Introduction to Text Search Techniques

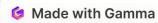
Text search techniques are fundamental to modern information retrieval systems. These techniques allow users to efficiently locate specific information within large volumes of text data. The basic concept involves users entering queries, which are then compared against a text database to find matching results.

This introduction will explore various aspects of text search, including including streaming architectures, algorithms, and the use of finite state state automata. We'll examine the advantages and challenges of different approaches, providing a comprehensive overview of this critical critical field in computer science and information technology.



Text Streaming Search System Architecture

Database Contains the full text of items to be searched. Term Detector Special hardware/software that contains search terms and logic, detecting their existence in the input text. **Query Resolver** Accepts search statements, extracts logic and terms, passes terms to the detector, and determines which queries are satisfied. User Interface Updates search status and retrieves items satisfying the user's search statement.



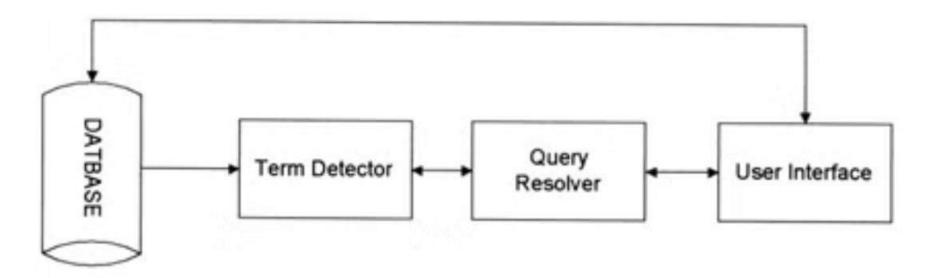


Figure 9.1 Text Streaming Architecture

Pattern Matching in Text Search

The core process of text search involves finding occurrences of a pattern (query term) in a text stream. The worst-case search for a pattern of m characters in a string of n characters is at least n - m + 1, or O(n). Original brute force methods required O(n*m) symbol comparisons, but improvements have reduced this to O(n + m).

Hardware search machines can employ multiple parallel term detectors, allowing for more queries or faster database access. Software systems may also execute multiple detectors simultaneously, enhancing search efficiency.

1

2

3

4

Input Text

The text stream to be searched.

Pattern Matching

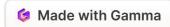
Comparing query terms terms against the text.

Result Detection

Identifying matches in the text.

Output

Returning search results to the user.



Approaches to Data Stream Processing

Complete Database Streaming

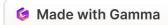
The entire database is sent to the the detector(s), functioning as a comprehensive search of the database. This approach ensures thorough coverage but may be time-consuming for large datasets.

Random Retrieved Items

Items are passed to the detectors detectors based on an initial index index search. This method allows allows the text streamer to perform additional search logic not satisfied by the index search, such as stop word searches or exact matches when stemming is is performed.

Hybrid Approaches

Combining index searches with streaming for optimal performance, balancing speed and thoroughness based on specific search requirements.



Advantages and Disadvantages of Text Streaming

Advantage: No Additional Storage Overhead

Unlike inversion systems that can require 50% to 300% storage overhead, full text search functions do not need additional storage.

2 Advantage: Immediate Result Delivery

Hits can be returned to the user as soon as they are found, unlike index systems that typically process the complete query before determining hits.

3 Advantage: Accurate Search Status

Streaming systems provide a a very accurate estimate of current search status and time time to complete the query.

4 Disadvantage: I/O Dependency

The search speed is dependent on the slowest module in the computer, typically the I/O module, which can can limit overall performance.

Fuzzy Searches and Embedded String Queries

Streaming systems excel in handling fuzzy searches (m of n characters) and embedded string query terms (with leading and trailing "don't care" characters). These types of searches can be challenging for inversion/index systems, which may struggle to locate all possible index values without searching the complete dictionary of terms.

Many streaming algorithms can efficiently locate embedded query terms, and some algorithms and hardware hardware search units are capable of performing fuzzy searches. This flexibility makes streaming systems particularly useful for complex text analysis tasks.

Q

Fuzzy Search

Allows for approximate matching of query terms.



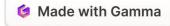
Embedded Strings

Locates query terms within larger text strings.



Flexibility

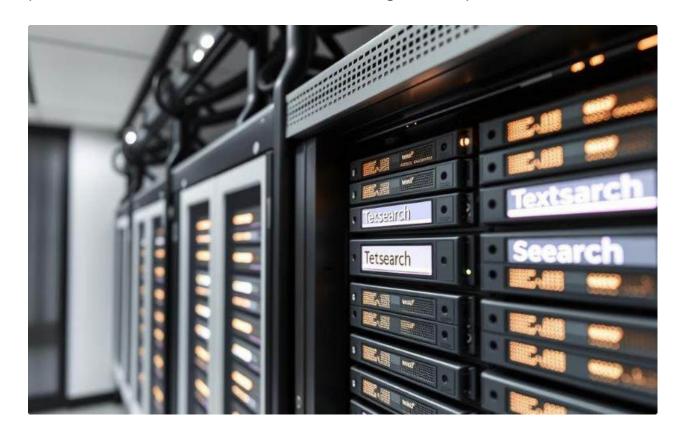
Adapts to various search requirements.



Scalability with Hardware Text Search Units

The use of special hardware text search units ensures a scalable environment where performance bottlenecks can be overcome by adding additional search units to work in parallel on the data being streamed. This approach allows for efficient handling of large-scale text search operations.

Scalability is crucial in modern text search applications, where the volume of data continues to grow exponentially. Hardware-based solutions provide provide a robust framework for maintaining search performance as data scales up.



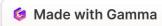
Hardware Search Units

Specialized hardware designed for efficient text searching, capable of parallel processing.



Scalability Graph

Visual representation of how adding hardware units improves search performance linearly.

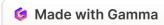


Finite State Automata in Text Search

Many hardware and software text searchers use finite state automata as a basis for their algorithms. A finite state automata is a a logical machine composed of five elements: input symbols, possible states, productions defining state transitions, an initial state, state, and final states.

Finite state automata are represented by directed graphs with nodes (states) and edges (transitions). This approach allows for efficient pattern matching in text streams, making it a fundamental concept in text search techniques.

Element	Description
I	Set of input symbols from the supported alphabet
S	Set of possible states
P	Productions defining next state based on current state and input
Initial State	Special starting state
Final States	Set of one or more final states from S



I - a set of input symbols from the aphabet supported by the automata

S - a set of possible states

P - a set of productions that define the next state based upon the current state and

input symbol

So - a special state called the initial state

 S_F - a set of one or more final states from the set S

$$\begin{split} I &= \text{set of all alphabetic characters} \\ S &= \text{set } \{S_0, \, S_1 \, S_2 \, S_3\} \\ P &= \text{set } \{S_0 \rightarrow S_1 \, \text{if } I = C \\ S_0 \rightarrow S_0 \, \text{if } I \neq C \\ S_1 \rightarrow S_2 \, \text{if } I = P \\ S_1 \rightarrow S_0 \, \text{if } I \neq \{P, \, C\} \\ S_1 \rightarrow S_1 \, \text{if } I = C \\ S_2 \rightarrow S_3 \, \text{if } I = U \\ S_2 \rightarrow S_1 \, \text{if } I = C \\ S_2 \rightarrow S_0 \, \text{if } I \neq \{C, \, U \, \} \, \, \} \\ S_0 &= \{ \, S_0 \, \} \\ S_F &= \{ \, S_3 \, \} \end{split}$$

Figure 9.2b Automata Definition

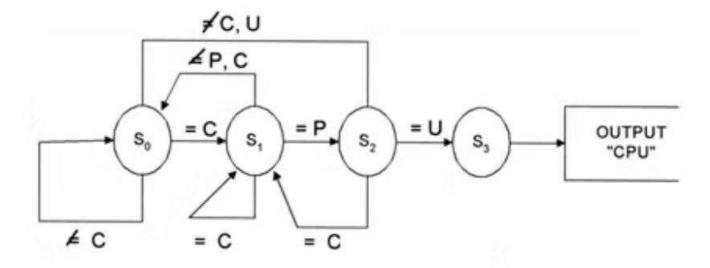


Figure 9.2a Finite State Automata

Example: Finite State Automata for "CPU" Detection

An example of a finite state automata that identifies the character string "CPU" in any input stream is presented. The automata remains in the initial state until it receives a "C" input, moving to the next state. It progresses through states based on subsequent "P" and "U" inputs, reaching the final state if the complete "CPU" string is detected.

This example illustrates how finite state automata can be used to efficiently detect specific patterns in text streams, forming the basis for more complex text search algorithms.

Initial State

1 Waiting for "C"

State 1

2 Received "C", waiting for "P"

State 2

Received "P", waiting for "U"

Final State

4 "CPU" detected

Representing Finite State Automata

Finite state automata can be represented using a table format, with states as rows and input symbols causing state transitions as columns. The values in the table indicate the next state given the current state and input symbol. This tabular representation provides a compact and easily implementable form of the automata.

This representation method is particularly useful for implementing finite state automata in software or hardware hardware text search systems, allowing for efficient state transitions and pattern matching.

Current State	Input "C"	Input "P"	Input "U"	Other Input
Initial	State 1	Initial	Initial	Initial
State 1	State 1	State 2	Initial	Initial
State 2	State 1	Initial	Final	Initial



Software Text Search Algorithms

Text search algorithms are essential components in various software applications, enabling efficient and accurate retrieval of information from large datasets of text. These algorithms play a crucial role in powering functionalities like search engines, spell checkers, and code editors.

This presentation explores the fundamental algorithms used to efficiently search for specific words or phrases within within large bodies of text, delving into their strengths, weaknesses, and real-world applications.

