

Bringing Big Data to Chinese Politics: Learning From the People's Daily Corpus

Matthew P. Robertson

Australian National University

January 06, 2022

Significance of People's Daily corpus

- Official mouthpiece of the Central Committee of the Chinese Communist Party
- Perhaps most authoritative public source of information about Chinese politics
- Published daily since 1946 - over 2 million articles
- Used by CCP to communicate official policy both to rank-and-file and the public
- Placement, emphasis, and construction of stories and headlines is an indicator of their salience. Given that PD is the voice of the Party, this is an indicator of the significance the Party attaches to topics
- Editor-in-chief is at minister-rank; editorials are often written on instruction of leaders and subject to extensive review

Prior art

- PD corpus has been used to a variety of topics:
 - Foreign policy (Kivimäki 2015)
 - Public health (Dong, Chang, and Chen 2008)
 - Propaganda (G. Wu 1994)
 - Transgender issues and homosexuality (Zhang 2014; Huang 2018)
 - People with disabilities (Ye and Zeldes 2020)
 - Climate change (Wang 2018)
 - Revival of Confucianism (S. Wu 2014)
 - Drug control (Liang and Lu 2013)
 - Military metaphors in describing disease (Yang 2021)
 - Death penalty (Trevaskes, Smith & Robertson, forthcoming)

Two of the biggest uses of the corpus are (1) machine learning project to predict policy shifts, outbreaks, and crackdowns, and (2) a study of the shifting meaning of bureaucratism.

Only those two seem to have featured the full corpus.

Prior art

- Policy Change Index: Chan, Julian Tszkin, and Weifeng Zhong. 2018. “Reading China: Predicting Policy Change with Machine Learning.” AEI Economics Working Paper, No. 2018-11

Policy Change Index (PCI)

A Machine Learning Framework to Predict Policy Changes

The Policy Change Index (PCI) is a series of open-source machine learning projects that predict authoritarian regimes' major policy moves by "reading" their propaganda publications. The first three projects inducted into the series are about China's policies and based on its official newspaper — the *People's Daily*.

1. **PCI-China**: predicts China's policy changes from 1951 Q1 to the present.
2. **PCI-Crackdown**: predicts how close in time the 2019-20 Hong Kong protests are to a Tiananmen-like crackdown by China.
3. **PCI-Outbreak**: measures the severity of an epidemic outbreak in China, such as COVID-19.

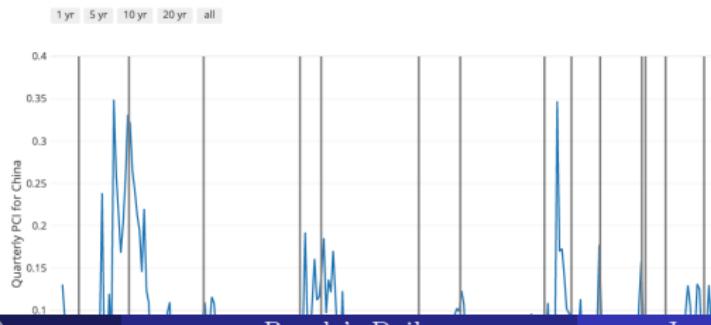
PCI-China

PCI-Crackdown

PCI-Outbreak

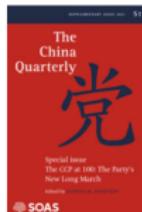
PCI-China and major events in China, 1951 Q1 to 2021 Q2

A spike in the PCI-China signals a major policy change, while a vertical bar marks the ground truth of the change labeled by the event. The PCI-China often spikes months before policy changes take place, validating the index's predictive power. Click [here](#) to learn more about how it works.



