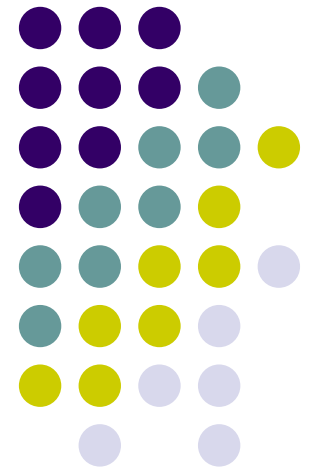


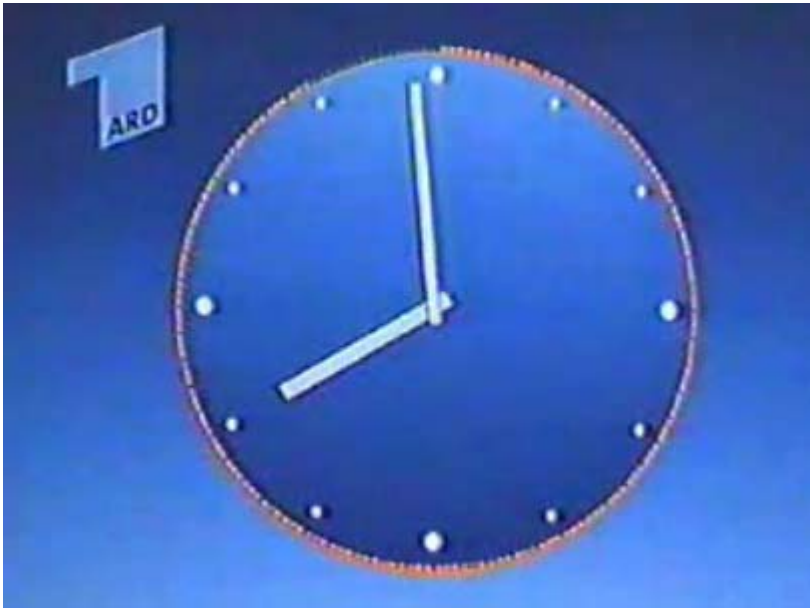
ConceptDoppler: A Weather Tracker for Internet Censorship

Daniel Zinn

Joint work with Jedidiah R. Crandall, Michael Byrd,
Earl Barr, and Rich East



Censorship is Not New

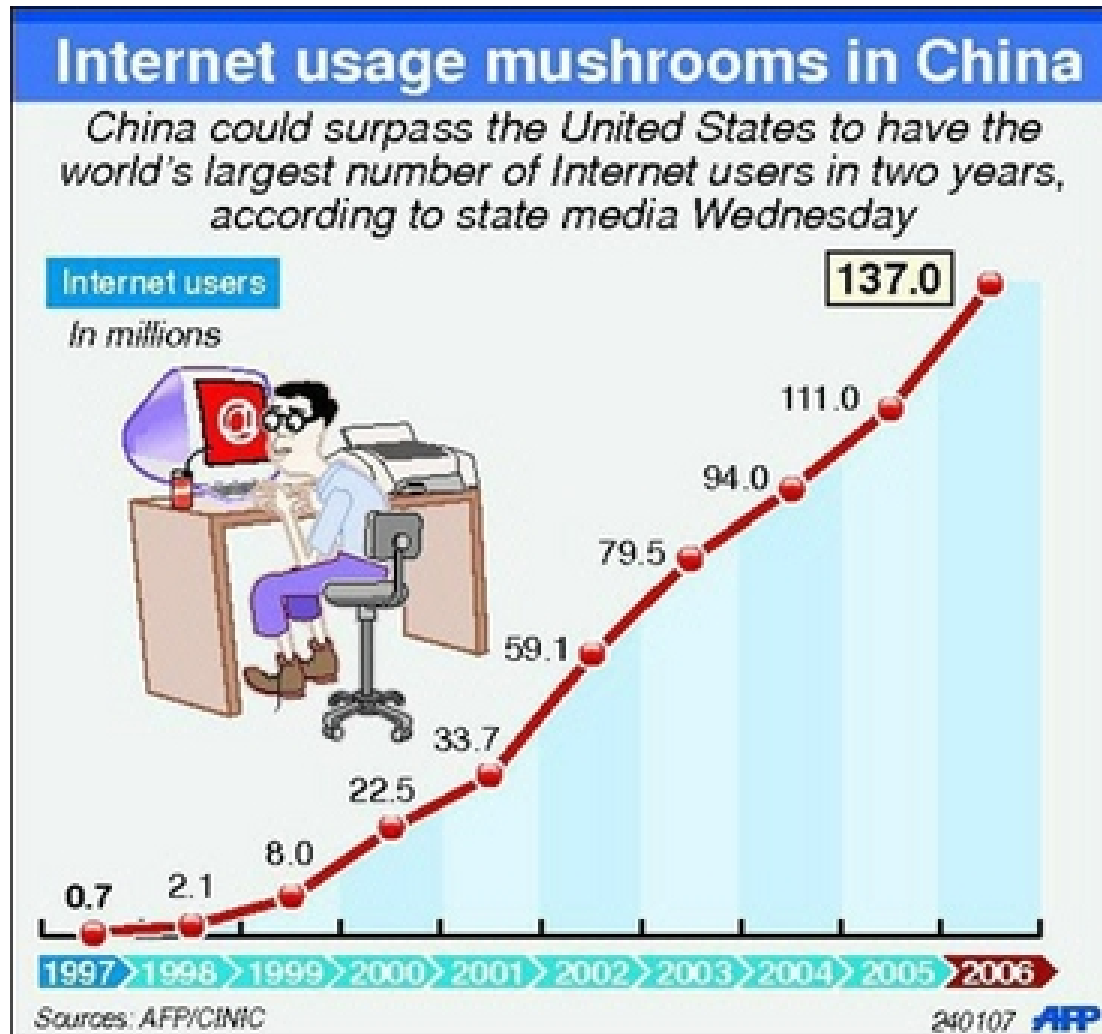
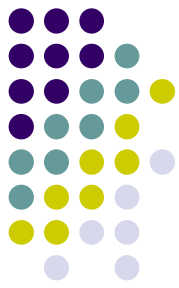


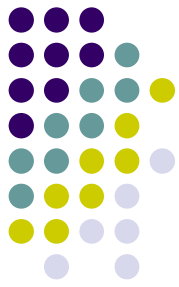
Tagesschau
Western Germany



Aktuelle Kamera
Eastern Germany

China's Internet Usage will Probably Surpass the US Soon





Internet Censorship in China

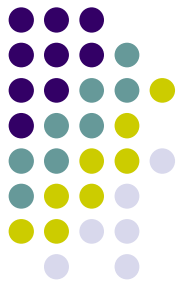
Called the “Great Firewall of China,” or “Golden Shield”

- IP address blocking
- DNS redirection
- Legal restrictions
- etc...
- Keyword filtering
 - Blog servers, chat, HTTP traffic



All probing was performed from outside of China

Why is Keyword Filtering Interesting?



- Chinese government claims to be targeting pornography and sedition
- The keywords provide insights into what material the government is targeting with censorship, *e.g.*
 - 专政机关 --- Dictatorship organs
 - 希特勒 (Hitler), and 我的奋斗 (Mein Kampf)
 - 多维尔 --- Deauville, a town in France



