

# Topic trends on the Hungarian Internet - news and academic web presence

---

György Kampis, László Gulyás,

Zsolt Jurányi, **Sándor Soós**

Lorand Eötvös University, Budapest, Hungary









# Pilot Study

- News Portals
  - Public discourse
- Academic Institutions
  - Academic discourse

Started in 2013.



# Tools

- Heritix crawler
  - Modified and configured to our purposes
  - Applied by several earlier national archive projects
  - Including at the British Library
- Dell T710 server
  - 2x4 core Xeon E5520
  - 48 GB RAM, 2 TB HDD



# Sources

- Two major news venues
  - Hírarchívum.hu, Index.hu
  - Containing: 155 different online news venues
- Home pages of 48 academic institutions
  - All universities
  - Institutions of the Hungarian Academy of Science
  - Full list is available at <http://hungarianscience.org/>

Hungarian Science.org

THE RESEARCH PERFORMANCE OF HUNGARY: INSTITUTIONAL PATTERNS AND COMPETENCES



# Content Stored

- Sites are revisited
  - ~once in 2 weeks
- Mainly text and video files stored
  - Excluded: exe, gz, iso, jar, mp3, ogg, ppt, rar, waw, xls, xlsx and zip
- No recreation of archived content (long term goal)



# Big Data is Small (for Hungary)

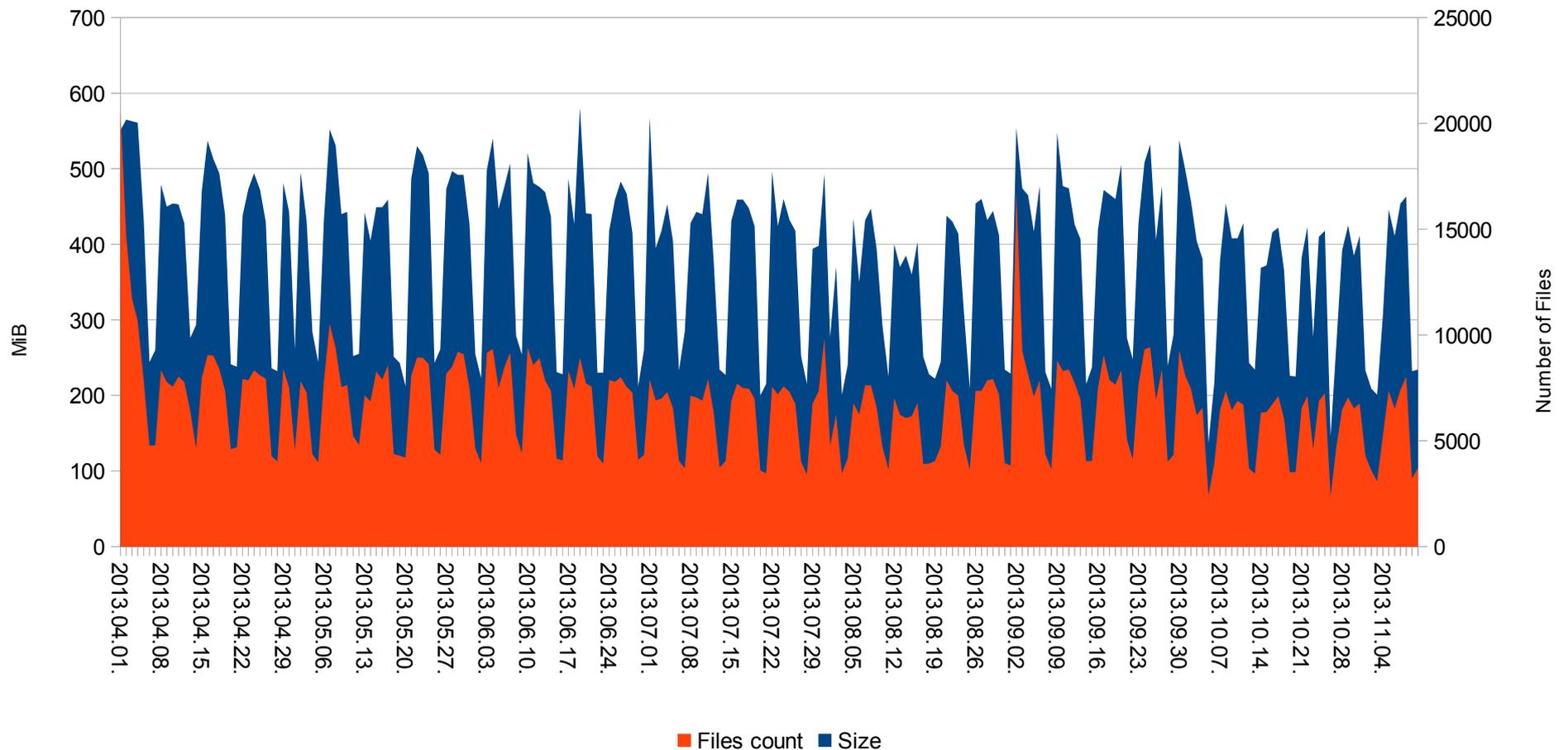
- Overall Quantity is Small
  - Update frequency is also modest
- News (one copy)
  - Size: ~ 600 MiB
  - Num of Files: ~8000 files
- Academia (one copy)
  - Size: 53 Gb
- Yet, distributions are long-tailed



