

# Tweeting Under Pressure: Analyzing Trending Topics and Evolving Word Choice on Sina Weibo

**Le Chen**

Chi Zhang

Christo Wilson

College of Computer and Information Science  
Northeastern University

# Censorship in China

- Great Firewall of China (GFW)



# Censorship on Weibo

- What will be censored?



- How is censorship implemented?

THE WORST  
THING ABOUT  
CENSORSHIP  
IS [REDACTED]



# Current Understanding of Censorship on Weibo

- Scale of censorship<sup>[1][2]</sup>
  - 16% of censored tweets vs. 0.01% of censored tweets
- Velocity of censorship<sup>[3]</sup>
  - Focus on 3K of Chinese celebrities
  - 30% of tweets are deleted within 30 minutes

抱歉，此微博不适宜对外公开。如需帮助，请联系客服。 <http://t.cn/z0D6ZaQ>

Sorry, this tweet is inappropriate to be public...

[1] BAMMAN, D., O'CONNOR, B., AND SMITH, N. A. Censorship and Deletion Practices in Chinese Social Media. First Monday 17, 3 (2012).

[2] WA FU, K., HONG CHAN, C., AND CHAU, M. Assessing censorship on microblogs in china: Discriminatory keyword analysis and the real-name registration policy. IEEE Internet Computing 17, 3 (2013).

[3] ZHU, T., PHIPPS, D., PRIDGEN, A., CRANDALL, J. R., AND WALLACH, D. S. The velocity of censorship: High-fidelity detection of microblog post deletions. In Proc. of USENIX Security (2013).

# Impact of Censorship

- Information dissemination
  - Extensively studied on Twitter
  - Missing: adversarial forces
- Chilling effect?
  - Dissuade people from discussing about censored topics
- Do users adapt?
  - Change linguistic conventions or behaviors

# Outline

- Motivation
- Methodology of data gathering
  - Data gathering: API or DIY
  - Crawler design and data collection
  - Validation of the method
- Topics analysis
- Word usage on Weibo

# Sina Weibo Overview

- Launched in August 2009
- By December 2012,  $\approx$  500 million users
- Over 4.6 million users and 100 million tweets daily

Profile image  
and  
stats

Username



functionalPL

<http://weibo.com/u/3102867863>

一句话介绍一下自己吧，让别人更了解你

已婚 | 海外 美国 | 毕业于北京邮电大学 | 标签

编辑个人资料

71  
关注

16  
粉丝

143  
微博

我的主页

微博

个人资料

关注/粉丝

相册

微博 | 原创 | 图片

雄起勒哦 //@微博小秘书 转发微博

@青年说

【#中关村牵手贵阳#】中关村管委会与贵阳市政府9月8日签订战略合作框架协议，贵州省委书记赵克志、北京市政协主席吉林为“中关村贵阳科技园”揭牌，生态文明与自主创新两个国家示范正式结盟。期间，举行中关村企业家贵阳峰会，联想、京东等一批高新企业签约106个项目、投资总额

A retweet with original  
tweet embedded

今天吃了这次在天朝最难吃的一顿 什么西红柿牛肉面 就三片火锅牛肉 俩圣女果 泥马酸甜甜全是味精 配套的7个饺子全部粘在一起 扯一个破一个 这叫昨天才吃完干锅火锅的我情何以堪 看来一直长不胖是有原因的@gxz\_99 tmd太难吃了

+加标签

7月27日 11:11 来自iPhone客户端

阅读(355) | 转发 | 收藏 | 评论(13)

An  
original  
tweet

新浪微博社区管理中心举报处理大厅，欢迎查阅！

同时转发到我的微博

评论

全部 | 共同评论 | 认证用户 | 关注的人

共13条



xuexuexue: 回复@functionalPL 哈哈，你一点都不不是哪种看中钱的人 (8月6日 07:37)

删除 | 查看对话 | 回复



functionalPL: 回复@xuexuexue 我是那种看重钱的人么？反正不送钱的那桌就只有西红柿鸡蛋盖饭吃 配套一瓶矿泉水 (8月6日 01:07)