Prior art

- Ding, I., & Thompson-Brusstar, M. (2021). The Anti-Bureaucratic Ghost in China's Bureaucratic Machine. *The China Quarterly*, 248(S1), 116-140



The China Quarterly

Article contents

- Abstract
- The Spirit of the Ghost
- The Churning of the Machine
- Dead Letter or Living Faith?
- Conclusion
- Footnotes
- References

The Anti-Bureaucratic Ghost in China's Bureaucratic Machine

Published online by Cambridge University Press: 10 November 2021

Iza Ding  and Michael Thompson-Brusstar 

Show author details 

Article

Figures Metrics



Save PDF



Share



Cite



Rights & Permissions

Abstract

The Chinese Communist Party's (CCP) ideology, rooted in its foundational struggles, explicitly denounces "bureaucratism" (*guanliaozhuyi*) as an intrinsic ailment of bureaucracy. Yet while the revolutionary Party has blasted bureaucratism, its revolutionary regime has had to find a way to coexist with bureaucracy, which is a requisite for effective governance. An anti-bureaucratic ghost thus dwells in the machinery of China's bureaucratic state. We analyse the CCP's anti-bureaucratism through two steps. First, we perform a historical analysis of the Party's anti-bureaucratic ideology, teasing out its substance and emphasizing its roots in and departures from European Marxism and Leninism. Second, we trace both the continuity and evolution in the Party's anti-bureaucratic rhetoric, taking an interactive approach that combines close reading with computational analysis of the [entire corpus of the People's Daily \(1947–2020\)](#). We find striking endurance as well as subtle shifts in the substance of the CCP's anti-bureaucratic ideology. We show that bureaucratism is an umbrella term that expresses the revolutionary Party's anxiety about losing its popular legitimacy. Yet the substance of the Party's concern evolved from commandism and revisionism under Mao, to corruption and

This project's contribution

The corpus...

- does not seem to have been used to study elite politics and regime power dynamics
- has not been made widely available in a format that is accessible to other researchers

Primary motivation of project is to demonstrate the utility of the dataset while making it widely accessible for further work.

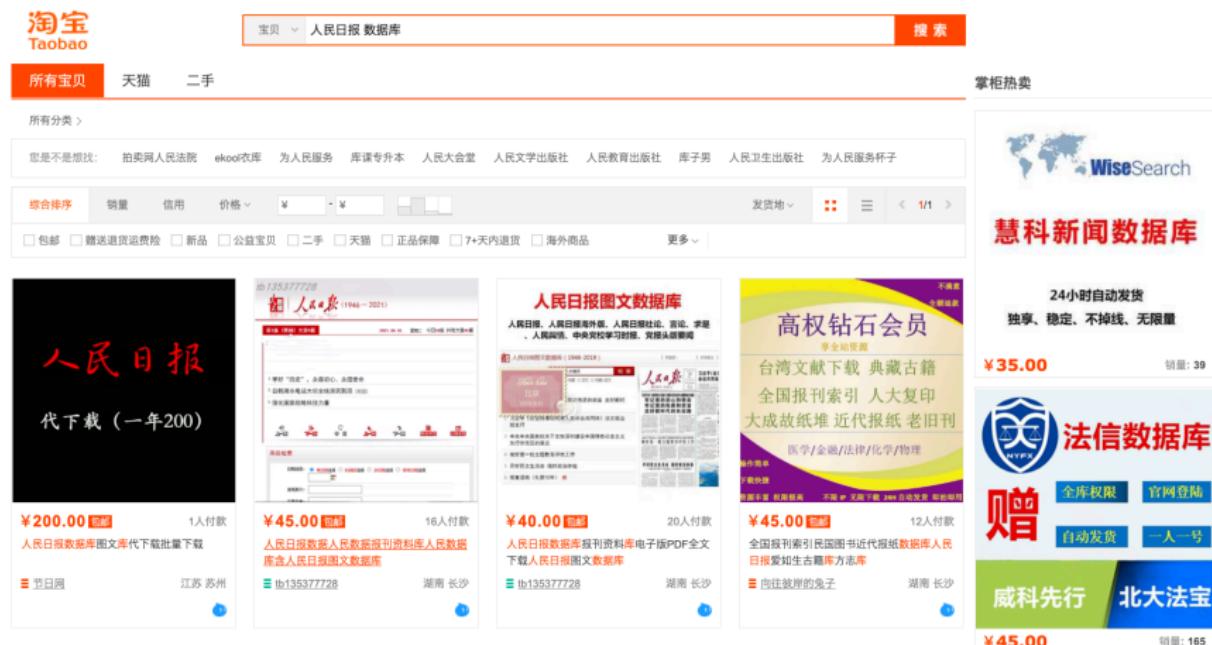
PD corpus is a rich source of data and other researchers are likely to come up with questions I've never thought of.