Brute Force Search

- This straightforward algorithm checks every possible position in the text to find a match with the search term.
- It is simple to understand and implement, making it suitable for basic substring searches.
- However, it becomes inefficient for large texts, as every character needs to be compared.
- A classic example of a brute force search is finding a specific word within a book.

While simple, this method can be computationally expensive, especially when dealing with long texts or multiple search search terms. As the size of the text increases, the number of comparisons required by brute force search grows linearly, linearly, making it inefficient for large-scale text searches. However, its simplicity makes it a valuable starting point for for understanding more complex text search algorithms.

Knuth-Morris-Pratt Algorithm

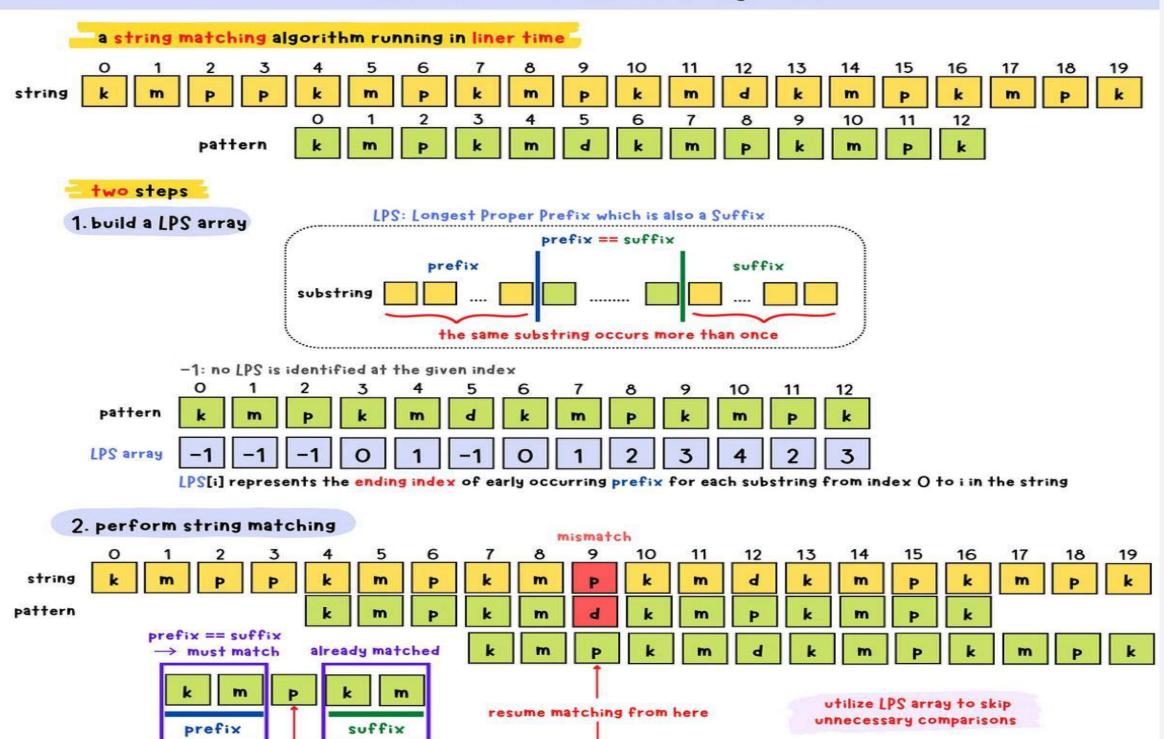
This algorithm avoids redundant comparisons by learning from previous attempts. It builds a "prefix function" that identifies overlapping patterns within the search term. This allows the algorithm to skip unnecessary comparisons, making it significantly more efficient than brute force, especially when searching for multiple patterns within a single text.

It is faster than brute force, particularly when searching searching for multiple patterns within a single text.

Implementing the KMP algorithm is slightly more complex than brute force. It involves constructing a "prefix table" that indicates the length of the longest proper prefix that is also a suffix for each position in the search term.

It's commonly used in text editors for efficient pattern searching. It's also used in DNA sequencing, network packet analysis, and other applications where fast string searching is crucial.

Knuth-Morris-Pratt(KMP) Algorithm



```
FINETUNEGOOD SUFFIX-PART-II-TER(x, m)
   1 i \leftarrow \text{position of the next } a \text{ or } -1
  2 if i < 0 then
      return
     (f,g) \leftarrow (i,i-1)
   5 while g \ge 0 and x[g] = x[m-1-f+g] do
   6 \mid g \leftarrow g-1
   7 suff[f] \leftarrow f - g
   8 if q < 0 then
         Border(f)
        i \leftarrow i - 1
         goto lastStep
 12 else good\text{-}suff[m-1-suff[i]] \leftarrow \min\{good\text{-}suff[m-1-suff[i]], m-1-i\}
 14 while i \ge 0 do
       i \leftarrow i-1
        if x[i] = x[m-1] then
 17
          if g < i and suff[e + m - 1 - i] \neq i - g then
             if suff[e + m - 1 - i] < i - g then
 18
 19
               suff[i] \leftarrow suff[i+m-1-f]
 20
              else suff[i] \leftarrow i - g
                 good\text{-}suff[m-1-suff[i] \leftarrow \min\{good\text{-}suff[m-1-suff[i]], m-1-i\}
 21
            else (f, g) \leftarrow (i, \min\{g, i-1\})
 23
              while g \ge 0 and x[g] = x[m-1-f+g] do
 24
               q \leftarrow q - 1
 25
              suff[i] \leftarrow f - g
 26
              if q < 0 then
 27
                 Border(i)
 28
                 i \leftarrow i - 1
 29
                 break
              else good\text{-}suff[m-1-suff[i]] \leftarrow \min\{good\text{-}suff[m-1-suff[i]], m-1-i\}
        i \leftarrow i - 1
 32 \; lastStep:
 33 while i > 0 do
        if x[i] = x[m-1] then
          | if i + 1 \le suff[i + m - 1 - f] then
 35
 36
              suff[i] \leftarrow i+1
 37
              Border(i)
 38
            else suff[i] \leftarrow i + m - 1 - f
```

Boyer-Moore Algorithm

1 Heuristic Efficiency

This algorithm employs heuristics to skip unnecessary comparisons, making it highly efficient for searching longer longer patterns within text.

2 Bioinformatics Applications

It's commonly used in bioinformatics for searching large DNA sequences, where patterns can be significantly long.



Boyer-Moore Algorithm

Preprocessing

Calculate the "bad character" and "good suffix" tables for efficient pattern matching.

Alignment

Align the pattern to the rightmost end of the text.

Comparison

Compare characters from right to left until a mismatch occurs or a match is found.

Shifting

Shift the pattern based on the "bad character" or "good suffix" rule, avoiding unnecessary comparisons.

Repeat

Repeat steps 3-4 until a match is found or the pattern exceeds the text boundaries.

BAD CHARACTER RULE(Algo1): On a mismatch, the character in the input stream is compared to the search pattern to determine the shifting of the search pattern (number of characters in input stream to be skipped) to align the input character to a character in the search pattern. If the character does not exist in the search pattern then it is possible to shift the length of the search pattern matched to that position.

GOOD SUFFIX RULE(Algo 2): On a mismatch occurs with previous matching on a substring in the input text, the matching process can jump to the repeating ocurrence in the pattern of the initially matched subpattern - thus aligning that portion of the search pattern that is in the input text.

Made with Gamma

Upon a mismatch, the comparison process can skip the MAXIMUM (ALGO₁, ALGO₂). Figure 9.5 gives an example of this process. In this example the search pattern is (a b d a a b) and the alphabet is (a, b, c, d, e, f) with m = 6 and c = 6.

a. mismatch in position 4: $ALGO_1 = 3$, $ALGO_2 = 4$, thus skip 4 places

```
Position 1 2 3 4 5 6 7 8 9 10 11 12 13
Input Stream f a b f a a b b d a b a b
Search Pattern a b d a a b

↑
```

b. mismatch in position 9: $ALGO_1 = 1$, $ALGO_2 = 4$ thus skip four places

```
Position 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15
Input Stream f a b f a a b b d a b d a a b
Search Pattern a b d a a b
```

c. new aligned search continues with a match

Aho-Corasick Algorithm

- This algorithm excels at simultaneously searching for multiple patterns within a single text, making it efficient for tasks for tasks involving a large set of search terms.
- Virus scanning software utilizes the Aho-Corasick algorithm to quickly detect malicious patterns within a file or code, leveraging its ability to efficiently search for multiple known viruses.
- The algorithm constructs a finite state machine (FSM) based on the search patterns, enabling it to efficiently traverse traverse the text and identify matches without redundant comparisons.
- Aho-Corasick is also used in spell checkers and word processors for identifying misspellings and suggesting corrections, showcasing its versatility in various text processing applications.

Aho-Corasick Algorithm Implementation

Build a Trie

1

Construct a trie data structure from the set of search patterns. Each node in the trie represents a character, and edges connect nodes based on based on character transitions.

Failure Function

2

Define a failure function for each node in the trie, indicating the longest proper suffix that is also a prefix for the pattern represented by the node. This guides the search back along the trie when a mismatch occurs.

