

Scalable and Cost-Effective Deduplication: Leveraging Algorithms and LLMs

SponsorMotion

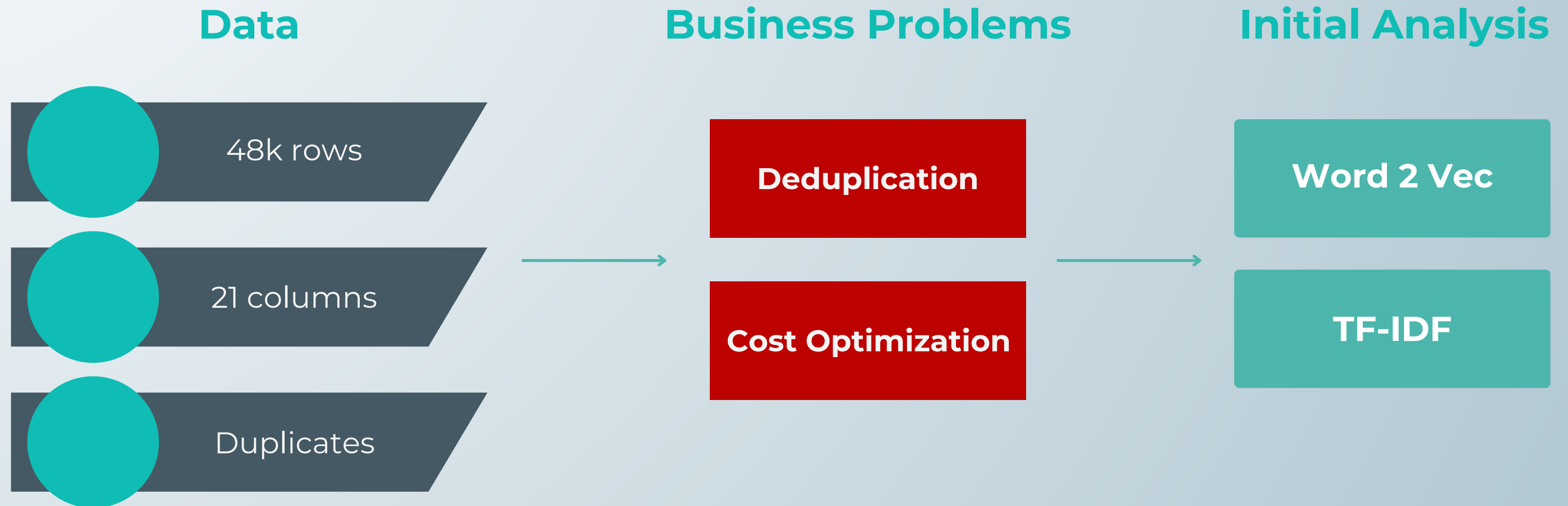
Rohan Chaudhary, Sarmad Kahut, Valentina Torres



Table of Contents

1. Business Problem & Data and Initial Analysis
2. Research and Initial Approach & Final Approach and Method
3. Key Findings & Results
4. Cost Optimization
5. Business Impact & Implications
6. Known limitations and avenues for further improvement
7. Summary & Conclusions

Business Problem, Data and Initial Analysis



Deduplication Algorithm; Python Notebook Workflow



Dropped all rows with no name, no state, no start date. Data went from 48k to 45k



Removed special characters and blank spaces, all names were changed to lowercase



Picked 6 states with most records to test our approach (CA, MA, NY, TX, FL, and NV)



Found exact duplicates, with same name and same start date (+/- 1 day)



Kept duplicate records with longest summaries and dropped the rest



Applied fuzzy matching to remaining records and applied 3 thresholds (75, 80, and 85)



Added 'Human Verification' column to show rows needing attention (Yes; similar name but different start dates, No; similar names and date)