Technical challenges in collection

- Recent data (last ~12 months) is publicly available
- Historical data behind subscription/paywall
- Rate limits on access volume
- Browser emulation is slow & RAM heavy
- 2m articles at a theoretical 10k per day = 200 days

But...

- Of course it's all on taobao



宝贝 人民日报 数据库 搜索

所有宝贝 天猫 二手

掌柜热卖

是不是想找: 拍卖网人民法院 ekoo衣库 为人民服务 库课专升本 人民大会堂 人民文学出版社 人民教育出版社 库子男 人民卫生出版社 为人民服务杯子

综合排序 销量 信用 价格 ¥ ~ ¥ 发货地~ 更多 < 1/1 >

更多~

包邮 赠送退货运费险 新品 公益宝贝 二手 天猫 正品保障 7天内退货 海外商品

人民日报 (1949-2021) 代下载 (一年200) ￥200.00 1人付款 人民日报数据库 图文库代下载批量下载 节日园 江苏 苏州

人民日报图文数据库 ￥45.00 16人付款 人民日报数据库人民数据报刊资料库人民数据库含人民日报图文数据库 tb135377728 湖南 长沙

人民日报图文数据库 ￥40.00 20人付款 人民日报数据库报刊资料库电子版PDF全文下载人民日报图文数据库 tb135377728 湖南 长沙

人民日报 (1949-2021) 代下载 (一年200) ￥45.00 12人付款 全国报刊索引民国图书近代报纸数据库人民日报数据库含人民日报古籍库 赠自动发货 一元一号 向往彼岸的兔子 湖南 长沙

高权钻石会员 ￥35.00 39 销量: 39 24小时自动发货 独享、稳定、不掉线、无限量 台湾文献下载 典藏古籍 全国报刊索引 人大复印 大成故纸堆 近代报纸 老旧刊 医学/金融/法律/化学/物理

慧科新闻数据库 ￥35.00 39 销量: 39 24小时自动发货 独享、稳定、不掉线、无限量 台湾文献下载 典藏古籍 全国报刊索引 人大复印 大成故纸堆 近代报纸 老旧刊 医学/金融/法律/化学/物理

法信数据库 ￥45.00 165 销量: 165 全库权限 官网登陆 赠 自动发货 一元一号 威科先行 北大法宝 ￥45.00 165 销量: 165 全库权限 官网登陆 赠 自动发货 一元一号

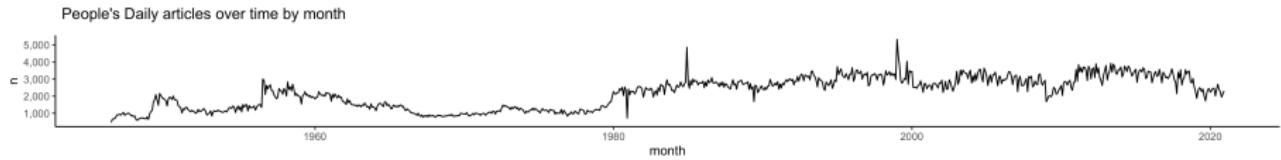
Technical challenges in organization & analysis I