Outline

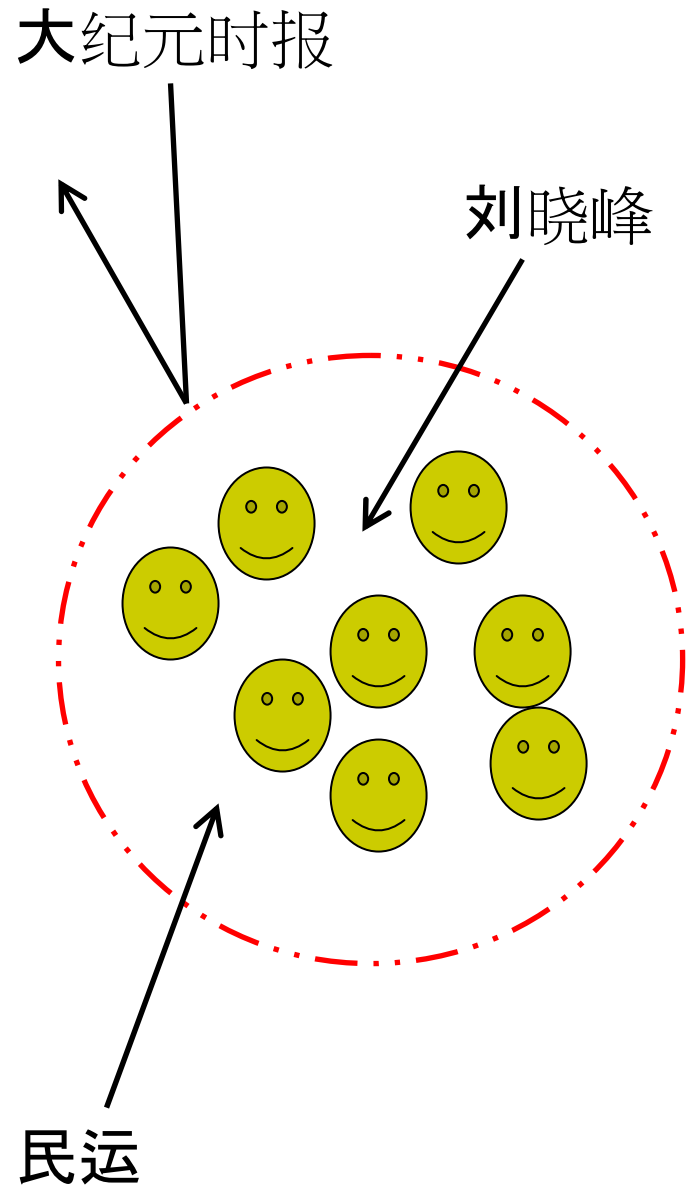
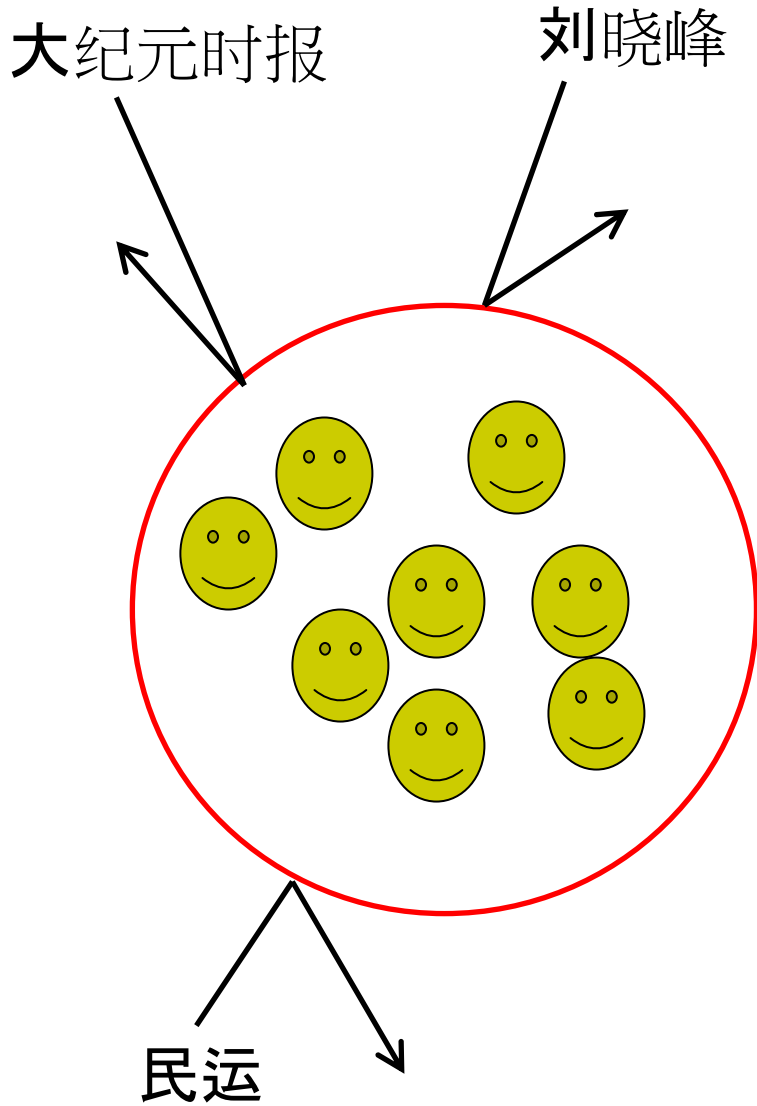
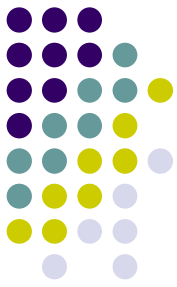
- Firewall or Something Else?
 - Where are filtering routers?
 - Who is doing filtering?
 - How reliable is filtering?
- Blocked Words
 - Which words to select?
 - Which words are blocked?
- Imprecise Filtering
 - What implications does keyword filtering have?



Outline

- Firewall or Something Else?
 - Where are filtering routers?
 - Who is doing filtering?
 - How reliable is filtering?
- Blocked Words
 - Which words to select?
 - Which words are blocked?
- Imprecise Filtering
 - What implications does keyword filtering have?

Firewall?

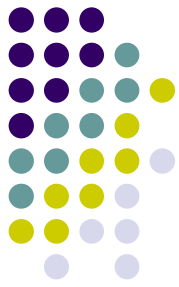




Where Are Filtering Routers

Different opinions about where censorship occurs:

- In three big centers in Beijing, Guangzhou, and Shanghai
- At the border
- Throughout the country's backbone
- At a local level
- An amalgam of the above

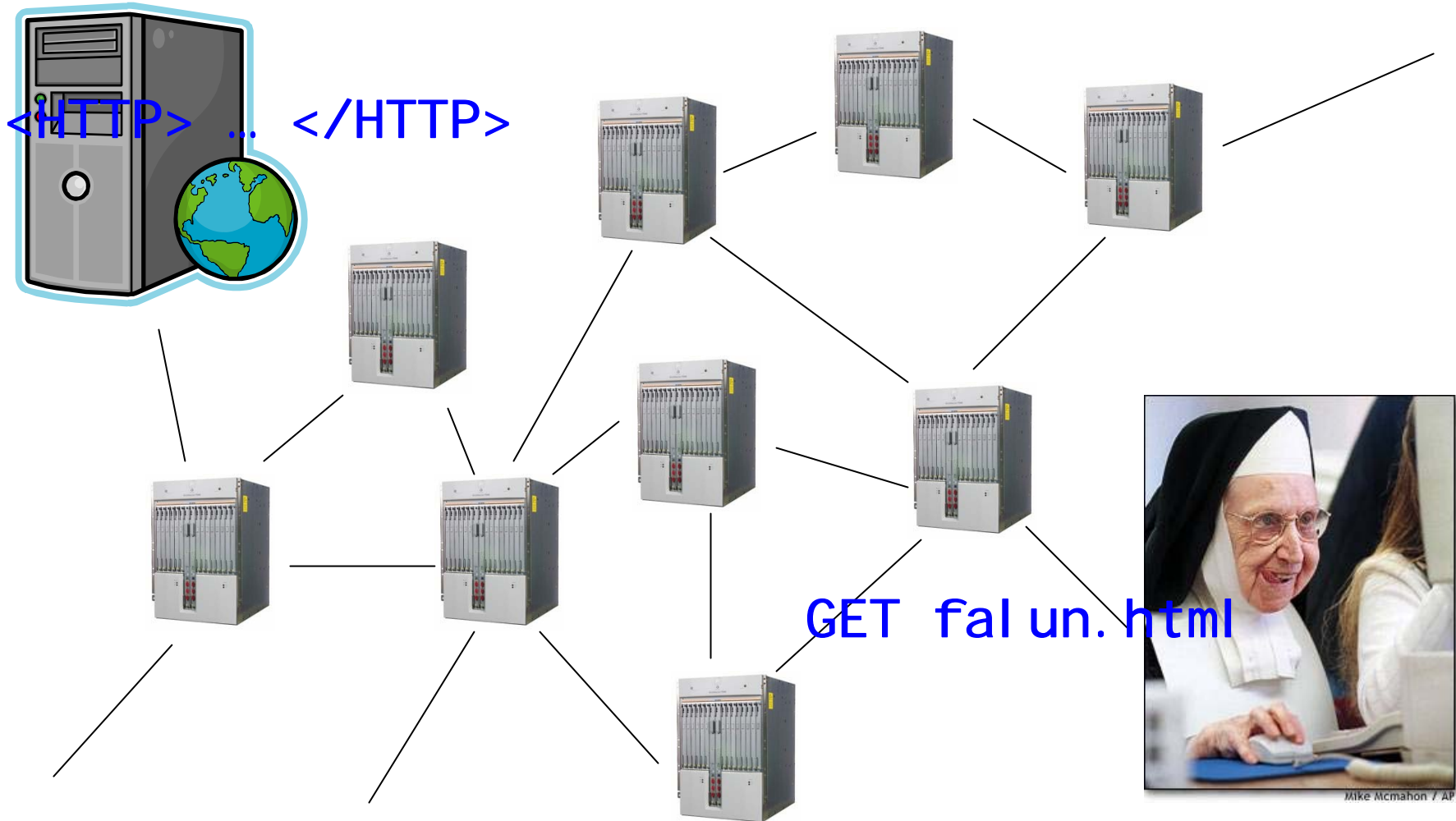
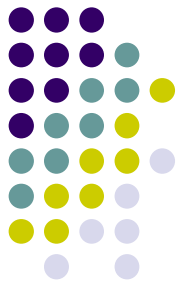