Final Approach and Method

Skip to Main Content
All Events
SGO ANNUAL MEETING ON WOMEN,ÂØS CANCER
SGO Annual Meeting on Women,ÂØS Cancer
Registration
Hotel Information
Schedule
Accreditation Information
Exhibits And Sponsorships
Exhibitor Information
Key Concepts
Destination Information & Activities
Industry Supported Symposia
Promotional Resources
Request Meeting Space
Conflict Of Interest
Media & Press
Mar 25 - Mar 28, 2023 Tampa, Florida
Tampa Florida
SGO 2023 ANNUAL MEETING ON WOMEN,ÂØS CANCER
Join us March 25,ÂØ 28, 2023, at the Tampa Convention Center in Tampa, FL.
The SGO Annual Meeting on Women,ÂØS Cancer-ÂØE is the premier educational and scientific event for those who treat and care for women with gynecologic cancer.
Since 1969, women,ÂØs health care professionals have convened at the Society of Gynecologic Oncology (SGO),ÂØs Annual Meeting on Women,ÂØs Cancer to discuss the latest science in the field, receive educational programming, and to network.
Members of the entire gynecologic cancer care team who provide treatment and care in the areas of chemotherapy, radiation therapy, surgery, and palliative care attend the SGO Annual Meeting. Gynecologic oncologists make up a large population of the attendees, along with medical oncologists, pathologists, radiation oncologists, hematologists, surgical oncologists, obstetrician/gynecologists, nurses, physician assistants, fellows in training, residents, and pharmacists.
The SGO Annual Meeting also brings together exhibitors from a variety of medical device and service companies. An array of state-of-the-art products and services geared towards members of the gynecologic cancer care team are on display in the exhibition area. There are numerous sponsorship opportunities as well for various events, services and takeaway items at the SGO Annual Meeting.
*Full abstract details posted the day of scheduled presentation.
REGISTER NOW BOOK HOTEL
VIEW SCHEDULE & ABSTRACTS
Abstract titles can be found in the full schedule by clicking into the sessions. Full abstract details posted the day of scheduled presentation.
2023 Program Committee
The Program Committee for the SGO 2023 Annual Meeting on Women,ÂØs Cancer is co-chaired by Dineo Khabele, MD, from Washington University in St. Louis, and Paul DiSilvestro, MD, from the Women & Infants Hospital of Brown University.
CO-CHAIRS
Dineo Khabele, MD
Paul DiSilvestro, MD
STEERING COMMITTEE
Floor Backes, MD
Leslie Boyd, MD
Lee-May Chan, MD
Bill Cliby, MD
Bradley Carr, MD
Andrea Hagemann, MD
Emily Hill, MD
Neil Horowitz, MD
Adrienne Mallen, MD
Alexander Olawaaye, MD
Dmitriy Zamarin, MD, PhD
PROGRAM COMMITTEE
Victoria Bae-Jump, MD, PhD
Joyce Barlin, MD
Evelyn Cantillo, MD, MPH
Allan Coorens, MD, FRCS
Keichi Fujiwara, MD, PhD
Stephanie Gaillard, MD, PhD
Sophia George, PhD
Lindsay Hasselwander, MSN
Gloria Huang, MD
Marilyn Huang, MD
Tilley Jenkins Vogel, MD
Andrea Jewell, MD
Emily Ko, MD, MSCR
Elizabeth Lokich, MD
Jyoti Mayadev, MD
Alexander Melamed, MD, MPH
Navya Nair, MD
Michaela Onstad Grimsfelder, MD, MPH
Neil Shipper, MD
Monica Prasad Hayes, MD
Angelica Rodrigues, PhD
Mian Shahzad, MD, PhD
Kristy Ward, MD
Annelise Wilhite, MD
Israel Zigelboim, MD
Annie Ellis (Patient Advocate)
Jubilee Brown, MD (Board Liaison)
THANK YOU FOR YOUR INTEREST IN EXHIBITING AT THE SGO 2023 ANNUAL MEETING!
More than 2,000 leaders in the field of gynecologic oncology convene every year at the SGO Annual Meeting on Women,ÂØs Cancer-ÂØE. Attendees come from the entire gynecologic cancer care team who provide treatment and care to cancer patients, including gynecologic oncologists, medical oncologists, pathologists, radiation oncologists, hematologists, surgical oncologists, obstetrician/gynecologists, nurses, physician assistants, fellows in training, residents and social workers.
There is no better time to reach the key decision-makers in gynecologic oncology care than at the SGO Annual Meeting. Interested in Sponsoring or Exhibiting? Please contact Jessie Par'ic, Manager, Industry Relations and Business Development.
Not an SGO member?
Sign up to receive notifications about SGO meetings, education and membership.
Fields marked with an * are required
FIRST NAME *
LAST NAME *
EMAIL *
PROFESSION
SIGN UP
Advertisement
Advertise With Us
Future Dates
MARCH 16,ÂØ MARCH 19, 2024
San Diego, CA
MARCH 15,ÂØ MARCH 18, 2025
Seattle, WA
Meetings & Events
SGO is the premier source for gynecologic oncology education. Several meetings are hosted by SGO throughout the year for professional development.
VIEW UPCOMING