- Some practical challenges to wrangling so much data; flat files don't cut it
- Speed, convenience, robustness: It needs to be put in a database structure
- Went with SQLite; compressed 25gb (2m html files) to 3gb
- Streaming cloud backups with Litestream
- Took days to parse the html and insert the records:
 - parallel processing
 - unfamiliarity with SQLite API (locks, write-ahead log, timeouts, everything)
 - constant crashes, used cron to stack jobs one after another
- Extensive validation required just to know if you did what you wanted to do

Technical challenges in organization & analysis II

- Protip: use data.table when dealing with millions of rows
- Protip 2: databases are good. They allow joins, filters and manipulations out of memory
- Whole corpus yet to be tokenized
- Following are simple dictionary counts. But also normalized by volume of text

Shape of dataset



Comparing Xi to Mao

- Dataset useful for illuminating debates on prominence of Xi Jinping versus Mao Zedong



SPECIAL REPORT

Xi Jinping Wants to Become the New Mao

The old Mao almost destroyed China.



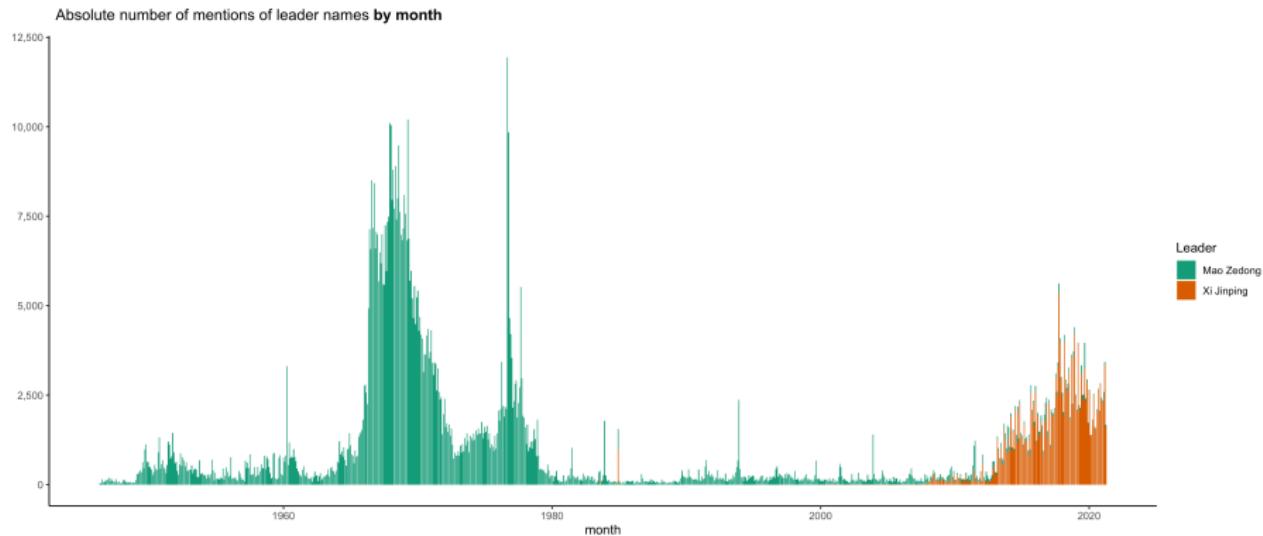
by DOUG BANDOW

May 26, 2020, 12:58 PM

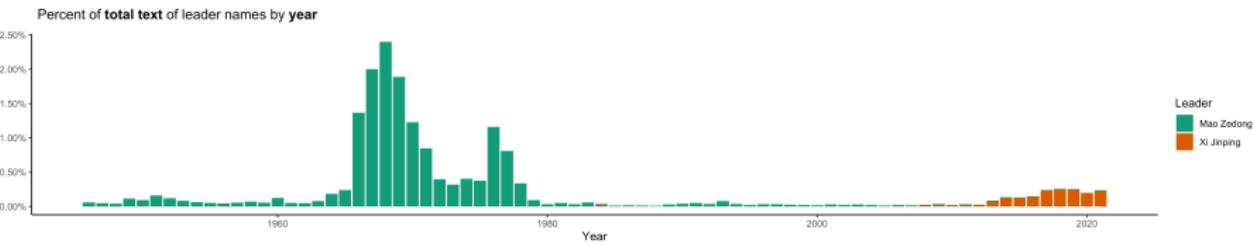


Comparing Xi to Mao

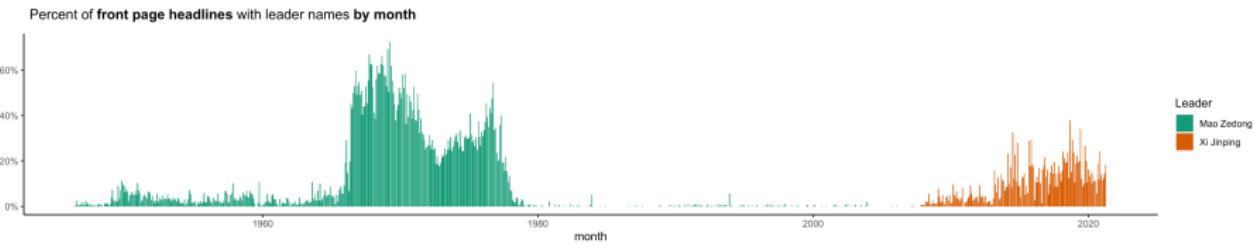
- Is Xi really “the new Mao?”
- Data shows that, at least in People’s Daily ink, Mao eclipses Xi; though Xi does eclipse predecessors



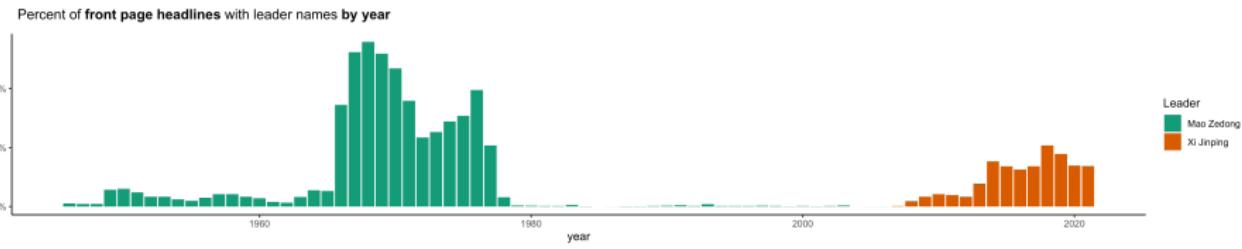
Comparing Xi to Mao



Comparing Xi to Mao



Comparing Xi to Mao



Comparing Xi to Mao

本条数据为：人民日报1986.05.26第1版

《社會研究》編輯組

前言

在工农兵群众和学生为主的运动中，不断地涌现出了大批一批批的先进人物。他们坚决听毛主席的话，把毛主席的指示当作自己行动的指南。树立全心全意为人民服务的思想，提高觉悟，改进工作，在阶级斗争、生产斗争和科学实验的革命实践中，取得了一个又一个的胜利。毛泽东思想武装起来的工农兵群众，是攻无不克，战无不胜的伟大力量。

工农兵群众对伟毛主席著作的根本态度，就是带着深厚的无产阶级感情读毛主席的书，就是坚决把毛主席的书当作各项工作的最高指示，真正做到“毛主席怎样说的，我就怎样做”。他们把毛主席的书比作“粮食”、“武器”、“方向盘”和“改革革命的指路明灯”。他们说：“毛主席的话就是最大的真理，毛主席的话就是真理。”“千条万条毛主席的话是第一条”，“千难万难有了毛泽东思想就不难”。“又一次一次地读毛主席的书，就一次又一次地认识毛主席思想的正确和伟大，又一次一次地加深了对毛主席著作的阶级感情”。他们带着浓厚兴趣学、灵活运用，特别是在“用”字上狠下工夫，掌握得用得得心应手。