Filtering With Forged RSTs

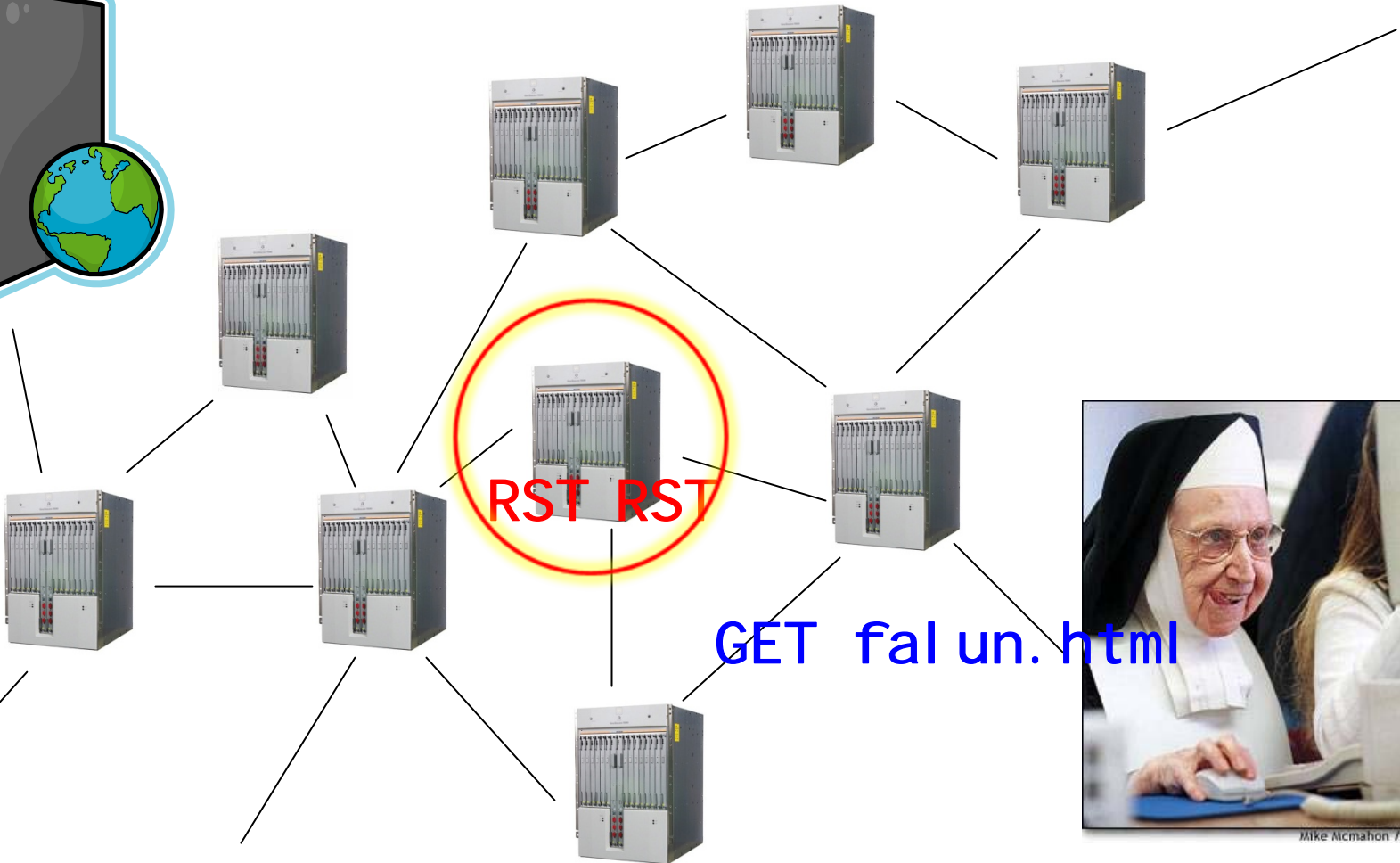
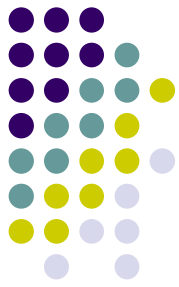
- Clayton et al., 2006.
- Comcast also uses forged RSTs

Example

Dissident Nuns on the Net

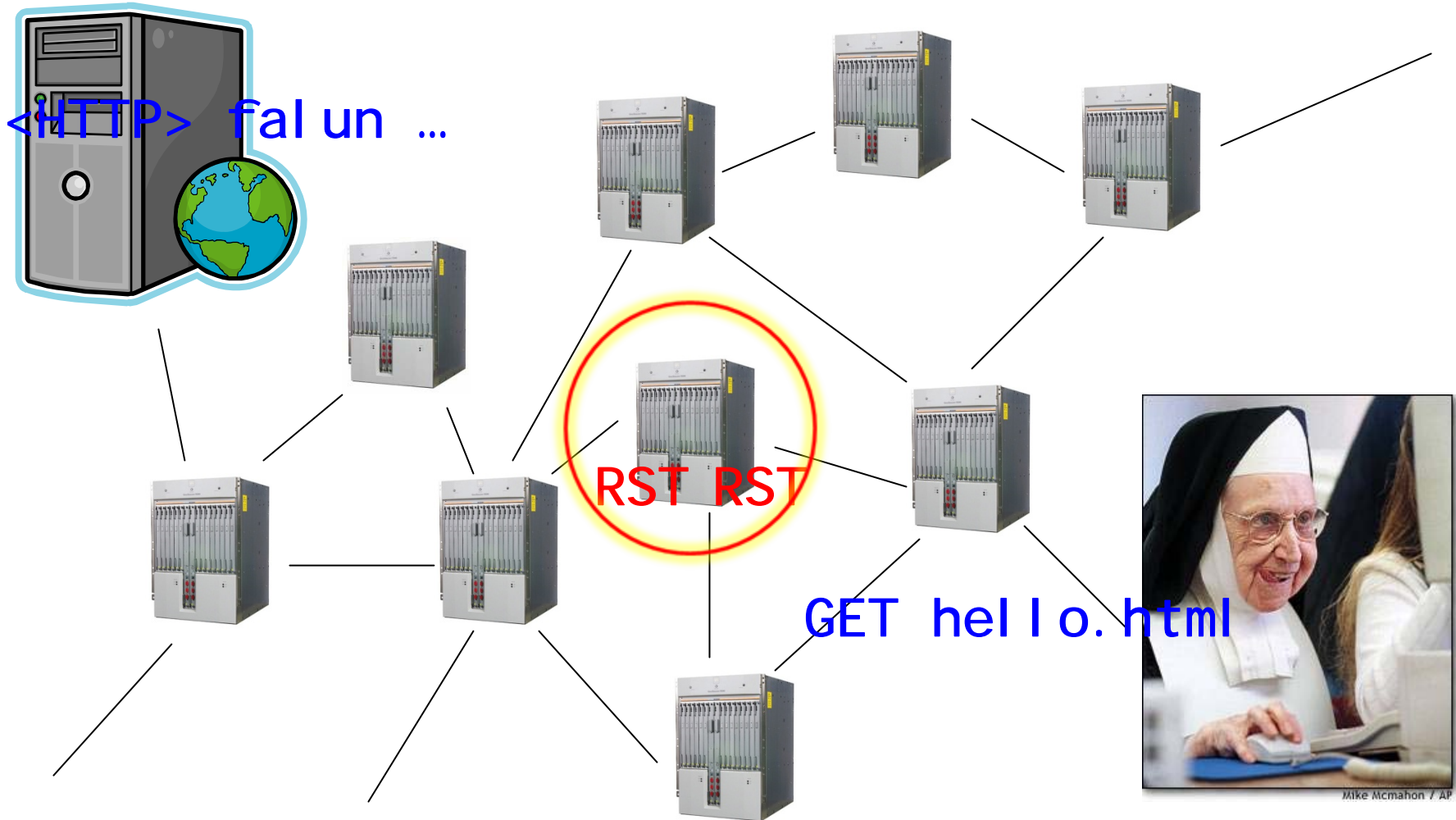
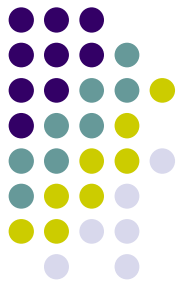


Censorship of HTML GET Requests

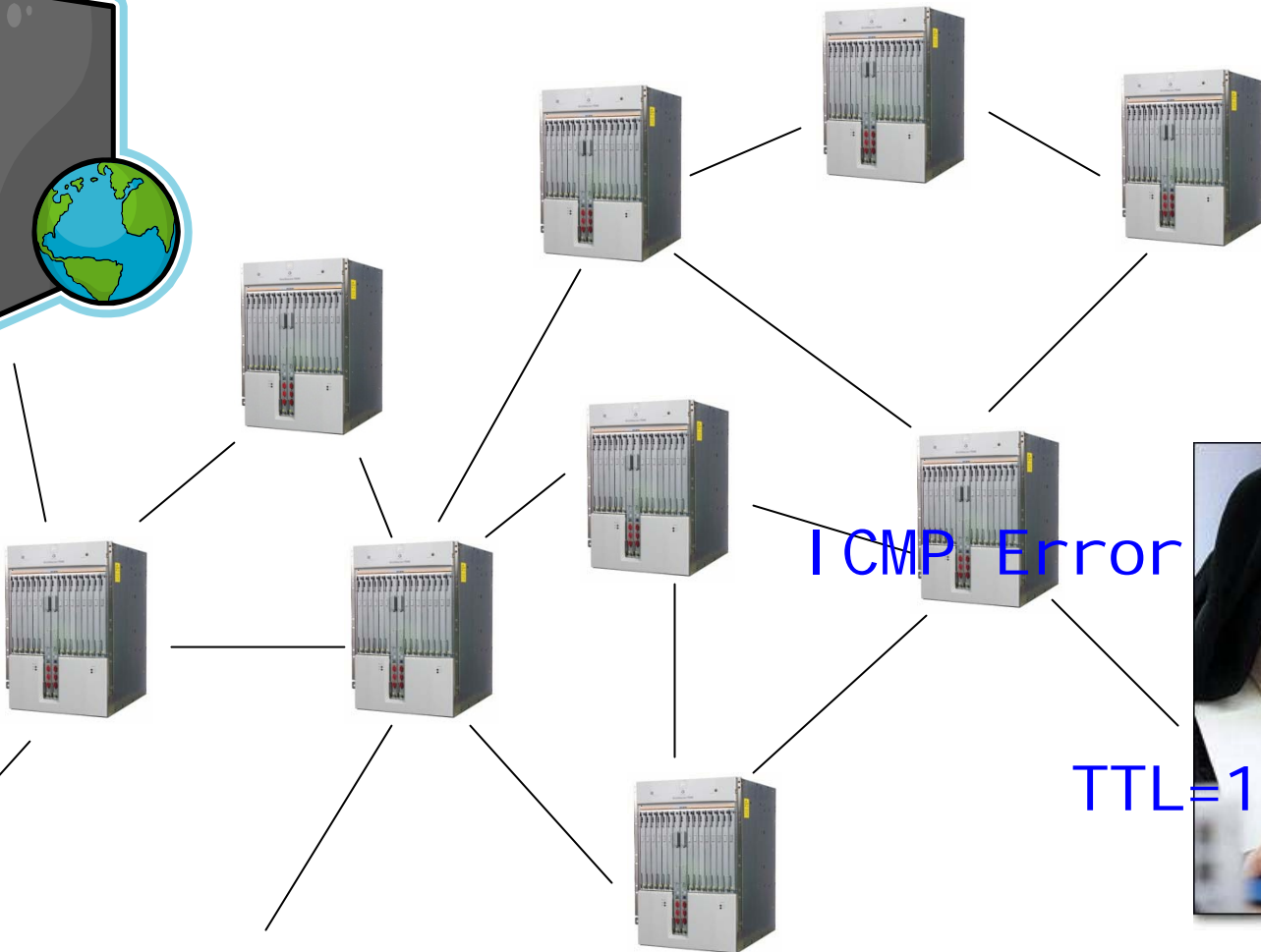
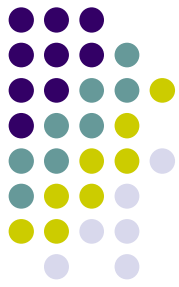


