

Exploratory Search Missions for TREC Topics

EuroHCIR 2013

Dublin, 1 August 2013

Martin Potthast

Matthias Hagen

Michael Völske

Benno Stein

Bauhaus-Universität Weimar

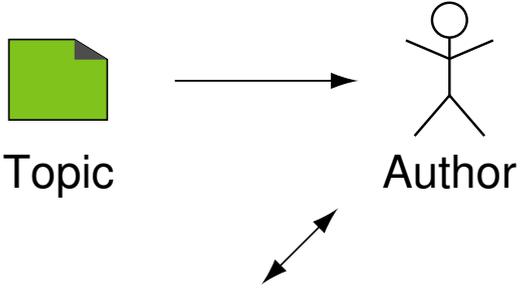
www.webis.de

Background

- ❑ Dataset for studying text reuse / plagiarism from web sources
- ❑ PAN workshop at CLEF
- ❑ Previous datasets: real plagiarism or algorithmically generated
- ❑ This time: crowdsourcing to freelance authors

Background

Corpus Construction



Editor

Barack Obama's Family
plagiarized by John Doe

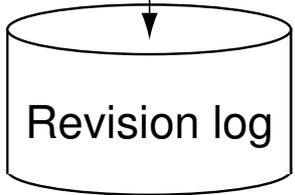
The Family of Barack Obama is an extended clan of African American, Indonesian, and Kenyan heritage. They are best known through the political career of Barack Obama, the current President of the United States. His immediate family is the First Family of the United States. The first First Family of African American descent in the United States youngest to enter the White House since the Kennedys. Obama's young family harks back to days of Camelot.

<http://webis15.medien.uni-weimar.de/chatnoir/clueweb?id=100011709993&tok>

In what follows, we give a detailed overview of Barack Obama's Far light on himself, his immediate and extended family, including maternal relations. Moreover, we give insights into the relations of Michelle Obama's wife, as well as some distant relations of both.

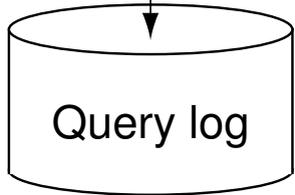
Barack Obama

Barack Hussein Obama II is the 44th and current President of the United States, the first African American to hold the office. Obama was the junior United States Senator from Illinois from 2005 until he resigned following his election as President.



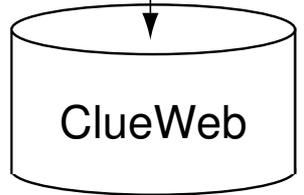
ChatNoir SE

A black silhouette of a cat sitting and looking to the right, with several short lines representing whiskers.