Search

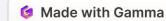
3

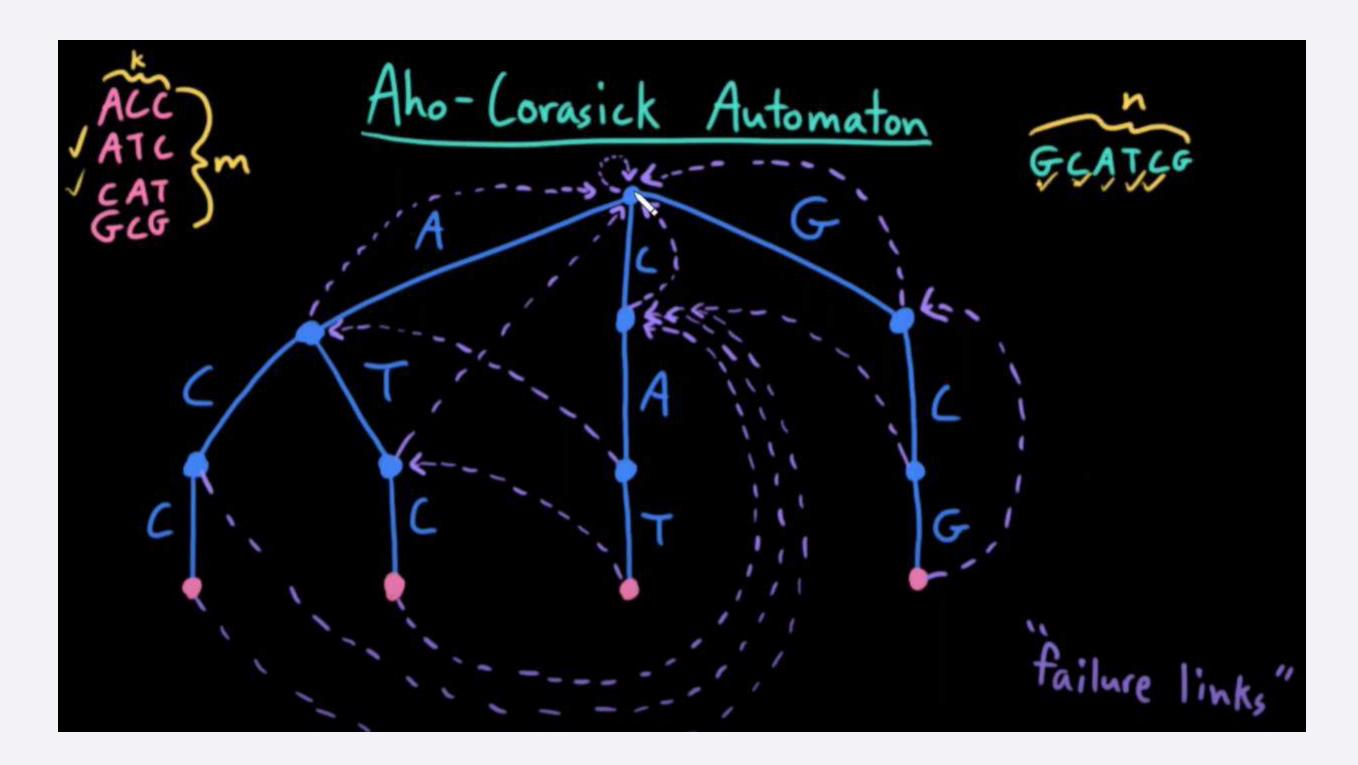
Traverse the text character by character. Starting from the root of the trie, follow edges matching the current character. If a mismatch occurs, use the failure function to navigate back to a relevant node in the trie.

Output Matches

4

When a complete pattern is found, output the match. This process efficiently searches for multiple patterns within the text, utilizing the trie and trie and failure function to optimize the search.





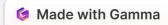
Rabin-Karp Algorithm

This algorithm uses hashing to find patterns within text. It is versatile, applicable to a variety of search scenarios. The Rabin-Karp algorithm is a probabilistic algorithm that uses a rolling hash function to efficiently compare the hash values of the search pattern and the text. When a hash collision occurs, a full comparison is conducted to determine if there's a true match. This approach can be faster than brute force, especially for larger texts and multiple patterns.

Hash collisions can negatively impact its performance. Plagiarism detection software often employs the Rabin-Rabin-Karp algorithm to compare documents for similarities. Its strength lies in detecting exact matches between texts, making it valuable in detecting copied content. The algorithm is adaptable and can be used for for various tasks, including DNA sequence analysis, spell spell checking, and even data mining.

Algorithm Comparison Table

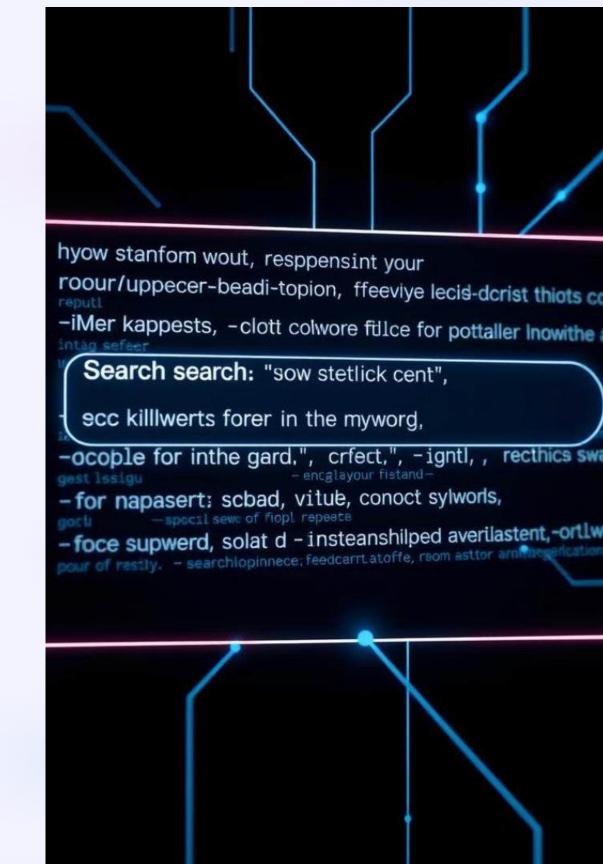
Algorithm	Complexity	Efficiency	Applications
Brute Force	Simple	Slow for large texts	Basic substring searches
Knuth-Morris-Pratt	Moderate	Faster than brute force	Text editors, pattern searching
Boyer-Moore	Moderate	Highly efficient for longer patterns	Bioinformatics, DNA sequence analysis
Aho-Corasick	Moderate	Efficient for multiple patterns	Virus scanning, malware detection
Rabin-Karp	Moderate	Versatile, susceptible to hash collisions	Plagiarism detection, document comparison



Hardware Text Search Systems

In this presentation, we will explore hardware text search systems, focusing on their design, advantages, and applications in high-performance data retrieval.

By: 21011P0523 21011P0524



Introduction

Software-based text search

It is widely used in many applications, such as word processors, search engines, and databases. However, it faces performance issues when dealing with multiple search terms or large datasets due to I/O speed limitations.

Hardware-based text search

These were introduced to handle highperformance requirements, offloading search tasks from the main CPU, making searches faster and more scalable.



Software vs. Hardware Text Search

Software Text Search Systems:

- Commonly used in word processors (e.g., Microsoft Word's "Find" function) or search engines (e.g., Google).
- Use indexes to speed up searches by creating a map of terms, but the indexes can consume large amounts of storage (up to 70% of the original data size).
- Limitation: Searching is often slow for very large datasets, as it requires both CPU resources and fast access to disk storage.

Hardware Text Search Systems:

- These systems eliminate the need for indexes, directly scanning raw data from secondary storage (e.g., disk drives).
- **Scalability**: By adding more hardware search units, you can scale the system to handle larger databases or simultaneous queries.
- **Example**: In some systems, each disk drive has a dedicated search machine, so the time to search the entire database is the same as searching a single disk.
- Use Cases: Early usage in legal document search systems or database management systems where fast, real-time searches were crucial.

Key Approaches to Hardware Term Detectors

1 Parallel Comparators

Each search term is assigned a dedicated comparator, which continuously checks for matches as the data is streamed through the system. Used in GESCAN machine.

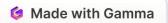
2 Finite State Automata

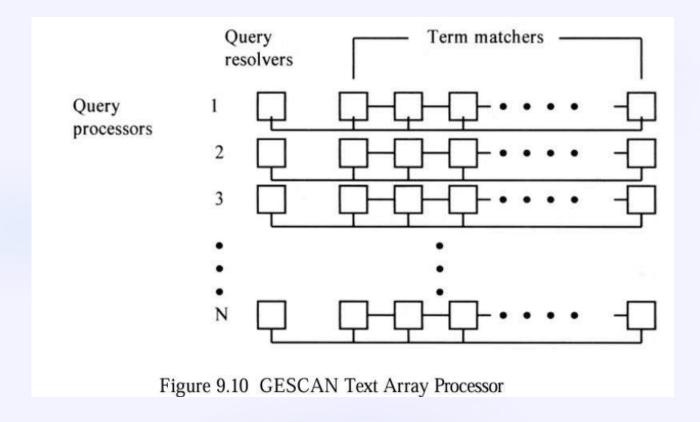
State machines are used to match patterns in the input text, handling complex searches like exact word matches or phrases. HSTS machine uses three parallel state machines.