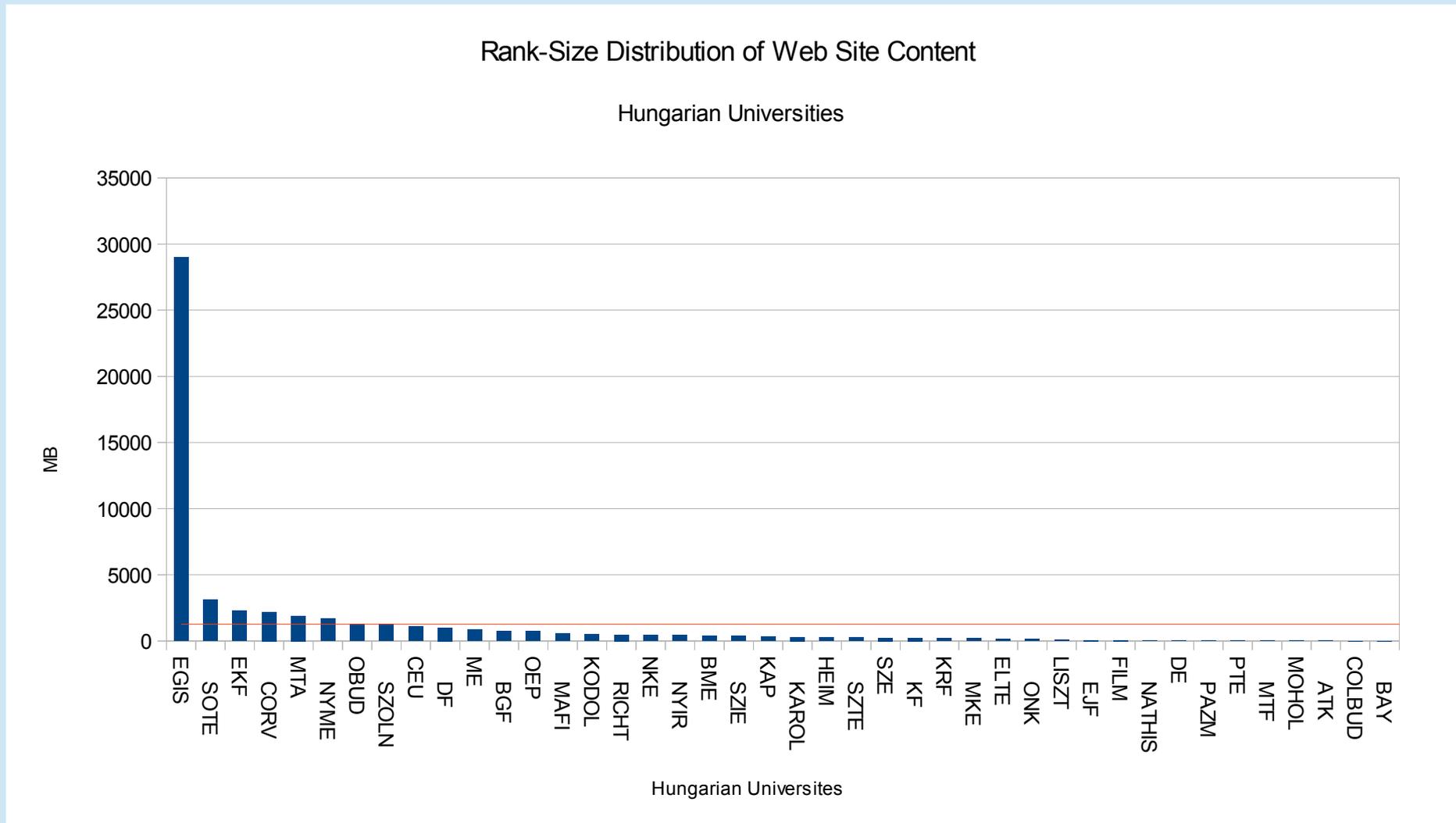
# Size of New Content on News Sites

Size of New