = [chunk.text for chunk in new_chunks]

                    # Encode in batch for efficiency
                    embeddings = self.embedding_model.encode(texts, convert_to_tensor=
    True)

                    # Assign embeddings back to chunks
                    for i, chunk in enumerate(new_chunks):
                        chunk.set_embedding(embeddings[i])

                except Exception as e:
                    logger.error(f"Error creating chunk embeddings: {str(e)}")

        return new_chunks

    def _chunk_by_sections(self, text: str, url: str) -> List[SemanticChunk]:
        """Chunk text by markdown section markers."""
        chunks = []

        # Split by section headers
        section_pattern = r'(^|\n)#+\s+.+?($|\n)'
        sections = re.split(section_pattern, text)

        # Group headers with content
        i = 0
        while i < len(sections):
            if i+2 < len(sections) and re.match(r'(^|\n)#+\s+', sections[i+1]):
                # This is a section header followed by content
                header = sections[i+1].strip()
                content = sections[i+2].strip()

                if content:
                    chunk = SemanticChunk(header + "\n" + content, url, len(chunks))
                    chunks.append(chunk)

            i += 1

        # If no chunks were created, treat the whole text as one chunk
        if not chunks and text:
            chunks.append(SemanticChunk(text, url, 0))

        return chunks

    def _chunk_by_paragraphs(self, text: str, url: str) -> List[SemanticChunk]:
        """Chunk text by paragraphs with some overlap."""
        chunks = []
```

```python
1824             # Split into paragraphs first
1825             paragraphs = []
1826
1827             # Check if text has natural paragraph breaks
1828             if '\n\n' in text:
1829                 paragraphs = [p.strip() for p in text.split('\n\n') if p.strip()]
1830             else:
1831                 # Try to identify paragraphs by looking for sentence boundaries
1832                 sentences = self.tokenize_fn(text)
1833
1834                 # Group sentences into paragraphs (simple method: ~5 sentences per paragraph
1835                 current_para = []
1836                 for sentence in sentences:
1837                     current_para.append(sentence)
1838                     if len(current_para) >= 5:
1839                         paragraphs.append(' '.join(current_para))
1840                         current_para = []
1841
1842                 # Add the last paragraph if any sentences remain
1843                 if current_para:
1844                     paragraphs.append(' '.join(current_para))
1845
1846             # Group paragraphs into chunks with reasonable sizes
1847             current_chunk = []
1848             current_length = 0
1849             target_length = 1000
1850
1851             for paragraph in paragraphs:
1852                 para_length = len(paragraph)
1853
1854                 if current_length > 0 and current_length + para_length > target_length:
1855                     chunk_text = ' '.join(current_chunk)
1856                     chunks.append(SemanticChunk(chunk_text, url, len(chunks)))
1857
1858                     # Start a new chunk with some overlap
1859                     if current_chunk:
1860                         current_chunk = [current_chunk[-1], paragraph]
1861                         current_length = len(current_chunk[-1]) + para_length
1862                     else:
1863                         current_chunk = [paragraph]
1864                         current_length = para_length
1865                 else:
1866                     # Add paragraph to current chunk
1867                     current_chunk.append(paragraph)
1868                     current_length += para_length
1869
1870             # Add the last chunk if there's anything left
1871             if current_chunk:
1872                 chunk_text = ' '.join(current_chunk)
1873                 chunks.append(SemanticChunk(chunk_text, url, len(chunks)))
1874
1875             return chunks
1876
1877     def _split_large_chunk(self, text: str, url: str, base_index: int) -> List[
     SemanticChunk]:
1878         """Split a large chunk into smaller ones with overlapping content."""
1879         sub_chunks = []
1880
1881         # Try to split on sentence boundaries
1882         sentences = self.tokenize_fn(text)
1883
1884         if not sentences:
1885             # If tokenization fails, just split by character count
1886             chunk_size = 1000
1887             overlap = 100
1888
1889             for i in range(0, len(text), chunk_size - overlap):
1890                 end = min(i + chunk_size, len(text))
1891                 if end - i < 200:
1892                     break
1893
1894                 sub_text = text[i:end]
```

```python
                idx = base_index * 100 + len(sub_chunks)
                sub_chunks.append(SemanticChunk(sub_text, url, idx))
        else:
            # Split by sentences with overlap
            chunk_size = 10
            overlap = 2

            for i in range(0, len(sentences), chunk_size - overlap):
                end = min(i + chunk_size, len(sentences))
                if end - i < 3:
                    break

                sub_text = ' '.join(sentences[i:end])
                idx = base_index * 100 + len(sub_chunks)
                sub_chunks.append(SemanticChunk(sub_text, url, idx))

        return sub_chunks

    def _analyze_chunk_entities(self, chunks: List[SemanticChunk]) -> None:
        """Extract entities from each chunk."""
        for chunk in chunks:
            entities = self.resource_manager.analyze_entities(chunk.text)
            for entity in entities:
                chunk.add_entity(entity)

#----------------------------------------------------------------------
# Knowledge Base for RAG
#----------------------------------------------------------------------

class KnowledgeBase:
    """Knowledge base for retrieval augmented generation."""

    def __init__(self, use_gpu: bool = True):
        """Initialize knowledge base."""
        self.use_gpu = use_gpu
        self.device = torch.device("cuda" if use_gpu and torch.cuda.is_available() else
    "cpu")

        # Storage for documents and chunks
        self.documents = {}
        self.chunks = []

        # For vector search
        self.embeddings = None
        self.chunk_ids = []
        self.embedding_dim = 0

        # Text search index
        self.chunk_text_index = {}  # term -> [chunk indices]

    def add_document(self, doc: SemanticDocument) -> None:
        """Add a document to the knowledge base."""
        if not doc or not doc.text:
            return

        # Add document
        self.documents[doc.url] = doc

        # Add chunks
        for chunk in doc.chunks:
            self.chunks.append(chunk)

            # Add to text index
            for term in set(chunk.text.lower().split()):
                if len(term) > 3:
                    if term not in self.chunk_text_index:
                        self.chunk_text_index[term] = []
                    self.chunk_text_index[term].append(len(self.chunks) - 1)

        # Rebuild search index if needed
        self.rebuild_search_index()

    def rebuild_search_index(self) -> None:
```

46

```
1967            """Rebuild the vector search index."""
1968            if not self.chunks:
1969                return
1970
1971            # Check if chunks have embeddings
1972            if not hasattr(self.chunks[0], 'embedding') or self.chunks[0].embedding is None:
1973                logger.warning("Chunks don't have embeddings, can't build search index")
1974                return
1975
1976            try:
1977                # Collect embeddings
1978                all_embeddings = []
1979                self.chunk_ids = []
1980
1981                for i, chunk in enumerate(self.chunks):
1982                    if chunk.embedding is not None:
1983                        all_embeddings.append(chunk.embedding)
1984                        self.chunk_ids.append(i)
1985
1986                if not all_embeddings:
1987                    logger.warning("No valid embeddings found in chunks")
1988                    return
1989
1990                # Convert to tensor
1991                if isinstance(all_embeddings[0], list):
1992                    self.embeddings = torch.tensor(all_embeddings)
1993                else:
1994                    self.embeddings = torch.stack(all_embeddings)
1995
1996                # Move to device
1997                if self.use_gpu:
1998                    self.embeddings = self.embeddings.to(self.device)
1999
2000                # Get embedding dimension
2001                self.embedding_dim = self.embeddings.shape[1]
2002
2003                logger.info(f"Built search index with {len(self.chunk_ids)} chunks,
       dimension: {self.embedding_dim}")
2004
2005            except Exception as e:
2006                logger.error(f"Error building search index: {str(e)}")
2007                self.embeddings = None
2008                self.chunk_ids = []
2009
2010        def search(self, query: str, embedding_model=None, top_k: int = 5) -> List[Dict]:
2011            """Search for relevant chunks using hybrid retrieval."""
2012            if not self.chunks:
2013                return []
2014
2015            # Check if we have embeddings
2016            if self.embeddings is None:
2017                # Fallback to text search
2018                return self._text_search(query, top_k)
2019
2020            # Check if we have an embedding model
2021            if embedding_model is None:
2022                logger.warning("No embedding model provided, falling back to text search")
2023                return self._text_search(query, top_k)
2024
2025            try:
2026                # Get query embedding
2027                query_embedding = embedding_model.encode(query, convert_to_tensor=True)
2028
2029                # Move to same device as index
2030                if self.use_gpu:
2031                    query_embedding = query_embedding.to(self.device)
2032
2033                # Calculate similarity scores
2034                similarity = torch.matmul(self.embeddings, query_embedding.unsqueeze(1)).
       squeeze(1)
2035
2036                # Get top-k chunks
2037                if len(similarity) <= top_k:
```

```python
                    top_indices = torch.argsort(similarity, descending=True)
                else:
                    top_indices = torch.topk(similarity, k=top_k).indices

                # Convert to Python list
                top_indices = top_indices.cpu().tolist()

                # Get the actual chunk indices
                chunk_indices = [self.chunk_ids[i] for i in top_indices]

                # Get the scores
                scores = [similarity[i].item() for i in top_indices]

                # Create results
                results = []
                for idx, score in zip(chunk_indices, scores):
                    chunk = self.chunks[idx]
                    results.append({
                        'chunk': chunk,
                        'score': score,
                        'doc_url': chunk.doc_url
                    })

                # Hybrid re-ranking: boost scores of chunks matching query terms
                query_terms = set(query.lower().split())
                for result in results:
                    # Check if chunk contains query terms
                    chunk_text = result['chunk'].text.lower()
                    matching_terms = sum(1 for term in query_terms if term in chunk_text)

                    # Boost score based on term matches (small boost to preserve vector
    similarity ordering)
                    result['score'] += matching_terms * 0.05

                # Re-sort by adjusted scores
                results.sort(key=lambda x: x['score'], reverse=True)

                return results[:top_k]

        except Exception as e:
            logger.error(f"Error in vector search: {str(e)}")
            # Fallback to text search
            return self._text_search(query, top_k)

    def _text_search(self, query: str, top_k: int = 5) -> List[Dict]:
        """Fallback text-based search."""
        if not self.chunks:
            return []

        # Simple term matching
        query_terms = set(query.lower().split())

        # Score each chunk
        scored_chunks = []

        for i, chunk in enumerate(self.chunks):
            # Count matching terms
            chunk_text = chunk.text.lower()
            matching_terms = sum(1 for term in query_terms if term in chunk_text)

            # Add additional score for exact phrases
            exact_matches = 0
            for size in range(2, min(5, len(query_terms) + 1)):
                for j in range(len(query_terms) - size + 1):
                    phrase = ' '.join(list(query_terms)[j:j+size])
                    if phrase in chunk_text:
                        exact_matches += 1

            # Calculate score
            score = matching_terms + exact_matches * 2

            if score > 0:
                scored_chunks.append((i, score))
```

```python
2110
2111          # Sort by score
2112          scored_chunks.sort(key=lambda x: x[1], reverse=True)
2113
2114          # Get top-k results
2115          results = []
2116          for idx, score in scored_chunks[:top_k]:
2117              chunk = self.chunks[idx]
2118              results.append({
2119                  'chunk': chunk,
2120                  'score': score,
2121                  'doc_url': chunk.doc_url
2122              })
2123
2124          return results
2125
2126      def get_chunks_by_url(self, url: str) -> List[SemanticChunk]:
2127          """Get all chunks for a specific URL."""
2128          return [chunk for chunk in self.chunks if chunk.doc_url == url]
2129
2130      def get_document(self, url: str) -> Optional[SemanticDocument]:
2131          """Get a document by URL."""
2132          return self.documents.get(url)
2133
2134      def get_stats(self) -> Dict:
2135          """Get knowledge base statistics."""
2136          return {
2137              'documents': len(self.documents),
2138              'chunks': len(self.chunks),
2139              'indexed_chunks': len(self.chunk_ids) if self.embeddings is not None else 0,
2140              'indexed_terms': len(self.chunk_text_index),
2141              'embedding_dim': self.embedding_dim
2142          }
2143
2144      def save(self, filepath: str) -> None:
2145          """Save knowledge base to file."""
2146          try:
2147              # Prepare for serialization
2148              data = {
2149                  'documents': {},
2150                  'chunks': [],
2151                  'chunk_text_index': self.chunk_text_index
2152              }
2153
2154              # Serialize documents (without embeddings to save space)
2155              for url, doc in self.documents.items():
2156                  data['documents'][url] = {
2157                      'text': doc.text,
2158                      'title': doc.title,
2159                      'topics': doc.topics,
2160                      'entities': doc.entities
2161                  }
2162
2163              # Serialize chunks (with embeddings)
2164              for chunk in self.chunks:
2165                  chunk_data = chunk.to_dict()
2166                  data['chunks'].append(chunk_data)
2167
2168              # Save to file
2169              with open(filepath, 'w', encoding='utf-8') as f:
2170                  json.dump(data, f)
2171
2172              logger.info(f"Knowledge base saved to {filepath}")
2173
2174          except Exception as e:
2175              logger.error(f"Error saving knowledge base: {str(e)}")
2176
2177      def load(self, filepath: str) -> bool:
2178          """Load knowledge base from file."""
2179          try:
2180              # Load from file
2181              with open(filepath, 'r', encoding='utf-8') as f:
2182                  data = json.load(f)
```

```python
                # Clear existing data
                self.documents = {}
                self.chunks = []
                self.chunk_text_index = {}
                self.embeddings = None
                self.chunk_ids = []

                # Load documents
                for url, doc_data in data['documents'].items():
                    doc = SemanticDocument(doc_data['text'], url, doc_data['title'])
                    doc.topics = doc_data['topics']
                    doc.entities = doc_data['entities']
                    self.documents[url] = doc

                # Load chunks
                for chunk_data in data['chunks']:
                    chunk = SemanticChunk(
                        chunk_data['text'],
                        chunk_data['doc_url'],
                        chunk_data['index']
                    )
                    chunk.entities = chunk_data['entities']

                    # Load embedding if present
                    if 'embedding' in chunk_data and chunk_data['embedding']:
                        if isinstance(chunk_data['embedding'], list):
                            chunk.embedding = torch.tensor(chunk_data['embedding'])
                        else:
                            # Handle string or other formats
                            logger.warning(f"Unexpected embedding format for chunk {chunk.
    index}")

                    self.chunks.append(chunk)

                # Load chunk text index
                self.chunk_text_index = data.get('chunk_text_index', {})

                # Rebuild search index
                self.rebuild_search_index()

                logger.info(f"Knowledge base loaded from {filepath}: {len(self.documents)}
    documents, {len(self.chunks)} chunks")
                return True

        except Exception as e:
            logger.error(f"Error loading knowledge base: {str(e)}")
            return False

#----------------------------------------------------------------------
# QA Generator with RAG and CoT
#----------------------------------------------------------------------

class QAGenerator:
    """Advanced QA pair generator using RAG and Chain-of-Thought with enhanced retrieval
     capabilities."""