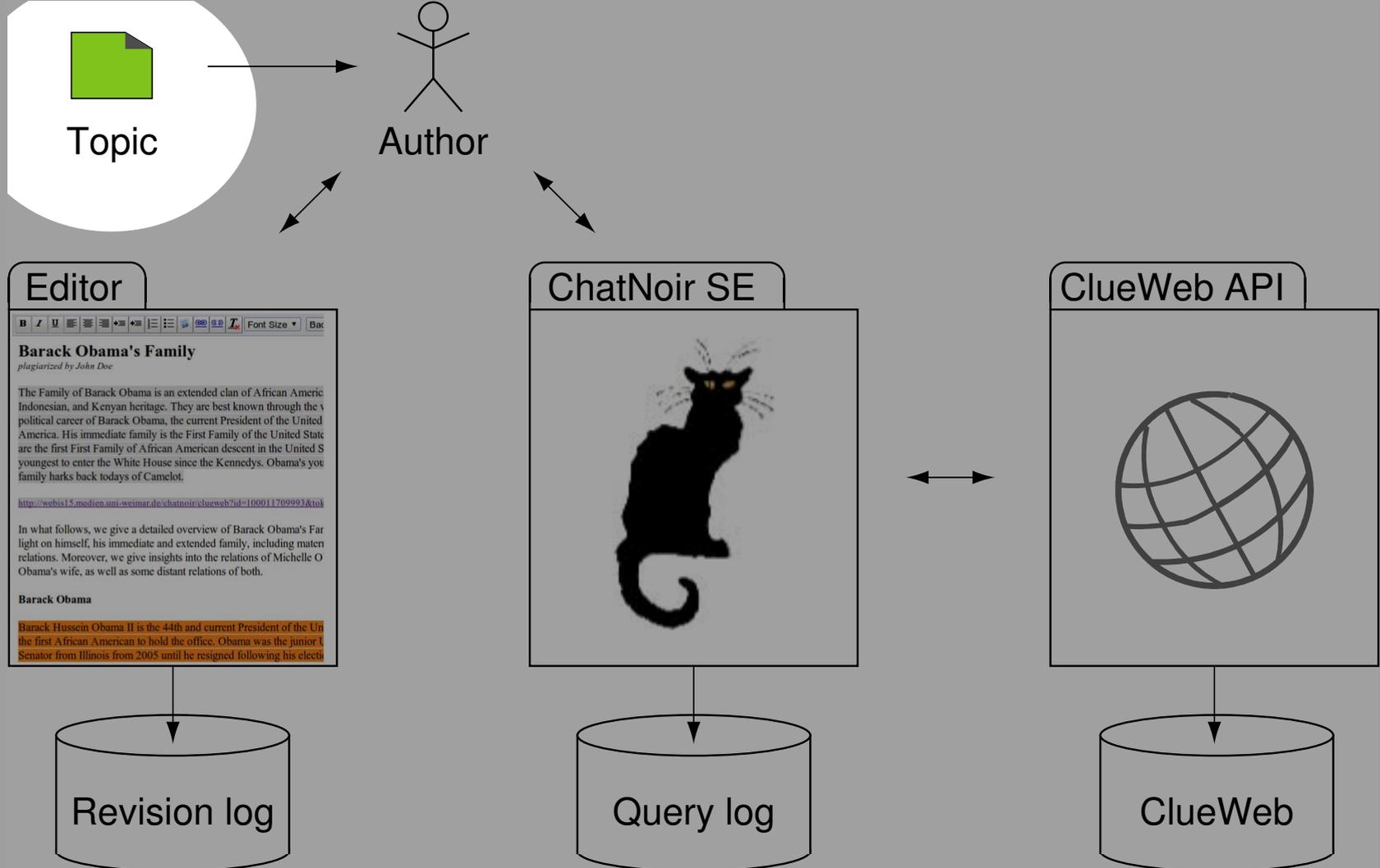
ClueWeb API

A wireframe globe icon representing a web corpus.



Background

Topics



Background

Topics

Example topic:

Obama's family.

Write about President Barack Obama's family history, including genealogy, national origins, places and dates of birth, etc. Where did Barack Obama's parents and grandparents come from? Also include a brief biography of Obama's mother.

Original topic 001 of the TREC Web Track 2009:

Query. obama family tree

Description. Find information on President Barack Obama's family history, including genealogy, national origins, places and dates of birth, etc.

Sub-topic 1. Find the TIME magazine photo essay "Barack Obama's Family Tree."

Sub-topic 2. Where did Barack Obama's parents and grandparents come from?

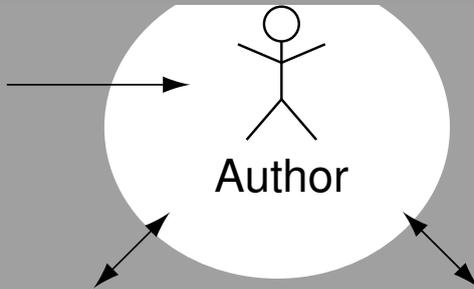
Sub-topic 3. Find biographical information on Barack Obama's mother.