删除 | 查看对话 | 回复



xuexuexue: 回复@functionalPL 反正请哦我也不送钱， (7月30日 00:26)

删除 | 查看对话 | 回复

Comments

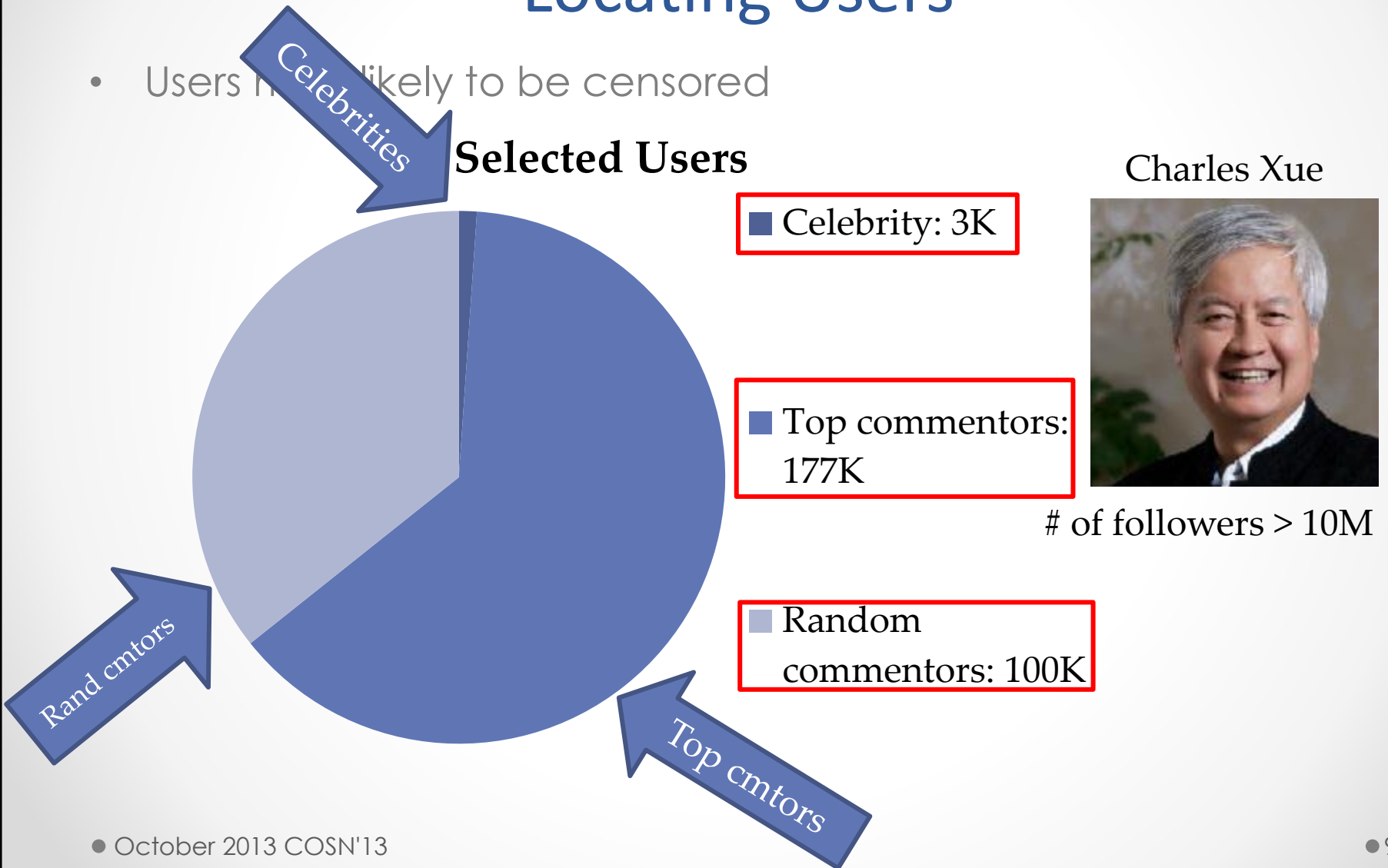
# Data Gathering Challenges

- Unlike Twitter, Weibo removes data
- Censorship happens in minutes
  - Frequently revisit users to keep track of censorship
- Strict rate-limit on requests