Mike McMahon / AP

Censorship of HTML Responses



Locating Filtering Routers



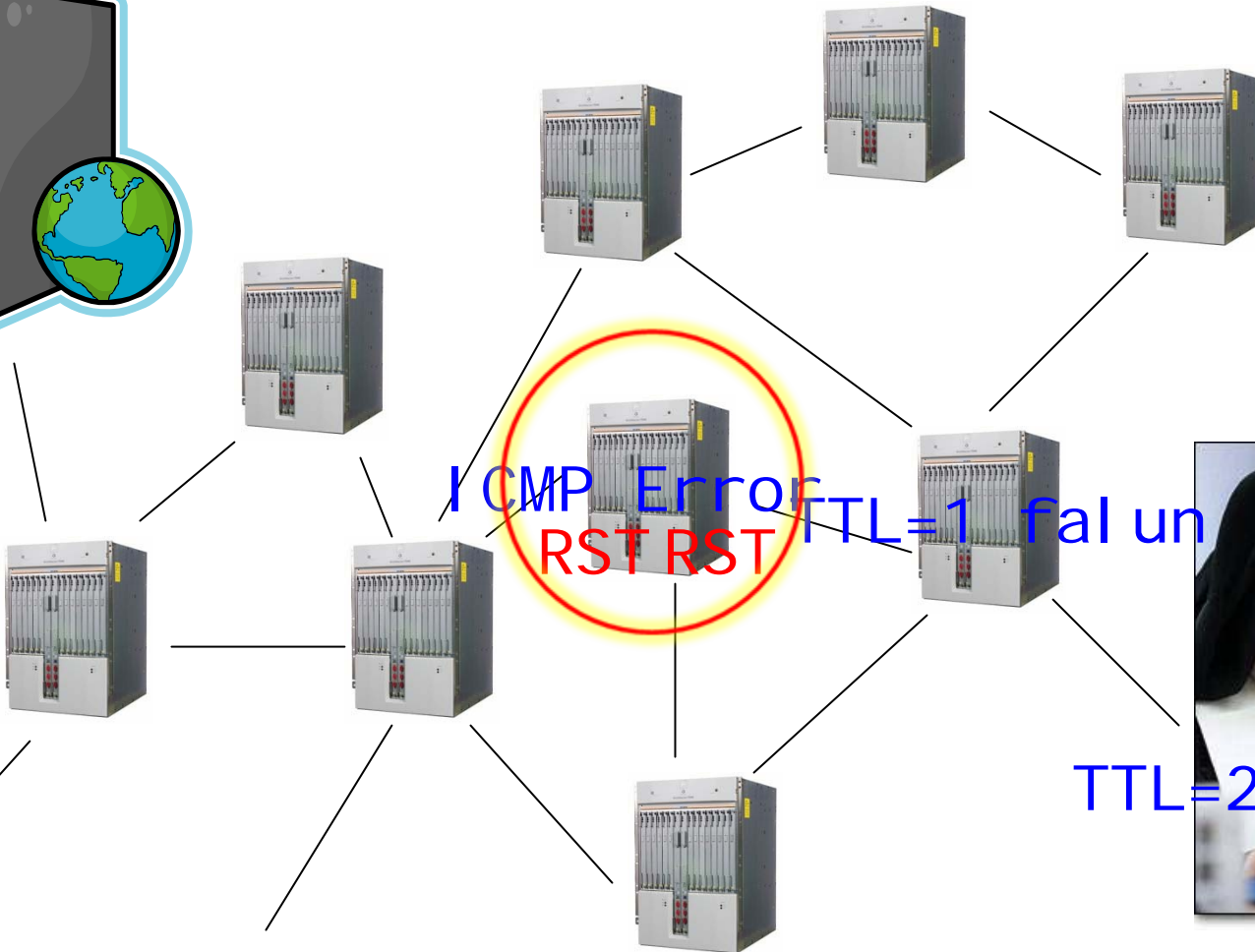
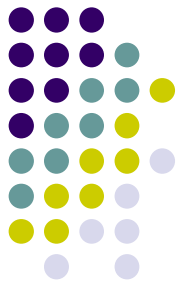
ICMP Error

TTL=1 fail un



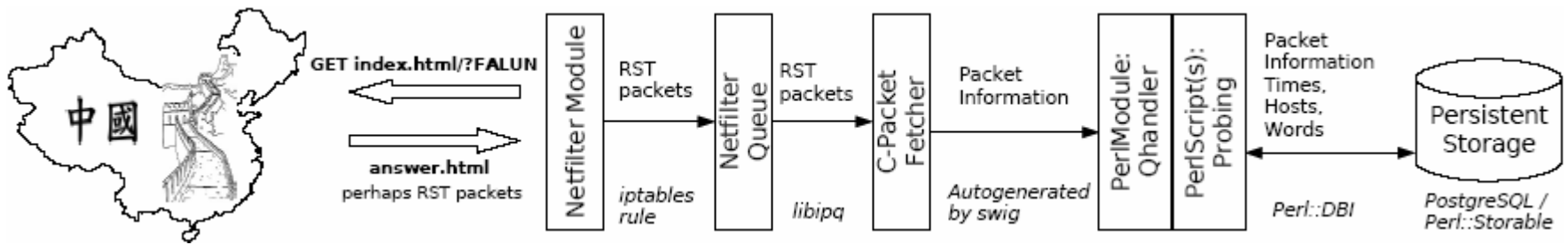
Mike McMahon / AP

Locating Filtering Routers



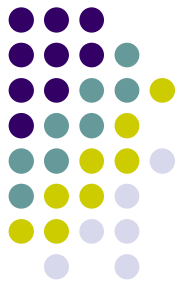
Mike McMahon / AP

ConceptDoppler Framework



- Netfilter (iptables) to capture packets
- Queue module to handle packets over to user-space
- Own TCP stack implementation
- Scapy for constructing custom packets
- Storing packets in PostgreSQL database
- Scapy stored procedures in DB

Experimental Setup

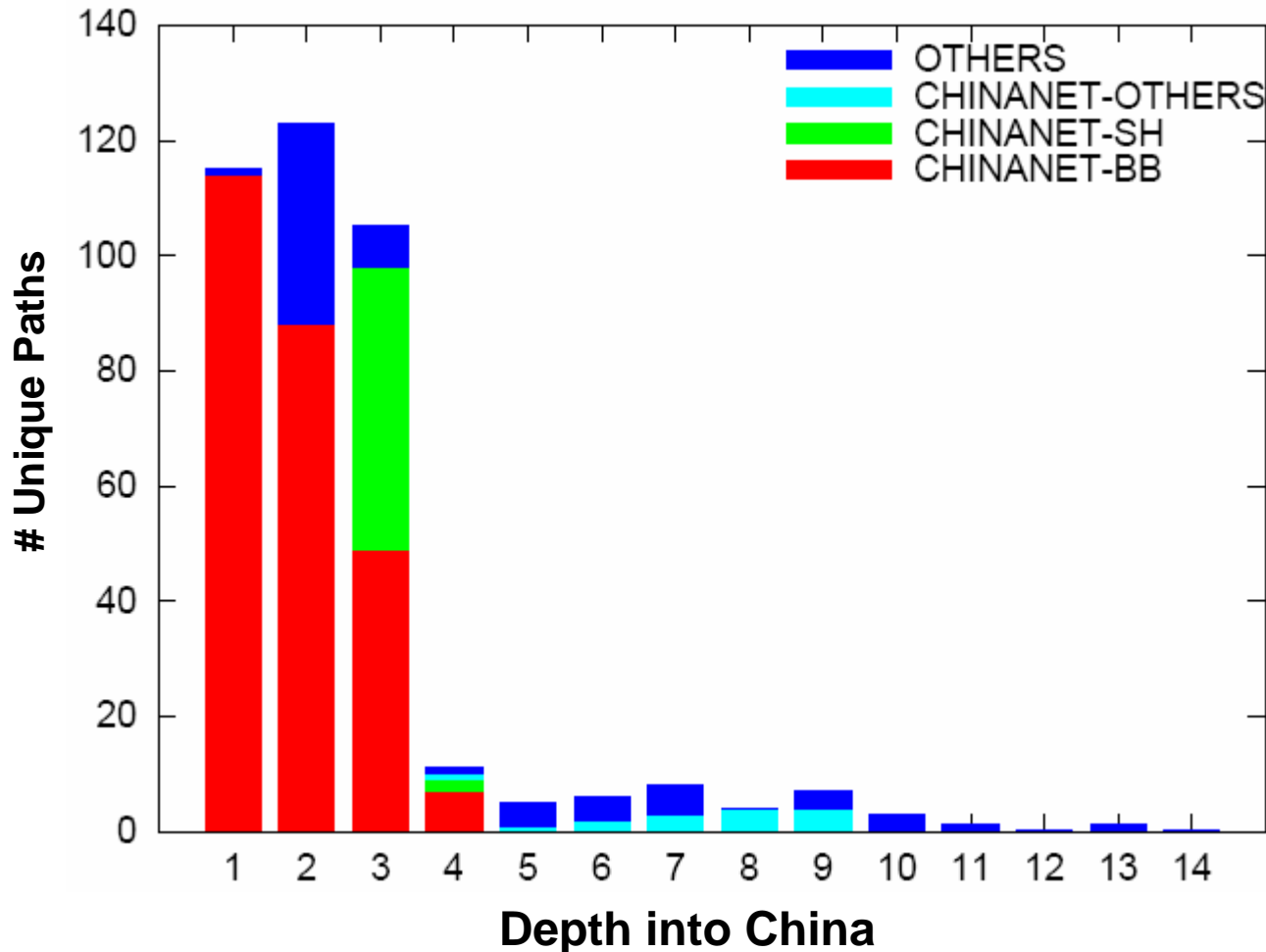


- Google “site:.cn” to find random destination sites in China
- Performed TTL-Modulation Experiment
 - Traceroute immediately before blocking test
 - Whois to query ISPs
- Probed over a two-week period
- Result:
Where are the GFC routers? Which ISP?