Background

Authors



Topic



Author

Editor

Barack Obama's Family
plagiarized by John Doe

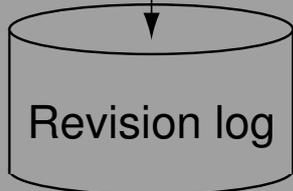
The Family of Barack Obama is an extended clan of African American, Indonesian, and Kenyan heritage. They are best known through the political career of Barack Obama, the current President of the United States. His immediate family is the First Family of the United States. The first First Family of African American descent in the United States youngest to enter the White House since the Kennedys. Obama's young family harks back to days of Camelot.

<http://webis15.medien.uni-weimar.de/chatnoir/clueweb?id=100011709993&to>

In what follows, we give a detailed overview of Barack Obama's Far light on himself, his immediate and extended family, including maternal relations. Moreover, we give insights into the relations of Michelle Obama's wife, as well as some distant relations of both.

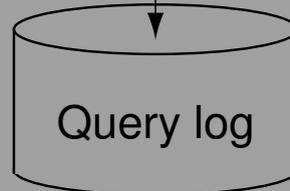
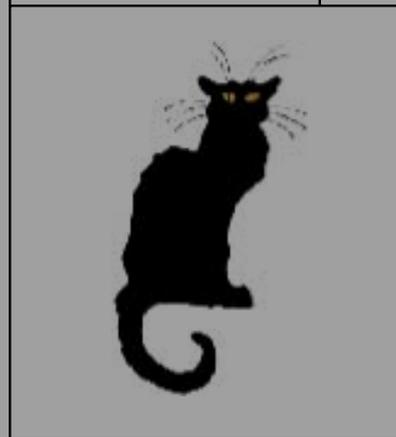
Barack Obama

Barack Hussein Obama II is the 44th and current President of the United States, the first African American to hold the office. Obama was the junior U.S. Senator from Illinois from 2005 until he resigned following his election as President.



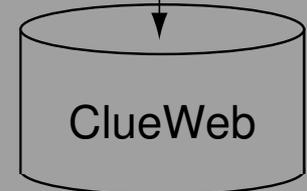
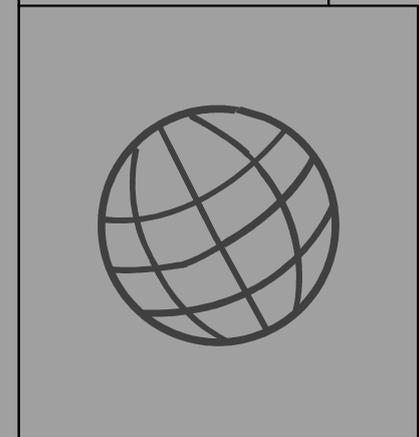
Revision log