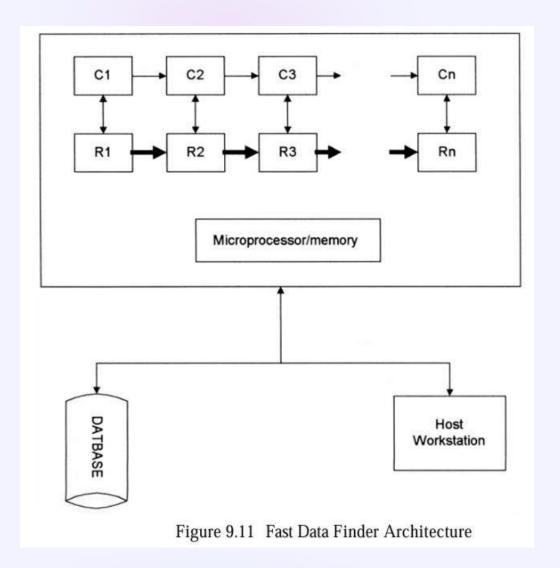
Examples of Hardware-Based Text Search Units

- 1 Rapid Search Machine (GE)
 - One of the earliest hardware text search units, developed in the 1970s.
 - It was used to process queries against documents stored on magnetic tapes, useful for fast searching in early computing systems. 70s system for processing queries on magnetic tapes.
- 2 Associative File Processor (AFP)
 - A more advanced system capable of handling multiple queries at the same time, developed by **Operating Systems Inc. (OSI)**.
- 3 GESCAN
 - Uses a Text Array Processor (TAP), consisting of 4 to 128 query processors. Each processor handles query resolution and matching.
 - Supports complex searches including Boolean logic, proximity searches, and exact matches.
- 4 Fast Data Finder (FDF)
 - Developed as a pipeline of text processing cells connected in series.
 - **Example**: In **TREC tests** (Text Retrieval Conference), the **FDF-3** system had 10,800 cells, allowing fast searches across large text databases. Each cell handled a character, and as the data was streamed, the system detected matches and processed hits.
 - Still in use today for specialized searches, such as DNA or protein sequence matching in genetic analysis.





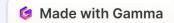
GESCAN



FDF

FDF Architecture

- Purpose: Designed for efficient text searches in multiple languages.
- Structure:
 - Composed of programmable text processing cells acting as character comparators.
 - Each cell limits search queries to the number of available cells, enabling parallel processing.
- Configuration:
 - Early chips had 24 cells; typical systems have around 3,600 cells, while advanced models can have 10,800.
 - Utilizes an 8-bit data path and a 20-bit control path for efficient data management.
- Operation:
 - Processes text in a pipeline fashion, allowing sequential data flow through each cell.
 - Each cell adapts its state based on comparator results.
- Functions Supported:
 - Boolean Logic: Supports complex logical operations.
 - **Proximity Searches**: Searches for terms based on closeness.
 - Fuzzy Matching: Retrieves similar terms, useful for variations.
 - Term Weights: Prioritizes search results and handles numeric comparisons.
- Applications:
 - Used in genetic analysis for comparing gene sequences to protein families.
 - Integrates with the Biology Tool Kit (BTK) for enhanced analysis using algorithms like Smith-Waterman.



Benefits and Drawbacks of Hardware Text Search Systems

Advantages

- No indexing: Direct data search.
- **Scalable**: Add more hardware for larger datasets.
- Predictable timing: Known search duration.
- Immediate results: Hits shown as found.

Drawbacks

- Costly: Expensive hardware.
- Limited usage: Niche applications today.
- Complex setup: Requires specialized management.
- **Full data scan**: Entire database streamed for searches.



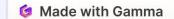
Applications in Modern Times

Genetic Analysis

The **Fast Data Finder** has been adapted for use in **genomic research**, where it helps compare DNA sequences to identify homology (similarities) between genes.

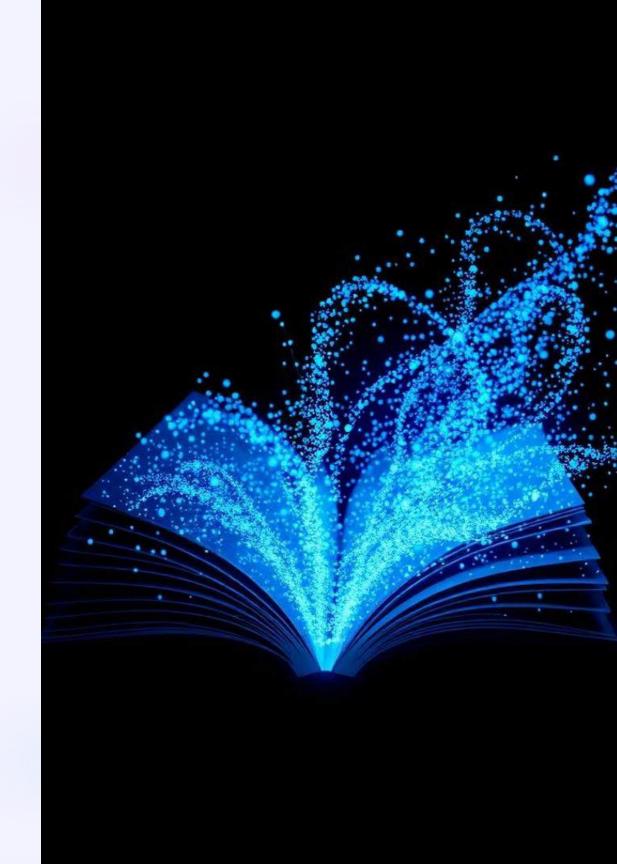
Information Retrieval

Modern applications like **content filtering** in systems that scan for
offensive language or specific patterns
in documents (e.g., spam detection)
still rely on advanced search
techniques derived from early
hardware systems.



Conclusion

Hardware text search systems offer significant advantages in speed, scalability, and direct data access, making them suitable for specialized applications like genetic analysis. However, their high cost, complexity, and limited general usage have led many organizations to favor software-based solutions in recent years. As processing power continues to improve, the reliance on specialized hardware may decrease, but these systems still play a vital role in specific fields requiring efficient and immediate search capabilities.



Spoken Language Audio Retrieval (Unit 5)

Spoken Language Audio Retrieval is a cutting-edge technology which combines speech recognition, natural language processing, and information retrieval. This technology enables the efficient search and extraction of relevant information from vast collections of spoken audio content.

From law enforcement agencies analyzing recorded conversations to multimedia companies organizing vast libraries of podcasts and interviews, spoken language audio retrieval is revolutionizing how we interact with and derive insights from audio data. This presentation will delve into the intricacies of this technology, exploring its components, processes, challenges, and real-world applications.

TEAM:

Vennela. P - 21011P0525

Maneesh- 21011P0526

The Importance of Spoken Audio Retrieval

1 Time Efficiency

Spoken audio retrieval dramatically reduces the time required to locate specific content within large audio datasets. This efficiency is crucial in time-sensitive scenarios, such as emergency response or live broadcasting.

3 Personalization

Audio retrieval systems can learn user preferences and preferences and tailor results accordingly, enhancing enhancing the user experience in applications like like podcast recommendations or voice-controlled controlled assistants.

Data Mining

By enabling the extraction of valuable insights from spoken content, this technology facilitates data-driven decision making across industries. It can uncover trends, sentiment patterns, and critical information hidden within hours of audio.

4 Accessibility

These systems make audio content more accessible to accessible to individuals with hearing impairments by impairments by providing accurate transcriptions and transcriptions and summaries of spoken content. content.

Core Components of Audio Retrieval Systems

Automatic Speech Recognition (ASR)

ASR technology converts converts spoken language language into text, forming forming the foundation of of audio retrieval systems. systems. Modern ASR employs deep learning models to achieve high accuracy across various accents and languages.

Keyword Spotting (KWS)

KWS algorithms identify specific terms or phrases within audio streams, enabling rapid search capabilities. This component is crucial for real-time applications and can work independently of full transcription.

Speaker Identification

This component uses voice voice biometrics to recognize individual speakers, facilitating speaker-based search and and segmentation of multimulti-speaker audio content.

Natural Language Processing (NLP)

NLP techniques analyze the the context, semantics, and and intent behind spoken spoken words, enabling more nuanced understanding and retrieval retrieval of audio content. content.

The Audio Retrieval Process

1 _____ Audio Ingestion and Preparation

The process begins with the collection and digitization of audio content from various sources.

2 Preprocessing

Audio signals undergo preprocessing to enhance quality. This includes noise reduction, speaker diarization, and segmentation of segmentation of continuous audio streams into manageable chunks.

3 ____ Transcription

ASR technology converts the preprocessed audio into text transcripts. Advanced systems may generate time-aligned transcripts for aligned transcripts for precise audio-text synchronization.

Indexing and Querying

The transcribed text and associated metadata are indexed for efficient searching. Query processing involves interpreting user queries interpreting user queries and matching them against the indexed content.

Result Filtering and Presentation

Retrieved results are filtered for relevance and accuracy. The system then presents the results to the user, often with options to play options to play the corresponding audio segments.

Key Technologies in Audio Retrieval



ASR Tools

Industry-leading ASR solutions like Google Speech-to-Text and OpenAI's Whisper Whisper provide highly accurate transcription capabilities, supporting multiple multiple languages and accents.



Machine Learning Models

Advanced ML models, including transformers and LSTM networks, enhance transcription accuracy and enable context-aware audio analysis for improved retrieval performance.



NLP Tools

NLP frameworks such as spaCy and NLTK facilitate semantic analysis, entity recognition, and sentiment analysis of transcribed audio content.



Database Solutions

Specialized databases like Elasticsearch offer high-performance indexing and querying and querying capabilities optimized for large-scale audio retrieval applications.



Challenges in Spoken Language Audio Retrieval

Speech-to-Text Accuracy

Variations in accents, dialects, and speaking styles can can significantly impact transcription accuracy.

Background noise and poor audio quality further complicate this challenge, requiring robust noise reduction and adaptation techniques.

Data Privacy and Security

Handling sensitive audio data, such as personal conversations or confidential business meetings, raises raises significant privacy concerns. Implementing robust robust encryption, access controls, and data anonymization techniques is crucial.

Context Understanding

Interpreting nuances like sarcasm, irony, and cultural cultural references remains a significant challenge for AI for AI systems. Misinterpretation of context can lead to lead to inaccurate retrieval results and misclassification of misclassification of audio content.

Scalability and Performance

Processing and indexing vast amounts of audio data in data in real-time presents significant computational computational challenges. Optimizing algorithms and and leveraging distributed computing architectures are architectures are essential for large-scale deployments. deployments.