Hops into China Where Filtering Occurs

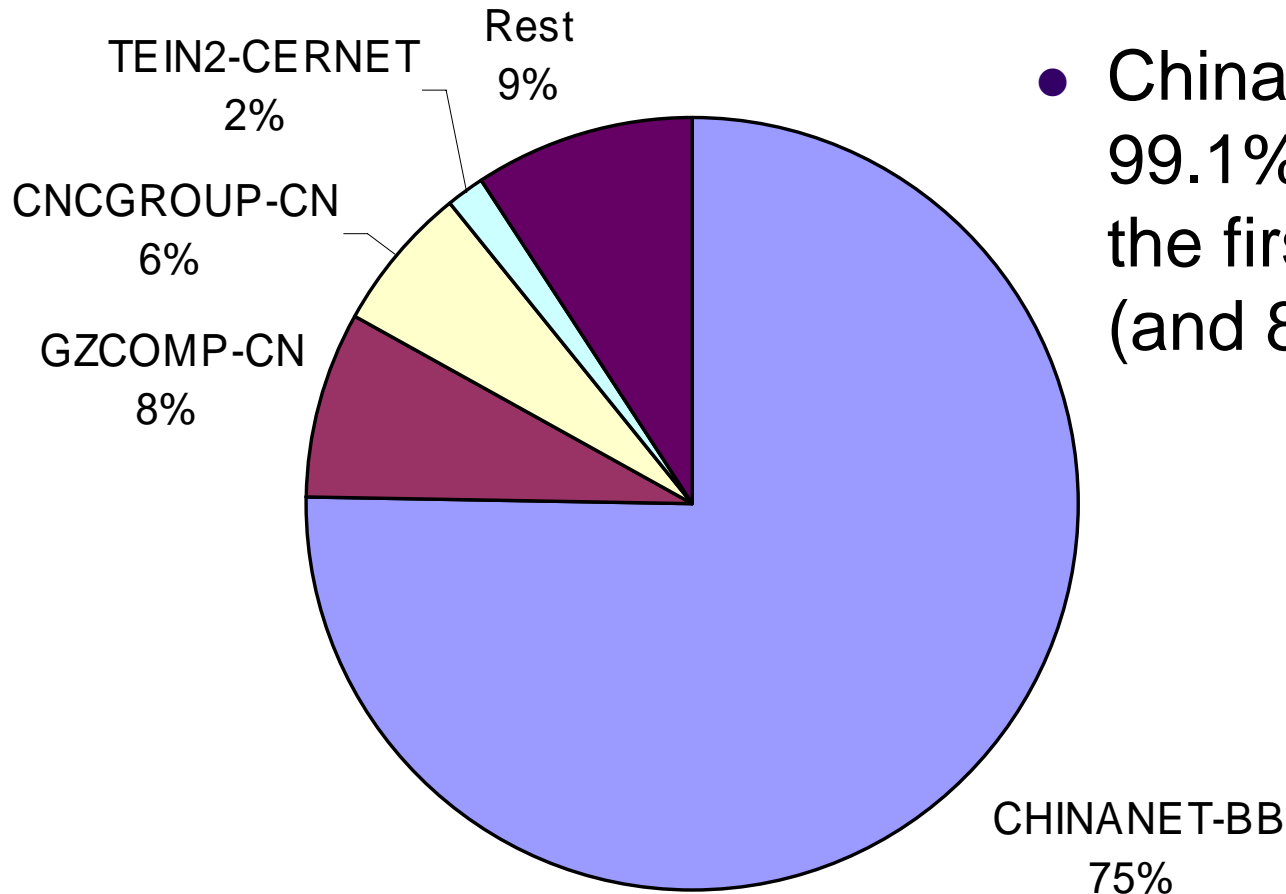
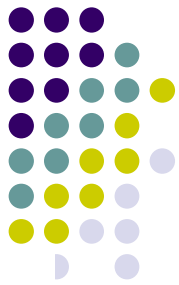


Blocked Paths

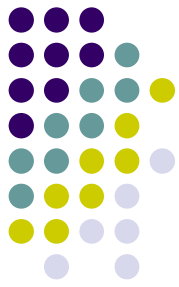


28% of paths were never filtered over two weeks of probing

First Hops



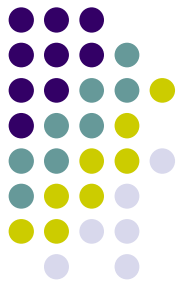
- ChinaNET performed 99.1% of all filtering at the first hop (and 83% of all filtering)



Outline

- Firewall or Something Else?
 - Where are filtering routers?
 - Who is doing filtering?
 - How reliable is filtering?
- Blocked Words
 - Which words to select?
 - Which words are blocked?
- Imprecise Filtering
 - What implications does keyword filtering have?

Slipping Words Through - Diurnal Pattern



Repeat

```
While "Falun" is not blocked
```

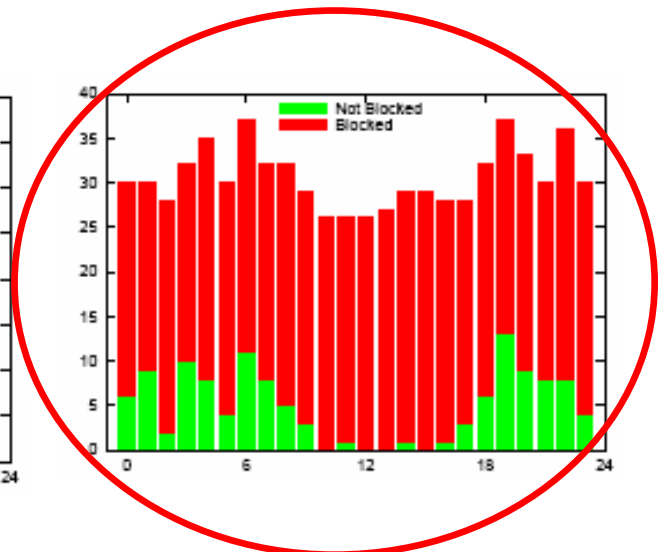
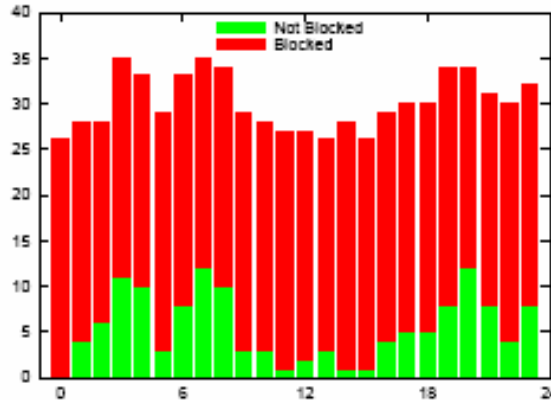
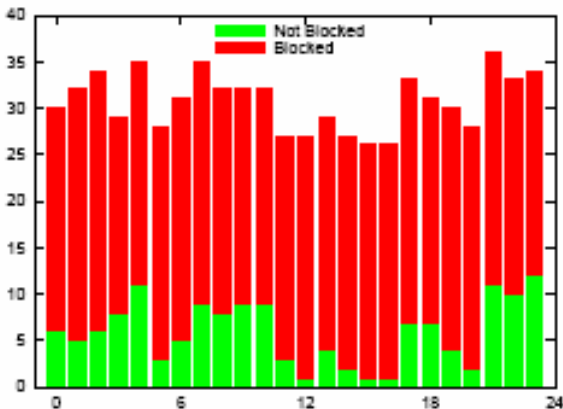
```
  green++
```

```
red++
```