Text Analysis
using LLM

Data Columns

Name

SGO Annual Meeting on Women

Start Date

3/25/2023
12:00:00 AM

Summary

The SGO Annual Meeting on Women,ÂØs Cancer-ÂØE is the premier educational and scientific event for those who treat and care for women with gynecologic cancer. Members of the entire gynecologic cancer care team who provide treatment and care in the areas of chemotherapy, radiation therapy, surgery, and palliative care attend the SGO Annual Meeting. The conference also includes exhibitors from a variety of medical device and service companies.

State

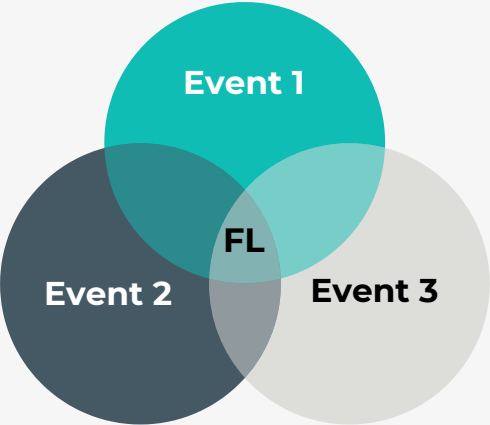
FL

Post Preprocessing

sgoannualmeet
ngonwomen

2023/03/25

=====



Key Findings: Deduplication

Performance Assessment

- Used 3 threshold level
- Six states with the highest number of records
- Actual duplicates = Duplicates identified + FN - FP
- % of Actual duplicates identified = Duplicates identified / Actual duplicates

Accuracy Results

- Higher thresholds = more false negatives
- **Cost function:** estimates the cost associated with each threshold
- **Test dataset:**
 - 15% actual duplicate records to be identified
 - 75% threshold identified 12% of those duplicates

Similarity Threshold	Total Records	Duplicates Identified	False Negatives	False Positives	Actual Duplicates	% of Actual Duplicates Identified
75	5646	680	223	65	838	81.15%
80	5646	530	288	35	783	67.69%
85	5646	426	336	26	736	57.88%