    def __init__(self,
                 resource_manager,
                 model_manager: ModelManager,
                 knowledge_base: KnowledgeBase,
                 use_gpu: bool = True):
        """Initialize QA generator with necessary components."""
        self.resource_manager = resource_manager
        self.model_manager = model_manager
        self.knowledge_base = knowledge_base
        self.use_gpu = use_gpu
        self.device = torch.device("cuda" if use_gpu and torch.cuda.is_available() else
    "cpu")

        # Track loaded models
        self.embedding_model = None
        self.question_generator = None
```

```python
        self.answer_generator = None
        self.qg_tokenizer = None
        self.ag_tokenizer = None

        # For answer generation quality
        self.fact_checker = None
        self.qa_evaluator = None

        # For enhanced retrieval
        self.reranker = None
        self.reranker_tokenizer = None

    def load_models(self):
        """Load necessary models for QA generation with enhanced retrieval capabilities.
        """
        try:
            # 1. Load embedding model for retrieval
            if self.embedding_model is None:
                # Try to use knowledge base's existing model first
                if hasattr(self.knowledge_base, 'embedding_model') and self.
    knowledge_base.embedding_model:
                    self.embedding_model = self.knowledge_base.embedding_model
                    logger.info("Using knowledge base's embedding model")
                else:
                    # Otherwise load our own
                    logger.info("Loading embedding model")
                    self.embedding_model = self.model_manager.load_model(
                        "qa_embeddings",
                        SentenceTransformer,
                        "sentence-transformers/all-mpnet-base-v2"
                    )

            # 2. Load question generation model
            if self.question_generator is None:
                logger.info("Loading question generation model")
                try:
                    # Try to load the XXL model first
                    logger.info("Attempting to load Flan-T5-XXL model")
                    self.question_generator = self.model_manager.load_model(
                        "question_generator",
                        T5ForConditionalGeneration,
                        "google/flan-t5-xxl",
                        quantize=True
                    )
                    self.qg_tokenizer = AutoTokenizer.from_pretrained("google/flan-t5-
    xxl")
                    logger.info("Successfully loaded Flan-T5-XXL model")
                except Exception as e:
                    logger.warning(f"Failed to load T5-XXL model: {e}")
                    try:
                        # Fall back to XL model
                        logger.info("Falling back to Flan-T5-XL model")
                        self.question_generator = self.model_manager.load_model(
                            "question_generator",
                            T5ForConditionalGeneration,
                            "google/flan-t5-xl"
                        )
                        self.qg_tokenizer = AutoTokenizer.from_pretrained("google/flan-
    t5-xl")
                        logger.info("Successfully loaded Flan-T5-XL model")
                    except Exception as e:
                        logger.warning(f"Failed to load T5-XL model: {e}")
                        # Fall back to an even smaller model
                        try:
                            logger.info("Falling back to Flan-T5-Large model")
                            self.question_generator = self.model_manager.load_model(
                                "question_generator",
                                T5ForConditionalGeneration,
                                "google/flan-t5-large"
                            )
                            self.qg_tokenizer = AutoTokenizer.from_pretrained("google/
    flan-t5-large")
                            logger.info("Successfully loaded Flan-T5-Large model")
```

```
                    except Exception as e:
                        logger.warning(f"Failed to load T5-Large model: {e}")
                        # Final fallback
                        logger.info("Falling back to Flan-T5-Base model")
                        self.question_generator = self.model_manager.load_model(
                            "question_generator",
                            T5ForConditionalGeneration,
                            "google/flan-t5-base"
                        )
                        self.qg_tokenizer = AutoTokenizer.from_pretrained("google/
    flan-t5-base")
                        logger.info("Successfully loaded Flan-T5-Base model")

        # 3. Load answer generation model (same model can be used for efficiency)
        if self.answer_generator is None:
            logger.info("Loading answer generation model")
            # Re-use the question generator model if it's suitable
            if self.question_generator and isinstance(self.question_generator,
    T5ForConditionalGeneration):
                self.answer_generator = self.question_generator
                self.ag_tokenizer = self.qg_tokenizer
                logger.info("Reusing question generation model for answer generation
    ")
            else:
                # Only executed if question generator failed or is not T5
                logger.warning("Need to load separate answer generator model")
                # Try the same cascade of models
                try:
                    self.answer_generator = self.model_manager.load_model(
                        "answer_generator",
                        T5ForConditionalGeneration,
                        "google/flan-t5-xxl",
                        quantize=True
                    )
                    self.ag_tokenizer = AutoTokenizer.from_pretrained("google/flan-
    t5-xxl")
                except Exception as e:
                    logger.warning(f"Failed to load XXL answer generator: {e}")
                    self.answer_generator = self.model_manager.load_model(
                        "answer_generator",
                        T5ForConditionalGeneration,
                        "google/flan-t5-base"
                    )
                    self.ag_tokenizer = AutoTokenizer.from_pretrained("google/flan-
    t5-base")

        # 4. Optional fact checking model for quality control
        try:
            logger.info("Loading fact checking model")
            self.fact_checker = self.model_manager.load_model(
                "fact_checker",
                AutoModelForSequenceClassification,
                "vectara/hallucination_evaluation_model",
                trust_remote_code=True
            )
            # Load the tokenizer for this fact-checker
            self.fact_checker_tokenizer = AutoTokenizer.from_pretrained(
                "vectara/hallucination_evaluation_model",
                trust_remote_code=True
            )
            logger.info("Successfully loaded fact_checker")
        except Exception as e:
            logger.warning(f"Failed to load fact checker: {e}")
            self.fact_checker = None
            self.fact_checker_tokenizer = None

        # 5. Load cross-encoder re-ranker for enhanced retrieval (higher quality
    than QA evaluator)
        try:
            logger.info("Loading cross-encoder re-ranker model (L-12)")
            self.reranker = self.model_manager.load_model(
                "cross_encoder_reranker",
                CrossEncoder,
```

```
2387                    "cross-encoder/ms-marco-MiniLM-L-12-v2",
2388                    max_length=512
2389                )
2390                logger.info("Successfully loaded cross-encoder re-ranker (L-12)")
2391            except Exception as e:
2392                logger.warning(f"Failed to load L-12 cross-encoder re-ranker: {e}")
2393                # Try a smaller re-ranker model
2394                try:
2395                    logger.info("Loading smaller cross-encoder model as fallback")
2396                    self.reranker = self.model_manager.load_model(
2397                        "cross_encoder_reranker",
2398                        CrossEncoder,
2399                        "cross-encoder/ms-marco-MiniLM-L-6-v2",
2400                        max_length=512
2401                    )
2402                    logger.info("Successfully loaded cross-encoder re-ranker (L-6)")
2403                except Exception as e:
2404                    logger.error(f"Failed to load any re-ranker: {e}")
2405                    self.reranker = None
2406
2407            # 6. Load QA evaluator for quality assessment (if not already loaded as
       reranker)
2408            if self.qa_evaluator is None:
2409                try:
2410                    # Directly load MS MARCO as the evaluator
2411                    self.qa_evaluator = self.model_manager.load_model(
2412                        "qa_evaluator",
2413                        AutoModelForSequenceClassification,
2414                        "cross-encoder/ms-marco-MiniLM-L-6-v2"
2415                    )
2416                    self.qa_tokenizer = AutoTokenizer.from_pretrained("cross-encoder/ms-
       marco-MiniLM-L-6-v2")
2417                except Exception as e:
2418                    logger.warning(f"Failed to load QA evaluator: {e}")
2419                    self.qa_evaluator = None
2420                    self.qa_tokenizer = None
2421
2422        except Exception as e:
2423            logger.error(f"Error loading QA generation models: {str(e)}")
2424            return False
2425
2426    def identify_topic_type(self, topic: Dict) -> str:
2427        """Identify the type of a topic for question generation context."""
2428        topic_text = topic["text"].lower()
2429
2430        # Check for explicit type indicators in the text
2431        type_indicators = {
2432            "SERVICE": ["service", "center", "office", "desk", "support"],
2433            "PROGRAM": ["program", "initiative", "project", "series", "system"],
2434            "RESOURCE": ["resource", "tool", "material", "guide", "handbook"],
2435            "WELLNESS": ["wellness", "health", "medical", "counseling", "therapy", "
       wellbeing"],
2436            "SUPPORT": ["support", "help", "assistance", "aid", "advising"],
2437            "LOCATION": ["hall", "building", "center", "campus", "laboratory", "library"
       ],
2438            "CONTACT": ["contact", "email", "phone", "reach", "connect"],
2439            "ELIGIBILITY": ["eligible", "qualify", "requirement", "criteria"],
2440            "ACADEMIC": ["academic", "class", "course", "study", "learning", "education"
       ],
2441            "FINANCIAL": ["financial", "money", "fund", "payment", "cost", "expense", "
       scholarship"]
2442        }
2443
2444        # Check if the topic text contains any type indicators
2445        for type_name, indicators in type_indicators.items():
2446            if any(indicator in topic_text for indicator in indicators):
2447                return type_name
2448
2449        # If no clear indicators, try to infer from entities
2450        entity_type_mapping = {
2451            "ORG": "SERVICE",
2452            "GPE": "LOCATION",
2453            "PERSON": "CONTACT",
```

```python
                "DATE": "PROGRAM",
                "MONEY": "FINANCIAL"
        }

        if "entities" in topic and topic["entities"]:
            for entity in topic["entities"]:
                if entity["type"] in entity_type_mapping:
                    return entity_type_mapping[entity["type"]]

        return "GENERAL"

    def generate_questions_from_documents(self, urls: List[str], max_questions_per_url:
    int = 10) -> List[Dict]:
        """Generate diverse questions from documents."""
        all_questions = []

        for url in urls:
            # Get all chunks for this URL
            chunks = self.knowledge_base.get_chunks_by_url(url)
            if not chunks:
                logger.warning(f"No chunks found for URL: {url}")
                continue

            # Get document if available
            document = self.knowledge_base.get_document(url)

            # Extract document topics if available or analyze chunks
            topics = []
            if document and document.topics:
                topics = document.topics
            else:
                # Analyze chunks to extract topics
                all_text = " ".join([chunk.text for chunk in chunks])
                topics = self.resource_manager.extract_key_phrases(all_text, top_n=15)

            # Generate questions for each important topic
            url_questions = []

            for topic in topics:
                # Identify topic type for context
                topic_type = self.identify_topic_type(topic)

                # Generate questions for this topic
                topic_questions = self.generate_questions_for_topic(
                    topic["text"],
                    topic_type,
                    chunks,
                    max_questions=max(2, max_questions_per_url // len(topics))
                )

                # Add metadata to questions
                for question in topic_questions:
                    question["topic"] = topic["text"]
                    question["topic_type"] = topic_type
                    question["source_url"] = url
                    url_questions.append(question)

            # Generate general questions about the document
            general_questions = self.generate_general_questions(chunks)
            for question in general_questions:
                question["topic"] = "General"
                question["topic_type"] = "GENERAL"
                question["source_url"] = url
                url_questions.append(question)

            # Limit to max questions per URL
            if len(url_questions) > max_questions_per_url:
                # Sort by quality score if available, otherwise keep first ones
                if all("quality_score" in q for q in url_questions):
                    url_questions.sort(key=lambda x: x["quality_score"], reverse=True)
                url_questions = url_questions[:max_questions_per_url]

            all_questions.extend(url_questions)
```

```python
2526
2527            return all_questions
2528
2529    def generate_questions_for_topic(self,
2530                                     topic: str,
2531                                     topic_type: str,
2532                                     chunks: List[SemanticChunk],
2533                                     max_questions: int = 3) -> List[Dict]:
2534        """Generate questions for a specific topic using neural models."""
2535        questions = []
2536
2537        # Primary approach: Use neural models to generate contextually-appropriate
       questions
2538        if self.question_generator and self.qg_tokenizer:
2539            neural_questions = self._generate_questions_neural(topic, topic_type, chunks
       )
2540            questions.extend(neural_questions)
2541
2542        # If we still don't have enough questions, use rule-based generation as fallback
2543        if len(questions) < max_questions:
2544            rule_based_questions = self._generate_questions_rule_based(topic, topic_type
       , chunks)
2545            questions.extend(rule_based_questions)
2546
2547        # Remove duplicates
2548        unique_questions = self._deduplicate_questions(questions)
2549
2550        # Ensure all questions end with a question mark
2551        for q in unique_questions:
2552            if not q["text"].endswith("?"):
2553                q["text"] = q["text"] + "?"
2554
2555        # Return top questions limited by max_questions
2556        return unique_questions[:max_questions]
2557
2558    def _generate_questions_neural(self, topic: str, topic_type: str, chunks: List[
       SemanticChunk]) -> List[Dict]:
2559        """Generate questions using neural models with enhanced context awareness."""
2560        if not self.question_generator or not self.qg_tokenizer:
2561            return []
2562
2563        questions = []
2564
2565        try:
2566            # Extract relevant content about this topic from the chunks
2567            topic_content = self._extract_topic_content(topic, chunks)
2568            if not topic_content:
2569                return []
2570
2571            # Create a variety of prompts that encourage diverse, natural question
       generation
2572            prompts = []
2573
2574            # General prompt for contextual questions
2575            prompts.append(
2576                f"Based on this content about {topic}, generate a natural, informative
       question that a university student might ask:\n\n"
2577                f"Content: {topic_content}\n\n"
2578                f"Question:"
2579            )
2580
2581            # Prompt for specific question types based on content analysis
2582            if topic_type == "SERVICE" or topic_type == "PROGRAM" or topic_type == "
       RESOURCE":
2583                prompts.append(
2584                    f"Create an informative question exploring what {topic} is and how
       it can benefit students:\n\n"
2585                    f"Content: {topic_content}\n\n"
2586                    f"Question:"
2587                )
2588
2589                prompts.append(
```

```python
                    f"Generate a question about accessing or utilizing {topic} at the
    university:\n\n"
                    f"Content: {topic_content}\n\n"
                    f"Question:"
                )

            elif topic_type == "LOCATION":
                prompts.append(
                    f"Generate a question about where to find {topic} and what services
    are available there:\n\n"
                    f"Content: {topic_content}\n\n"
                    f"Question:"
                )

            elif topic_type == "WELLNESS" or topic_type == "SUPPORT":
                prompts.append(
                    f"Create a question asking how {topic} supports student wellbeing or
     success:\n\n"
                    f"Content: {topic_content}\n\n"
                    f"Question:"
                )

            elif topic_type == "FINANCIAL":
                prompts.append(
                    f"Generate a question about financial aspects of {topic} that would
    be relevant to students:\n\n"
                    f"Content: {topic_content}\n\n"
                    f"Question:"
                )