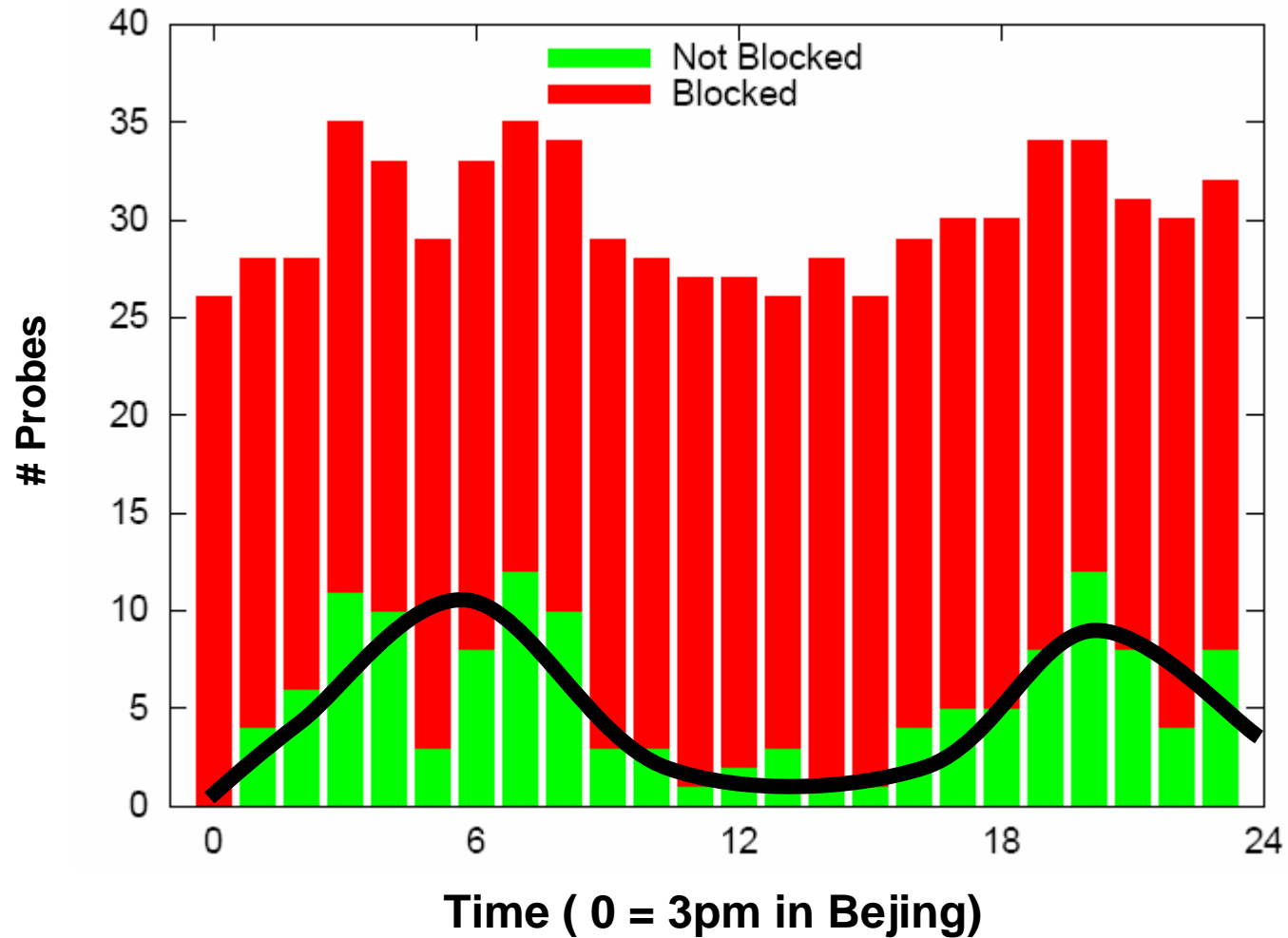
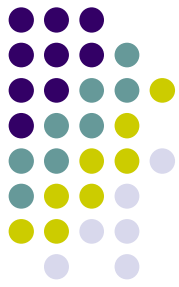
```
While "Test" is blocked
```

```
  wait
```

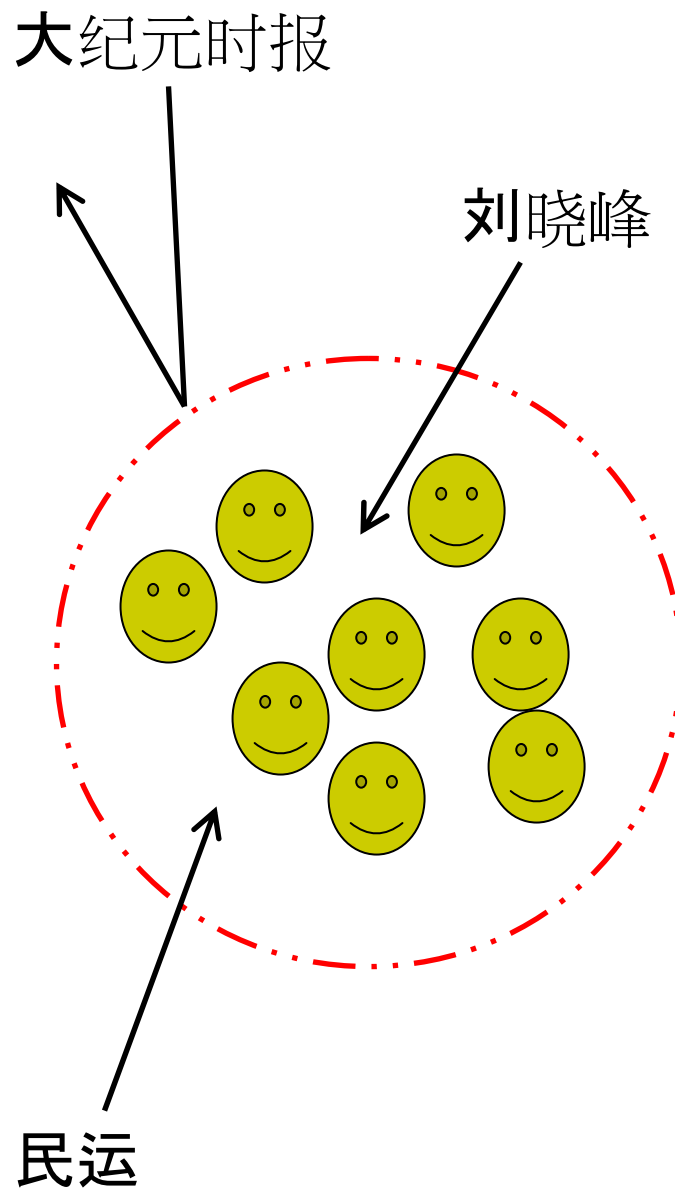
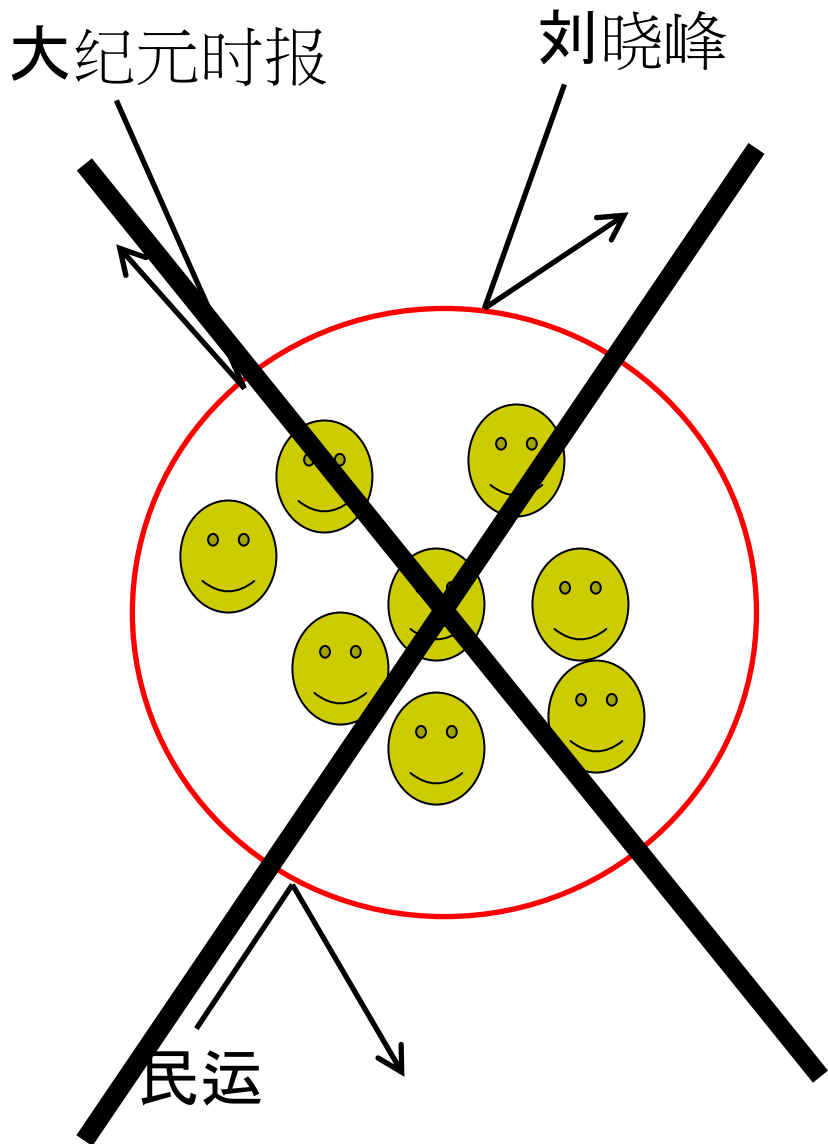
Forever



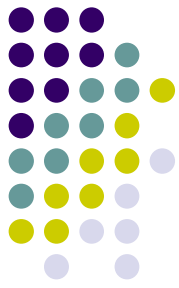
Slipping Words Through - Diurnal Pattern



Firewall?



Panopticon!