$$\text{Cost (fn)} = 3 \times \text{FN} + 1 \times \text{FP}$$

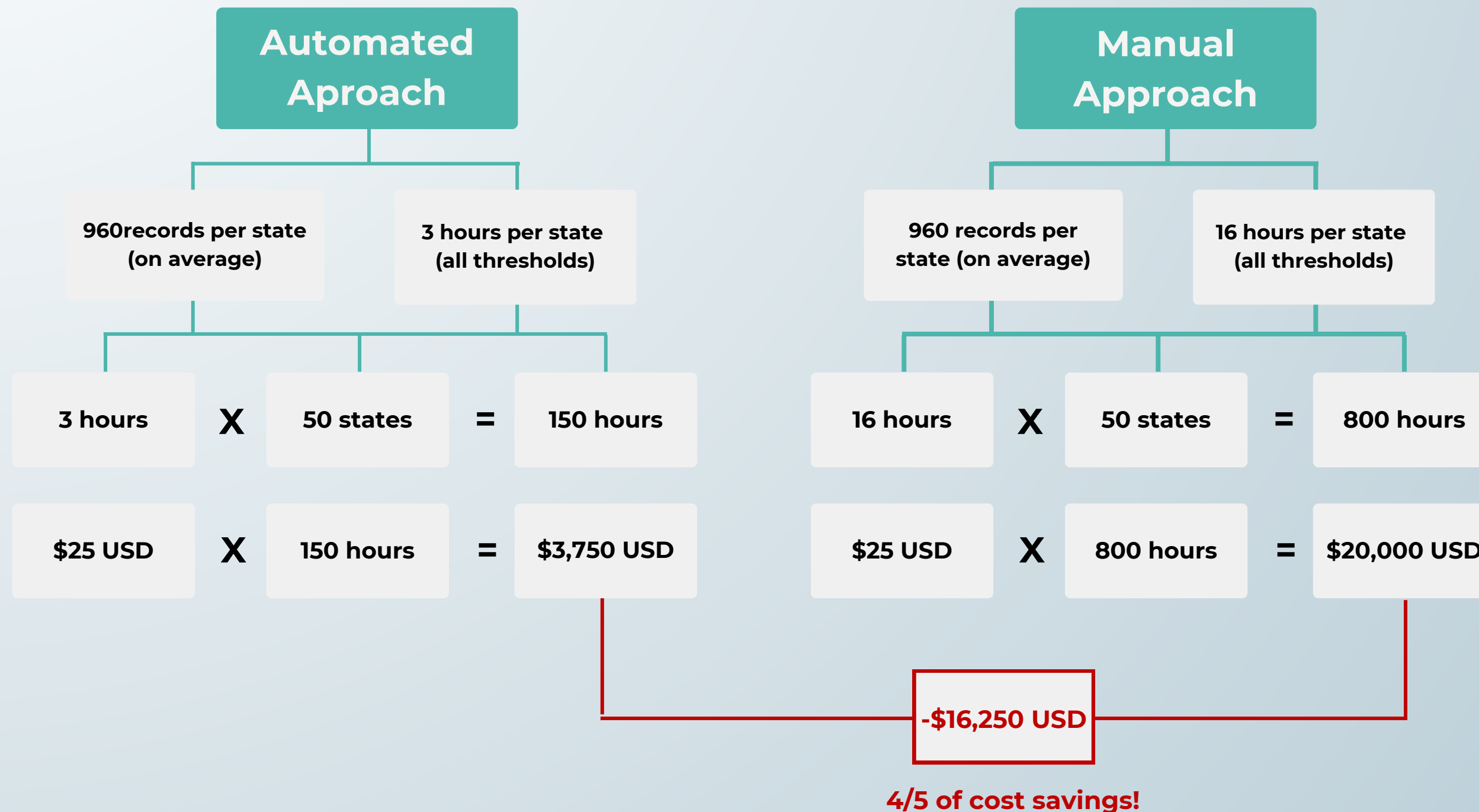
- Cost (Threshold: **75**) = $3 \times 223 + 1 \times 65 = 669 + 65 = \mathbf{734}$
- Cost (Threshold: **80**) = $3 \times 288 + 1 \times 35 = 864 + 35 = \mathbf{899}$
- Cost (Threshold: **85**) = $3 \times 336 + 1 \times 26 = 1008 + 26 = \mathbf{1034}$