对德毛主席著作采取什么态度。这不是一般的学风态度问题，而是阶级立场问题，是世界观的问题。

毛主席的书，是天才地、全面地、系统地、创造性地发展了马克思列宁主义，是对中国革命和建设经验的科学总结，是对国际共产主义运动正反两方面经验的科学总结，是被实践反复证明了的伟大的革命的科学真理。毛泽东思想是当代马克思列宁主义的顶峰，是最高最活的马克思列宁主义。是战胜现代修正主义和资产阶级思想的庞大武器。是不是听毛主席的话，是不是跟着毛主席指引的共产主义方向前进，是不是坚决遵循毛主席的指示办事，是不是把毛主席的书当作各项工作的最高指示，这是革命还是反革命的分水岭。是马克思列宁主义还是修正主义的分水岭。

“谁主生杀的书，听谁生杀的话，照谁生杀的指示办事”。这是关系到巩固无产阶级专政，防止资本主义复辟的大问题，是关系到我们党和国家的命运和前途的大问题，也是关系到世界革命前途的大问题。

一切的工作中，不论在什么时候和什么情况下，我们都要坚决地按照毛主席的指示办事。凡是毛主席指示的，我们就一定要坚决执行，坚决照办；凡是违背毛主席指示的，我们就一定要坚决抵制，坚决反对。

工农兵群众活学活用毛泽东著作，高举毛泽东思想的伟大红旗，积极参加社会主义文化大革命，他们不仅是阶级斗争、生产斗争、科学实验三大革命运动的主力军，而且也是社会主义文化大革命的主力军。他们活学活用毛泽东著作，为我们树立了光辉的榜样。知识分子、学术工作者必须认真地学好工农兵学习，深入到工农兵群众中去。同他们一起生活、一起斗争，改造自己的思想感情，真正站到无产阶级的立场上来。真正树立无产阶级世界观，学习工农兵群众对毛泽东著作的态度和学习方法。只有这样，才能把毛泽东思想学到手，在全党全军全国各族人民中很好地贯彻执行。