ChatNoir SE



Query log

ClueWeb API

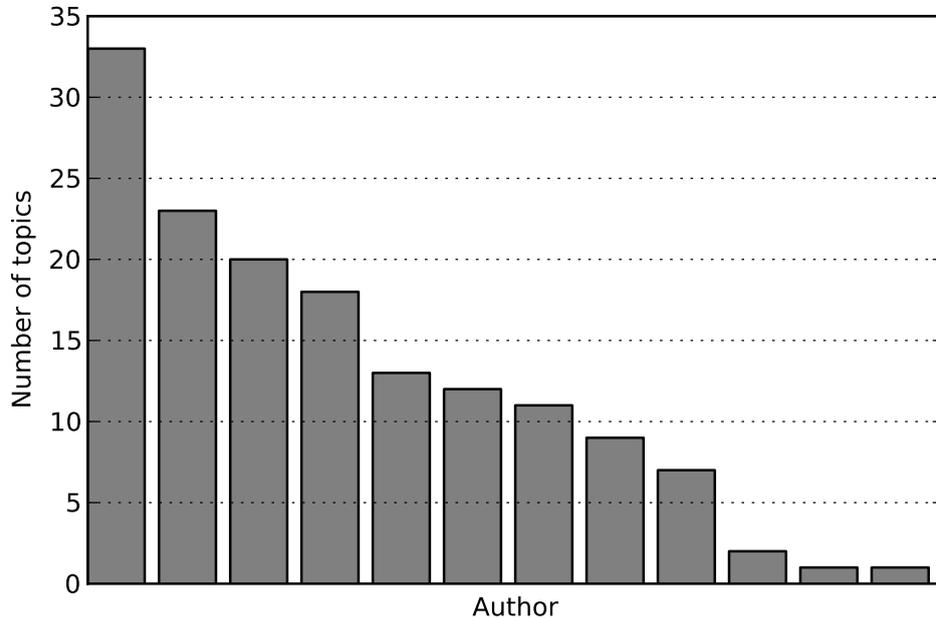


ClueWeb



Background

Authors

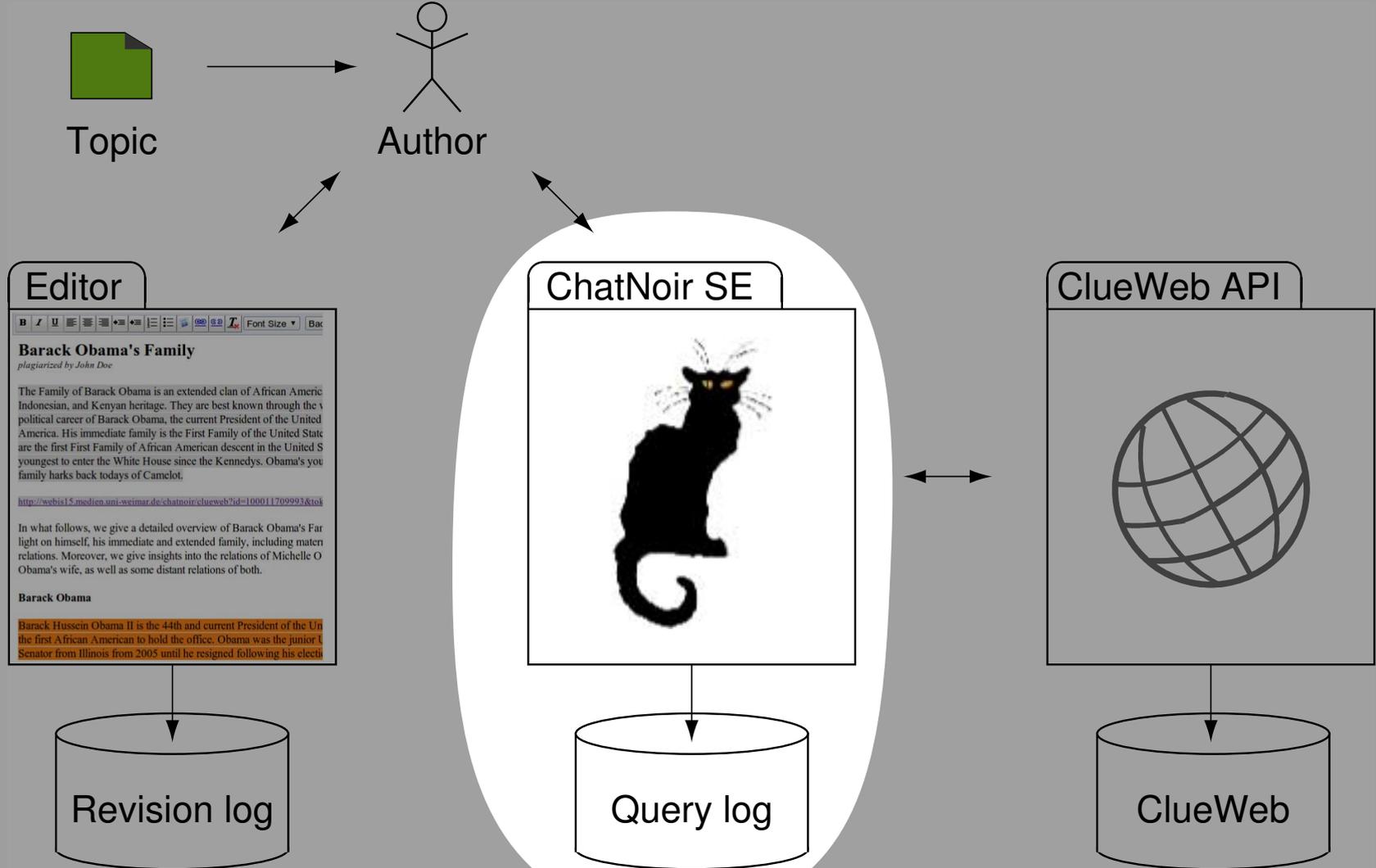


Author Demographics	
Age (Median)	37
Years Writing (Median)	8
<i>Academic degree</i>	
Postgrad	33%
Undergrad	25%
<i>English</i>	
Native	67%
Second Language	33%

- ❑ 12 Professional writers hired on oDesk [www.odesk.com]
- ❑ Fluent English speakers with several years writing experience