Content Downloaded

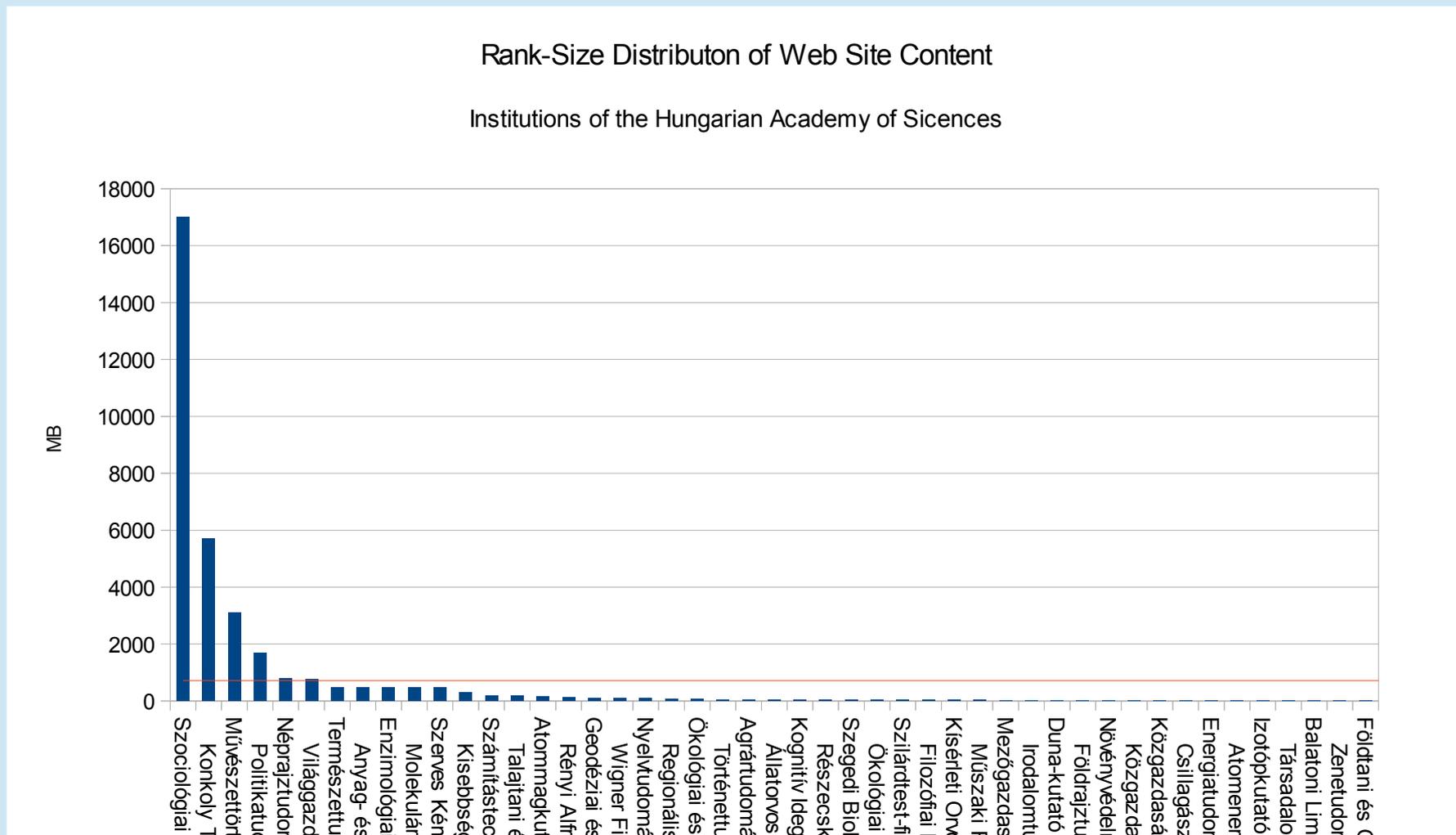
at each new archiving day