Where

FN = Number of False Negatives

FP = Number of False Positives

Cost Optimization



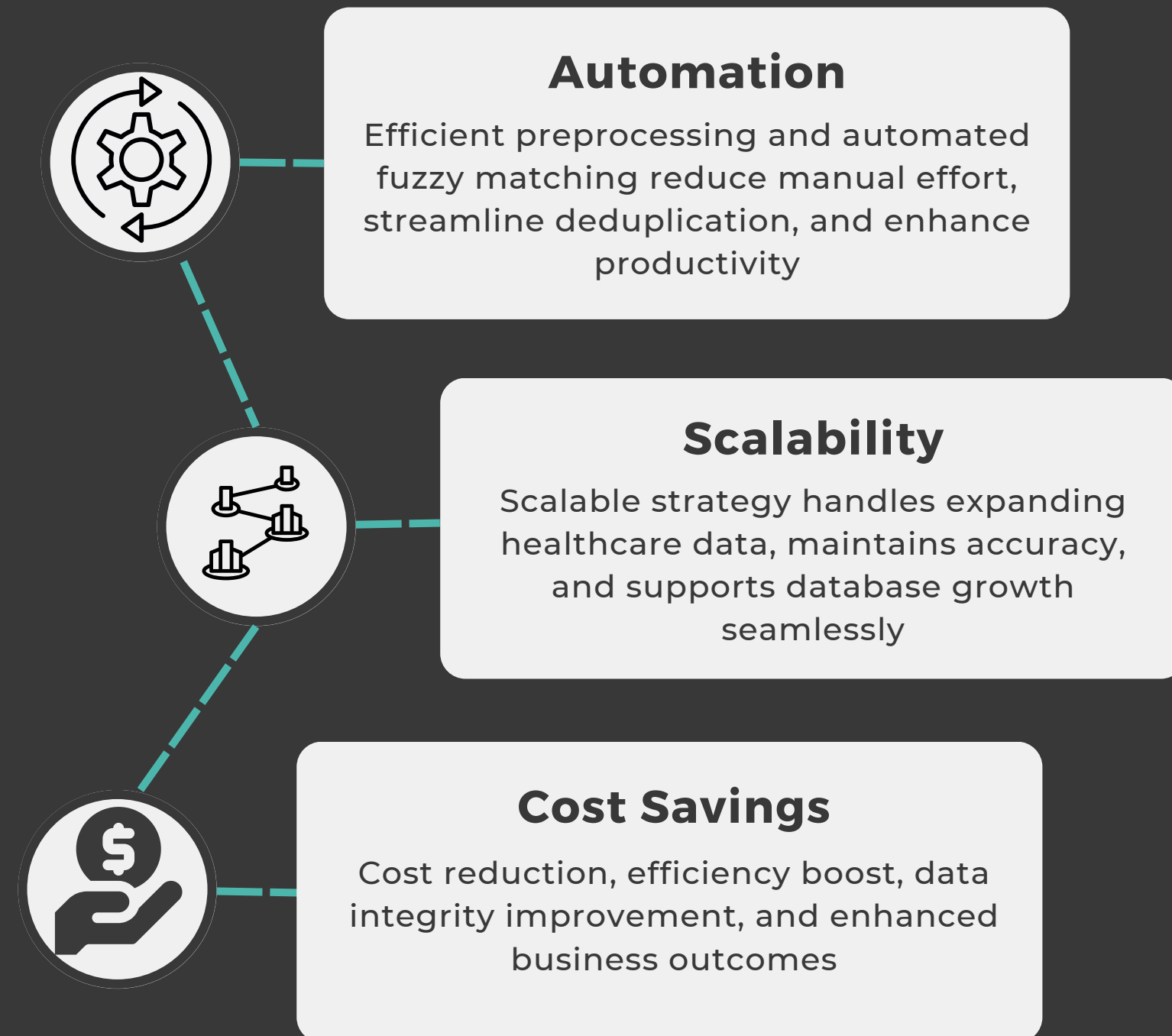
Dynamic Model Selection

- Choose LLMs (GPT-J, GPT-4, ChatGPT) based on:
 - Query needs
 - Balancing accuracy
 - Budget for data extraction
- Pilot Project: Compare GPT-3.5 and frugal LLMs (GPT-J, ChatGPT) on 100 healthcare URLs for accuracy, cost, and data quality insights.