二、主要的表是各項工作的量能顯示

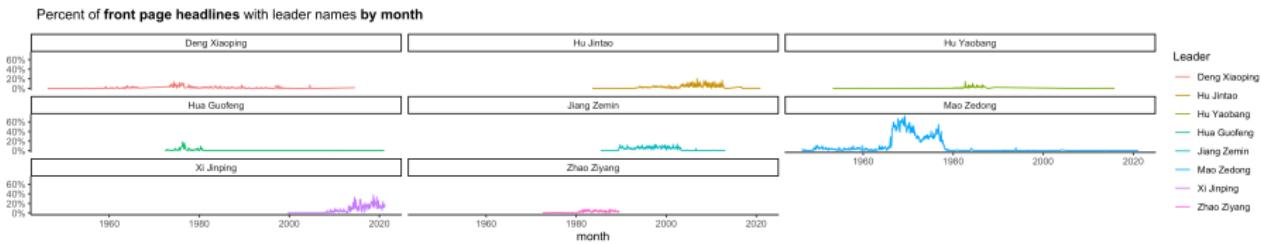
编者按：工农兵群众怀着浓厚的工农阶级感情，把毛主席的书当作各项工作的最高指导。

“千篇万篇有了毛泽东思想就不难”。

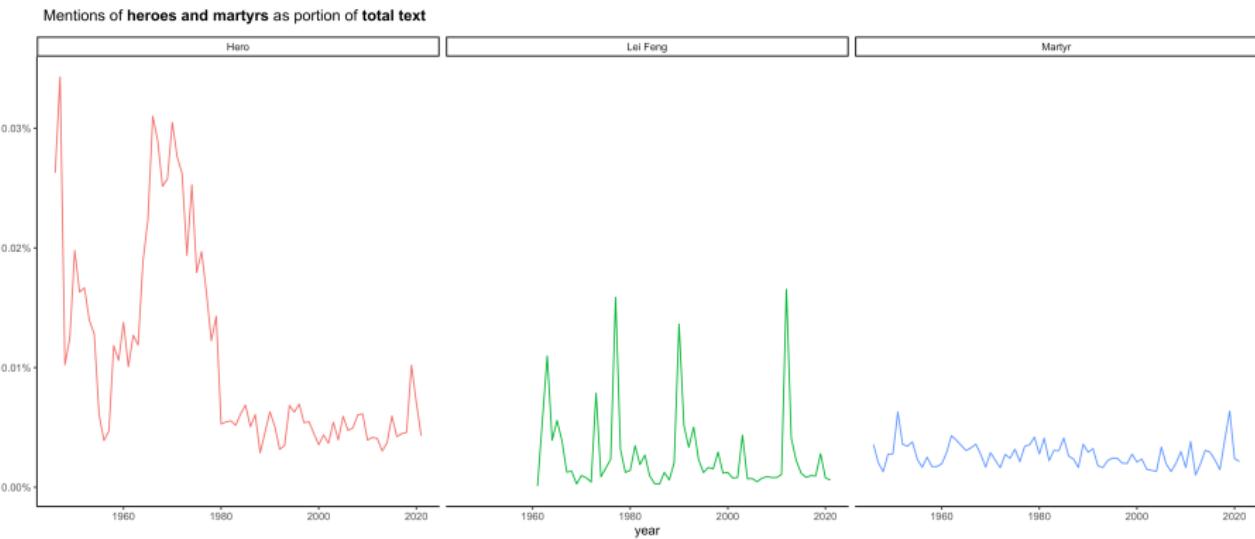
毛泽东思想是我们党最伟大的旗帜。这是中国几十年的革命经验总结了的。是工农兵群众这些生动的、形象的语言，是他们对我国几十年来革命经验的概括。也是他们在当前三大革命斗争的实践中，活学活用毛泽东思想的切身体会。

有志者事竟成 反戈一击
(解放军某部班长熊双全)

Comparing leaders

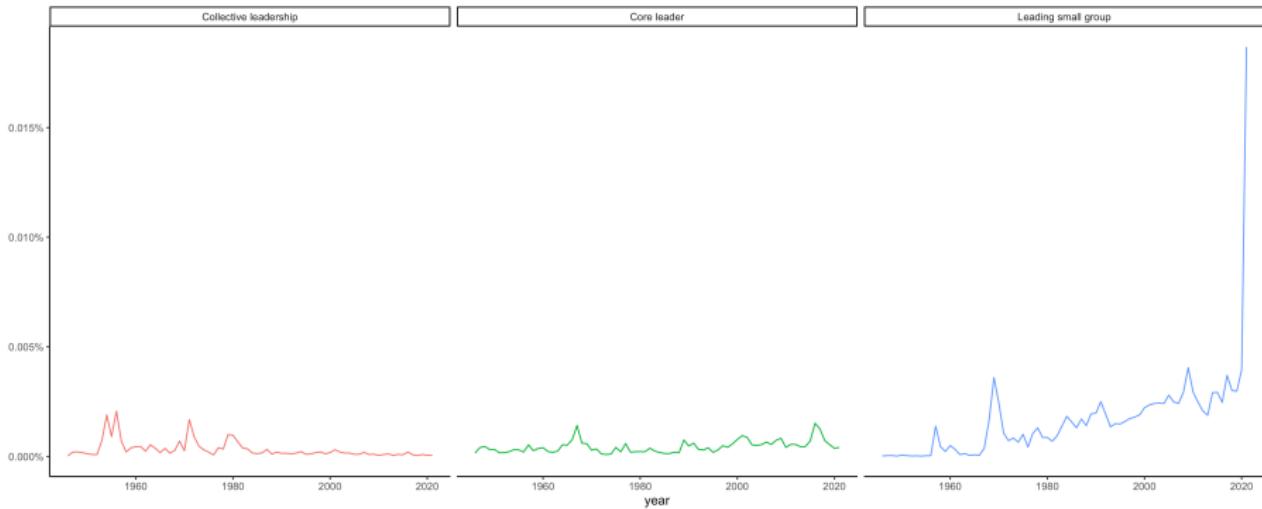


Thematic analyses: heroism, martyrdom, collective leadership I



Thematic analyses: heroism, martyrdom, collective leadership II

Mention of **party leadership structure** as portion of **total text**



Purges

- Nope.
- Data insufficient to say anything useful.