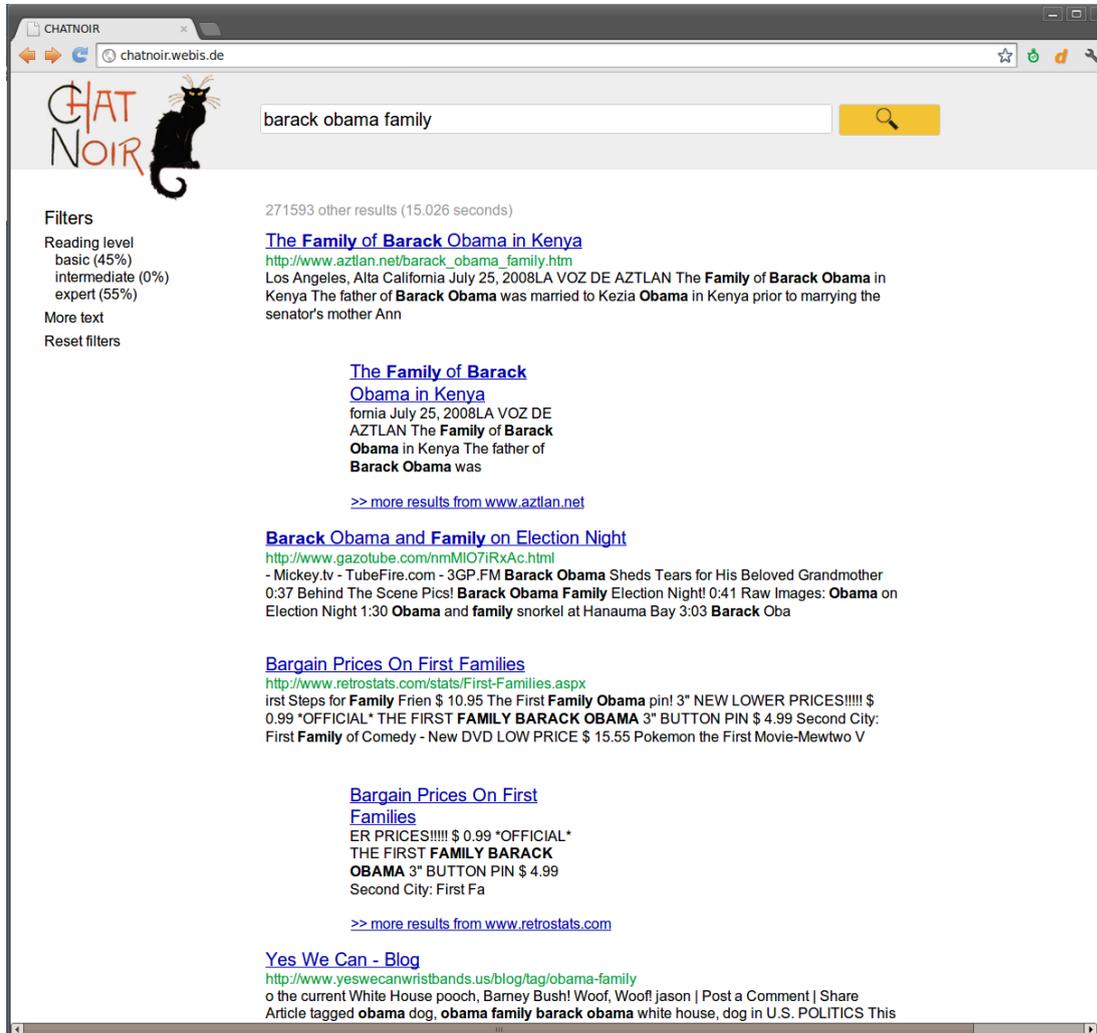
Background

Search engine and query log



Background

Clueweb search engine



The screenshot shows a web browser window with the address bar at 'chatnoir.webis.de'. The search bar contains the text 'barack obama family'. The search results are displayed on the right side of the page, with a search filter on the left. The search results include several links and snippets of text related to Barack Obama's family, including his time in Kenya and election night events.

CHATNOIR

chatnoir.webis.de

barack obama family

271593 other results (15.026 seconds)

[The Family of Barack Obama in Kenya](#)
http://www.aztlan.net/barack_obama_family.htm
Los Angeles, Alta California July 25, 2008LA VOZ DE AZTLAN The **Family of Barack Obama** in Kenya The father of **Barack Obama** was married to Kezia **Obama** in Kenya prior to marrying the senator's mother Ann

[The Family of Barack Obama in Kenya](#)
formia July 25, 2008LA VOZ DE AZTLAN The **Family of Barack Obama** in Kenya The father of **Barack Obama** was

[Barack Obama and Family on Election Night](#)
<http://www.gazolube.com/nmMIO7iRxAc.html>
- Mickey.tv - TubeFire.com - 3GP FM **Barack Obama** Sheds Tears for His Beloved Grandmother 0:37 Behind The Scene Pics! **Barack Obama Family** Election Night! 0:41 Raw Images: **Obama** on Election Night 1:30 **Obama** and **family** snorkel at Hanauma Bay 3:03 **Barack Oba**

[Bargain Prices On First Families](#)
<http://www.retrostats.com/stats/First-Families.aspx>