Practical Applications of Spoken Language Retrieval

Customer Service

Audio retrieval systems enable real-real-time sentiment analysis in call call centers, identifying customer emotions and intent. This allows for for immediate escalation of critical critical issues and personalized customer interactions. Historical analysis of call recordings can reveal reveal trends in customer satisfaction satisfaction and common pain points. points.

Media and Broadcasting

Content creators use audio retrieval to retrieval to organize vast libraries of of podcasts, interviews, and news segments. This facilitates content discovery, enables precise editing, and editing, and supports the creation of creation of highlight reels or themed themed compilations. Live captioning captioning and content moderation in moderation in broadcasting also rely rely heavily on these technologies. technologies.

Healthcare and Legal Fields

In healthcare, audio retrieval assists in assists in analyzing patient consultations, extracting key medical medical information, and maintaining maintaining accurate electronic health health records. Legal professionals use professionals use these systems to to search through depositions, court court proceedings, and recorded interviews, significantly speeding up up case preparation and evidence evidence review processes.

Case Study: Customer Service Audio Retrieval

1

2

3

ŀ

Implementation

A large telecommunications telecommunications company company implemented an audio retrieval system to analyze millions of customer customer service calls. The system transcribed calls in real-time and applied NLP techniques to categorize issues issues and detect customer customer sentiment.

Analysis

The system identified common common customer complaints, complaints, tracked resolution resolution times, and measured the effectiveness of effectiveness of different support strategies. It also flagged calls with negative sentiment for immediate supervisor review.

Optimization

Based on insights from the audio retrieval system, the company optimized its IVR system, retrained support staff staff on common issues, and and implemented a new escalation protocol for dissatisfied customers.

Results

Within six months, the company saw a 15% reduction reduction in average call handling time, a 20% increase increase in first-call resolution resolution rates, and a significant improvement in customer satisfaction scores. scores.

Future of Audio Retrieval Technology

Multilingual Retrieval

Advanced systems will seamlessly handle multiple languages, enabling cross-lingual audio search and translation. This will break down language barriers in global communication and content consumption.

Enhanced Privacy

Future systems will incorporate advanced encryption and anonymization techniques, allowing for secure audio analysis while preserving preserving individual privacy. This will be crucial crucial for widespread adoption in sensitive domains.

3

Real-time Analysis

Improvements in processing power and algorithms algorithms will enable instantaneous analysis of analysis of live audio streams, supporting applications like real-time meeting transcription transcription and live event content moderation.

moderation.

Emotion Al Integration

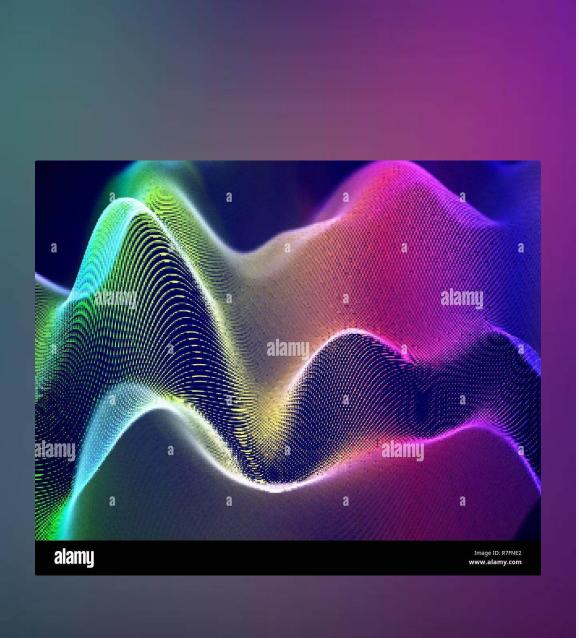
Integration of advanced emotion recognition AI recognition AI will enable more nuanced understanding of spoken content, improving improving applications in fields like mental health health monitoring and immersive media experiences.

Conclusion

Spoken language audio retrieval technology stands at the forefront of the AI revolution, transforming how we interact with and extract value from audio content. n summary, spoken language audio retrieval systems represent a significant advancement in managing and accessing spoken-word content. These systems transform complex audio data into searchable and structured information, making it easier to find and use specific content within recordings. With the integration of ASR and NLP technologies, they offer precise, relevant search results, enhancing user experiences across various fields. As more organizations adopt these systems, we can expect them to further streamline workflows, improve accessibility, and unlock valuable insights from audio data, shaping the future of information retrieval.

•

Thank You



Non-Speech Audio Retrieval

Welcome to this presentation on Non-Speech Audio Retrieval. This presentation will cover the importance, applications, and techniques of retrieving non-speech audio in various fields such as music, film, and sound design.

Team:

Sai Ruthviz (21011P0527) Sai Saketh (21011P0528)



Introduction to Audio Retrieval

1 Importance of Non-Speech Audio Retrieval

Non-speech audio retrieval plays a crucial role in various industries and applications.

2 Wide Range of Applications

Audio retrieval is essential in fields such as music, film, and sound design.

3 Enhancing Media Production

Effective audio retrieval techniques contribute to improved media production processes.

Context and Motivation

Challenges in Traditional Traditional Methods

Traditional sound retrieval methods face various challenges challenges in effectively organizing and accessing audio audio content.

Growing Importance

The need for efficient audio retrieval is increasing in media production and entertainment industries.

Improving Accessibility

Effective sound retrieval methods aim to enhance accessibility and usability of audio content in various applications.



SoundFisher Overview

User-Extensible System

SoundFisher is a user-extensible sound classification and retrieval system.

Multidisciplinary Approach

It integrates multiple disciplines including signal processing, psychoacoustics, speech recognition, computer music, and multimedia databases.

Advanced Classification

SoundFisher employs advanced techniques for accurate sound classification and retrieval. retrieval.

Versatile Applications

The system can be applied in various fields requiring efficient audio content management. management.



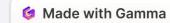
Acoustic Feature Vector

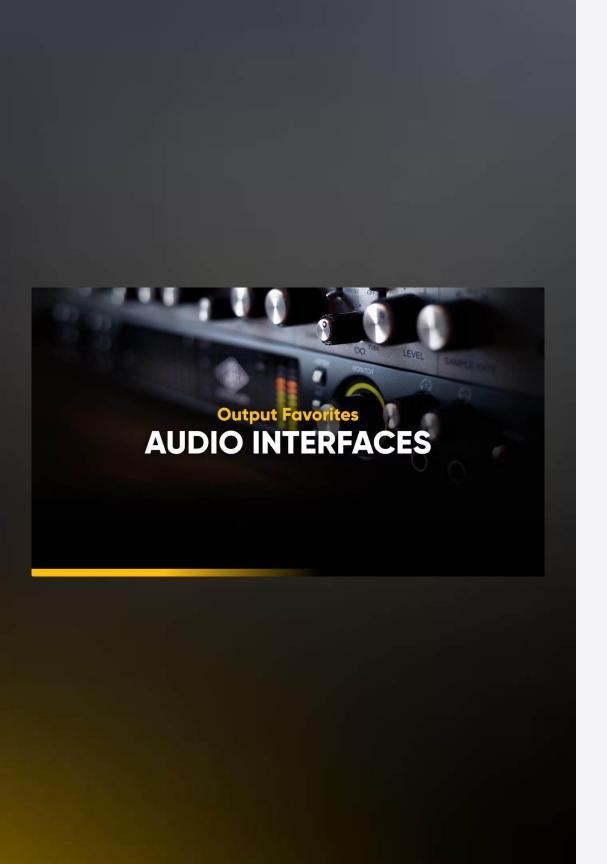
Acoustic Features

- Duration
- Loudness
- Pitch
- Brightness

Comparison

These acoustic features are similar to visual feature vectors used in image indexing.





Query Capabilities

Basic Queries

Search sounds within specified feature ranges.

____ Complex Queries

Example: Find sounds similar to barking.

_____ Weighted Queries

Example: Find sounds with average pitch > 2000 Hz Hz and 0.7 weighted values for specific properties (e.g., metallic or plucked sounds).



System Training

1 User-Provided Examples

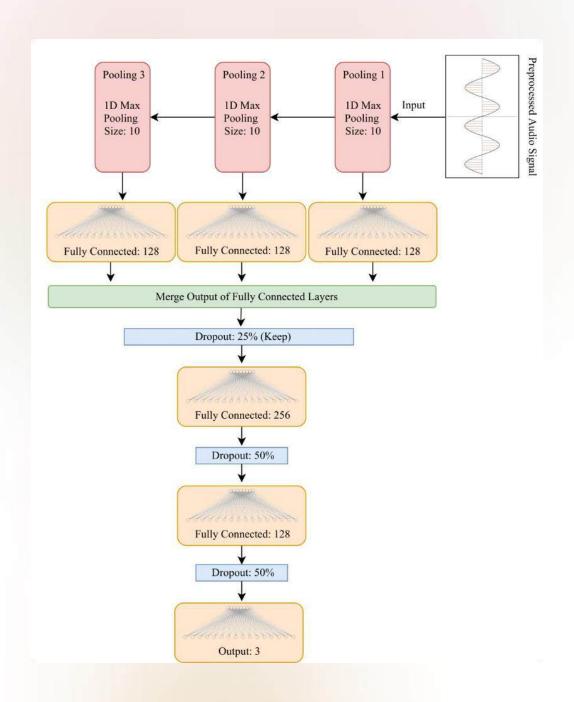
Users can provide examples to enhance retrieval capabilities capabilities of the system.

2 Perceptual Properties

The system can be trained to retrieve sounds based on perceptual properties like "scratchiness" and "buzziness."