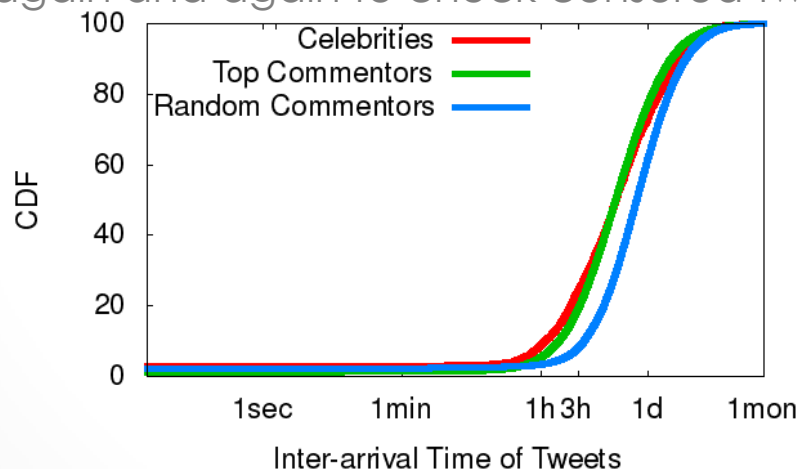
# Locating Users

- Users most likely to be censored



# Crawler Design

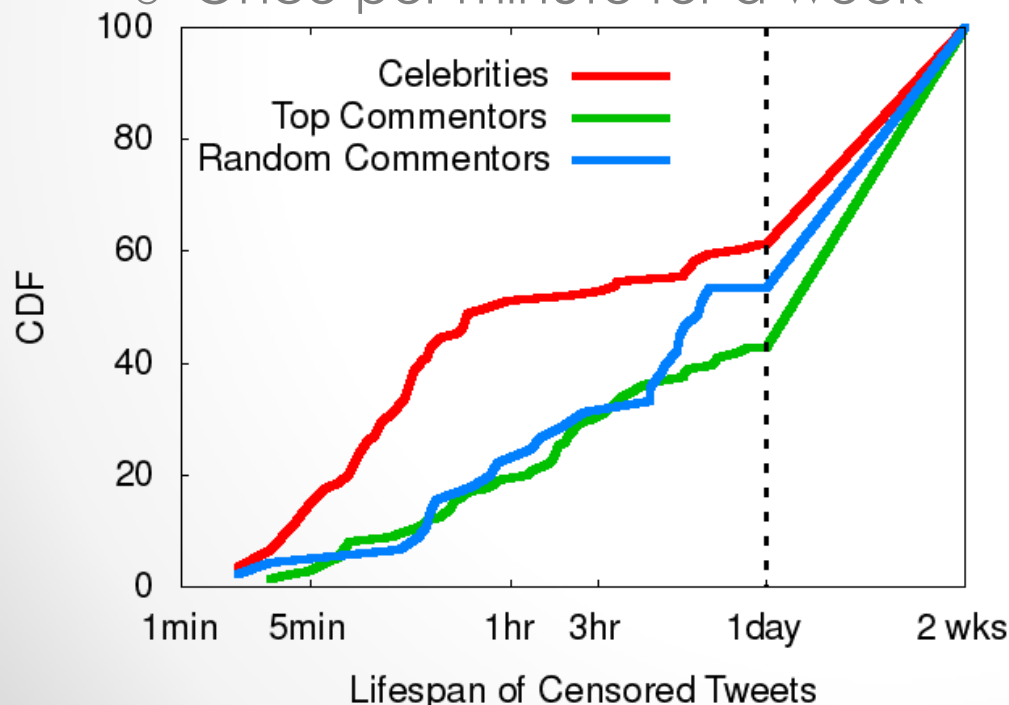
- Design challenges
  - Different users tweet at different speed
  - Censorship happens shortly after tweets are posted
- Solution
  - Crawl the users at different speed
  - Revisit again and again to check censored tweets