irst Steps for **Family** Frien \$ 10.95 The First **Family Obama** pin! 3" NEW LOWER PRICES!!!! \$ 0.99 *OFFICIAL* THE FIRST **FAMILY BARACK OBAMA** 3" BUTTON PIN \$ 4.99 Second City: First **Family** of Comedy - New DVD LOW PRICE \$ 15.55 Pokemon the First Movie-Mewtwo V

[Bargain Prices On First Families](#)
ER PRICES!!!! \$ 0.99 *OFFICIAL*
THE FIRST **FAMILY BARACK OBAMA** 3" BUTTON PIN \$ 4.99
Second City: First Fa

[Yes We Can - Blog](#)
<http://www.yeswecanwristbands.us/blog/tag/obama-family>
o the current White House pooch, Barney Bush! Woof, Woof! jason | Post a Comment | Share
Article tagged **obama** dog, **obama family** **barack obama** white house, dog in U.S. POLITICS This

[chatnoir.webis.de]

Background

Clueweb search engine

The screenshot shows a web browser window with the URL 'chatnoir.webis.de'. The search bar contains the text 'barack obama family'. Below the search bar, there is a search button with a magnifying glass icon. The search results are displayed in a list format. The first result is titled 'The Family of Barack Obama in Kenya' and includes a URL 'http://www.aztlan.net/barack_obama_family'. The second result is titled 'The Family of Barack Obama in Kenya' and includes a URL 'http://www.gazolube.com/nmMIO7iRxAc.htm'. The third result is titled 'Bargain Prices On First Families' and includes a URL 'http://www.retrostats.com/stats/First-Families'. The fourth result is titled 'Bargain Prices On First Families' and includes a URL 'http://www.retrostats.com'. The fifth result is titled 'Yes We Can - Blog' and includes a URL 'http://www.yeswecanwristbands.us/blog/tag/obama-family'. The interface also includes a 'Filters' section on the left with options for 'Reading level' (basic, intermediate, expert) and 'More text'.