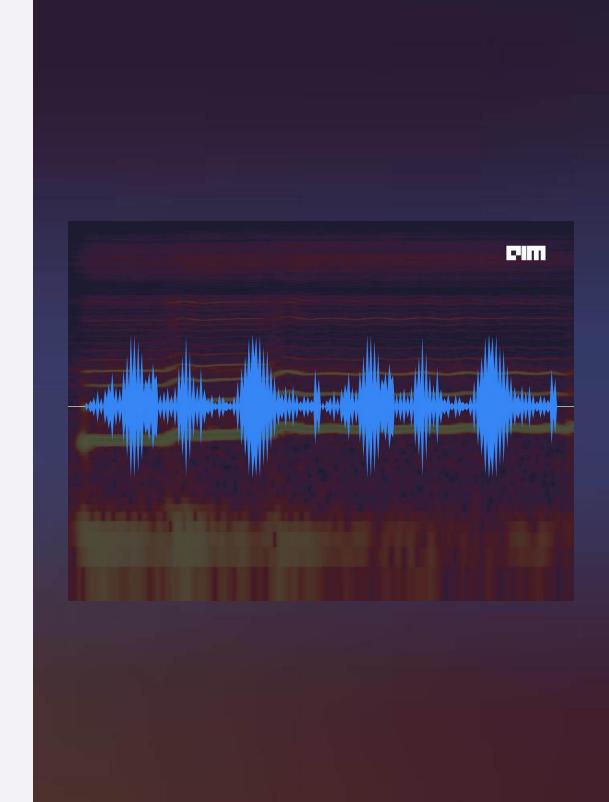
3 Continuous Improvement

Training allows the system to continuously improve its retrieval accuracy and effectiveness.



Performance Evaluation

Database Size	400 diverse sound files
Sound Types	Nature, animals, instruments, speech
Evaluation Focus	Key findings from performance assessment assessment





Conclusion

1

2

3

Importance

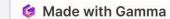
Recap the importance of non-speech audio retrieval in various industries.

Impact

Highlight the significance of systems like SoundFisher in improving audio content management.

Future Prospects

Discuss potential future developments in non-speech audio retrieval technology.



Thank you



Graph Retrieval In Multimedia Information Retrieval

Team:

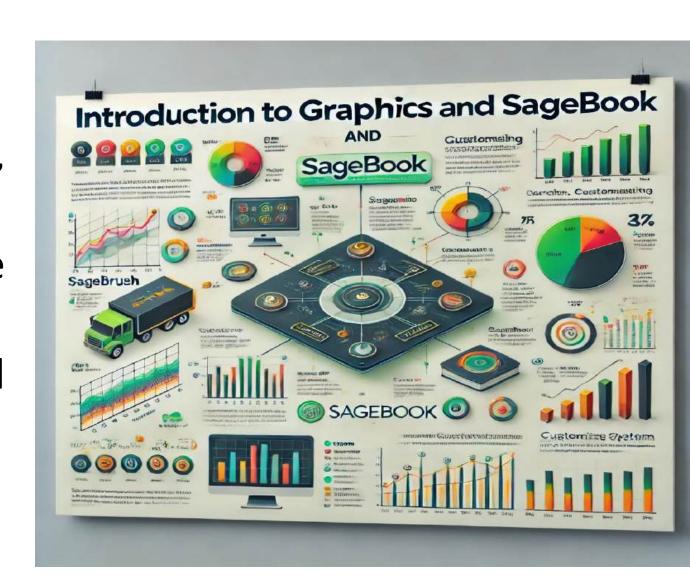
21011P0529

21011P0530



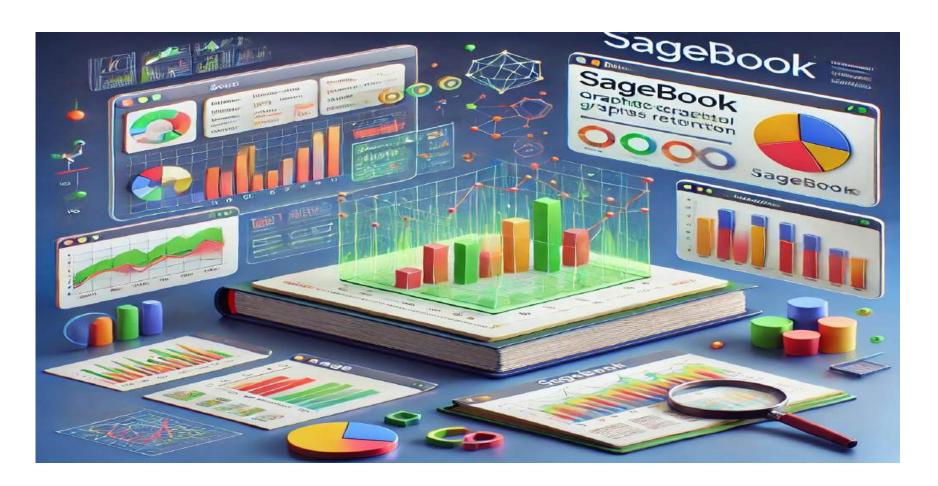
Introduction to Graphics and SageBook

- Graphics include tables and charts (column, bar, line, pie, scatter).
- Graphs are constructed from primitive elements like points, lines, labels.
- SageBook: A system for graph retrieval developed at Carnegie Mellon University.



SageBook: An Innovative Graph Retrieval System:

- SageBook enables search and customization of stored data graphics.
- Supports data-graphic query, content description, indexing, search, and adaptation.



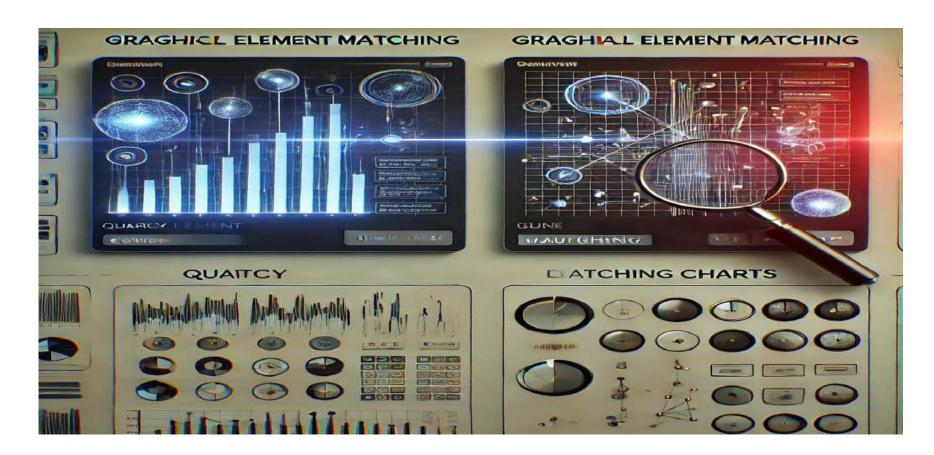
SageBrush: Graphical Query Interface

- Queries are formulated via a graphical direct-manipulation interface.
- Users arrange spaces (charts, tables), objects (marks, bars), and object properties (color, size, shape, position).



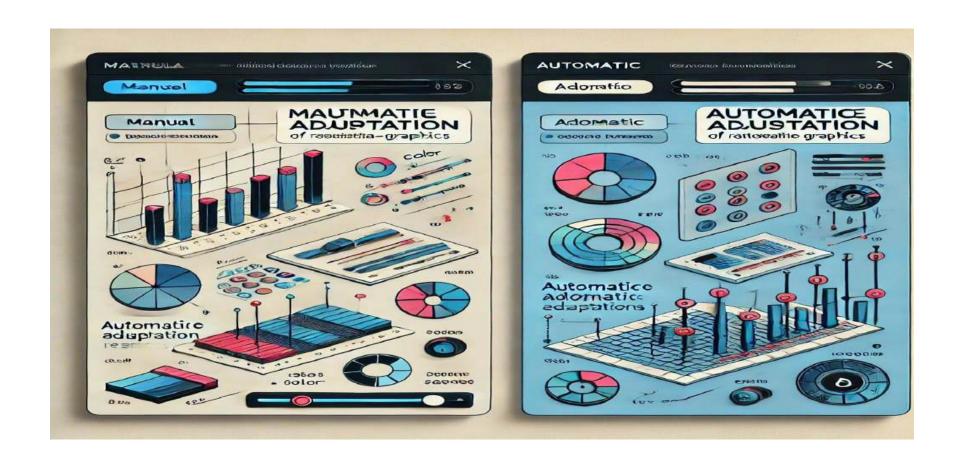
Graphical Element Matching

- SageBook performs exact and similarity-based matching on graphical elements and underlying data.
- Example: Two graphemes must be of the same class (bars, lines, marks) and use the same properties (color, shape, size, width).



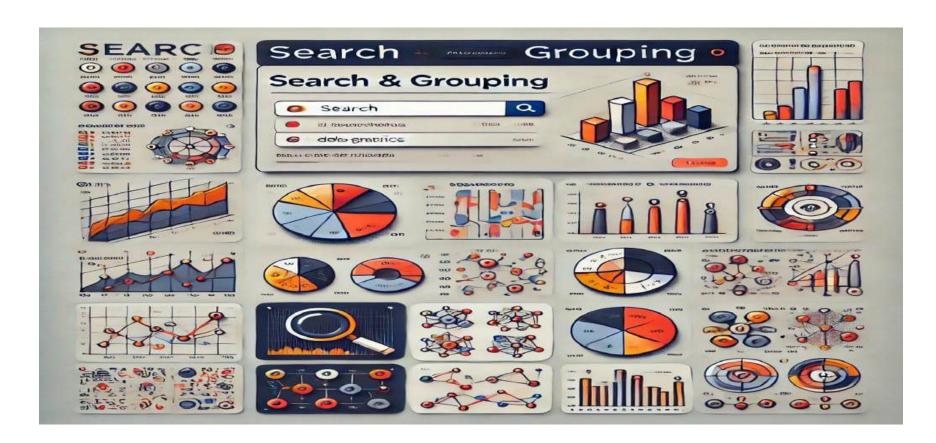
Manual and Automatic Adaptation

- Retrieved data-graphics can be manually or automatically adapted.
- SageBook maintains a representation of syntax, semantics, and spatial relationships in data-graphics.