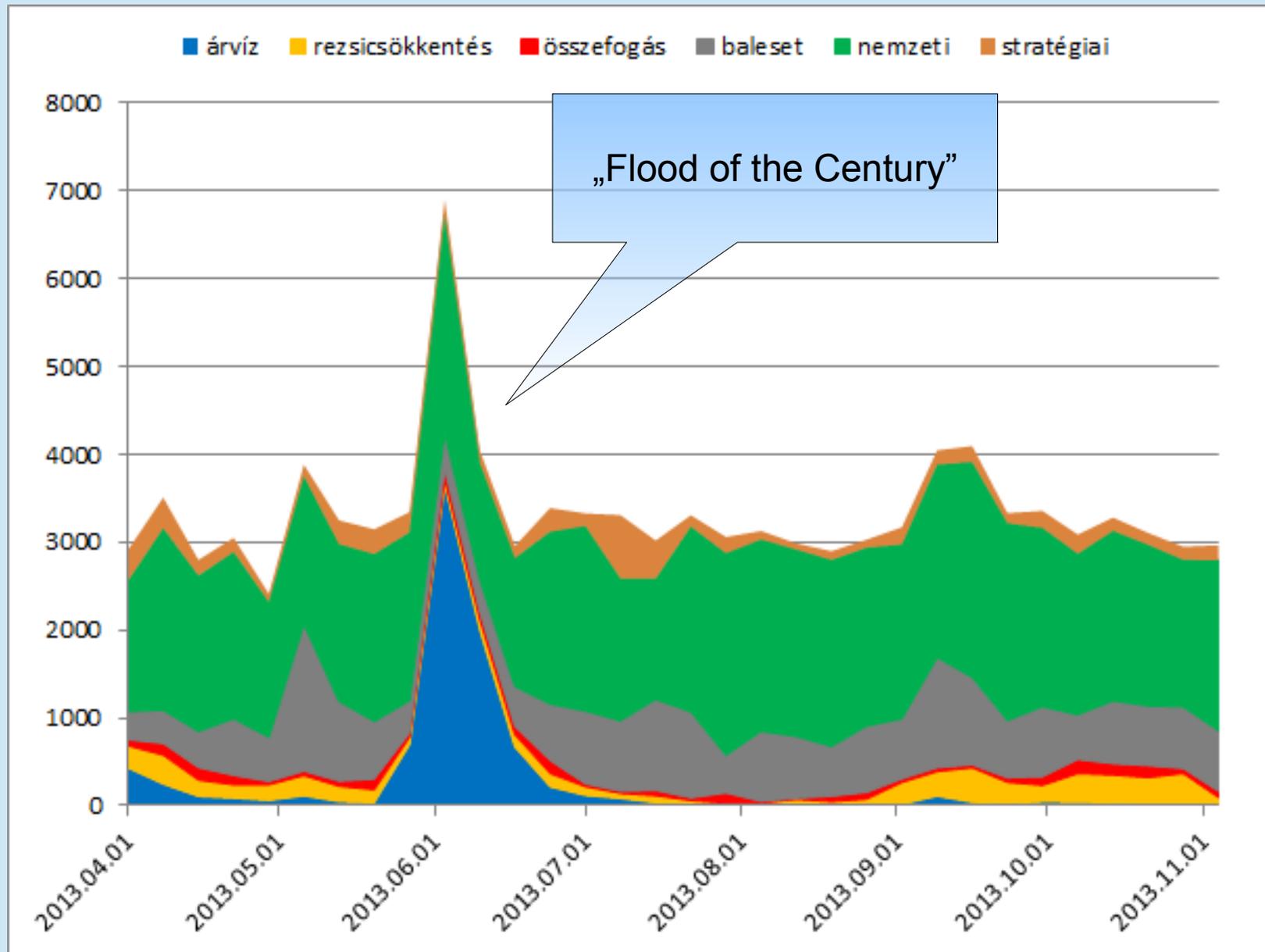
# Rank-Size Distribution of University Web Sites



