

# CS533 – Information Retrieval Systems

*SINCAP: SIMPLE NEWS' CONTENT  
eXTRACTION  
APPLICATION*



Gülsüm Ece BIÇAKCI  
Abdurrahman YAŞAR

# Outline

- Introduction
- Methodology
- Algorithm
- Infrastructure
- Experiments
- Results
- Conclusion
- Future Work



# Introduction

- Usage of internet increases
- Newspapers have lots of noise
- Aim: erase the noise and get the content with an application



# Methodology

- Parsing HTML documents
  - Removing unnecessary parts
  - Get the header and related part
- Get the content
- Send it to database
- Showing the data with SINCAP app



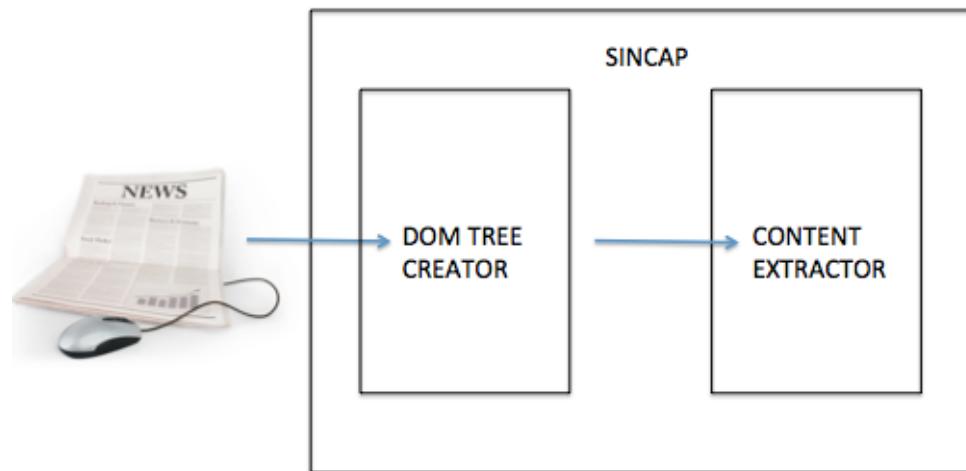
# Algorithm

- Get URLs of news
- Parsing their HTML docs
  - Creating DOM tree
- Find published date
- Look for title
- Get content and header



# Infrastructure

- Java, php, SQLite



# Experiments

- With several newspapers
- News from at most 3 days



# Newspapers...

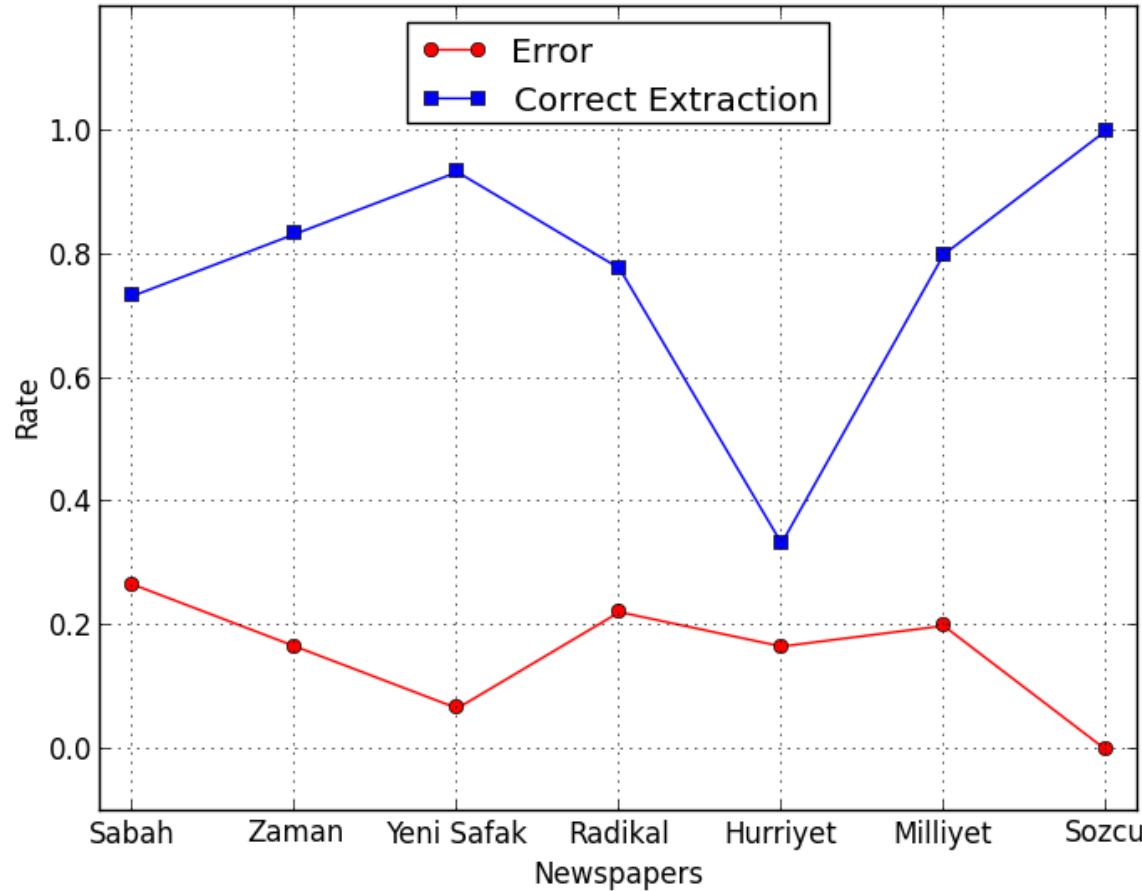


# Results

- Get good results in well structured news
  - Sabah, Yeni Şafak, Milliyet
- Avg results in news have many banners



# Results



# Conclusion

- Solving a problem with Information Retrieval perspective
- Doing a real life application



# Future Work

- Extend the project
- Mobile application form



# Demo...



# Thank You...

...he could trust

...with how his racing  
pastured by the media that  
he will be remembered.  
Funny people, folk.

...in the chapter concluding  
with the exhilarating Grand  
Prix, for example, and his  
spat with his dangerous  
team-mate, his driving is so  
so positive.

ment focused on the 100,000  
prisoners released every year."