            # Add a chain-of-thought prompt to generate more sophisticated questions
            prompts.append(
                f"Based on this content, create a thoughtful question about {topic}:\n\n
    "
                f"Content: {topic_content}\n\n"
                f"Step 1: Identify the most important information about {topic}.\n"
                f"Step 2: Consider what students would want to know about {topic}.\n"
                f"Step 3: Formulate a clear, specific question.\n"
                f"Question:"
            )

            # Add a prompt for eligibility or requirements if relevant
            if "eligibility" in topic_content.lower() or "requirement" in topic_content.
    lower() or "qualify" in topic_content.lower():
                prompts.append(
                    f"Create a question about eligibility or requirements for {topic}:\n
    \n"
                    f"Content: {topic_content}\n\n"
                    f"Question:"
                )

            # Add a prompt for process-related questions if relevant
            if "process" in topic_content.lower() or "step" in topic_content.lower() or
    "procedure" in topic_content.lower():
                prompts.append(
                    f"Generate a question about the process or steps involved with {
    topic}:\n\n"
                    f"Content: {topic_content}\n\n"
                    f"Question:"
                )

            # Generate questions from each prompt
            generated_questions = []

            for prompt in prompts:
                try:
                    # Tokenize prompt
                    inputs = self.qg_tokenizer(prompt, return_tensors="pt", truncation=
    True, max_length=1024)
                    inputs = {k: v.to(self.device) for k, v in inputs.items()}

                    # Generate with diverse sampling parameters
                    outputs = self.question_generator.generate(
```

```python
                        **inputs,
                        max_length=128,
                        num_return_sequences=2,
                        do_sample=True,
                        temperature=0.8,
                        top_p=0.9,
                        no_repeat_ngram_size=3
                    )

                    # Decode outputs
                    for output in outputs:
                        question_text = self.qg_tokenizer.decode(output,
    skip_special_tokens=True)

                        # Clean up question
                        question_text = self._clean_question(question_text)

                        if question_text:
                            generated_questions.append(question_text)

                except Exception as e:
                    logger.error(f"Error generating question from prompt: {str(e)}")
                    continue

            # Process and add generated questions
            for question_text in generated_questions:
                questions.append({
                    "text": question_text,
                    "source": "neural",
                    "quality_score": 0.8
                })

            return questions

        except Exception as e:
            logger.error(f"Error in neural question generation: {str(e)}")
            return []

    def _extract_topic_content(self, topic: str, chunks: List[SemanticChunk]) -> str:
        """Extract content relevant to a topic from chunks."""
        if not chunks:
            return ""

        # Find chunks that mention the topic
        topic_lower = topic.lower()
        relevant_chunks = []

        for chunk in chunks:
            if topic_lower in chunk.text.lower():
                relevant_chunks.append((chunk, 2))  # Direct mention gets higher score
            else:
                # Check for partial matches (topic terms)
                topic_terms = set(topic_lower.split())
                if len(topic_terms) > 1:  # Only check multi-word topics
                    matches = sum(1 for term in topic_terms if term in chunk.text.lower
    ())
                    if matches >= len(topic_terms) // 2:  # At least half the terms
    match
                        relevant_chunks.append((chunk, 1))  # Partial match gets lower
    score

        # Sort by relevance score
        relevant_chunks.sort(key=lambda x: x[1], reverse=True)

        # Extract text from most relevant chunks (limit length)
        text_parts = []
        total_length = 0
        max_length = 1000

        for chunk, _ in relevant_chunks:
            if total_length + len(chunk.text) > max_length:
                # If adding this chunk would exceed max length, just take enough to
    reach max
```

```
2721              remaining = max_length - total_length
2722              if remaining > 100:  # Only add if we can get a meaningful amount
2723                  text_parts.append(chunk.text[:remaining])
2724              break

2726          text_parts.append(chunk.text)
2727          total_length += len(chunk.text)

2729          if total_length >= max_length:
2730              break

2732      # Combine text parts
2733      return " ".join(text_parts)

2735  def _generate_questions_rule_based(self, topic: str, topic_type: str, chunks: List[
      SemanticChunk]) -> List[Dict]:
2736      """Generate questions using rule-based approaches when neural generation fails.
      """
2737      questions = []

2739      # Extract key sentences from chunks that mention the topic
2740      topic_lower = topic.lower()
2741      topic_sentences = []

2743      tokenize_fn = self.resource_manager.get_fallback_tokenize()

2745      for chunk in chunks:
2746          if topic_lower in chunk.text.lower():
2747              # Get sentences from this chunk
2748              sentences = tokenize_fn(chunk.text)

2750              # Find sentences that mention the topic
2751              for sentence in sentences:
2752                  if topic_lower in sentence.lower():
2753                      topic_sentences.append(sentence)

2755      # Generate questions from key sentences and content patterns
2756      if topic_sentences:
2757          for sentence in topic_sentences[:3]:  # Limit to first few sentences
2758              # Try to identify sentence type and generate appropriate question
2759              if re.search(r'(is|are|was|were|will be)', sentence.lower()):
2760                  # Definition/description sentence
2761                  questions.append({
2762                      "text": f"What is {topic} and what does it offer?",
2763                      "source": "rule_based",
2764                      "quality_score": 0.65
2765                  })
2766              elif re.search(r'(can|could|may|might|should)', sentence.lower()):
2767                  # Capability/possibility sentence
2768                  questions.append({
2769                      "text": f"How can students use {topic}?",
2770                      "source": "rule_based",
2771                      "quality_score": 0.65
2772                  })
2773              elif re.search(r'(located|found|available|offered|provided)', sentence.
      lower()):
2774                  # Location/availability sentence
2775                  questions.append({
2776                      "text": f"Where can students access {topic}?",
2777                      "source": "rule_based",
2778                      "quality_score": 0.65
2779                  })
2780              elif re.search(r'(eligible|qualify|qualifies|requirement)', sentence.
      lower()):
2781                  # Eligibility sentence
2782                  questions.append({
2783                      "text": f"Who is eligible for {topic}?",
2784                      "source": "rule_based",
2785                      "quality_score": 0.65
2786                  })
2787      else:
2788          # If no topic-specific sentences found, generate generic questions based on
      topic type
```

```python
            if topic_type == "SERVICE":
                questions.append({
                    "text": f"What services does {topic} provide?",
                    "source": "rule_based",
                    "quality_score": 0.6
                })
            elif topic_type == "LOCATION":
                questions.append({
                    "text": f"Where is {topic} located on campus?",
                    "source": "rule_based",
                    "quality_score": 0.6
                })
            elif topic_type == "PROGRAM":
                questions.append({
                    "text": f"What is the purpose of the {topic} program?",
                    "source": "rule_based",
                    "quality_score": 0.6
                })
            else:
                # Generic fallback
                questions.append({
                    "text": f"What information is available about {topic}?",
                    "source": "rule_based",
                    "quality_score": 0.6
                })

        return questions

    def _deduplicate_questions(self, questions: List[Dict]) -> List[Dict]:
        """Remove duplicate and nearly-duplicate questions."""
        if not questions:
            return []

        unique_questions = []
        question_texts = set()

        # First, sort by quality score (highest first)
        questions.sort(key=lambda x: x.get("quality_score", 0), reverse=True)

        for question in questions:
            # Normalize question text
            text = question["text"].lower().strip()

            # Skip if exact duplicate
            if text in question_texts:
                continue

            # Check for near-duplicates
            is_duplicate = False
            for existing in unique_questions:
                if self._questions_are_similar(text, existing["text"].lower()):
                    is_duplicate = True
                    break

            if not is_duplicate:
                question_texts.add(text)
                unique_questions.append(question)

        return unique_questions

    def _questions_are_similar(self, q1: str, q2: str) -> bool:
        """Check if two questions are semantically similar."""
        # Method 1: Jaccard similarity on words
        words1 = set(q1.split())
        words2 = set(q2.split())

        if not words1 or not words2:
            return False

        jaccard = len(words1.intersection(words2)) / len(words1.union(words2))

        # Questions with high word overlap are likely similar
        if jaccard > 0.7:
```

```
2862                return True
2863
2864            # Method 2: Check edit distance for short questions
2865            if len(q1) < 50 and len(q2) < 50:
2866                edit_distance = self._levenshtein_distance(q1, q2)
2867                if edit_distance / max(len(q1), len(q2)) < 0.3:
2868                    return True
2869
2870            return False
2871
2872        def _levenshtein_distance(self, s1: str, s2: str) -> int:
2873            """Calculate the Levenshtein distance between two strings."""
2874            if len(s1) < len(s2):
2875                return self._levenshtein_distance(s2, s1)
2876
2877            if len(s2) == 0:
2878                return len(s1)
2879
2880            previous_row = range(len(s2) + 1)
2881            for i, c1 in enumerate(s1):
2882                current_row = [i + 1]
2883                for j, c2 in enumerate(s2):
2884                    insertions = previous_row[j + 1] + 1
2885                    deletions = current_row[j] + 1
2886                    substitutions = previous_row[j] + (c1 != c2)
2887                    current_row.append(min(insertions, deletions, substitutions))
2888                previous_row = current_row
2889
2890            return previous_row[-1]
2891
2892        def _clean_question(self, question: str) -> str:
2893            """Clean and normalize a generated question."""
2894            # Remove any prompt leftovers
2895            question = re.sub(r'^(Question:|Q:|Step \d+:)', '', question).strip()
2896
2897            # Ensure first letter is capitalized
2898            if question and not question[0].isupper():
2899                question = question[0].upper() + question[1:]
2900
2901            # Ensure question ends with question mark
2902            if question and not question.endswith('?'):
2903                question = question + '?'
2904
2905            # Remove repetitive question words at start
2906            question = re.sub(r'^(What|How|Why|Where|When|Who)\s+(is|are|can|does|do|did)\s
        +\1\s+', r'\1 \2 ', question, flags=re.IGNORECASE)
2907
2908            return question
2909
2910        def generate_general_questions(self, chunks: List[SemanticChunk]) -> List[Dict]:
2911            """Generate general questions about a document using the model."""
2912            questions = []
2913
2914            # If we have a neural model, generate contextual general questions
2915            if self.question_generator and self.qg_tokenizer and chunks:
2916                try:
2917                    # Combine chunks into a summary
2918                    summary = ""
2919                    total_length = 0
2920                    for chunk in chunks:
2921                        if total_length > 1500:
2922                            break
2923                        summary += chunk.text + " "
2924                        total_length += len(chunk.text)
2925
2926                    # Create multiple prompts for diverse general questions
2927                    general_prompts = [
2928                        f"Based on this content, generate a question that would help someone
        understand the main purpose of this information:\n\nContent: {summary[:1000]}\n\
        nQuestion:",
2929
2930                        f"Create a question that would help a student find key resources
        described in this content:\n\nContent: {summary[500:1500]}\n\nQuestion:",
```

```python
                    f"Generate a question about how students can access the services
mentioned in this content:\n\nContent: {summary[:1000]}\n\nQuestion:",

                    f"Create a question about who students should contact for the
services described in this content:\n\nContent: {summary[500:1500]}\n\nQuestion:"
                ]

                # Generate questions from each prompt
                for prompt in general_prompts:
                    inputs = self.qg_tokenizer(prompt, return_tensors="pt", truncation=
True, max_length=1024)
                    inputs = {k: v.to(self.device) for k, v in inputs.items()}

                    outputs = self.question_generator.generate(
                        **inputs,
                        max_length=128,
                        num_return_sequences=1,
                        do_sample=True,
                        temperature=0.7
                    )

                    # Process generated question
                    question_text = self.qg_tokenizer.decode(outputs[0],
skip_special_tokens=True)
                    question_text = self._clean_question(question_text)

                    if question_text:
                        questions.append({
                            "text": question_text,
                            "source": "neural_general",
                            "quality_score": 0.8
                        })

        except Exception as e:
            logger.error(f"Error generating general neural questions: {str(e)}")
            # Fall back to basic questions if neural generation fails
            questions.append({
                "text": "What services are described on this page?",
                "source": "fallback_general",
                "quality_score": 0.6
            })
            questions.append({
                "text": "How can students access the resources mentioned here?",
                "source": "fallback_general",
                "quality_score": 0.6
            })

    # Deduplicate and return
    return self._deduplicate_questions(questions)[:4]  # Limit to 4 general
questions

def _generate_adaptive_questions(self, urls: List[str]) -> List[Dict]:
    """Generate context-aware adaptive questions for limited content."""
    adaptive_questions = []

    # Get all chunks across all URLs for context
    all_chunks = []
    for url in urls:
        chunks = self.knowledge_base.get_chunks_by_url(url)
        all_chunks.extend(chunks)

    if not all_chunks:
        return []

    # Extract any text content we can find
    all_text = " ".join([chunk.text for chunk in all_chunks])

    # Try neural generation first if model is available
    if self.question_generator and self.qg_tokenizer and all_text:
        try:
            # Extract representative sample of the text
            sample_text = all_text[:2000]  # Take the first 2000 chars as a sample
```

```python
                # Create adaptive prompts based on available content
                adaptive_prompts = [
                    f"Based on this limited information, generate a general question
    that would be appropriate regardless of the specific details:\n\nContent: {
    sample_text}\n\nQuestion:",

                    f"Create a question asking what resources or services are available
    based on this information:\n\nContent: {sample_text}\n\nQuestion:",

                    f"Generate a question about how to find more information about the
    topics mentioned here:\n\nContent: {sample_text}\n\nQuestion:"
                ]

                # Look for specific content patterns and create relevant prompts
                if "contact" in all_text.lower() or "email" in all_text.lower() or "
    phone" in all_text.lower():
                    adaptive_prompts.append(
                        f"Generate a question about how to contact or reach out for the
    services mentioned:\n\nContent: {sample_text}\n\nQuestion:"
                    )

                if "location" in all_text.lower() or "building" in all_text.lower() or "
    office" in all_text.lower():
                    adaptive_prompts.append(
                        f"Create a question about where to find the services or offices
    mentioned:\n\nContent: {sample_text}\n\nQuestion:"
                    )

                # Generate questions from prompts
                for prompt in adaptive_prompts:
                    inputs = self.qg_tokenizer(prompt, return_tensors="pt", truncation=
    True, max_length=1024)
                    inputs = {k: v.to(self.device) for k, v in inputs.items()}

                    outputs = self.question_generator.generate(
                        **inputs,
                        max_length=128,
                        num_return_sequences=1,
                        do_sample=True,
                        temperature=0.7
                    )

                    question_text = self.qg_tokenizer.decode(outputs[0],
    skip_special_tokens=True)
                    question_text = self._clean_question(question_text)

                    if question_text:
                        adaptive_questions.append({
                            "text": question_text,
                            "source": "adaptive_neural",
                            "topic": "General",
                            "topic_type": "GENERAL",
                            "quality_score": 0.75,
                            "source_url": urls[0] if urls else ""
                        })

        except Exception as e:
            logger.error(f"Error generating adaptive neural questions: {str(e)}")
            # Fall back to rule-based questions below