# Rank-Size Distribution of Hungarian Academy of Sciences Web Sites



# Topics Trends in News



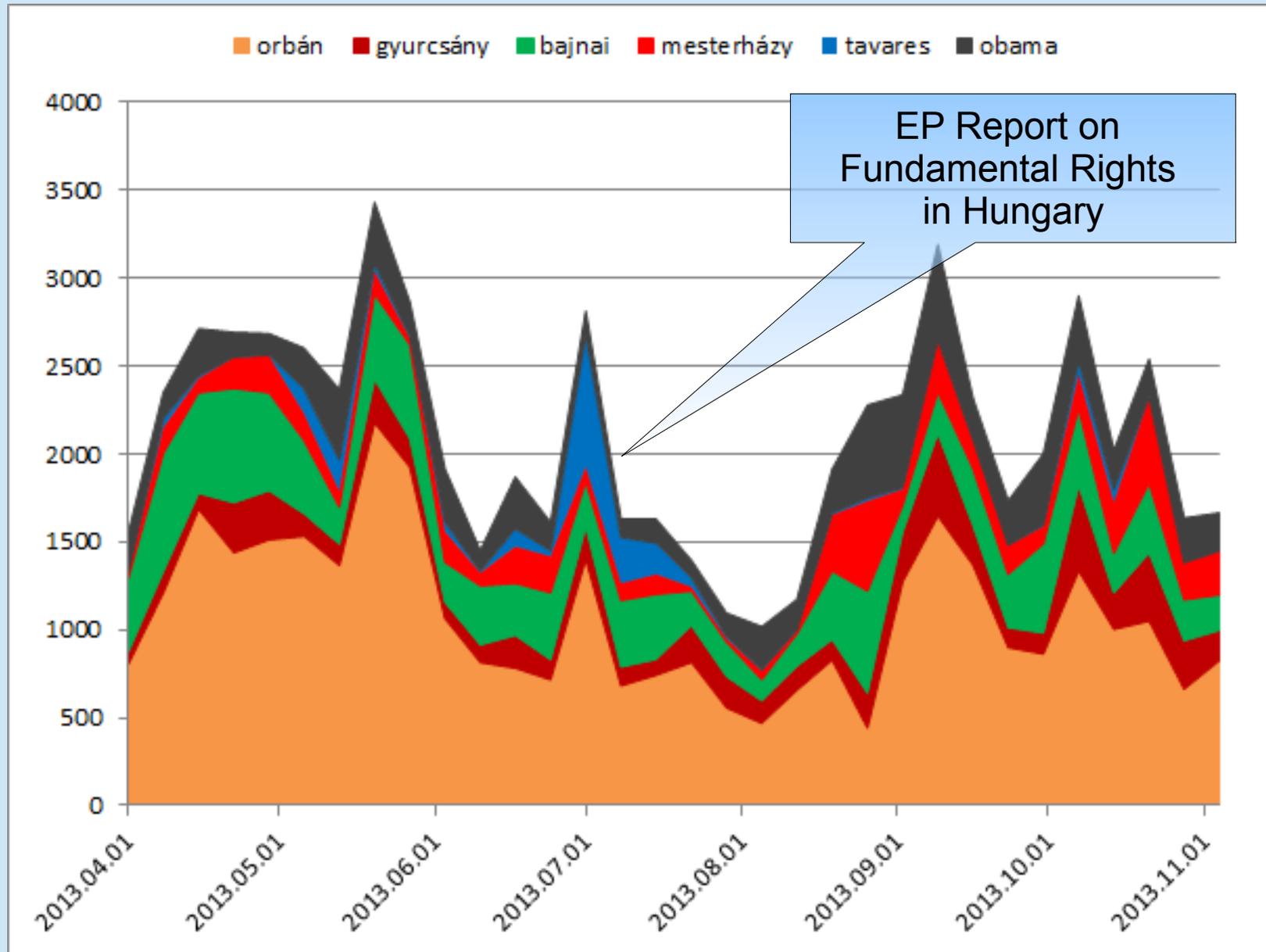
Motivation

Materials and Methods

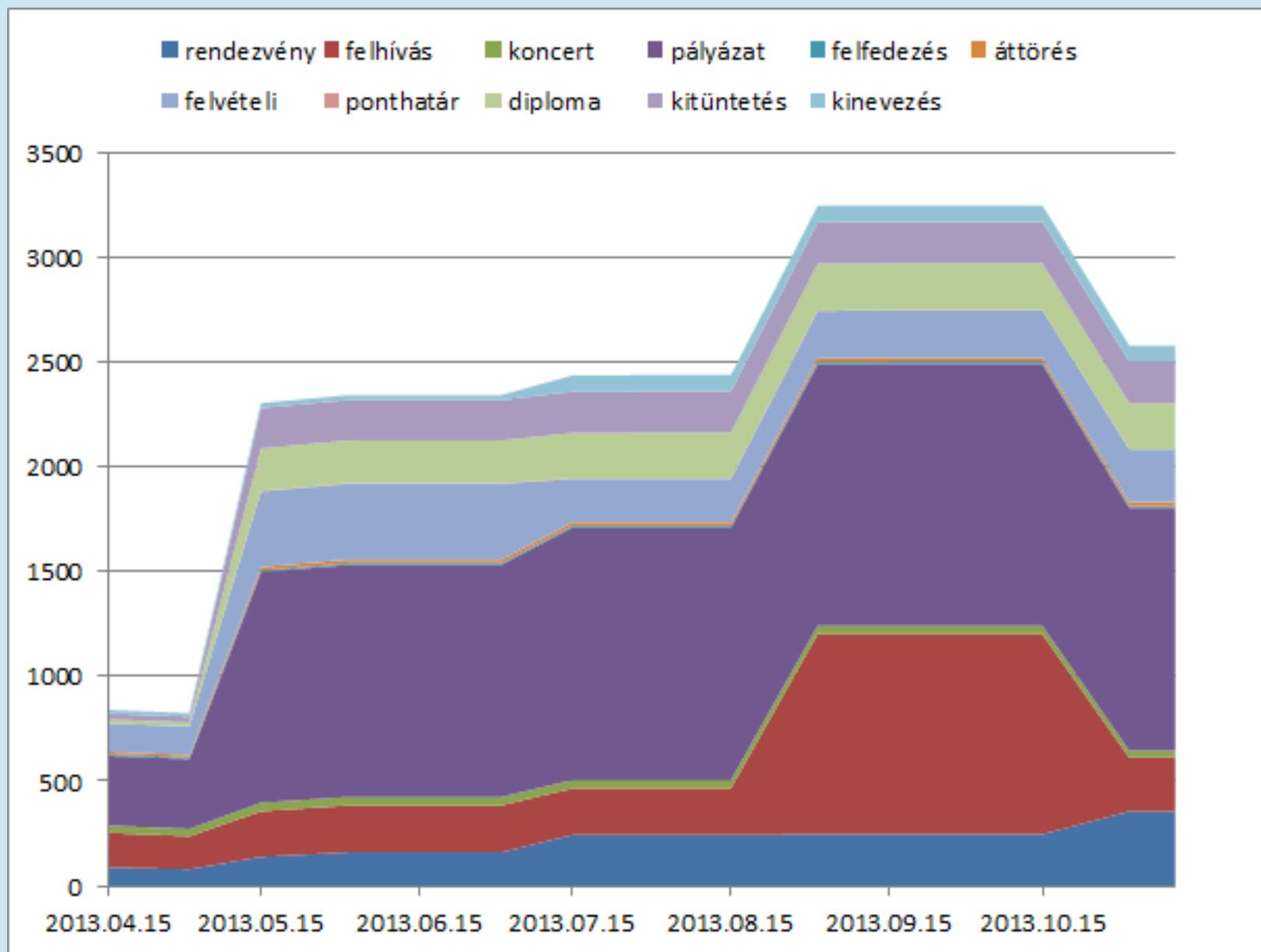


Message

# Topics Trends in News



# Topics Trends in Academia



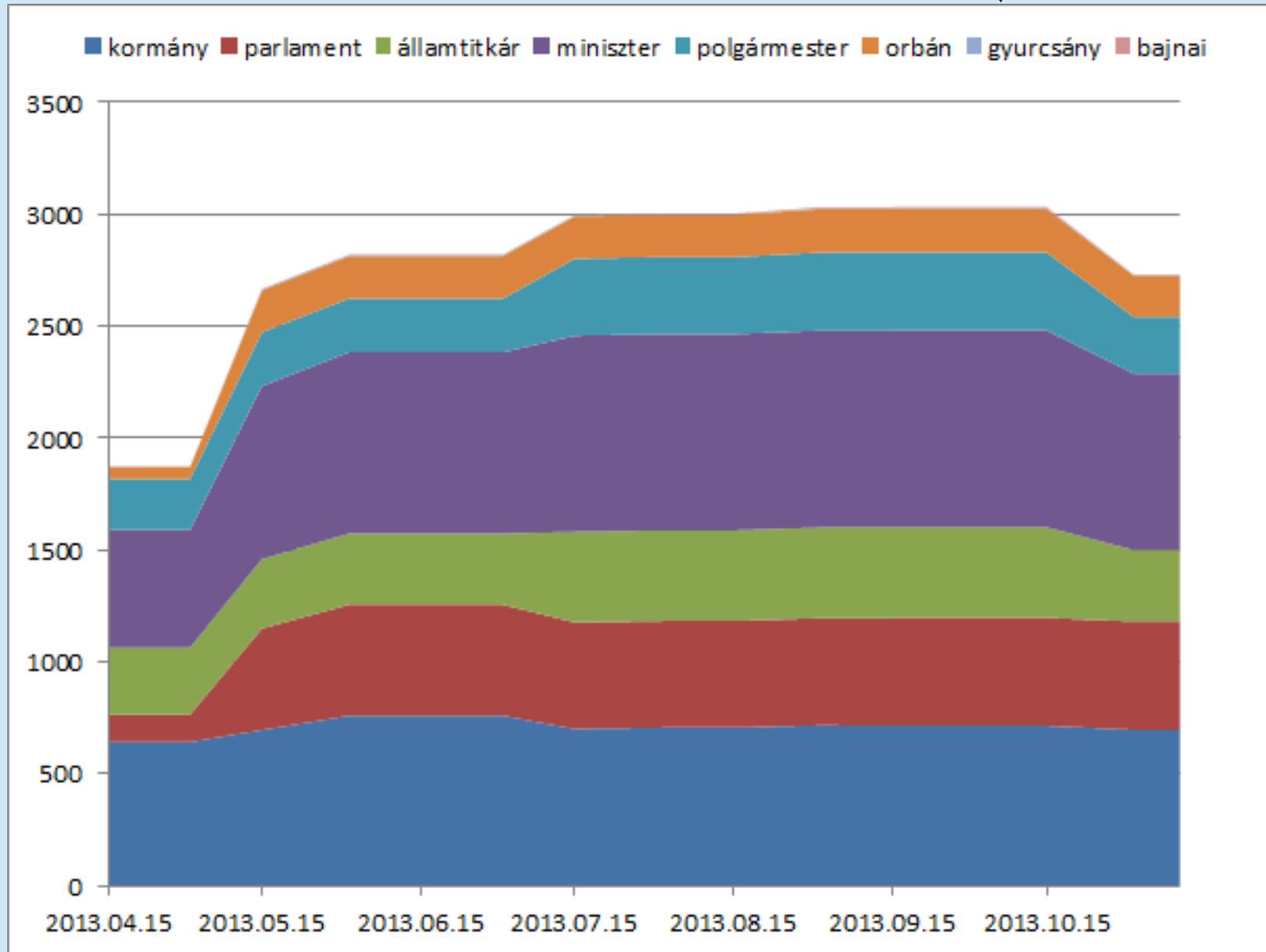
Motivation

Materials and Methods



Message

# Topics Trends in Academia (political)



Motivation

Materials and Methods



Message

# Summary and Conclusions

- A pilot study in Hungary targeting public internet content of
  - Public news portals
  - Academic research institutions
- Tools and methods
  - In continuous development
- "Big Data" is unexpectedly small
  - But has long tails



# Future Works

- More sophisticated longitudinal analysis
  - More data
  - Seasonal corrections
  - Etc.
- Change rate analysis
  - Speed
  - Correlation between (academic) web site activity and scientific output / success



# Acknowledgements

- This work was partially supported by the European Union and the European Social Fund through project (TAMOP-4.2.2.C-11/1/KONV- 2012-0013).



- Building on technology from



- With contributions from



- Zsolt Jurányi, Balázs Bálint, Attila Pálmai

Hungarian Science .org

THE RESEARCH PERFORMANCE OF HUNGARY: INSTITUTIONAL PATTERNS AND COMPETENCES