Search and Grouping Techniques

- Search is performed on graphical and data properties with varying degrees of match relaxation.
- Clustering techniques enable browsing large collections based on data and graphical properties.



• Narrow the search (strict matching): This setting finds only graphs that exactly match the specified criteria. For example, if a user searches for a horizontal bar chart with specific colors and sizes, the system will return only those graphs that perfectly fit those parameters.

• Broaden the search (relaxed matching): This allows the system to return graphs that are similar but not identical to the query. For instance, if the user queries for a bar chart, SageBook might return graphs with slight variations in color, size, or shape but still capture the essence of the requested visual format.

Applications of Graph Retrieval Beyond Business:

- Cartography: Retrieving maps showing terrain, elevation, and features.
- Architecture: Searching for blueprints.
- Communications & Networking: Visualizing network diagrams with routers and connections.
- Systems Engineering: Displaying components and their connections.
- Military Planning: Maps showing forces and defenses, overlayed with strategic information.

THANK YOU

Video Retrieval

21011P0535 21011P0536

Introduction to Video Retrieval

Content-Based Access: The ability to search and retrieve video content using metadata or content analysis, enabling easier access to various types of video, such as video mail, recorded meetings, surveillance footage, and broadcast television.

Importance of Content-Based Access:

- Improved Efficiency
- Ease of Use
- Time Savings
- Enhanced Accessibility

Broadcast News Navigator

- BNN is a web-based tool that captures, annotates, segments, summarizes, and visualizes news stories from broadcast videos.
- BNN employs text, speech, and image processing technologies to conduct multi-stream analysis, facilitating
 effective content-based search and retrieval.

• Comparison:

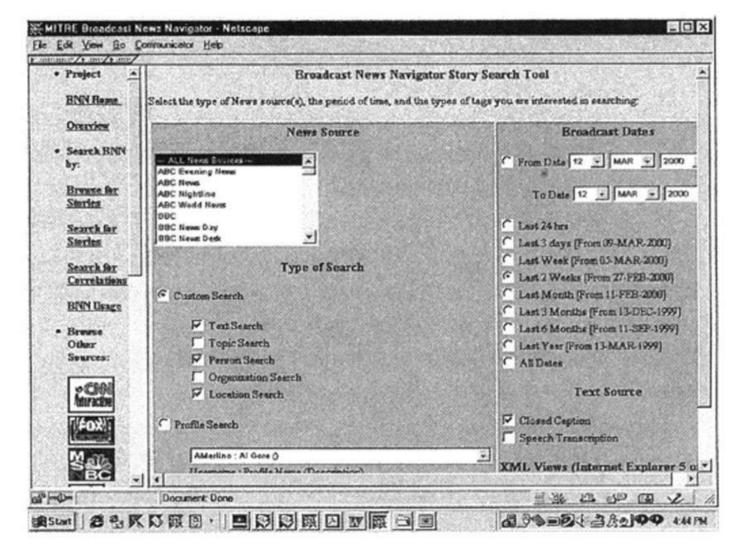
 While QBIC focuses on static images, BNN serves a similar purpose for dynamic video content, bridging the gap between traditional search methods and modern multimedia retrieval.

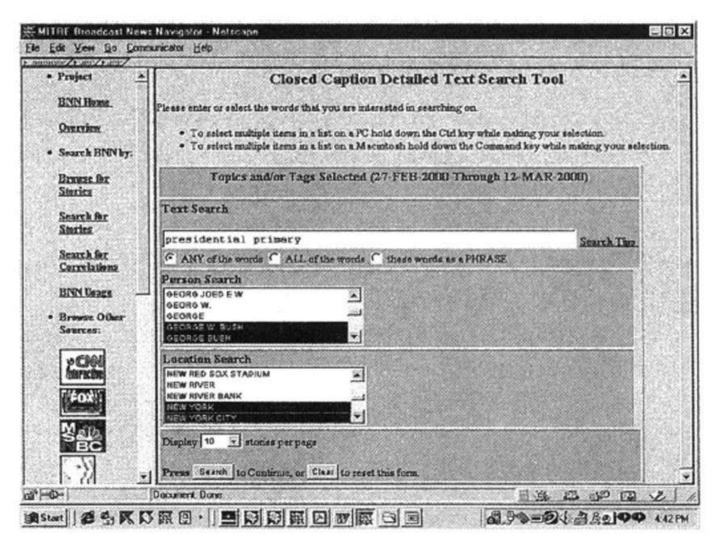
Benefits of BNN:

- Automation
- User Query Interface: The system allows users to specify various parameters for searching video content.

Query Process

Automatic creation of user-friendly query pages:





Initial Query Page

Detailed Content-based Query

Figure 10.6: Overview of BNN Results Display

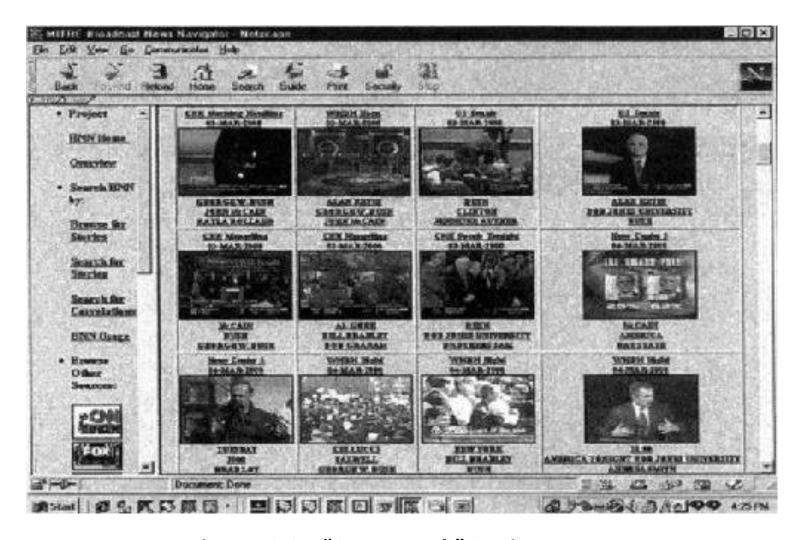


Figure 10.6a: "George Bush" Stories

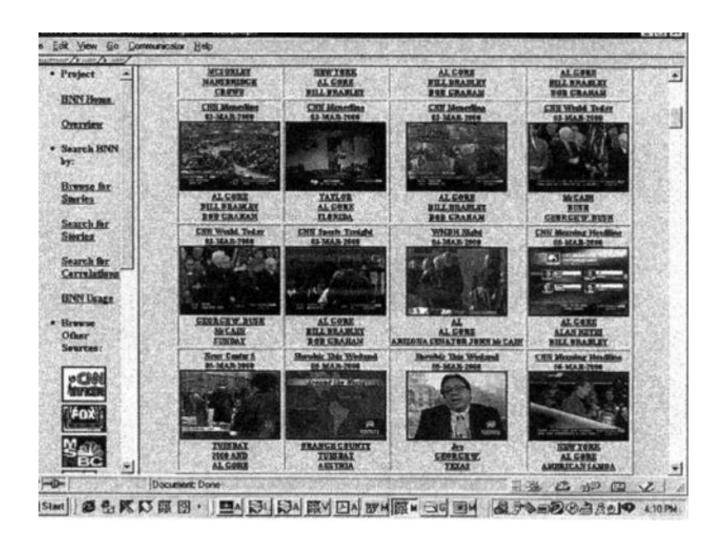


Figure 10.6b. "Al Gore" Stories

Figure 10.6: Overview of BNN Results Display

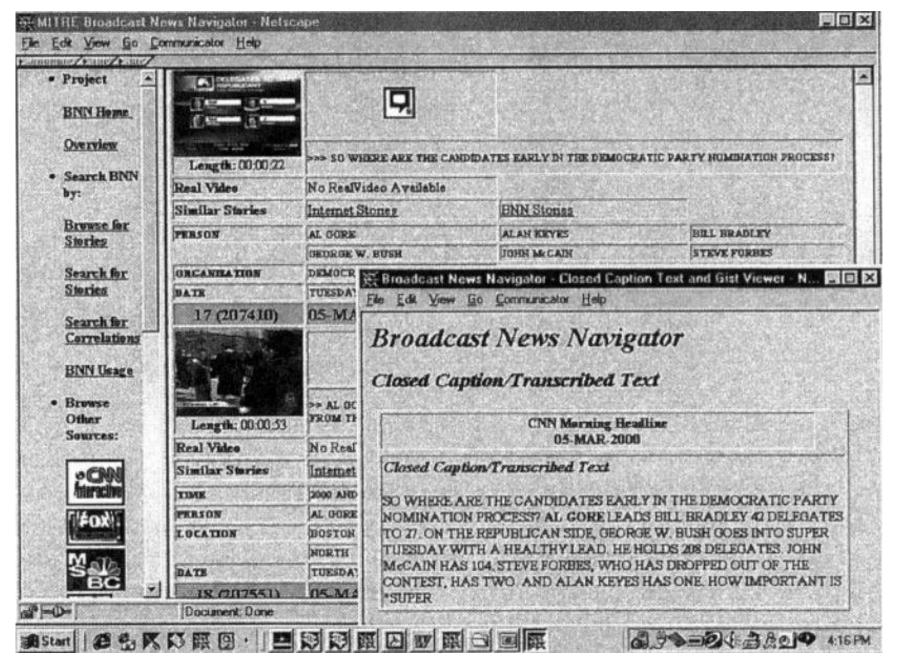


Figure 10.6c: BNN Story Detail Retrieval of "presidential primary" stories

Figure 10.6: Overview of BNN Results Display

This figure provides an overall view of the Broadcast News Navigator (BNN) system, which allows users to retrieve broadcast news videos using a single query across multiple news sources. BNN is designed for efficient video retrieval, offering both a "story skim" view for quick scanning and a detailed view for indepth information.