    # If we have no questions yet or too few, add some fallback questions
    if len(adaptive_questions) < 3:
        fallback_questions = [
            {
                "text": "What information is provided on this page?",
                "source": "adaptive_fallback",
                "topic": "General",
                "topic_type": "GENERAL",
                "quality_score": 0.6,
                "source_url": urls[0] if urls else ""
            },
            {
```

```python
                    "text": "What services or resources are described here?",
                    "source": "adaptive_fallback",
                    "topic": "Services",
                    "topic_type": "SERVICE",
                    "quality_score": 0.6,
                    "source_url": urls[0] if urls else ""
                },
                {
                    "text": "How can students get more information about what's
    mentioned here?",
                    "source": "adaptive_fallback",
                    "topic": "Information",
                    "topic_type": "GENERAL",
                    "quality_score": 0.6,
                    "source_url": urls[0] if urls else ""
                }
            ]

            adaptive_questions.extend(fallback_questions)

        # Extract meaningful terms we can use for additional questions
        keywords = set()
        meaningful_phrases = [
            "student", "university", "campus", "service", "resource", "support",
            "wellness", "health", "academic", "financial", "career", "housing",
            "registration", "advising", "tutoring", "counseling", "aid", "scholarship"
        ]

        for phrase in meaningful_phrases:
            if phrase in all_text.lower():
                keywords.add(phrase)

        # Add keyword-based questions if we found any
        for keyword in list(keywords)[:3]:  # Limit to 3 keywords
            keyword_question = {
                "text": f"What information is provided about {keyword} resources or
    services?",
                "source": "adaptive_keyword",
                "topic": keyword.title(),
                "topic_type": "KEYWORD",
                "quality_score": 0.65,
                "source_url": urls[0] if urls else ""
            }
            adaptive_questions.append(keyword_question)

        # Deduplicate questions
        unique_questions = self._deduplicate_questions(adaptive_questions)

        return unique_questions

    def generate_answers(self, questions: List[Dict]) -> List[Dict]:
        """Generate answers for a list of questions using RAG."""
        qa_pairs = []

        # Process each question
        for question in tqdm(questions, desc="Generating answers"):
            try:
                # Get question text and metadata
                question_text = question["text"]
                source_url = question.get("source_url", "")

                # 1. Retrieve relevant chunks
                relevant_chunks = self._retrieve_context(question_text, source_url)

                # 2. Generate answer
                answer = self._generate_answer(question_text, relevant_chunks)

                # 3. Evaluate answer quality
                scores = self._evaluate_answer(question_text, answer, relevant_chunks)

                # 4. Create QA pair with metadata
                qa_pair = {
                    "question": question_text,
```

```python
                    "answer": answer,
                    "source_url": source_url,
                    "topic": question.get("topic", ""),
                    "topic_type": question.get("topic_type", ""),
                    "scores": scores
                }

                qa_pairs.append(qa_pair)

            except Exception as e:
                logger.error(f"Error generating answer for question '{question['text']}': {str(e)}")

        return qa_pairs

    def _fallback_search(self, question: str, source_url: str = "", top_k: int = 10) ->
    List[Dict]:
        """
        Fallback search method for when the primary search returns too few results.
        Uses more lenient matching and keyword-based approaches.
        """
        fallback_results = []

        # Extract key terms from the question (excluding stopwords)
        question_terms = [term.lower() for term in question.split()
                          if term.lower() not in self.resource_manager.stopwords]

        # Get all chunks
        all_chunks = []
        if source_url:
            # Prioritize chunks from source URL
            source_chunks = self.knowledge_base.get_chunks_by_url(source_url)
            all_chunks.extend([(chunk, 2.0) for chunk in source_chunks])  # Higher
    weight for source chunks

        # Add other chunks with lower weight
        other_chunks = [chunk for chunk in self.knowledge_base.chunks
                       if not source_url or chunk.doc_url != source_url]
        all_chunks.extend([(chunk, 1.0) for chunk in other_chunks])

        # Score chunks based on term overlap
        scored_chunks = []
        for chunk, base_weight in all_chunks:
            chunk_text = chunk.text.lower()

            # Count matching terms
            matches = sum(1 for term in question_terms if term in chunk_text)
            if matches > 0:
                # Normalize by total terms and apply base weight
                score = (matches / len(question_terms)) * base_weight
                scored_chunks.append({
                    'chunk': chunk,
                    'score': score,
                    'doc_url': chunk.doc_url
                })

        # Sort by score and take top_k
        scored_chunks.sort(key=lambda x: x['score'], reverse=True)
        fallback_results = scored_chunks[:top_k]

        logger.info(f"Fallback search found {len(fallback_results)} additional chunks")
        return fallback_results

    def _retrieve_context(self, question: str, source_url: str = "") -> List[Dict]:
        """
        Context retrieval that searches across the entire knowledge base
        and re-ranks results using a cross-encoder.
        """
        # Initial parameters for retrieval
        initial_top_k = 60
        final_top_k = 30

        # Store retrieved chunks with metadata
```

```
3203          retrieved_chunks = []

3204
3205          # STEP 1: Retrieve chunks from the entire knowledge base
3206          global_results = self.knowledge_base.search(question, self.embedding_model,
        top_k=initial_top_k)
3207
3208          # STEP 2: Process results and mark source-specific chunks
3209          for result in global_results:
3210              result['is_source'] = (result['doc_url'] == source_url)
3211              retrieved_chunks.append(result)

3212
3213          # STEP 3: If we have very few results, try to increase retrieval scope
3214          if len(retrieved_chunks) < 3:
3215              logger.warning(f"Retrieved only {len(retrieved_chunks)} chunks, increasing
        search scope")
3216              # Try with more chunks and lower similarity threshold
3217              additional_results = self._fallback_search(question, source_url, top_k=10)
3218
3219              # Add non-duplicate chunks
3220              seen_chunk_ids = {id(chunk['chunk']) for chunk in retrieved_chunks}
3221              for result in additional_results:
3222                  chunk_id = id(result['chunk'])
3223                  if chunk_id not in seen_chunk_ids:
3224                      result['is_source'] = (result['doc_url'] == source_url)
3225                      retrieved_chunks.append(result)
3226                      seen_chunk_ids.add(chunk_id)

3227
3228          # STEP 4: Re-rank using cross-encoder if available
3229          if self.reranker and len(retrieved_chunks) > 1:
3230              try:
3231                  # Prepare query-passage pairs for re-ranking
3232                  pairs = [(question, chunk['chunk'].text) for chunk in retrieved_chunks]
3233
3234                  # Get cross-encoder scores
3235                  cross_scores = self.reranker.predict(pairs)
3236
3237                  # Update scores with cross-encoder scores
3238                  for i, score in enumerate(cross_scores):
3239                      retrieved_chunks[i]['cross_score'] = float(score)
3240
3241                      # Combine vector similarity with cross-encoder score (weighted)
3242                      vector_score = retrieved_chunks[i]['score']
3243                      if isinstance(vector_score, torch.Tensor):
3244                          vector_score = vector_score.item()
3245
3246                      # Final score: 0.3 * vector_score + 0.7 * cross_score
3247                      retrieved_chunks[i]['final_score'] = 0.3 * vector_score + 0.7 *
        float(score)
3248
3249                  # Sort by final score
3250                  retrieved_chunks.sort(key=lambda x: x.get('final_score', 0), reverse=
        True)
3251
3252                  logger.info(f"Re-ranked {len(retrieved_chunks)} chunks using cross-
        encoder")
3253
3254              except Exception as e:
3255                  logger.error(f"Error in cross-encoder re-ranking: {str(e)}")
3256                  # Fall back to original scores
3257                  retrieved_chunks.sort(key=lambda x: x['score'], reverse=True)
3258          else:
3259              # If no re-ranker, sort by original scores
3260              retrieved_chunks.sort(key=lambda x: x['score'], reverse=True)

3261
3262          # STEP 5: Apply diversity selection to ensure representation from different
        documents
3263          diverse_chunks = self._select_diverse_chunks(retrieved_chunks, final_top_k)

3264
3265          # Log the final number of chunks used for context
3266          logger.info(f"Using {len(diverse_chunks)} diverse chunks for answering: '{
        question}'")

3267
3268          return diverse_chunks
```

```python
      def _select_diverse_chunks(self, chunks: List[Dict], max_chunks: int) -> List[Dict]:
          """
          Select a diverse set of chunks using a greedy algorithm that balances
          relevance and diversity across documents.
          """
          if len(chunks) <= max_chunks:
              return chunks

          # Track URLs and chunks already selected
          selected_chunks = []
          selected_urls = set()
          remaining_chunks = chunks.copy()

          # STEP 1: First select the highest scoring chunk overall
          if remaining_chunks:
              best_chunk = max(remaining_chunks, key=lambda x: x.get('final_score', x['
      score']))
              selected_chunks.append(best_chunk)
              selected_urls.add(best_chunk['doc_url'])
              remaining_chunks.remove(best_chunk)

          # STEP 2: Prioritize source URL chunks if any
          source_chunks = [c for c in remaining_chunks if c.get('is_source', False)]
          if source_chunks:
              # Take the best source chunk
              best_source = max(source_chunks, key=lambda x: x.get('final_score', x['score
      ']))
              if best_source not in selected_chunks:
                  selected_chunks.append(best_source)
                  remaining_chunks.remove(best_source)

          # STEP 3: Select chunks with a mix of relevance and diversity
          while len(selected_chunks) < max_chunks and remaining_chunks:
              # Calculate diversity bonus for each chunk
              for chunk in remaining_chunks:
                  # If this chunk is from a new URL, give it a diversity bonus
                  diversity_bonus = 0.3 if chunk['doc_url'] not in selected_urls else 0.0

                  # Calculate adjusted score with diversity bonus
                  base_score = chunk.get('final_score', chunk['score'])
                  chunk['adjusted_score'] = base_score + diversity_bonus

              # Select the chunk with the highest adjusted score
              best_chunk = max(remaining_chunks, key=lambda x: x['adjusted_score'])
              selected_chunks.append(best_chunk)
              selected_urls.add(best_chunk['doc_url'])
              remaining_chunks.remove(best_chunk)

          return selected_chunks

      def _is_just_disclaimer(self, text: str) -> bool:
          """Check if the answer is just a disclaimer without actual content."""
          disclaimer_patterns = [
              r"^I don't have (enough|sufficient) information to answer this question\.?\s
      *$",
              r"^There is not enough (context|information|data) (provided|available|given)
       to answer this question\.?\s*$",
              r"^Based on the (provided|given|available) (context|information), I cannot
      answer this question\.?\s*$"
          ]

          # Check if the text matches any disclaimer pattern
          for pattern in disclaimer_patterns:
              if re.match(pattern, text.strip()):
                  return True

          # Check length and disclaimer ratio
          words = text.split()
          if len(words) < 20 and "don't have" in text.lower():
              return True

          return False
```

```python
    def _generate_answer(self, question: str, context_chunks: List[Dict]) -> str:
        """Generate an answer using retrieved context with improved prompting."""
        if not context_chunks:
            return "I don't have enough information to answer this question."

        if not self.answer_generator or not self.ag_tokenizer:
            return "Answer generation model not available."

        try:
            # Get merged context with improved formatting
            context_text = self._merge_context_advanced(question, context_chunks)

            # Prompt with clear instructions
            prompt = (
                f"You are a helpful assistant answering a question based on provided "
                f"information sources. "
                f"Your task is to synthesize a complete, accurate answer using ONLY the "
                f"information in the context below. "
                f"Maintain a confident, direct tone and NEVER say 'I don't have enough "
                f"information' if you can provide "
                f"any relevant details from the context.\n\n"

                f"If the information is incomplete, simply share what IS available in "
                f"the context. "
                f"If the context doesn't address the question at all, ONLY THEN state "
                f"that "
                f"you don't have the specific information requested.\n\n"

                f"CONTEXT:\n{context_text}\n\n"

                f"INSTRUCTIONS FOR ANSWERING:\n"
                f"1. Read the context carefully and identify all relevant information\n"
                f"2. Synthesize the information into a coherent, complete answer\n"
                f"3. If information is partial, provide what's available without "
                f"disclaimers\n"
                f"4. Ensure your answer is fully supported by the context\n"
                f"5. Write in complete sentences with proper formatting\n\n"

                f"QUESTION: {question}\n\n"
                f"ANSWER:"
            )

            # Tokenize prompt with increased max length to handle larger context
            inputs = self.ag_tokenizer(prompt, return_tensors="pt", truncation=True,
max_length=3072)
            inputs = {k: v.to(self.device) for k, v in inputs.items()}


            outputs = self.answer_generator.generate(
                **inputs,
                max_length=1000,
                min_length=50,
                num_beams=4,
                num_beam_groups=1,
                num_return_sequences=2,
                diversity_penalty=0.0,
                do_sample=False,
                temperature=1.0,
                top_p=1.0,
                no_repeat_ngram_size=3,
                length_penalty=1.5,
                early_stopping=True
            )

            # Process candidates
            candidates = []
            for output in outputs:
                answer_text = self.ag_tokenizer.decode(output, skip_special_tokens=True)
                candidates.append(answer_text)

            # Select the best answer
            if candidates:
```

```python
                # Filter out candidates that are just disclaimers
                valid_candidates = [c for c in candidates if not self.
_is_just_disclaimer(c)]

                if valid_candidates:
                    # Choose the most substantive answer
                    best_answer = max(valid_candidates, key=lambda x: len(x) - 10 * x.
count("I don't have"))
                else:
                    best_answer = candidates[0]  # Fallback to first answer

                # Apply formatting
                return self._format_answer(best_answer)
            else:
                return "I couldn't generate an answer based on the available information
."

        except Exception as e:
            logger.error(f"Error generating answer: {str(e)}")
            return "I encountered an error while generating the answer."

    def _clean_chunk_text(self, text: str) -> str:
        """Clean chunk text of artifacts and normalize formatting."""
        # Remove HTML artifacts
        text = re.sub(r' |&amp;|&lt;|&gt;|&quot;', ' ', text)

        # Fix ellipsis and other punctuation
        text = re.sub(r'\.{2,}', '. ', text)

        # Normalize whitespace
        text = re.sub(r'\s+', ' ', text)

        # Fix newlines
        text = re.sub(r'\r\n|\r', '\n', text)
        text = re.sub(r'\n{3,}', '\n\n', text)

        # Remove artifacts common in the data
        text = re.sub(r'rn\s', ' ', text)
        text = re.sub(r'\.s\.', '.', text)

        # Fix broken markdown headers
        text = re.sub(r'#\s+', '## ', text)

        return text.strip()

    def _is_text_too_similar(self, text1: str, text2: str) -> bool:
        """Check if two texts are too similar to both include."""
        # For very short texts, use exact matching
        if len(text1) < 100 or len(text2) < 100:
            return text1 in text2 or text2 in text1