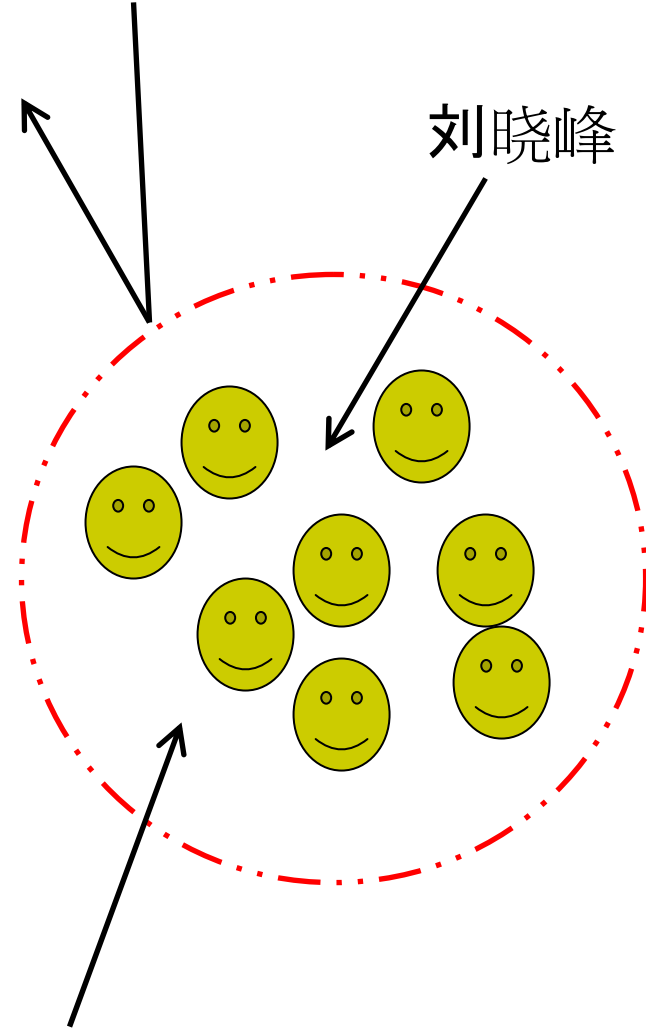


- Imperfect filtering
- Not strictly at the border
- Promotes self-censorship
- Good enough
- Defeating a Panopticon is different than defeating a firewall

大纪元时报

刘晓峰

民运





Outline

- Firewall or Something Else?
 - Where are filtering routers?
 - Who is doing filtering?
 - How reliable is filtering?
- **Blocked Words**
 - Which words to select?
 - Which words are blocked?
- Imprecise Filtering
 - What implications does keyword filtering have?

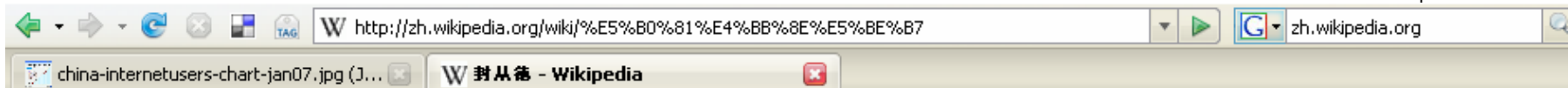
Latent Semantic Analysis (LSA)



- Deerwester et al., 1988
- Document summary technique to find relationships between documents and words
- Based on co-occurrence of words in a collection of documents

What to use as corpus?

Chinese Version of Wikipedia!



維基百科
自由的百科全書

导航

- [首页](#)
- [社区主页](#)
- [新闻动态](#)
- [最近更改](#)
- [特色内容](#)
- [随机条目](#)

帮助

- [帮助](#)
- [互助客栈](#)
- [询问处](#)
- [繁简转换](#)
- [所有页面](#)
- [联系我们](#)
- [资助维基百科](#)

搜索

[登录](#)或[创建](#)

[条目](#) [讨论](#) [编辑本页](#) [0](#) [历史](#) [不转换](#) ▼

若您来自中国大陆，并能顺利浏览维基百科，请[登入](#)后至[状况回報](#)。（注意：若未[注册](#)或[登入](#)，您的IP地址会被显示。）

封从德

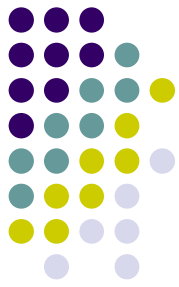
维基百科，自由的百科全书

[[编](#)]

封从德（1966年－），中国四川人。1989年六四天安门事件的学生领袖之一。原就读于北京大学。1989年间与另一学**柴玲**为夫妻关系，学潮遭到镇压后二人共同前往**法国**避难，后**柴玲**又辗转到了**美国**，封从德则继续留在**法国**读完了博士并和**柴玲**离婚。

封从德于1982年入读**北京大学**，1986年保送**北京大学**遥感所研究**计算机识别卫星图象专家系统**，1988年获**高级程序员**。1989年5月获**波士顿大学**五年博士学位奖学金。六四事件期间当选**北大筹委会**常委、**高自联**主席，及任**绝食团**和**广场指**总指挥。六四事件后在国内逃亡的十个月中感悟**传统文化**价值，遂弃理从文，出国后转入**法国实用高等研究院**宗教历史。1995年获硕士文凭，1996年通过博士候选资格答辩。1997年在**荷兰莱顿大学**研究**儒家学说**，1999年回**法国**撰写博士论**道教**和**中医**，2003年7月获得博士学位。以设计**网络数据库**为业。现主编《[六四档案](#)》。兼任《[公民议政](#)》编委。

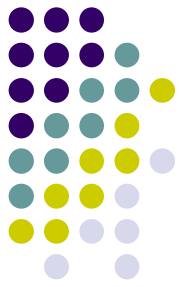
7个分类：[1966年出生](#) | [在世人物](#) | [四川人](#) | [中国持不同政见者](#) | [八九民运学生领袖](#) | [封姓](#) | [六四人物](#)



LSA of Chinese Wikipedia

	d_1	d_2	d_3	\cdots	d_n
t_1	2	0	2	\vdots	0
t_2	0	2	1	\vdots	0
t_3	5	0	0	\vdots	0
t_4	0	0	0	\vdots	4
t_5	1	0	0	\vdots	0
t_6	0	0	1	\vdots	0
\vdots	\cdots	\cdots	\cdots	\ddots	\cdots
t_m	0	1	0	\vdots	0

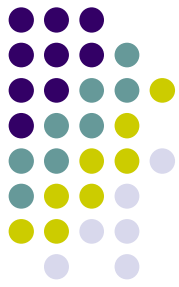
- $n=94863$ documents and $m=942033$ terms
- tf-idf weighting
- Matrix probably has rank r where $k < r < n < m$
- Implicit assumption that Wikipedia authors add additive Gaussian noise
- SVD and rank reduction to rank k



10 + 2 Seed Concepts

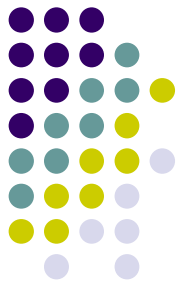
Bootstrap concept	Translation
六四事件	June_4th_events
高智晟	Gao_Zhisheng
赵紫阳	Zhao_Ziyang
选举	Election
红色恐怖	Red_Terror
大纪元时报	Epoch_Times
李洪志	Li_Hongzhi
台 (衛國)	Taiwan
东突厥斯坦	East_Turkistan
海啸	Tsunami
第二次世界大战	World_War_II
德国	Germany

Words correlated with 六四事件 – June 4th Events

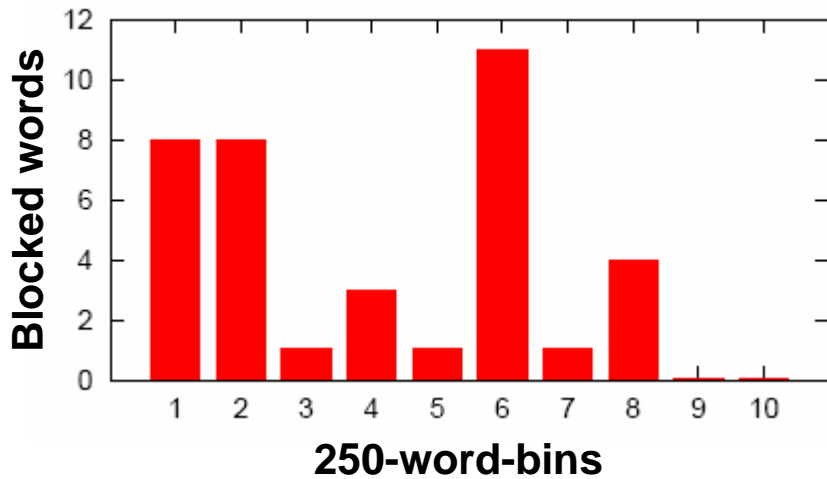


- 1 : 六四事件 – June 4th Events
- 2 : 重庆高家花园嘉陵江大桥 - Chongqing high family garden Jialing River bridge
- 3 : 樂提羌渠 - Yu Fulo (related to Chinese Eastern Han Dynasty)
- 4 : 李建良 - Li Jianliang
- 5 : 美丽岛事件 - Gaoxiong event (violent political event 1979) ←
- 6 : 赵紫阳 - Zhao Ziyang (Name, related to China travel logistics)
- 7 : 統戰部 - United front activities department
- 8 : 陈炳德 - Chen Bingde
- 9 : 洛杉磯安那罕天使歷任經營者與總教練 - Los Angeles Angels of Anaheim ... ←
- 10 : 李铁林 - Li Tielin (Government official)
- 11 : 邓力群 - Deng Liqun (Chinese politician) ←
- 12 : 中国政治 - Chinese politics
- 13 : 中共十四大 - The Chinese Communist Party 14th ...
- 14 : 改革开放 - Reform and open policy
- 15 : 报禁 - The newspaper endures
- to 2500

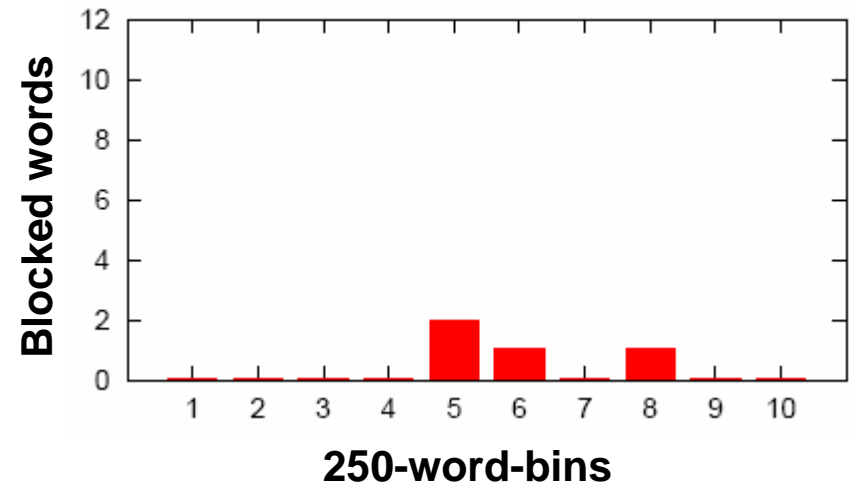
Efficient Probing



Epoch Times



Random Words

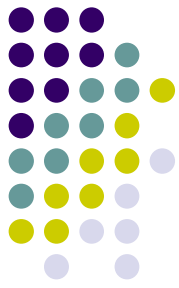


37

vs.