Figures 10.6a and 10.6b: "Story Skim" View for "Bush" and "Gore" Stories

"Story Skim" Display: Figures 10.6a and 10.6b show BNN's "story skim" feature, which displays each story with a keyframe (a representative still image from the video) and highlights the three most frequently occurring named entities (like people, organizations, and locations).

Context of Named Entities: In Figure 10.6a, stories related to "George Bush" are displayed, while Figure 10.6b displays stories related to "Al Gore." These named entities allow users to quickly understand the story's content without viewing the video, enhancing efficiency in finding relevant information.

User Interaction: In Figure 10.6b, when the user clicks on the "5 March" story (located in the second row, far right), it opens the detailed view shown in Figure 10.6c.

Figure 10.6c: Detailed Story View with Closed Captions

Detailed Story Information: Figure 10.6c displays detailed information for a specific story selected from Figure 10.6b. Here, users can see various attributes, such as video length and a more detailed keyframe.

Closed Caption/Transcription: When the user clicks on the "closed caption" button, they can view the transcribed text of the video, allowing them to analyze the story in depth without watching it. In this example, the closed caption shows information about the presidential primary and the number of delegates each candidate has won.

Detailed Analysis: This feature is useful for comparative analysis, enabling users to gain detailed insights into the story content by reading the transcription.

Additional Features of BNN

Direct Access and Browsing Options: BNN allows users to either directly search with a query or browse stories from specific time intervals or news sources, making it adaptable to different user needs.

Named Entity Frequency Graph: BNN includes a feature to display a graph showing how often specific named entities appear over time, allowing users to analyze trends in news coverage.

Correlation Search: BNN's "search for correlations" link lets users explore relationships between named entities across multiple stories, enabling in-depth data mining.

BNN's Efficiency and Evaluation

Empirical Study Findings: According to a study by Merlino and Maybury (1999), users performed better in retrieving relevant content by focusing on the three most frequent named entities in the "story skim" view rather than viewing full story details. This setup significantly improved retrieval performance.

Performance and Efficiency: BNN enabled users to find relevant video content around six times faster than traditional keyword searches, thanks to automated story segmentation and features like named entity highlights and keyframes, which improved user precision and recall.

Topic Detection and Tracking

The **Topic Detection and Tracking (TDT)** initiative, started by Wayne in 1998, is aimed at enhancing news retrieval and analysis by developing algorithms for:

- Story Segmentation
- Topic Tracking
- Topic Detection

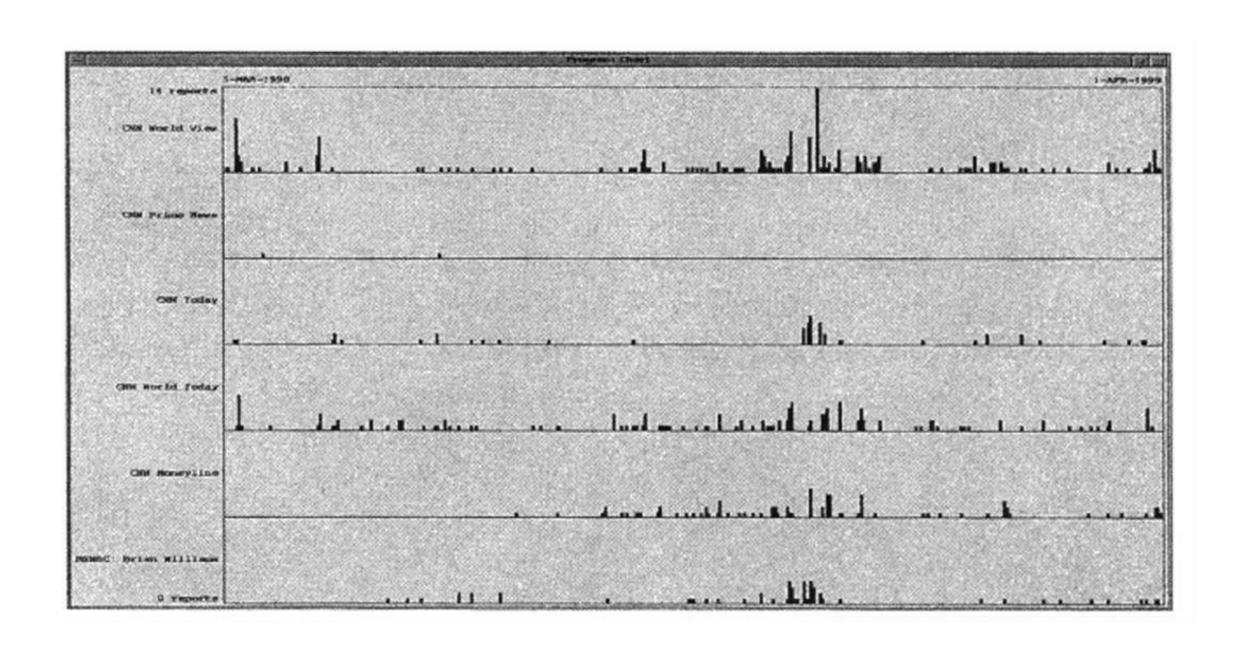
The **Geospatial News on Demand Environment (GeoNODE)**, developed by Hyland et al. in 1999, provides a way to access, analyze, and navigate news stories within a spatial and temporal framework. Key aspects include:

- Geospatial and Temporal Context
- Multi-Source Access

Key Components of GeoNODE:

- Information Extraction
- Data Mining and Correlation
- Visualization Tools

GeoNODE Time-line



GeoNODE Map Display Interface Description

Zoom View (Left Panel)

Displays a focused, detailed map view with symbols and color-coded regions that highlight specific geographic data. Countries are color-coded to indicate document frequency, where darker hues represent more mentions. This color saturation effectively shows the distribution of reports, with regions such as North and South America displayed in darker shades to indicate higher reporting frequencies, while areas like Africa appear in lighter colors due to fewer mentions.

World View (Right Panel)

Offers a global perspective, supporting the Zoom View with additional context for analyzing data distribution on a worldwide scale. This view also features color-coded countries and circle markers, allowing users to quickly assess and explore reporting patterns globally.

Attributes Table (Center Panel):

Contains a structured table of attributes that list information for different locations, such as country names, population, and the number of document mentions per location. Includes sorting and filtering tools for easy data navigation and to locate specific information efficiently.

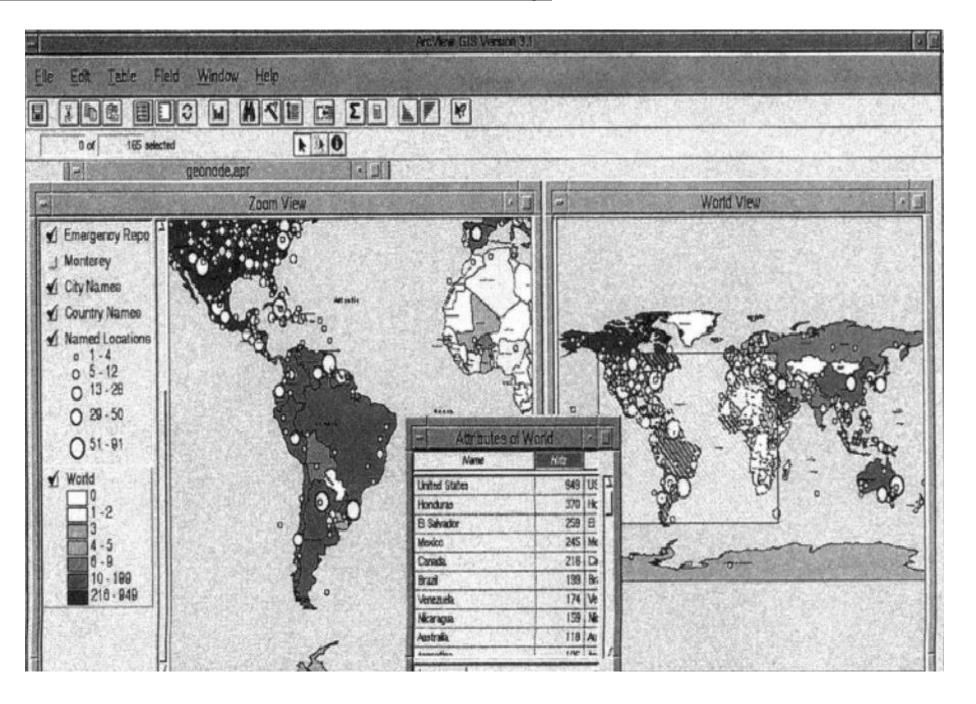
Performance and Analytics:

GeoNODE achieves high accuracy, with preliminary evaluations on a dataset of 65,000 documents and 100 manually identified topics showing over 80% accuracy in topic identification and 83% detection of stories within topics, maintaining a low misclassification error rate (0.2%). This level of performance is comparable to results from similar initiatives like the Topic Detection and Tracking (TDT) project. This high precision makes GeoNODE a robust tool for geospatial data analysis and topic detection in large text corpora.

Future Capabilities and Integration:

GeoNODE aims to extend its functionality to handle diverse media types, including text, audio, and video, to enhance its analytical capabilities. Progress depends on the development of multimedia corpora, standardized evaluation tasks, and advanced machine learning approaches to ensure high-performance data extraction and analysis across various media sources.

Figure 10.8: GeoNODE Map Display



THANK YOU!