        # For longer texts, use n-gram similarity
        words1 = text1.split()
        words2 = text2.split()

        # Create 3-grams
        def get_ngrams(words, n=3):
            return set(' '.join(words[i:i+n]) for i in range(len(words)-n+1))

        ngrams1 = get_ngrams(words1)
        ngrams2 = get_ngrams(words2)

        if not ngrams1 or not ngrams2:
            return False

        # Calculate Jaccard similarity
        intersection = len(ngrams1.intersection(ngrams2))
        union = len(ngrams1.union(ngrams2))

        # Higher threshold to avoid removing related but distinct content
        return intersection / union > 0.8

    def _merge_context_advanced(self, question: str, context_chunks: List[Dict]) -> str:
```

```python
        """Advanced context merging with clear source boundaries and improved structure.
"""
        if not context_chunks:
            return ""

        # Sort chunks by relevance
        sorted_chunks = sorted(context_chunks, key=lambda x: x.get('final_score', x['
score']), reverse=True)

        # Calculate maximum context size based on model
        if self.answer_generator and hasattr(self.answer_generator, 'config'):
            if hasattr(self.answer_generator.config, 'model_type'):
                model_path = getattr(self.answer_generator.config, '_name_or_path', '').
lower()
                if 'xxl' in model_path:
                    max_context_chars = 30000
                    logger.info("Using expanded context size for XXL model")
                else:
                    max_context_chars = 18000
            else:
                max_context_chars = 18000
        else:
            max_context_chars = 18000

        # Group chunks by source document
        doc_chunks = {}
        for chunk_data in sorted_chunks:
            chunk = chunk_data['chunk']
            doc_url = chunk.doc_url
            if doc_url not in doc_chunks:
                doc_chunks[doc_url] = []
            doc_chunks[doc_url].append((chunk, chunk_data.get('final_score', chunk_data[
'score'])))

        # Assemble context with clear document sections
        context_parts = []
        current_length = 0

        # First add chunks from the highest scoring documents
        for doc_url, chunks in sorted(doc_chunks.items(),
                                    key=lambda x: max([score for _, score in x[1]]),
                                    reverse=True):
            # Sort chunks within this document by score
            doc_chunks_sorted = sorted(chunks, key=lambda x: x[1], reverse=True)

            # Extract domain for reference
            domain = doc_url.replace('https://', '').replace('http://', '').split('/')
[0]
            path = doc_url.split('/')[-1] if '/' in doc_url else ''

            # Clean and combine text from this document
            doc_texts = []
            for chunk, _ in doc_chunks_sorted:
                # Clean and normalize text
                text = self._clean_chunk_text(chunk.text)

                # Check if this would exceed our max context
                if current_length + len(text) + 100 > max_context_chars:
                    # If we already have content, just stop adding more
                    if doc_texts or context_parts:
                        break
                    # If this is the first chunk, take a portion to fit
                    truncated = text[:max_context_chars - 200] + "..."
                    doc_texts.append(truncated)
                    current_length += len(truncated)
                    break

                # Add text if not too similar to existing content
                if not any(self._is_text_too_similar(text, existing) for existing in
doc_texts):
                    doc_texts.append(text)
                    current_length += len(text)
```

```python
                # Only add this document section if we have content
                if doc_texts:
                    # Create section header
                    section_header = f"DOCUMENT: {domain}/{path}"
                    section_content = "\n\n".join(doc_texts)
                    section = f"{'='*50}\n{section_header}\n{'='*50}\n{section_content}"
                    context_parts.append(section)

        # Join all sections with clear separation
        full_context = "\n\n" + "\n\n".join(context_parts)

        # Add helpful metadata at the beginning
        context_intro = (
            f"QUESTION: {question}\n\n"
            f"The following information comes from {len(doc_chunks)} different sources "
    about this topic. "
            f"Use this information to construct a complete, accurate answer."
        )

        full_context = context_intro + full_context

        return full_context

    def _split_into_semantic_units(self, text: str) -> List[str]:
        """Split text into semantic units (paragraphs or coherent sections)."""
        # First try to split by paragraph breaks
        if '\n\n' in text:
            segments = [seg.strip() for seg in text.split('\n\n') if seg.strip()]
            # Filter out very short segments and merge them with adjacent ones
            filtered_segments = []
            current_segment = ""

            for segment in segments:
                if len(segment) < 50:  # Short segment
                    current_segment += " " + segment
                else:
                    if current_segment:
                        filtered_segments.append(current_segment)
                        current_segment = segment
                    else:
                        current_segment = segment

            # Add the last segment if it exists
            if current_segment:
                filtered_segments.append(current_segment)

            return filtered_segments if filtered_segments else [text]

        # If no paragraph breaks, try to use sentence tokenization
        tokenize_fn = self.resource_manager.get_fallback_tokenize()
        sentences = tokenize_fn(text)

        if len(sentences) <= 3:
            # Text is already small enough
            return [text]

        # Group sentences into coherent segments
        segments = []
        current_segment = []
        for sentence in sentences:
            current_segment.append(sentence)
            if len(current_segment) >= 4:
                segments.append(" ".join(current_segment))
                current_segment = []

        # Add the last segment if it exists
        if current_segment:
            segments.append(" ".join(current_segment))

        return segments if segments else [text]

    def _eliminate_redundancy(self, blocks: List[Dict]) -> List[Dict]:
        """Detect and eliminate redundant content from blocks."""
```

```
3612          if len(blocks) <= 1:
3613              return blocks
3614
3615          # Sort blocks by score (highest first)
3616          sorted_blocks = sorted(blocks, key=lambda x: x['score'], reverse=True)
3617
3618          # Calculate similarity threshold
3619          similarity_threshold = 0.7
3620
3621          # Track which blocks are redundant
3622          is_redundant = [False] * len(sorted_blocks)
3623
3624          # For each high-scoring block, check if lower-scoring blocks are redundant
3625          for i in range(len(sorted_blocks) - 1):
3626              if is_redundant[i]:
3627                  continue
3628
3629              block_i = sorted_blocks[i]['text'].lower()
3630
3631              for j in range(i + 1, len(sorted_blocks)):
3632                  if is_redundant[j]:
3633                      continue
3634
3635                  block_j = sorted_blocks[j]['text'].lower()
3636
3637                  # Simple n-gram based similarity check
3638                  similarity = self._calculate_text_similarity(block_i, block_j)
3639
3640                  # Mark as redundant if similarity is high
3641                  if similarity > similarity_threshold:
3642                      # If blocks are from same chunk, always mark the lower-scoring one
3643                      if sorted_blocks[i]['chunk_id'] == sorted_blocks[j]['chunk_id']:
3644                          is_redundant[j] = True
3645                      else:
3646                          score_ratio = sorted_blocks[j]['score'] / sorted_blocks[i]['
     score']
3647                          if score_ratio > 0.85 and len(block_j) < len(block_i) * 0.8:
3648                              # Keep the shorter block if scores are comparable
3649                              continue
3650                          else:
3651                              is_redundant[j] = True
3652
3653          # Return non-redundant blocks
3654          return [block for i, block in enumerate(sorted_blocks) if not is_redundant[i]]
3655
3656      def _calculate_text_similarity(self, text1: str, text2: str) -> float:
3657          """Calculate similarity between two text blocks using n-grams."""
3658          # Use tf-idf weighted bigram similarity for a balance of efficiency and accuracy
3659
3660          # Get bigrams
3661          def get_bigrams(text):
3662              words = text.split()
3663              return set(" ".join(words[i:i+2]) for i in range(len(words)-1))
3664
3665          bigrams1 = get_bigrams(text1)
3666          bigrams2 = get_bigrams(text2)
3667
3668          # Handle empty sets
3669          if not bigrams1 or not bigrams2:
3670              return 0.0
3671
3672          # Calculate Jaccard similarity
3673          intersection = len(bigrams1.intersection(bigrams2))
3674          union = len(bigrams1.union(bigrams2))
3675
3676          return intersection / union
3677
3678      def _select_context_blocks(self, question: str, blocks: List[Dict]) -> List[Dict]:
3679          """
3680          Select context blocks to include, balancing relevance and diversity
3681          while managing context length.
3682          """
3683          if self.answer_generator and hasattr(self.answer_generator, 'config'):
```

```python
            if hasattr(self.answer_generator.config, 'model_type'):
                model_path = getattr(self.answer_generator.config, '_name_or_path', '').
    lower()
                if 'xxl' in model_path:
                    max_context_length = 3000
                elif 'xl' in model_path:
                    max_context_length = 2000
                elif 'large' in model_path:
                    max_context_length = 1500
                else:
                    max_context_length = 1000
            else:
                max_context_length = 1500
        else:
            max_context_length = 1500

        # If we have few blocks, include them all if they fit
        total_length = sum(len(block['text']) for block in blocks)
        if total_length <= max_context_length:
            return blocks

        # We need to be selective - first prioritize blocks from different documents
        doc_urls = set(block['doc_url'] for block in blocks)

        selected_blocks = []
        remaining_blocks = blocks.copy()
        current_length = 0

        # First phase: Select highest scoring block from each document
        for url in doc_urls:
            doc_blocks = [b for b in remaining_blocks if b['doc_url'] == url]
            if doc_blocks:
                best_block = max(doc_blocks, key=lambda x: x['score'])
                if current_length + len(best_block['text']) <= max_context_length:
                    selected_blocks.append(best_block)
                    remaining_blocks.remove(best_block)
                    current_length += len(best_block['text'])

        # Second phase: Select additional blocks based on score, with diminishing
    returns
        # for blocks from the same document
        remaining_blocks.sort(key=lambda x: x['score'], reverse=True)

        # Count blocks per document
        doc_counts = Counter(block['doc_url'] for block in selected_blocks)

        # Adjust scores based on document representation
        for block in remaining_blocks:
            # Apply penalty for documents that are already well-represented
            doc_count = doc_counts.get(block['doc_url'], 0)
            adjusted_score = block['score'] * (0.95 ** doc_count)  # Diminishing returns
            block['adjusted_score'] = adjusted_score

        # Re-sort with adjusted scores
        remaining_blocks.sort(key=lambda x: x.get('adjusted_score', x['score']), reverse
    =True)

        # Add blocks until we reach the length limit
        for block in remaining_blocks:
            if current_length + len(block['text']) <= max_context_length:
                selected_blocks.append(block)
                current_length += len(block['text'])
                # Update document count
                doc_counts[block['doc_url']] = doc_counts.get(block['doc_url'], 0) + 1
            else:
                # Try to fit as much relevant content as possible
                space_left = max_context_length - current_length
                if space_left > 200:  # Only add if we can fit something substantial
                    # Truncate the block to fit
                    truncated_text = block['text'][:space_left].rsplit('.', 1)[0] + '.'
                    if len(truncated_text) > 100:  # Only add if it's still meaningful
                        block['text'] = truncated_text
                        selected_blocks.append(block)
```

```python
3754                    break
3755
3756            return selected_blocks
3757
3758        def _assemble_context(self, blocks: List[Dict]) -> str:
3759            """Assemble the final context, organized by document source."""
3760            if not blocks:
3761                return ""
3762
3763            # Group blocks by document URL
3764            blocks_by_url = {}
3765            for block in blocks:
3766                url = block['doc_url']
3767                if url not in blocks_by_url:
3768                    blocks_by_url[url] = []
3769                blocks_by_url[url].append(block)
3770
3771            # Assemble context with source information
3772            context_parts = []
3773
3774            for url, url_blocks in blocks_by_url.items():
3775                # Sort blocks from same document by their original order
3776                url_blocks.sort(key=lambda x: (x['chunk_id'], x['segment_id']))
3777
3778                # Combine text from this document
3779                doc_text = " ".join(block['text'] for block in url_blocks)
3780
3781                # Add document source indicator for multi-document context
3782                if len(blocks_by_url) > 1:
3783                    # Extract domain for cleaner reference
3784                    domain = url.replace('https://', '').replace('http://', '').split('/')
        [0]
3785                    source_indicator = f"[From: {domain}] "
3786                    context_parts.append(source_indicator + doc_text)
3787                else:
3788                    context_parts.append(doc_text)
3789
3790            # Join all parts
3791            return "\n\n".join(context_parts)
3792
3793        def _prepare_context(self, question: str, chunks: List[Dict]) -> str:
3794            """
3795            Legacy method kept for compatibility.
3796            Now calls the advanced context merging function.
3797            """
3798            return self._merge_context_advanced(question, chunks)
3799
3800        def _extract_answer_from_cot(self, text: str) -> str:
3801            """Extract the final answer from a chain-of-thought generation."""
3802            # If text contains step markers, extract the part after the last step
3803            step_matches = list(re.finditer(r'Step \d+:', text))
3804            if step_matches:
3805                last_step_match = step_matches[-1]
3806                last_step_end = last_step_match.end()
3807
3808                # Find the next step or the end of text
3809                next_step_match = re.search(r'Step \d+:', text[last_step_end:])
3810                if next_step_match:
3811                    # Extract between last step and next step
3812                    answer_text = text[last_step_end:last_step_end + next_step_match.start()
        ].strip()
3813                else:
3814                    # Extract from last step to end
3815                    answer_text = text[last_step_end:].strip()
3816
3817                return answer_text
3818
3819            # If we see a clear "Answer:" marker
3820            answer_match = re.search(r'Answer:(.*?)(?:$|Step \d+:)', text, re.DOTALL)
3821            if answer_match:
3822                return answer_match.group(1).strip()
3823
3824            # If none of the above patterns match, just return the original text
```

```
3825        return text
3826
3827    def _select_best_answer(self, question: str, candidates: List[str], context: str) ->
         str:
3828        """Select the best answer from multiple candidates based on quality and
       consistency."""
3829        if not candidates:
3830            return ""
3831
3832        if len(candidates) == 1:
3833            return candidates[0]
3834
3835        # 1. Score candidates by length (prefer longer, more detailed answers)
3836        length_scores = []
3837        for candidate in candidates:
3838            # Normalize length (prefer answers between 50-200 characters)
3839            length = len(candidate)
3840            if length < 20:
3841                score = length / 20  # Penalize very short answers
3842            elif length < 50:
3843                score = 0.5 + (length - 20) / 60  # Ramp up to 0.5-1.0
3844            elif length <= 200:
3845                score = 1.0  # Ideal length
3846            else:
3847                score = 1.0 - (length - 200) / 800  # Gradually penalize very long
       answers
3848                score = max(0.5, score)  # Don't go below 0.5
3849
3850            length_scores.append(score)
3851
3852        # 2. Check factual consistency with context
3853        factual_scores = []
3854        for candidate in candidates:
3855            # Simple heuristic: count overlapping n-grams with context
3856            candidate_ngrams = self._get_ngrams(candidate.lower(), 2)
3857            context_ngrams = self._get_ngrams(context.lower(), 2)
3858
3859            overlap = len(candidate_ngrams.intersection(context_ngrams))
3860            total = len(candidate_ngrams)
3861
3862            if total == 0:
3863                factual_scores.append(0.0)
3864            else:
3865                factual_scores.append(min(1.0, overlap / total))
3866
3867        # 3. Combined scoring
3868        final_scores = []
3869        for i in range(len(candidates)):
3870            # Weight factual consistency more heavily
3871            score = 0.3 * length_scores[i] + 0.7 * factual_scores[i]
3872            final_scores.append(score)
3873
3874        # Return the candidate with the highest score
3875        best_idx = final_scores.index(max(final_scores))
3876        return candidates[best_idx]
3877
3878    def _get_ngrams(self, text: str, n: int) -> set:
3879        """Get n-grams from text."""
3880        words = text.split()
3881        ngrams = set()
3882        for i in range(len(words) - n + 1):
3883            ngram = ' '.join(words[i:i+n])
3884            ngrams.add(ngram)
3885        return ngrams
3886
3887    def _fix_structural_issues(self, text: str) -> str:
3888        """Fix structural issues in the answer."""
3889        # Remove any instruction artifacts
3890        text = re.sub(r'^(Answer:|ANSWER:|Step \d+:|To answer this question:)', '', text
       ).strip()
3891
3892        # Remove meta-commentary about the answering process
```