4

Blocked Words (122 discovered)



Pornography:

- 色情 --- Pornography
- 处女卖淫案 --- Virgin prostitution law case

Politics:

- 反人类罪 --- Crime against humanity
- 专政 --- Dictatorship (party), also 群众专政, 独裁, 一党专政, 专制
- 红色恐怖 --- Red Terror
- 六四事件 --- June 4th events (1989 Tiananmen Square protests)
- 藏独 --- Tibet Independence Movement

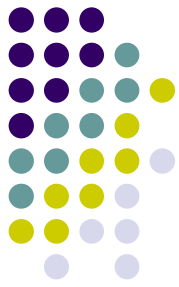
Others:

- 封杀 --- Block
- 桥头电厂 --- (Qinghai) Qiaotou power plant
- 卢多维克·阿里奥斯托 --- Ludovico Ariosto



Outline

- Firewall or Something Else?
 - Where are Filtering Routers?
 - Who is doing Filtering?
 - How Reliable is Filtering?
- Blocked Words
 - Which words to select?
 - Which words are blocked?
- Imprecise Filtering
 - What implications does keyword filtering have?



Imprecise Filtering

- Filtered are:

- 北莱茵-威斯特法伦
(Nordrhein-Westfalen – German state)
- 国际地质科学联合会
(International geological scientific federation)
- 卢多维克·阿里奥斯托
(Ludovico Ariosto – Italian Poet)

- Because:

法伦

(Sounds like *Falun* Gong)

学联

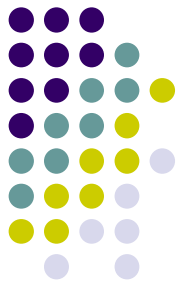
(student federation)

多维

(multidimensional)



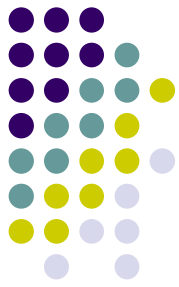
Keyword-based Censorship



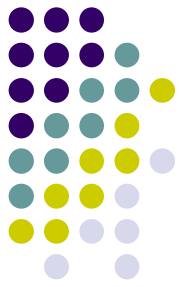
Censor the Wounded Knee Massacre in the Library of Congress

- Remove “Bury my Heart at Wounded Knee” and a few other select books?
- Remove every book containing the keyword “massacre” in its text?

~~Massacre~~



- Dante's *"Inferno"*
- *"The War of the Worlds"* by H. G. Wells
- *"King Richard III,"* and *"King Henry VI,"*
Shakespeare
- *"Adventures of Tom Sawyer,"* Mark Twain
- Jack London, *"Son of the Sun,"* *"The Acorn-planter,"* *"The House of Pride"*
- Thousands more

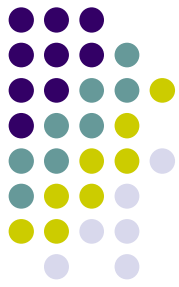


More Imprecision

- ~~Crime against humanity~~ *"The Economic Consequences of the Peace,"* John Maynard Keynes
- ~~Dictatorship~~ *The U.S. Constitution*
- ~~Suppression~~ *"Origin of Species,"* by Charles Darwin
- ~~Block~~ *"Computer Organization and Design,"* P. H.
- ~~Hitler~~ Virtually every book about World War II
- ~~Strike~~ *"White Fang," "The Sea Wolf,"* and *"The Call of the Wild,"* Jack London

Hypothetical?

Actually Blocked



屠杀	Massacre
反人类罪	Crime against humanity
专政 or 专制	Dictatorship
镇压	Suppression
封杀	Block
希特勒	Hitler
罢工	Strike

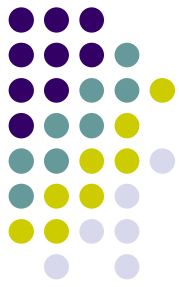


Future Work

- ConceptDoppler –
A Censorship Weather Report

What words are censored today?

- Track the blacklist over a period of time, to correlate with current events
 - Named entity extraction, online learning
- Scale up (bigger corpus, more words, advanced document summary techniques)



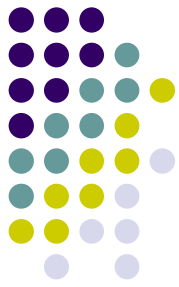
Future Work

- What are the effects of keyword filtering?
 - What content is being targeted?
 - What content is collateral damage due to imprecise filtering?
- Where *exactly* is filtering implemented?
 - More sources
 - Topological considerations
 - IP tunneling, IPv6, IXPs, ...



Conclusions

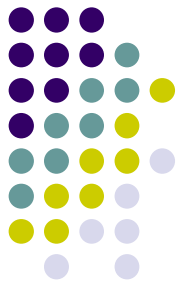
- Firewall vs. Panopticon
 - GFC implemented mostly at the borders by Chinanet, but also inner routers do filter
 - Filtering is **NOT** reliable:
 - Routes without GFC routers
 - Slip through during busy periods of the day
- Blocked words
 - Blocked more than pornography and sedition
 - LSA can help to increase probing efficiency
- Imprecise Filtering
 - You block a whole lot more than you probably want to



Thank You.

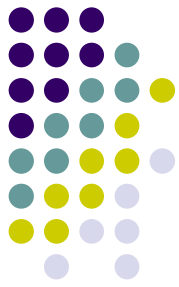
Questions?

<http://www.conceptdoppler.org>

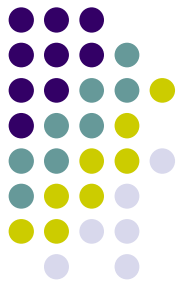


Unsponsored add: University of New Mexico CS dept. is hiring for 2 junior level positions and 1 senior level position.

Thanks, Jed + michael!



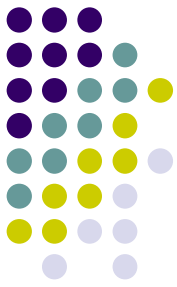
~~Crime against humanity~~



- *“The Economic Consequences of the Peace,”*
John Maynard Keynes
- Thousands more?

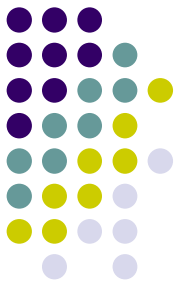
~~Dictatorship~~

- *The U.S. Constitution*
- Thousands more?



~~Traitor~~

- “*Fahrenheit 451*,” Ray Bradbury
- Thousands more?

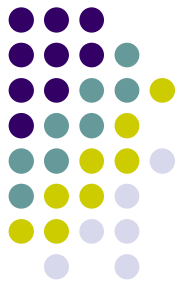


~~Suppression~~



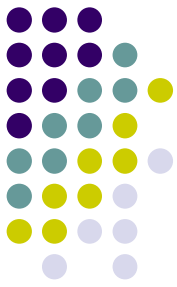
- “*Origin of Species*,” by Charles Darwin
- Thousands more?

~~Block~~



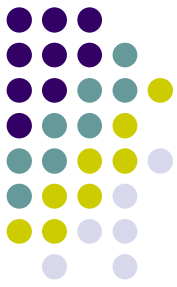
- “*Computer Organization and Design*,” Patterson and Hennessy
- “*Artificial Intelligence: 4th Edition*,” George F. Luger
- Millions more?

~~Hitler~~

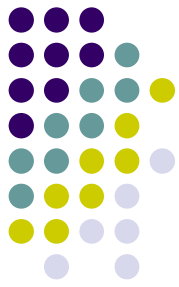


- Virtually every book about World War II

~~Strike~~

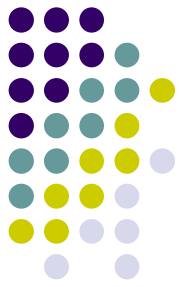


- “*White Fang*,” “*The Sea Wolf*,” and “*The Call of the Wild*,” Jack London
- Millions more?



Outline

- Firewall or Something Else?
 - Where are Filtering Routers?
 - Who is doing Filtering?
 - How Reliable is Filtering?
- Blocked Words
 - Which words to select?
 - Which words are blocked?
- Imprecise Filtering
 - What implications has keyword filtering?



Outline

- Implications of Imprecise Filtering
 - What are consequences of key-word-based filtering?
- Panopticon vs. Firewall
 - How is filtering implemented?
 - Where is filtering implemented?
 - How “reliable” is filtering?
- Blocking Words
 - How to efficiently discover blocked words?
 - What words are blocked?



Outline

- Implications of Imprecise Filtering
 - What are consequences of key-word-based filtering?
- Panopticon vs. Firewall
 - How is filtering implemented?
 - Where is filtering implemented?
 - How “reliable” is filtering?
- **Blocking Words**
 - How to efficiently discover blocked words?
 - What words are blocked?