To do

- Tokenize text (jieba/spacy)
- Productionize, update
- More analysis to prove out value of the corpus
- Word embeddings?
- Figure out safe way of making dataset accessible to other researchers
 - maybe just the document feature matrix public?
 - put the SQLite files on a cloud bucket and share the key per request?
 - copyright issues at play

Publication plans

- Not ambitious: Simple ‘dataset’ paper
- Two other exploratory analyses to prove out utility of the data:
 - Policy priorities: Topic models of front page articles (by decade?)
 - Foreign policy: sentiment analysis of discussion of other countries
- Potentially useful dataset as one building block in a large scale comparative analysis of policy in autocratic regimes

References I

- Dong, Dong, Tsan-Kuo Chang, and Dan Chen. 2008. “Reporting AIDS and the Invisible Victims in China: Official Knowledge as News in the People’s Daily, 1986–2002.” *Journal of Health Communication* 13 (4): 357–74. <https://doi.org/10.1080/10810730802063793>.
- Huang, Yixiong. 2018. “Media Representation of Tongxinglian in China: A Case Study of the People’s Daily.” *Journal of Homosexuality* 65 (3): 338–60. <https://doi.org/10.1080/00918369.2017.1317475>.
- Kivimäki, Timo. 2015. “People’s Daily and the Reality of the South China Sea Territorial Disputes.” *Asian Politics & Policy* 7 (2): 319–23. <https://doi.org/10.1111/aspp.12183>.
- Liang, Bin, and Hong Lu. 2013. “Discourses of Drug Problems and Drug Control in China: Reports in the People’s Daily, 1946–2009.” *China Information* 27 (3): 301–26. <https://doi.org/10.1177/0920203x13491387>.

References II

- Wang, Sidan. 2018. "Dynamic Constructed Climate Change Discourses and Discourse Networks Across Newspapers in China Around Three Critical Policy Moments: A Comparative Study of People's Daily, China Daily, and Southern Weekend." PhD thesis, University of Exeter. <https://ore.exeter.ac.uk/repository/handle/10871/33375>.
- Wu, Guoguang. 1994. "Command Communication: The Politics of Editorial Formulation in the People's Daily." *The China Quarterly*, no. 137: 194–211. <http://www.jstor.org/stable/655694>.
- Wu, Shufang. 2014. "The Revival of Confucianism and the CCP's Struggle for Cultural Leadership: A Content Analysis of the People's Daily, 2000–2009." *Journal of Contemporary China* 23 (89): 971–91. <https://doi.org/10.1080/10670564.2014.882624>.
- Yang, Zheng. 2021. "Military Metaphors in Contemporary Chinese Disease Coverage: A Case Study of the People's Daily, 1946–2019." *Chinese Journal of Communication* 14 (3): 259–77. <https://doi.org/10.1080/17544750.2020.1818593>.

References III

- Ye, Wen, and Geri Alumit Zeldes. 2020. “The Representation of People with Disabilities in an Official Newspaper in China: A Longitudinal Study of the People’s Daily from 2003 to 2013.” *Journal of Disability Policy Studies* 31 (1): 26–34.
<https://doi.org/10.1177/1044207319868783>.
- Zhang, Qing Fei. 2014. “Transgender Representation by the People’s Daily Since 1949.” *Sexuality & Culture* 18 (1): 180–95.
<https://doi.org/10.1007/s12119-013-9184-3>.