```python
        text = re.sub(r'(Based on the (provided|given|available) (context|information)
,?\s*)', '', text).strip()
        text = re.sub(r'(According to the (provided|given|available) (context|
information),?\s*)', '', text).strip()

        # Remove duplicate phrases at the beginning
        words = text.split()
        if len(words) > 10:
            first_5 = ' '.join(words[:5]).lower()
            for i in range(1, min(10, len(words) - 5)):
                next_5 = ' '.join(words[i:i+5]).lower()
                if first_5 == next_5:
                    text = ' '.join(words[i:])
                    break

        return text

    def _clean_formatting_artifacts(self, text: str) -> str:
        """Clean formatting artifacts from the answer."""
        # Fix newlines and spacing
        text = re.sub(r'\r\n|\r', '\n', text)
        text = re.sub(r'\n{2,}', '\n\n', text)
        text = re.sub(r'[ \t]+', ' ', text)

        # Fix ellipses and quotes
        text = re.sub(r'\.{2,}', '. ', text)
        text = re.sub(r'"{2,}', '"', text)
        text = re.sub(r"'{2,}", "'", text)

        # Fix common artifacts
        text = re.sub(r'\.s\.', '.', text)
        text = re.sub(r'\bn\s', ' ', text)
        text = re.sub(r'rn', '', text)

        # Fix list formatting
        text = re.sub(r'(\d+)\)\s*', r'\1. ', text)

        return text.strip()

    def _fix_truncated_ending(self, text: str) -> str:
        """Fix truncated endings to ensure complete sentences."""
        if not text:
            return text

        # If text doesn't end with sentence-ending punctuation
        if not re.search(r'[.!?]$', text):
            # Try to find the last complete sentence
            last_period = max(
                text.rfind('.'),
                text.rfind('!'),
                text.rfind('?')
            )

            if last_period > 0.7 * len(text):
                # If we have most of the content, trim to last complete sentence
                return text[:last_period+1]
            else:
                # Otherwise add a period to the existing text
                return text + "."

        return text

    def _fix_capitalization_punctuation(self, text: str) -> str:
        """Fix capitalization and punctuation in the answer."""
        # Capitalize first letter
        if text and text[0].isalpha() and not text[0].isupper():
            text = text[0].upper() + text[1:]

        # Remove filler conclusions
        conclusion_patterns = [
            r' I hope (this|that) (helps|answers your question)\.?$',
            r' (Please )?[Ll]et me know if you (have|need) (any|more|further|additional)
 (questions|information|details)\.?$',
```

```
3963              r' (Feel free to|Please) (ask|contact|reach out)( me)?( if you have| for)? (
       more|any other|additional) questions\.?$'
3964          ]
3965
3966          for pattern in conclusion_patterns:
3967              text = re.sub(pattern, '.', text)
3968
3969          return text
3970
3971      def _format_answer(self, answer: str) -> str:
3972          """Format the answer with improved handling of disclaimers and artifacts."""
3973          if not answer:
3974              return "I don't have specific information about this in the available
       sources."
3975
3976          # Extract meaningful content after disclaimers
3977          if answer.startswith("I don't have enough information"):
3978              # Look for content after the disclaimer
3979              parts = answer.split(".", 1)
3980              if len(parts) > 1 and len(parts[1].strip()) > 30:
3981                  # If substantial content follows the disclaimer, use it
3982                  answer = parts[1].strip()
3983              else:
3984                  # If there's nothing substantial, keep the disclaimer but improve it
3985                  return "Based on the available sources, I don't have specific
       information to answer this question completely."
3986
3987          # Fix answer structure issues
3988          answer = self._fix_structural_issues(answer)
3989
3990          # Clean up formatting artifacts
3991          answer = self._clean_formatting_artifacts(answer)
3992
3993          # Ensure the answer doesn't end mid-sentence
3994          answer = self._fix_truncated_ending(answer)
3995
3996          # Ensure proper capitalization and punctuation
3997          answer = self._fix_capitalization_punctuation(answer)
3998
3999          return answer
4000
4001      def _evaluate_answer(self, question: str, answer: str, context_chunks: List[Dict])
       -> Dict[str, float]:
4002          """Evaluate answer quality with enhanced checks for formatting issues and
       completeness."""
4003          # Convert context_chunks format for compatibility with the original method
4004          semantic_chunks = [chunk_data['chunk'] for chunk_data in context_chunks]
4005
4006          scores = {
4007              'relevance': 0.0,
4008              'factuality': 0.0,
4009              'completeness': 0.0,
4010              'formatting': 0.0,
4011              'overall': 0.0
4012          }
4013
4014          # Check for empty or very short answers
4015          if not answer or len(answer) < 15:
4016              return scores
4017
4018          # Check if we have limited context - adjust scoring if needed
4019          limited_context = len(context_chunks) <= 1
4020
4021          # 1. Relevance score - more lenient with limited content
4022          relevance = self._score_relevance(question, answer)
4023          if limited_context:
4024              # Boost relevance score for limited content
4025              relevance = min(1.0, relevance * 1.2)
4026          scores['relevance'] = relevance
4027
4028          # 2. Factuality score - adjust for limited content
4029          factuality = self._score_factuality(answer, semantic_chunks)
4030          if limited_context:
```

```
4031              # Set a minimum factuality score for limited content
4032              factuality = max(0.5, factuality)
4033          scores['factuality'] = factuality
4034
4035          # 3. Completeness score - with enhanced checks for truncation
4036          completeness = self._score_completeness(question, answer)
4037
4038          # Check for truncated answers or formatting issues
4039          formatting_score = 1.0
4040
4041          # Penalize answers that appear truncated
4042          if len(answer) > 20 and not answer.endswith(('.', '!', '?')):
4043              completeness *= 0.8
4044              formatting_score *= 0.7
4045
4046          # Penalize answers with "I don't have enough information" followed by content
4047          if "I don't have enough information" in answer and len(answer) > 70:
4048              completeness *= 0.9
4049              formatting_score *= 0.8
4050
4051          # Penalize answers with strange formatting artifacts
4052          formatting_artifacts = [
4053              r'\.{3,}',
4054              r'\brn\b',
4055              r'\.s\.',
4056              r'\n{3,}',
4057              r'\s{3,}'
4058          ]
4059
4060          for pattern in formatting_artifacts:
4061              if re.search(pattern, answer):
4062                  formatting_score *= 0.85
4063
4064          if limited_context:
4065              # Boost completeness for limited content
4066              completeness = min(1.0, completeness * 1.2)
4067
4068          scores['completeness'] = completeness
4069          scores['formatting'] = formatting_score
4070
4071          # 4. Calculate overall score with weighted combination
4072          if limited_context:
4073              # For limited content, weigh relevance more heavily
4074              overall = (0.35 * relevance + 0.25 * factuality + 0.25 * completeness + 0.15
      * formatting_score)
4075          else:
4076              # Standard weighting
4077              overall = (0.25 * relevance + 0.35 * factuality + 0.25 * completeness + 0.15
      * formatting_score)
4078
4079          scores['overall'] = overall
4080
4081          return scores
4082
4083      def _score_relevance(self, question: str, answer: str) -> float:
4084          """Score the relevance of an answer to the question."""
4085          # Method 1: Check for question terms in the answer
4086          question_terms = set(question.lower().split()) - self.resource_manager.stopwords
4087          answer_lower = answer.lower()
4088
4089          # Count matching terms
4090          matching_terms = sum(1 for term in question_terms if term in answer_lower)
4091          term_score = min(1.0, matching_terms / max(1, len(question_terms)))
4092
4093          # Method 2: Use semantic model if available
4094          if self.qa_evaluator:
4095              try:
4096                  # Use the model to score the pair
4097                  inputs = self.qa_tokenizer([question], [answer], return_tensors="pt",
      padding=True, truncation=True)
4098                  inputs = {k: v.to(self.device) for k, v in inputs.items()}
4099
4100                  outputs = self.qa_evaluator(**inputs)
```

```
4101                    logits = outputs.logits
4102
4103                    # Convert logits to probability
4104                    if logits.shape[1] > 1:
4105                        # Multi-class classification
4106                        probs = torch.softmax(logits, dim=1)
4107                        model_score = probs[0, 1].item()  # Assume second class is "relevant
       "
4108                    else:
4109                        # Binary classification
4110                        prob = torch.sigmoid(logits)
4111                        model_score = prob.item()
4112
4113                    # Combine model score with term score
4114                    return 0.7 * model_score + 0.3 * term_score
4115
4116            except Exception as e:
4117                logger.error(f"Error using QA evaluator: {str(e)}")
4118                return term_score
4119
4120        return term_score
4121
4122    def _score_factuality(self, answer: str, context_chunks: List[SemanticChunk]) ->
       float:
4123        """Score the factual consistency of the answer with the context."""
4124        if not context_chunks:
4125            return 0.5  # Neutral if no context
4126
4127        # Combine context
4128        context = " ".join(chunk.text for chunk in context_chunks)
4129
4130        # Use fact checking model if available
4131        if self.fact_checker and self.fact_checker_tokenizer:
4132            try:
4133                # Merge claim + partial context into a single string
4134                combined_text = f"Claim: {answer}\nContext: {context[:1000]}"
4135                inputs = self.fact_checker_tokenizer(
4136                    combined_text, return_tensors="pt", truncation=True, max_length=512
4137                )
4138                inputs = {k: v.to(self.device) for k, v in inputs.items()}
4139                outputs = self.fact_checker(**inputs)
4140
4141                # Check outputs.logits
4142                if hasattr(outputs, "logits"):
4143                    logits = outputs.logits
4144                    # If single-logit, interpret > 0 => more truthful
4145                    if logits.shape[-1] == 1:
4146                        prob = torch.sigmoid(logits).item()
4147                        return prob
4148                    else:
4149                        # Multi-logit => assume index 1 is "not hallucinated" or "
       truthful"
4150                        probs = torch.softmax(logits, dim=-1)
4151                        return probs[0, 1].item()
4152
4153                return 0.5
4154
4155            except Exception as e:
4156                logger.error(f"Error using fact checker: {str(e)}")
4157                # Fall back to n-gram overlap below
4158
4159        # If fact checker not available or it failed
4160        answer_ngrams = self._get_ngrams(answer.lower(), 2)
4161        context_ngrams = self._get_ngrams(context.lower(), 2)
4162        if not answer_ngrams:
4163            return 0.0
4164
4165        overlap = len(answer_ngrams.intersection(context_ngrams))
4166        score = overlap / len(answer_ngrams)
4167        return 0.2 + (score * 0.8)
4168
4169
4170    def _score_completeness(self, question: str, answer: str) -> float:
```

```
            """
        Score how completely the answer addresses all aspects of the question
        with enhanced detection of truncated or malformed answers.
        """
        # Base score
        score = 0.5

        # 1. Length-based scoring (very short answers are likely incomplete)
        words = len(answer.split())
        if words < 15:
            score -= 0.3
        elif words > 50:
            score += 0.1

        # 2. Check for truncated answers
        if not answer.endswith(('.', '!', '?')) and len(answer) > 20:
            score -= 0.15

        # 3. Detect formatting issues or artifacts
        if re.search(r'\brn\b|\.{3,}|\.s\.|\s{3,}', answer):
            score -= 0.1

        # 4. Check for "I don't have enough information" pattern
        if answer.startswith("I don't have enough information"):
            # If it's just that phrase or very little after it
            if len(answer) < 70:
                score -= 0.2
            else:
                # If it has substantial content after the disclaimer
                score -= 0.1

        # 5. Check for question type and expected answer elements
        question_lower = question.lower()

        # "What" questions typically define or explain something
        if question_lower.startswith("what"):
            if re.search(r'is|are|was|were', question_lower[:15]):
                # Definition question - answer should define the topic
                if re.search(r'(is|are|refers to|defined as|means)', answer.lower()
    [:50]):
                    score += 0.2

        # "How" questions explain a process or method
        elif question_lower.startswith("how"):
            # Process question - answer should include steps or a method
            if re.search(r'(first|second|then|next|finally|by|through)', answer.lower())
    :
                score += 0.2

        # "Where" questions should mention a location
        elif question_lower.startswith("where"):
            # Location question - answer should mention a place
            if re.search(r'(located|at|in|on|near|building|room|floor|campus)', answer.
    lower()):
                score += 0.2

        # "Who" questions should mention a person or organization
        elif question_lower.startswith("who"):
            # Person/org question - answer should mention a name or title
            if re.search(r'(staff|faculty|office|center|department|director|coordinator)
    ', answer.lower()):
                score += 0.2

        # Ensure score is in valid range
        return max(0.0, min(1.0, score))

    def generate_qa_pairs(self, urls: List[str], max_pairs_per_url: int = 10) -> List[
    Dict]:
        """Generate high-quality QA pairs for a list of URLs with adaptive strategies
    for limited content."""
        # Store original filter method for potential restoration
        if not hasattr(self, '_original_filter_qa_pairs'):
            self._original_filter_qa_pairs = self.filter_qa_pairs
```

```python
        # 1. Generate questions
        questions = self.generate_questions_from_documents(urls, max_pairs_per_url)

        # Log progress
        logger.info(f"Generated {len(questions)} questions for {len(urls)} URLs")

        # Check if we have very few questions - adapt strategy if needed
        if len(questions) < 5 and len(urls) > 0:
            logger.warning("Very few questions generated, using adaptive strategies")

            # Generate more general questions that require less specific content
            general_questions = self._generate_adaptive_questions(urls)

            # Add to questions list
            questions.extend(general_questions)
            logger.info(f"Added {len(general_questions)} adaptive questions, total: {len(questions)}")