Three buckets: one hour (5K), three hours (22K), daily (253K); crawled active users more often

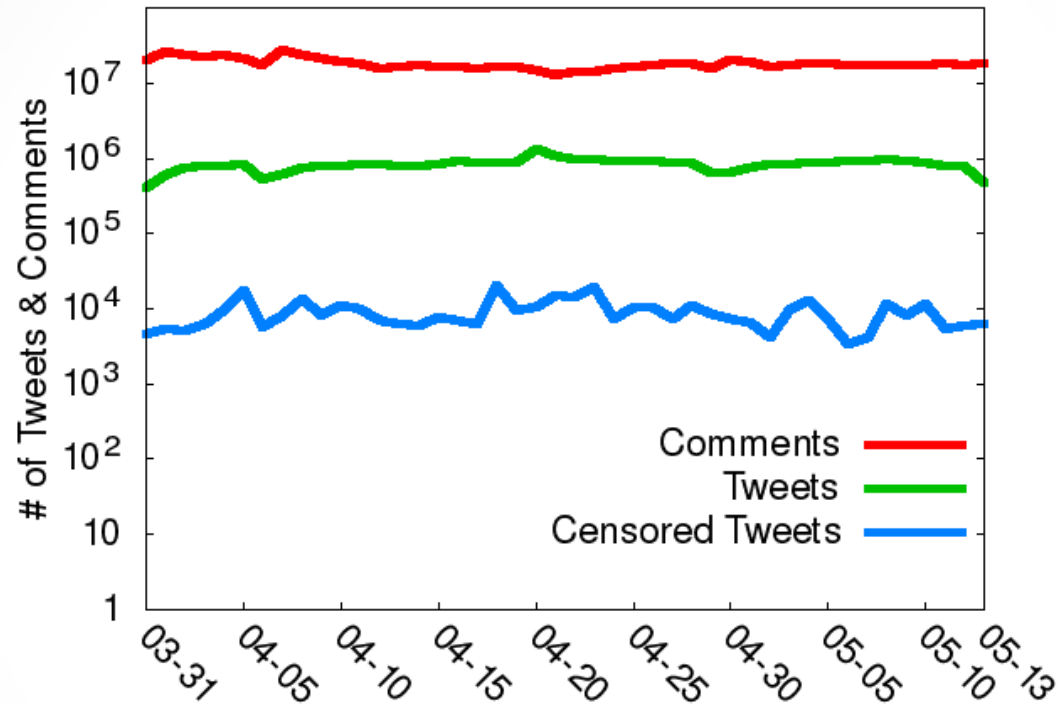
# Validation

- Crawl in hours but censorship happens in minutes
- Intensive crawling
  - 500 random users from each group
  - Once per minute for a week



1. Gathers > 50% censored tweets
2. > 99% of content on Weibo is comments and comments are never censored

# Daily Activity on Weibo



- 830K tweets; 18M comments; 9K censored tweets
- 1% of tweets are censored
  - Estimates from prior work: 0.01% --16%

# Outline

- Motivation
- Methodology of data gathering
- Topics analysis
  - Locating trending topics
  - Analysis of trending topics
- Word usage on Weibo

# Key Research Question

- What is the impact of censorship on discourse on Sina Weibo?
  - Impact of censorship in trending topics
  - Censored topics vs. uncensored topics

# Locating Trending Topics

- Topic Extraction
  - Latent Dirichlet Allocation (LDA)
  - LDA output: 1) topics; 2) words in each topic
- Challenges
  - Chinese does not include breaks between words
    - Solution: word segmentation using OpenCLAS
  - Many new words on Weibo
    - OpenCLAS has a small dictionary (104K)
    - Solution: Use the Sogou Pinyin Dictionary (including words from social media)
    - Sogou pinyin: crowd-sourced and keep updating

# Topic Validation

- Output of LDA: 300 topics
  - 37 trending topics corresponded to real-world events;
- Locating trending topics