- ❑ Source retrieval for essay writing
- ❑ Indexes ClueWeb09 [lemurproject.org/clueweb09]
- ❑ Fine-grained interaction logs

[chatnoir.webis.de]

Query logs

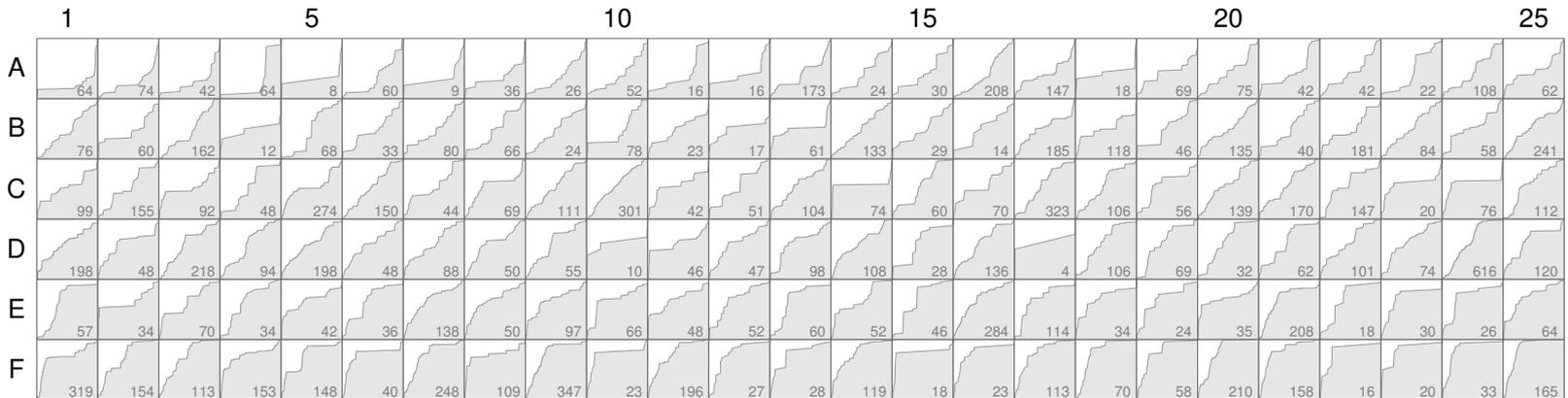
Basic stats

Corpus Characteristic	Distribution				Σ
	min	avg	max	stdev	
Queries					13 651
Queries / Topic	4	91.0	616	83.1	
Clicks					16 739
Clicks / Topic	12	111.6	443	80.3	
Clicks / Query	0	0.8	76	2.2	
Sessions					931
Sessions / Topic	1	12.3	149	18.9	
Days					201
Days / Topic	1	4.9	17	2.7	
Hours					2068
Hours / Writer	3	129.3	679	167.3	
Hours / Topic	3	7.5	10	2.5	