        # 2. Generate answers
        qa_pairs = self.generate_answers(questions)

        # Log progress
        logger.info(f"Generated {len(qa_pairs)} QA pairs")

        # 3. Filter by quality - adaptive threshold based on content quantity
        if len(qa_pairs) < 5:
            # Set lower quality threshold for limited content
            logger.warning("Few QA pairs generated, lowering quality threshold")
            filtered_pairs = self.filter_qa_pairs(qa_pairs, min_score=0.4)
        else:
            # Use standard threshold
            filtered_pairs = self.filter_qa_pairs(qa_pairs)

        # If we still don't have enough pairs, take the best regardless of threshold
        if len(filtered_pairs) < 3 and len(qa_pairs) > 0:
            logger.warning("Very few QA pairs after filtering, taking best available")
            qa_pairs.sort(key=lambda x: x.get("scores", {}).get("overall", 0), reverse=True)
            filtered_pairs = qa_pairs[:min(5, len(qa_pairs))]

        # Log results
        logger.info(f"Final QA pairs after filtering: {len(filtered_pairs)}")

        return filtered_pairs

    def filter_qa_pairs(self, qa_pairs: List[Dict], min_score: float = 0.5) -> List[Dict]:
        """Filter QA pairs by quality scores and remove duplicates, with adaptive
        threshold for limited content."""
        if not qa_pairs:
            return []

        # Determine if we have limited content
        limited_content = len(qa_pairs) < 5

        # Adjust minimum score based on content availability
        if limited_content:
            adjusted_min_score = min_score * 0.8  # 20% lower threshold for limited
content
            logger.info(f"Limited content detected, adjusting quality threshold to {adjusted_min_score:.2f}")
        else:
            adjusted_min_score = min_score

        # Filter by minimum quality score
        quality_pairs = [pair for pair in qa_pairs
                         if pair.get("scores", {}).get("overall", 0) >=
        adjusted_min_score]

        # If we have very few pairs after filtering, accept lower quality ones
        if len(quality_pairs) < 3 and len(qa_pairs) > 3:
            # Sort by quality and take top 3 regardless of threshold
```

```python
            sorted_pairs = sorted(qa_pairs, key=lambda x: x.get("scores", {}).get("
    overall", 0), reverse=True)
            quality_pairs = sorted_pairs[:3]
            logger.info(f"Few high-quality pairs, accepting top {len(quality_pairs)}
    pairs regardless of threshold")

        # Group by URL
        url_to_pairs = {}
        for pair in quality_pairs:
            url = pair.get("source_url", "")
            if url not in url_to_pairs:
                url_to_pairs[url] = []
            url_to_pairs[url].append(pair)

        # For each URL, deduplicate and select best pairs
        final_pairs = []

        for url, pairs in url_to_pairs.items():
            # Sort by overall score
            pairs.sort(key=lambda x: x.get("scores", {}).get("overall", 0), reverse=True
    )

            # Select unique pairs (avoid answer duplication)
            unique_pairs = []
            seen_answers = set()

            for pair in pairs:
                answer_key = self._get_answer_signature(pair["answer"])
                if answer_key not in seen_answers:
                    unique_pairs.append(pair)
                    seen_answers.add(answer_key)

            final_pairs.extend(unique_pairs)

        # Log summary of filtering
        logger.info(f"QA filtering: {len(qa_pairs)} original pairs -> {len(final_pairs)}
     final pairs")

        return final_pairs

    def _get_answer_signature(self, answer: str) -> str:
        """Create a signature for an answer to identify near-duplicates."""
        # Remove stopwords and normalize
        words = [w.lower() for w in answer.split() if w.lower() not in self.
    resource_manager.stopwords]

        # Sort to make order-independent
        words.sort()

        # Take first 10 words as signature
        signature = " ".join(words[:10])

        return signature

#---------------------------------------------------------------------
# Complete QA Generation System with Checkpointing
#---------------------------------------------------------------------

class QAGenerationSystem:
    """Complete QA generation system with checkpointing and evaluation."""

    def __init__(self,
                 base_url: str,
                 output_dir: str,
                 use_gpu: bool = True,
                 checkpoint_dir: str = "checkpoints"):
        """Initialize the QA generation system."""
        self.base_url = base_url
        self.output_dir = output_dir
        self.use_gpu = use_gpu
        self.checkpoint_dir = checkpoint_dir

        # Create output and checkpoint directories
```

```python
        os.makedirs(output_dir, exist_ok=True)
        os.makedirs(checkpoint_dir, exist_ok=True)

        # Create unique ID for this run based on base URL
        self.run_id = hashlib.md5(base_url.encode('utf-8')).hexdigest()

        # Initialize resource manager
        self.resource_manager = ResourceManager()

        # Initialize model manager
        self.model_manager = ModelManager(use_gpu=use_gpu)

        # Initialize knowledge base
        self.knowledge_base = KnowledgeBase(use_gpu=use_gpu)

        # Initialize crawler with longer delay (2.0 seconds instead of 1.0)
        self.crawler = PersistentCrawler(base_url, delay=2.0, checkpoint_dir=
    checkpoint_dir)

        # Set common paths that should exist on most sites
        self.crawler.priority_paths = [
            "/",
            "/index.html",
            "/about",
            "/contact",
            "/services",
            "/resources"
        ]

        # Initialize document processor
        self.doc_processor = DocumentProcessor(self.resource_manager, self.model_manager
    , use_gpu=use_gpu)

        # Initialize QA generator
        self.qa_generator = QAGenerator(self.resource_manager, self.model_manager, self.
    knowledge_base, use_gpu=use_gpu)

        # Track progress
        self.progress = {
            "resources_setup": False,
            "crawling": False,
            "document_processing": False,
            "knowledge_base": False,
            "qa_generation": False,
            "evaluation": False
        }

    def _get_progress_path(self) -> str:
        """Get the path for the progress checkpoint file."""
        return os.path.join(self.checkpoint_dir, f"progress_{self.run_id}.json")

    def _get_knowledge_base_path(self) -> str:
        """Get the path for the knowledge base checkpoint file."""
        return os.path.join(self.checkpoint_dir, f"kb_{self.run_id}.json")

    def _get_qa_pairs_path(self) -> str:
        """Get the path for the QA pairs checkpoint file."""
        return os.path.join(self.checkpoint_dir, f"qa_pairs_{self.run_id}.json")

    def save_progress(self) -> None:
        """Save current progress to checkpoint file."""
        try:
            progress_path = self._get_progress_path()

            # Save progress
            with open(progress_path, 'w', encoding='utf-8') as f:
                json.dump({
                    "base_url": self.base_url,
                    "progress": self.progress,
                    "timestamp": datetime.now().isoformat()
                }, f)

            logger.info(f"Progress saved to {progress_path}")
```

```python
4442
4443        except Exception as e:
4444            logger.error(f"Error saving progress: {str(e)}")
4445
4446    def load_progress(self) -> bool:
4447        """Load progress from checkpoint file."""
4448        try:
4449            progress_path = self._get_progress_path()
4450
4451            if not os.path.exists(progress_path):
4452                logger.info("No progress checkpoint found")
4453                return False
4454
4455            # Load progress
4456            with open(progress_path, 'r', encoding='utf-8') as f:
4457                data = json.load(f)
4458
4459            # Verify base URL
4460            if data.get("base_url") != self.base_url:
4461                logger.warning(f"Progress checkpoint is for a different URL: {data.get('
    base_url')}")
4462                return False
4463
4464            # Load progress
4465            self.progress = data.get("progress", {})
4466
4467            logger.info(f"Progress loaded from {progress_path}")
4468            return True
4469
4470        except Exception as e:
4471            logger.error(f"Error loading progress: {str(e)}")
4472            return False
4473
4474    def save_qa_pairs(self, qa_pairs: List[Dict]) -> None:
4475        """Save QA pairs to checkpoint file."""
4476        try:
4477            qa_path = self._get_qa_pairs_path()
4478
4479            # Save QA pairs
4480            with open(qa_path, 'w', encoding='utf-8') as f:
4481                json.dump(qa_pairs, f)
4482
4483            logger.info(f"QA pairs saved to {qa_path}")
4484
4485        except Exception as e:
4486            logger.error(f"Error saving QA pairs: {str(e)}")
4487
4488    def load_qa_pairs(self) -> List[Dict]:
4489        """Load QA pairs from checkpoint file."""
4490        try:
4491            qa_path = self._get_qa_pairs_path()
4492
4493            if not os.path.exists(qa_path):
4494                logger.info("No QA pairs checkpoint found")
4495                return []
4496
4497            # Load QA pairs
4498            with open(qa_path, 'r', encoding='utf-8') as f:
4499                qa_pairs = json.load(f)
4500
4501            logger.info(f"Loaded {len(qa_pairs)} QA pairs from {qa_path}")
4502            return qa_pairs
4503
4504        except Exception as e:
4505            logger.error(f"Error loading QA pairs: {str(e)}")
4506            return []
4507
4508    def setup_resources(self) -> bool:
4509        """Set up all required resources."""
4510        try:
4511            # Skip if already done
4512            if self.progress.get("resources_setup", False):
4513                logger.info("Resources already set up, skipping")
```

```
                        return True

            # Set up NLTK resources
            logger.info("Setting up NLTK resources")
            self.resource_manager.setup_nltk()

            # Set up spaCy
            logger.info("Setting up spaCy")
            self.resource_manager.setup_spacy()

            # Set up HuggingFace access
            logger.info("Setting up HuggingFace access")
            self.resource_manager.setup_huggingface_access()

            # Load document processing models
            logger.info("Loading document processing models")
            self.doc_processor.load_models()

            # Load QA generation models
            logger.info("Loading QA generation models")
            self.qa_generator.load_models()

            # Update progress
            self.progress["resources_setup"] = True
            self.save_progress()

            return True

        except Exception as e:
            logger.error(f"Error setting up resources: {str(e)}")
            return False

    def crawl_website(self, max_pages: int = 30) -> bool:
        """Crawl website and extract content."""
        try:
            # Skip if already done
            if self.progress.get("crawling", False):
                logger.info("Crawling already done, skipping")
                return True

            # Try to load checkpoint first
            checkpoint_loaded = self.crawler.load_checkpoint()

            if checkpoint_loaded:
                # Validate checkpoint
                self.crawler.validate_checkpoint()

                # If we have enough pages, skip crawling
                if len(self.crawler.content_cache) >= max_pages:
                    logger.info(f"Checkpoint loaded with {len(self.crawler.content_cache
    )} pages, skipping crawl")

                    # Update progress
                    self.progress["crawling"] = True
                    self.save_progress()

                    return True

                logger.info(f"Checkpoint loaded with {len(self.crawler.content_cache)}
    pages, continuing crawl")

            # Crawl website
            logger.info(f"Crawling website: {self.base_url}")
            page_contents = self.crawler.crawl(max_pages=max_pages)

            # Check if we have enough pages
            if not page_contents:
                logger.error("No pages crawled")
                return False

            # If we have very few pages, try to crawl again with different settings
            if len(page_contents) < 3:
```

```python
                logger.warning(f"Only {len(page_contents)} pages crawled, trying again
with different settings")

                # Try a different base URL (www. version or non-www version)
                parsed_url = urlparse(self.base_url)
                if parsed_url.netloc.startswith('www.'):
                    new_base = parsed_url.netloc[4:]
                else:
                    new_base = 'www.' + parsed_url.netloc

                new_url = f"{parsed_url.scheme}://{new_base}{parsed_url.path}"

                # Create a new crawler with the alternative URL
                alt_crawler = PersistentCrawler(new_url, delay=3.0, checkpoint_dir=self.
checkpoint_dir)
                alt_crawler.priority_paths = self.crawler.priority_paths

                # Try to crawl with the alternative URL
                logger.info(f"Trying alternative URL: {new_url}")
                alt_contents = alt_crawler.crawl(max_pages=max_pages)

                # Merge results if we found more pages
                if len(alt_contents) > len(page_contents):
                    logger.info(f"Alternative URL yielded more pages: {len(alt_contents)
}")
                    page_contents = alt_contents
                    self.crawler.content_cache.update(alt_contents)

            # Log results
            logger.info(f"Crawled {len(page_contents)} pages")

            # Save checkpoint
            self.crawler.save_checkpoint()

            # Update progress
            self.progress["crawling"] = True
            self.save_progress()

            return True

        except Exception as e:
            logger.error(f"Error crawling website: {str(e)}")
            return False

    def process_documents(self) -> bool:
        """Process crawled documents into a knowledge base."""
        try:
            # Skip if already done
            if self.progress.get("document_processing", False) and self.progress.get("
knowledge_base", False):
                logger.info("Document processing already done, skipping")

                # Try to load knowledge base
                kb_loaded = self.knowledge_base.load(self._get_knowledge_base_path())

                if kb_loaded:
                    logger.info(f"Knowledge base loaded with {len(self.knowledge_base.
chunks)} chunks")
                    return True
                else:
                    logger.warning("Failed to load knowledge base, reprocessing
documents")

            # Check if we have crawled documents
            if not self.progress.get("crawling", False):
                logger.warning("No crawled documents, run crawl_website first")
                return False

            # Get crawled content
            page_contents = self.crawler.content_cache

            if not page_contents:
                logger.error("No crawled content found")
```

```
4651                    return False
4652
4653            # Process each document
4654            logger.info(f"Processing {len(page_contents)} documents")
4655
4656            for url, content in tqdm(page_contents.items(), desc="Processing documents")
        :
4657                try:
4658                    # Extract title from content
4659                    title_match = re.search(r'TITLE: (.*?)(\n|$)', content)
4660                    title = title_match.group(1) if title_match else ""
4661
4662                    # Process document
4663                    doc = self.doc_processor.process_document(content, url, title)
4664
4665                    # Add to knowledge base
4666                    if doc and doc.chunks:
4667                        self.knowledge_base.add_document(doc)
4668
4669                except Exception as e:
4670                    logger.error(f"Error processing document {url}: {str(e)}")
4671
4672            # Log results
4673            kb_stats = self.knowledge_base.get_stats()
4674            logger.info(f"Processed {kb_stats['documents']} documents into {kb_stats['
        chunks']} chunks")
4675
4676            # Save knowledge base
4677            self.knowledge_base.save(self._get_knowledge_base_path())
4678
4679            # Update progress
4680            self.progress["document_processing"] = True
4681            self.progress["knowledge_base"] = True
4682            self.save_progress()
4683
4684            return True
4685
4686        except Exception as e:
4687            logger.error(f"Error processing documents: {str(e)}")
4688            return False
4689
4690
4691    def evaluate_qa_pairs(self, qa_pairs: List[Dict]) -> Dict:
4692        """Evaluate QA pairs quality metrics."""
4693        try:
4694            # Skip if already done
4695            if self.progress.get("evaluation", False):
4696                logger.info("Evaluation already done, skipping")
4697                return {}
4698
4699            if not qa_pairs:
4700                logger.warning("No QA pairs to evaluate")
4701                return {}
4702
4703            # Calculate various metrics
4704            metrics = {
4705                "total_pairs": len(qa_pairs),
4706                "avg_scores": {},
4707                "distribution": {},
4708                "topic_coverage": {},
4709                "url_coverage": {}
4710            }
4711
4712            # Average scores
4713            score_keys = list(qa_pairs[0].get("scores", {}).keys())
4714            for key in score_keys:
4715                avg_score = sum(pair["scores"].get(key, 0) for pair in qa_pairs) / max
        (1, len(qa_pairs))
4716                metrics["avg_scores"][key] = round(avg_score, 3)
4717
4718            # Score distribution
4719            metrics["distribution"] = {
```