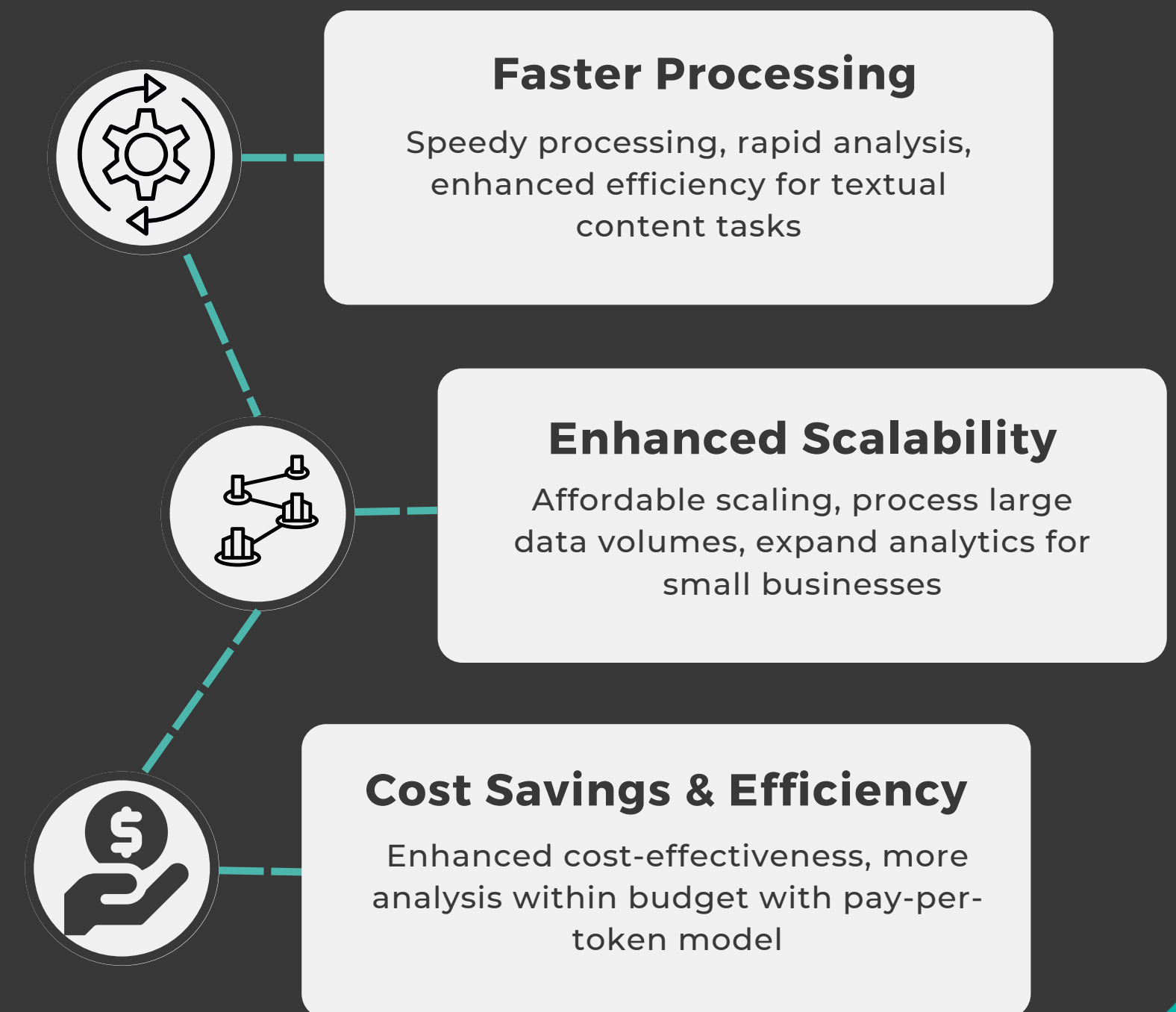


Business Impact & Implications

Deduplication



LLM Optimization



Limitations and Avenues for Further Improvement

Enhancing Data Quality for Accurate Analysis

- Implement filters to highlight ".org" URLs
- Optimize web scraping by eliminating null values in key columns



 www.ncchc.org

Correctional Health Care

Advancing User Experience through Taxonomy Labeling

- **Enhanced Categorization** of events (e.g., oncology, cardiology)
- **Improve user experience** through quick access and filtering based on categories

Balancing Business Benefits with Model Selection

- **Strategic Consideration:** Evaluate trade-offs in accuracy, response time, and comprehensiveness when choosing LLM sizes
- **Workflow Integration:** Integrate LLM cascading with minimal workflow disruption while assessing task suitability



Summary

Our Path of Achievements, Insights and Learning



**Automated Data
Quality Control and
Optimized Post-
Scape Filtering**



**Preprocessed textual
data and explored
similarity detection
techniques**



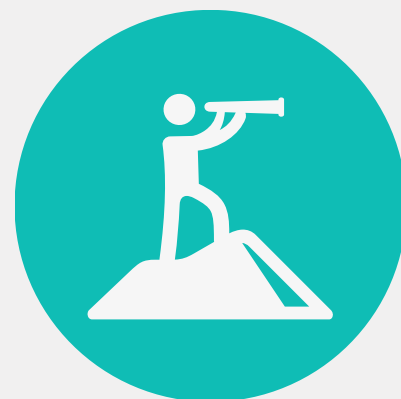
**Achieved 81.15%
Accuracy at 75%
Threshold**



**Developed a cost function
to assess results and
estimated savings through
Deduplication Process
Automation**



**Informed Decision-
Making: Pilot Project
for LLM Comparison:
GPT-3.5, GPT-J,
ChatGPT**



**Future Vision:
Streamlined
Operations and
Enhanced Navigation**



**Business impact and
deliverables aligning
with the objective of
the project**



**Balanced Trade-offs:
accuracy, response
time,
comprehensiveness
in LLM Usage**



the Final Project Report



Thank You for *your* time!