General	Real-world events	Censored
Economics	Boston marathon bombing	Lushan secretary
Food	Ya'an earthquake	Beijing taxi driver
English learning	Wechat charge	Bird flu

- Validation
  - All 11 trending topics listed on external sources



# High-level Overview

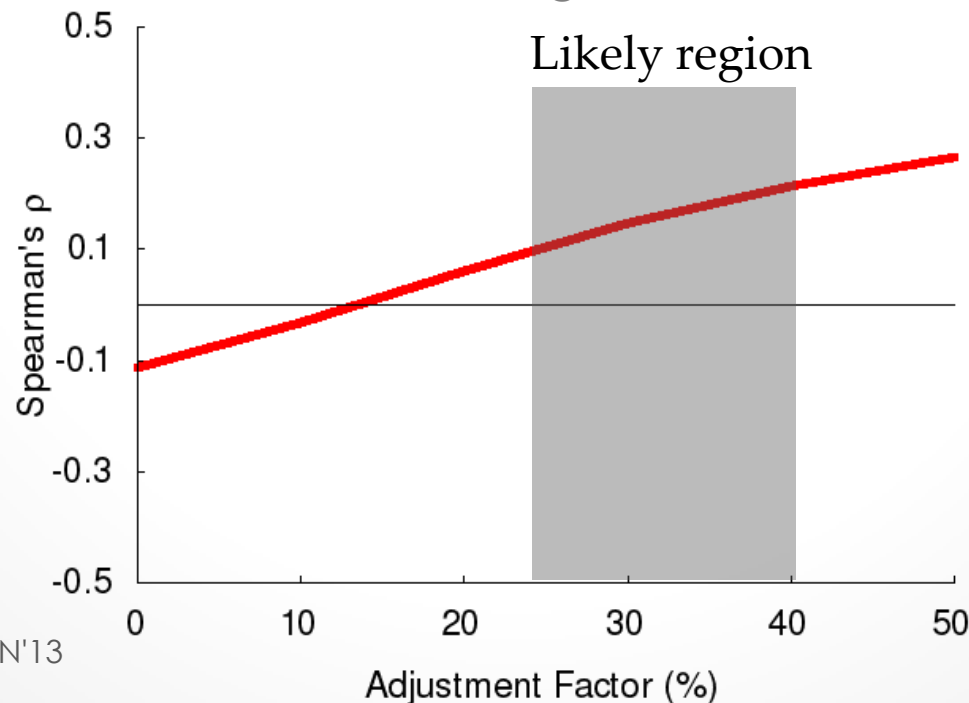
- 37 real-world events

Characteristics	Range	Average
Lifespan(days)	2 ~ 14	4.6
# of tweets	394 ~ 108K	19K
# of comments	538 ~ 3.1M	635K
% of censored tweets	0% ~ 82%	1.7%

- Most topics are not censored
  - 5 topics experienced >5% of censorship

# Impact of Censorship

- Does censorship have impact on user behavior?
- Correlation between censorship and # of tweets per user
  - Spearman's  $\rho = -0.1$
- But our crawler missed some censored tweets
  - Adjustment factor = % of missing tweets



## Key Research Question

- What is the impact of censorship on discourse on Sina Weibo?
- **Ans: We observe there is no chilling effect. Instead, censored topics see more active users tweeting more frequently.**

# Outline

- Motivation
- Methodology of data gathering
- Topics analysis
- Word usage on Weibo
  - Locating morphs
  - Morph usage and censorship
  - Dynamics of Morph usage

# Key Research Question

- Do Weibo users adapt in order to avoid censorship?
  - Keyword-based censorship leads to linguistic adaptation
  - To analyze adaptation, we focus on morphs