```
              "excellent": len([p for p in qa_pairs if p["scores"].get("overall", 0)
     >= 0.8]),
              "good": len([p for p in qa_pairs if 0.7 <= p["scores"].get("overall", 0)
      < 0.8]),
              "average": len([p for p in qa_pairs if 0.6 <= p["scores"].get("overall",
      0) < 0.7]),
              "below_avg": len([p for p in qa_pairs if p["scores"].get("overall", 0) <
      0.6])
          }

          # Topic coverage
          topics = {}
          for pair in qa_pairs:
              topic = pair.get("topic", "").lower()
              if topic and topic != "general":
                  topics[topic] = topics.get(topic, 0) + 1

          metrics["topic_coverage"] = dict(sorted(topics.items(), key=lambda x: x[1],
     reverse=True)[:10])

          # URL coverage
          urls = {}
          for pair in qa_pairs:
              url = pair.get("source_url", "")
              if url:
                  urls[url] = urls.get(url, 0) + 1

          metrics["url_coverage"] = dict(sorted(urls.items(), key=lambda x: x[1],
     reverse=True))

          # Save evaluation results
          eval_path = os.path.join(self.output_dir, "evaluation_metrics.json")
          with open(eval_path, 'w', encoding='utf-8') as f:
              json.dump(metrics, f, indent=2)

          logger.info(f"Evaluation metrics saved to {eval_path}")

          # Update progress
          self.progress["evaluation"] = True
          self.save_progress()

          return metrics

      except Exception as e:
          logger.error(f"Error evaluating QA pairs: {str(e)}")
          return {}

  def generate_qa_pairs(self, max_pairs_per_url: int = 10) -> List[Dict]:
      """Generate QA pairs from processed documents."""
      try:
          # Check if we can load from checkpoint
          if self.progress.get("qa_generation", False):
              logger.info("QA generation already done, loading from checkpoint")
              qa_pairs = self.load_qa_pairs()

              if qa_pairs:
                  logger.info(f"Loaded {len(qa_pairs)} QA pairs from checkpoint")
                  return qa_pairs
              else:
                  logger.warning("Failed to load QA pairs, regenerating")

          # Check if we have processed documents
          if not self.progress.get("document_processing", False) or not self.progress.
     get("knowledge_base", False):
              logger.warning("No processed documents, run process_documents first")
              return []

          # Check knowledge base
          kb_stats = self.knowledge_base.get_stats()
          if kb_stats["chunks"] == 0:
              logger.error("Knowledge base is empty")
              return []
```

```python
            # Get all URLs
            urls = list(set(chunk.doc_url for chunk in self.knowledge_base.chunks))

            if not urls:
                logger.error("No URLs found in knowledge base")
                return []

            # Store original filter function for potential reset
            if not hasattr(self.qa_generator, '_original_filter_qa_pairs'):
                self.qa_generator._original_filter_qa_pairs = self.qa_generator.
    filter_qa_pairs

            # Generate QA pairs
            logger.info(f"Generating QA pairs for {len(urls)} URLs")
            qa_pairs = self.qa_generator.generate_qa_pairs(urls, max_pairs_per_url)

            # Check if we have enough QA pairs
            if len(qa_pairs) < 5:
                logger.warning(f"Only {len(qa_pairs)} QA pairs generated, trying with
    lower quality threshold")

                # Lower quality threshold for limited content
                original_filter = self.qa_generator.filter_qa_pairs

                # Override with more lenient filter
                def lenient_filter(pairs, min_score=0.6):
                    return original_filter(pairs, min_score=0.4)

                # Apply the lenient filter
                self.qa_generator.filter_qa_pairs = lenient_filter

                # Try again
                qa_pairs = self.qa_generator.generate_qa_pairs(urls, max_pairs_per_url)
                logger.info(f"After adjustment: {len(qa_pairs)} QA pairs")

                # Restore original filter
                self.qa_generator.filter_qa_pairs = original_filter

            # Save QA pairs
            self.save_qa_pairs(qa_pairs)

            # Update progress
            self.progress["qa_generation"] = True
            self.save_progress()

            # Log results
            logger.info(f"Generated {len(qa_pairs)} QA pairs")

            return qa_pairs

        except Exception as e:
            logger.error(f"Error generating QA pairs: {str(e)}")
            return []

    def save_output(self, qa_pairs: List[Dict]) -> None:
        """Save final output in multiple formats."""
        try:
            if not qa_pairs:
                logger.warning("No QA pairs to save")
                return

            # 1. Save as JSON
            json_path = os.path.join(self.output_dir, "qa_pairs_final.json")
            with open(json_path, 'w', encoding='utf-8') as f:
                json.dump({"qa_pairs": qa_pairs}, f, indent=2)

            # 2. Save as CSV for easy viewing
            csv_path = os.path.join(self.output_dir, "qa_pairs_final.csv")
            with open(csv_path, 'w', encoding='utf-8', newline='') as f:
                # Write header
                f.write("Question,Answer,Source URL,Topic,Topic Type,Overall Score\n")

                # Write data
```

```
                 for pair in qa_pairs:
                     question = pair["question"].replace('"', '""')
                     answer = pair["answer"].replace('"', '""')
                     url = pair.get("source_url", "").replace('"', '""')
                     topic = pair.get("topic", "").replace('"', '""')
                     topic_type = pair.get("topic_type", "").replace('"', '""')
                     score = str(pair.get("scores", {}).get("overall", 0))

                     f.write(f'"{question}","{answer}","{url}","{topic}","{topic_type}",{
    score}\n')

             # 3. Save as Markdown for human reading
             md_path = os.path.join(self.output_dir, "qa_pairs_final.md")
             with open(md_path, 'w', encoding='utf-8') as f:
                 f.write(f"# QA Pairs for {self.base_url}\n\n")

                 # Group by URL
                 url_to_pairs = {}
                 for pair in qa_pairs:
                     url = pair.get("source_url", "Unknown")
                     if url not in url_to_pairs:
                         url_to_pairs[url] = []
                     url_to_pairs[url].append(pair)

                 # Write each URL's pairs
                 for url, pairs in url_to_pairs.items():
                     f.write(f"## {url}\n\n")

                     for i, pair in enumerate(pairs, 1):
                         f.write(f"### Q{i}: {pair['question']}\n\n")
                         f.write(f"{pair['answer']}\n\n")
                         f.write(f"*Topic: {pair.get('topic', 'N/A')} | "
                                 f"Type: {pair.get('topic_type', 'N/A')} | "
                                 f"Score: {pair.get('scores', {}).get('overall', 0):.2f}*\n
    \n")
                         f.write("---\n\n")

         logger.info(f"Output saved to {self.output_dir} in JSON, CSV, and Markdown
    formats")

     except Exception as e:
         logger.error(f"Error saving output: {str(e)}")

 def run(self, max_pages: int = 30, max_pairs_per_url: int = 10) -> List[Dict]:
     """Run the complete QA generation pipeline."""
     start_time = datetime.now()

     # 1. Load progress if available
     self.load_progress()

     try:
         # 2. Setup resources
         success = self.setup_resources()
         if not success:
             logger.error("Failed to set up resources")
             return []

         # 3. Crawl website
         success = self.crawl_website(max_pages)
         if not success:
             logger.error("Failed to crawl website")
             return []

         # 4. Process documents
         success = self.process_documents()
         if not success:
             logger.error("Failed to process documents")
             return []

         # 5. Generate QA pairs
         qa_pairs = self.generate_qa_pairs(max_pairs_per_url)
         if not qa_pairs:
             logger.error("Failed to generate QA pairs")
```

```
4927                    return []
4928
4929            # 6. Evaluate QA pairs
4930            self.evaluate_qa_pairs(qa_pairs)
4931
4932            # 7. Save output
4933            self.save_output(qa_pairs)
4934
4935            # 8. Log completion
4936            end_time = datetime.now()
4937            duration = end_time - start_time
4938
4939            logger.info(f"QA generation complete in {duration}")
4940            logger.info(f"Generated {len(qa_pairs)} QA pairs")
4941            logger.info(f"Output saved to {self.output_dir}")
4942
4943            return qa_pairs
4944
4945        except Exception as e:
4946            logger.error(f"Error running QA generation: {str(e)}")
4947            return []
4948        finally:
4949            # Clean up resources
4950            try:
4951                if hasattr(self, 'model_manager'):
4952                    self.model_manager.cleanup()
4953            except Exception as e:
4954                logger.error(f"Error cleaning up: {str(e)}")
4955
4956 #----------------------------------------------------------------------
4957 # Main Application Entry Point
4958 #----------------------------------------------------------------------
4959
4960 def parse_arguments():
4961     """Parse command line arguments."""
4962     parser = argparse.ArgumentParser(description="Generate QA pairs from a website")
4963
4964     parser.add_argument("--url", type=str, default="https://www.thrive.pitt.edu",
4965                         help="Base URL to crawl")
4966     parser.add_argument("--output", type=str, default="qa_output",
4967                         help="Output directory")
4968     parser.add_argument("--max-pages", type=int, default=30,
4969                         help="Maximum number of pages to crawl")
4970     parser.add_argument("--max-pairs", type=int, default=10,
4971                         help="Maximum number of QA pairs per page")
4972     parser.add_argument("--no-gpu", action="store_true",
4973                         help="Disable GPU usage")
4974     parser.add_argument("--checkpoint", type=str, default="checkpoints",
4975                         help="Checkpoint directory")
4976     parser.add_argument("--force", action="store_true",
4977                         help="Force regeneration even if checkpoints exist")
4978
4979     return parser.parse_args()
4980
4981 def main():
4982     args = parse_arguments()
4983
4984     print("=" * 80)
4985     print("Advanced QA Generation System with RAG, Neural Models and Semantic Processing
4986     ")
4986     print("=" * 80)
4987     print(f"URL: {args.url}")
4988     print(f"Output directory: {args.output}")
4989     print(f"Max pages: {args.max_pages}")
4990     print(f"Max QA pairs per page: {args.max_pairs}")
4991     print(f"GPU enabled: {not args.no_gpu}")
4992     print("=" * 80)
4993
4994     try:
4995         # Create system
4996         system = QAGenerationSystem(
4997             base_url=args.url,
4998             output_dir=args.output,
```

```
5999            use_gpu=not args.no_gpu,
5000            checkpoint_dir=args.checkpoint
5001        )
5002
5003        # If force option is set, clear progress
5004        if args.force:
5005            system.progress = {key: False for key in system.progress}
5006            system.save_progress()
5007            logger.info("Forced regeneration, cleared progress")
5008
5009
5010        qa_pairs = system.run(
5011            max_pages=args.max_pages,
5012            max_pairs_per_url=args.max_pairs
5013        )
5014
5015        print("\nGeneration Summary:")
5016        print(f"Generated {len(qa_pairs)} QA pairs")
5017        print(f"Output saved to {args.output}")
5018        print("=" * 80)
5019
5020        return 0
5021
5022    except KeyboardInterrupt:
5023        print("\nOperation cancelled by user")
5024        return 1
5025    except Exception as e:
5026        print(f"\nError: {str(e)}")
5027        logger.error(f"Fatal error: {str(e)}")
5028        traceback.print_exc()
5029        return 1
5030
5031 if __name__ == "__main__":
5032     sys.exit(main())
```

## 10 Code Usage Instructions

### 10.1 Prerequisites

Before running the QA generation system, ensure you have the following prerequisites:

- Python 3.8 or higher

- CUDA-compatible GPU (recommended) with appropriate drivers

- At least 8GB of RAM (16GB recommended)

- At least 5GB of free disk space

- Internet connection for downloading models and crawling websites

### 10.2 Installation

Install the required packages using pip:

```
1 pip install torch transformers beautifulsoup4 nltk spacy tqdm backoff
2     sentence-transformers requests numpy sklearn
3 python -m spacy download en_core_web_sm
```

## 10.3  Output Directory Structure

```
thrive_qa_generator.py ........................................ Main QA generation script
qa_output/ ................................................ Created automatically
    qa_pairs_final.json ........................................ QA pairs in JSON format
    qa_pairs_final.csv ......................................... QA pairs in CSV format
    qa_pairs_final.md ........................................ QA pairs in Markdown format
    evaluation_metrics.json .......................... Quality metrics for generated QA pairs
checkpoints/ ........................................ For storing crawl and model state
    content_*.json ............................................ Cached webpage content
    visited_*.json ............................................. Record of visited URLs
    failed_*.json .................................................. Failed URL attempts
    scores_*.json ................................................. Page quality scores
    kb_*.json ..................................................... Knowledge base data
    qa_pairs_*.json ........................................ Generated QA pair checkpoints
    progress_*.json ........................................... Pipeline progress tracking
debug_html/ ........................................ Optional: Created for HTML debugging
    page_0.html ............................................... Sample of downloaded pages
    page_1.html
```

**Note:** The asterisk (*) in filenames (e.g., `content_*.json`) represents a unique identifier generated for each crawling session, based on an MD5 hash of the website URL. This allows the system to maintain separate checkpoint files for different websites.

For example, when crawling `https://www.thrive.pitt.edu`, the system might create files like:

```
content_a7f92e3b4c.json
visited_a7f92e3b4c.json
kb_a7f92e3b4c.json
```

## 10.4  Command Structure

```
1  python thrive_qa_generator.py \
2      --url https://www.thrive.pitt.edu \
3      --output qa_output \
4      --max-pages 30 \
5      --max-pairs 10 \
6      --checkpoint checkpoints \
7      --no-gpu \
8      --force
```

## 10.5  Command-Line Arguments

| Argument | Default | Description |
|---|---|---|
| -url | https://www.thrive.pitt.edu | Base URL to crawl for generating QA pairs. |
| -output | qa_output | Directory to save the generated QA pairs and evaluation metrics. |
| -max-pages | 30 | Maximum number of pages to crawl from the website. |
| -max-pairs | 10 | Maximum number of QA pairs to generate per page. |
| -no-gpu | False | Flag to disable GPU usage and run only on CPU. |
| -checkpoint | checkpoints | Directory for saving and loading checkpoint data (crawled content, model states, etc.). |
| -force | False | Force regeneration of QA pairs even if checkpoints exist from previous runs. |

Table 5: Command-line arguments for the QA generation script.

## References

Ji, Z., Vu, M., Wang, X., and Neubig, G. (2023). A survey of hallucination in large language models. *arXiv preprint arXiv:2303.17085.*

OpenAI (2023). Gpt-4 system card.