Findings so far

Findings so far

Distribution of Queries Over Time

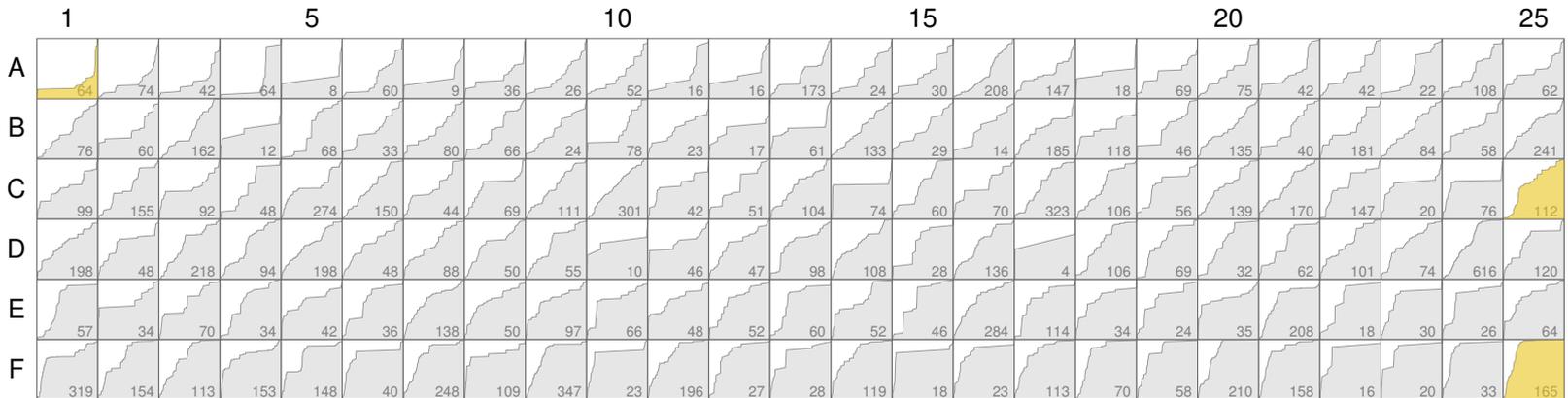


Distribution of queries over time.

- ❑ fraction of posed queries (y-axis) over elapsed time (x-axis) between the first query until essay completion
- ❑ each cell represents one of 150 essays
- ❑ the numbers denote the total amount of posed queries
- ❑ the cells are sorted by area under the curve

Findings so far

Distribution of Queries Over Time

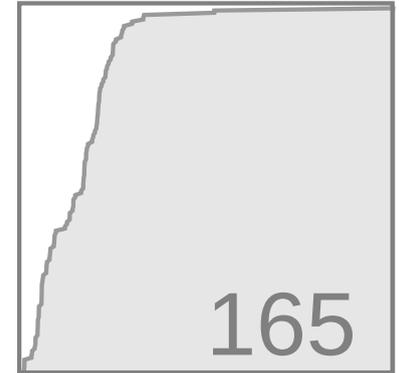
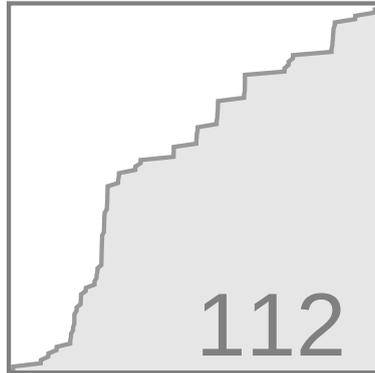
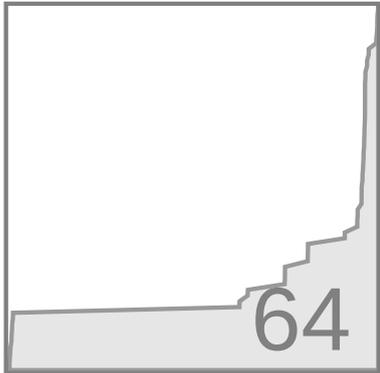


Distribution of queries over time.

- ❑ fraction of posed queries (y-axis) over elapsed time (x-axis) between the first query until essay completion
- ❑ each cell represents one of 150 essays
- ❑ the numbers denote the total amount of posed queries
- ❑ the cells are sorted by area under the curve

Findings so far

Distribution of Queries Over Time



Distribution of queries over time.

- ❑ fraction of posed queries (y-axis) over elapsed time (x-axis) between the first query until essay completion
- ❑ each cell represents one of 150 essays
- ❑ the numbers denote the total amount of posed queries
- ❑ the cells are sorted by area under the curve

Summary

1. Corpus of essay writing, focus on text reuse
2. 150 exploratory search missions for TREC web track topics
3. First appraisal indicates wide spectrum of search behavior
4. Data will be publicly available: <http://webis.de/research/corpora>

Summary

1. Corpus of essay writing, focus on text reuse
2. 150 exploratory search missions for TREC web track topics
3. First appraisal indicates wide spectrum of search behavior
4. Data will be publicly available: <http://webis.de/research/corpora>

Thank you for your attention!