# Locating Morphs

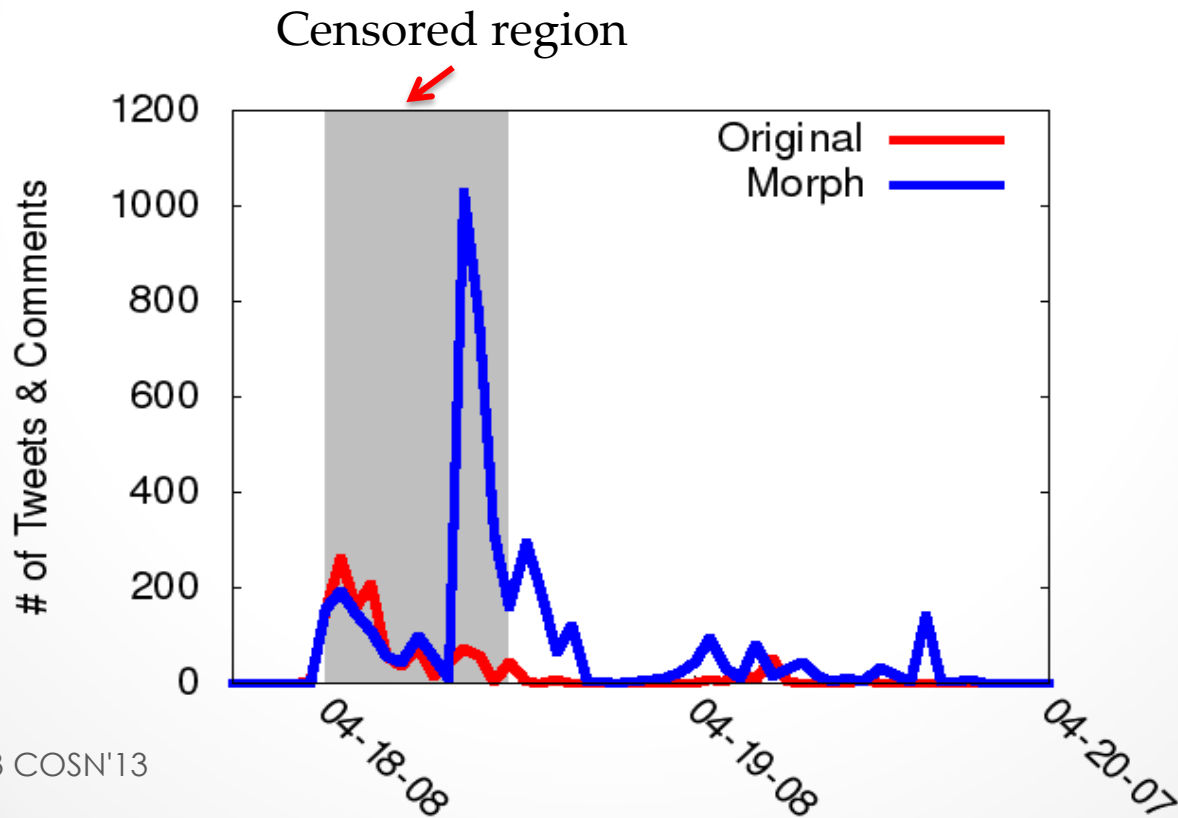
- What is a morph?

Type	Original	Literal Meaning	Morph	Literal Meaning
Homophone	和谐	Harmony	河蟹	River crab
Pinyin Initial	美眉	Pretty eyebrow	MM	N/A
Neologism	囧	Brightness	囧	Embarrassed
Anglicization	京温	Jingwen Mall	京wen	Jingwen Mall
Homograph	翻墙	Over the GFW	番羽土嗆	N/A

- Identify novel morphs
  - A morph from topic A used > 100 times prior to the start date of A was not invented during A, and is not a novel morph

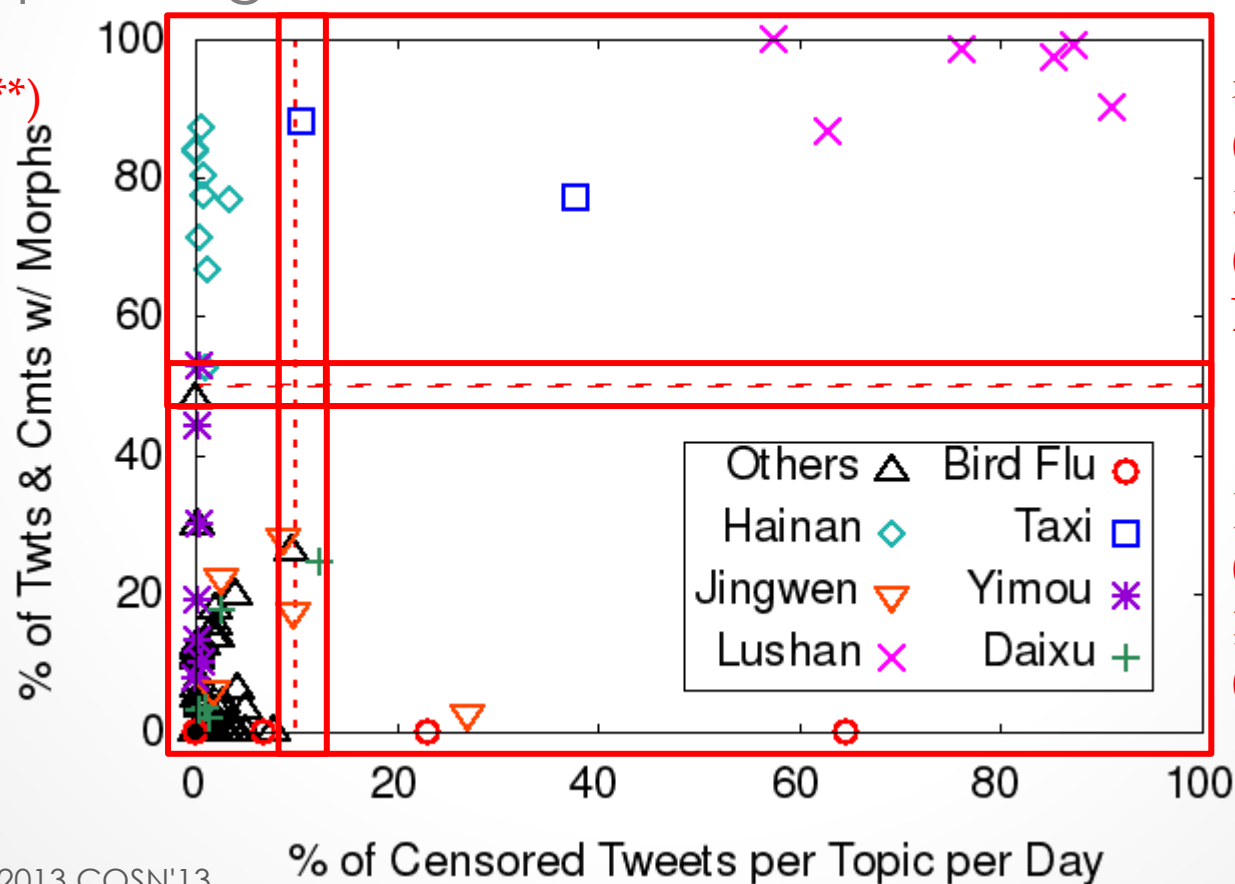
# Morph Usage Over Time

Topic	Story	Original	Morphs
Taxi	A rumor involving a taxi driver to meet President Xi.	郭立新(the driver's name)	郭师傅(Shifu Guo), 北京的哥(Beijing Taxi Driver)



# Morph Usage and Censorship

- Is there a relationship between censorship and morph usage?



表印哥  
(Bro Watchprint)  
北京的哥  
(Beijing Taxi Driver)

京wen  
(Jingwen Mall)  
戴某  
(Dai XXX)



# Key Research Question

- Do Weibo users adapt in order to avoid censorship?
- **Ans: We observe that as censorship becomes prominent, the morph usage becomes popular.**
- What are the dynamics of morph generation?
- **Please refer to the paper, or we can discuss after the talk.**

# Summary

- Study the impact of censorship on discourse and word choice on Weibo
- Crawled 280K users on an hourly basis for 44 days, gathering 839M tweets and comments
- There is positive relationship between censorship and user engagement
- Weibo users tend to introduce morphs to heavily censored topics
- Updating information dissemination models that take adversarial forces into account

# Questions?