

# 1 **Tracking the spread of novel coronavirus (2019-nCoV) based on big** 2 **data**

3 Zhao, Xumao<sup>[1]</sup>; Liu, Xiang<sup>[1]</sup>; Li, Xinhai<sup>[2,3]\*</sup>

4 <sup>[1]</sup> State Key Laboratory of Grassland Agro-Ecosystem/Institute of Innovation Ecology,  
5 Lanzhou University, Lanzhou, 730020, China

6 <sup>[2]</sup> Key Laboratory of Animal Ecology and Conservation Biology, Institute of Zoology,  
7 Chinese Academy of Sciences, Beijing, 100101, China

8 <sup>[3]</sup> University of Chinese Academy of Sciences, Yuquan Road, Beijing 100049, China

9 \*Correspondence to: [lixh@ioz.ac.cn](mailto:lixh@ioz.ac.cn) ; [zhaoxm@lzu.edu.cn](mailto:zhaoxm@lzu.edu.cn)

10 **Abstract:** The novel coronavirus (2019-nCoV) appeared in Wuhan in late 2019 have  
11 infected 34,598 people, and killed 723 among them until 8<sup>th</sup> February 2020. The new  
12 virus has spread to at least 316 cities (until 1<sup>st</sup> February 2020) in China. We used the  
13 traffic flow data from Baidu Map, and number of air passengers who left Wuhan from  
14 1<sup>st</sup> January to 26<sup>th</sup> January, to quantify the potential infectious people. We developed  
15 multiple linear models with local population and air passengers as predicted variables  
16 to explain the variance of confirmed cases in every city across China. We found the  
17 contribution of air passengers from Wuhan was decreasing gradually, but the effect of  
18 local population was increasing, indicating the trend of local transmission. However,  
19 the increase of local transmission is slow during the early stage of novel coronavirus,  
20 due to the super strict control measures carried out by government agents and  
21 communities.

22 **Key words:** Baidu Map, big data, quarantine, pandemic, traffic flow

## 23 **Introduction**

24 Since the first case of pneumonia named 2019 Novel Coronavirus (2019-nCoV)  
25 being identified in in Wuhan, China, a total of 34,598 confirmed infections were  
26 reported in the country until 8<sup>th</sup> February 2020, which caused 723 deaths (1). Guo et  
27 al (2020) have found that 2019-nCoVs was similar with bat coronaviruses (2), which  
28 occasionally transmitted to human due to wildlife consumption. Understanding the  
29 pattern of transmission characteristics of 2019-nCoV is important for effective

30 preventing and controlling the ongoing pandemic disease. The value of  $R_0$  (basic  
31 reproduction number) was estimated as 2.2, inferring a median size outbreak (3).  
32 However, based on the epidemic transmission model, the number of actual infections  
33 would be much larger than the number of confirmed cases (4).

34 At present, tracking the passengers from Wuhan in January 2020 is still the top  
35 task for preventing the further spread of novel coronavirus (2019-nCoV). To  
36 accurately estimate the risk of the novel coronavirus, we compiled the detailed daily  
37 traffic data outbound Wuhan from a big-data source, Baidu Map, before the lockdown  
38 of Hubei Province, in order to provide information for risk assessment of 2019-nCoV  
39 at the province level and the city level (Supp. Fig. 1).

## 40 **Methods**

41 The traffic flow data outbound Wuhan from 1<sup>st</sup> January to 26<sup>th</sup> January 2020 was  
42 downed from Baidu Map Huiyan platform (5). The number of air passengers from  
43 Wuhan from 30<sup>th</sup> December 2019 to 20 January 2020 was released by Aviationtalk (6).  
44 The time series data of confirmed 2019-nCoV cases from 10<sup>th</sup> January to 30 January  
45 2020 was obtained from People's daily-Dingxiangyuan (1), which was released by  
46 China National Health Commission.

47 We did Spearman correlation analysis for the daily traffic from Wuhan (from 1<sup>st</sup>  
48 January to 26<sup>th</sup> January) and the total traffic in this period with the number of  
49 confirmed cases (from 25<sup>th</sup> January to 30<sup>th</sup> January). To explain the variance of  
50 confirmed cases in all infected provinces, we developed multiple linear models  
51 including population, GDP, population density, and mean temperature as independent  
52 variables. All analysis was performed using R (version 3.6.2).

## 53 **Results**

54 From 20<sup>th</sup> December 2019 to 20 January 2020, 854,424 air passengers left  
55 Wuhan Tianhe Airport to 49 cities in China (Fig .1). From 1<sup>st</sup> to 26<sup>th</sup> January, about  
56 three million domestic passengers travelled from Wuhan to other cities. Among the  
57 passengers, a few thousands had been confirmed to infected by the novel coronavirus  
58 (Fig. 2). The distribution of air passengers from Wuhan to other cities in China had  
59 high correlation coefficients (0.71) with the number of confirmed infection cases in  
60 those cities on 22<sup>nd</sup> January. The correlation coefficient drops to 0.56 on 24<sup>th</sup> January.  
61 Then the number of confirmed infection cases was positive correlation with local

62 population size.

63 We used a multiple regression model to explain the variance of the number of  
64 cases in the infected cities. After model selection, only two variables remained, the  
65 number of passengers and local population (Fig. 3). Overall, the population of the  
66 provinces explains near half of the variance in the number of confirmed cases across  
67 34 provinces and province-level municipalities, whereas the number of passengers  
68 from Wuhan explained around 10% (Fig. 3).

69 Correlation coefficients of number of confirmed that the number of cases in the  
70 cities ( $n=97$ ) from 25<sup>th</sup> January to 30<sup>th</sup> January match the number of passengers from  
71 Wuhan during the period from 1<sup>st</sup> January to 26<sup>th</sup> January. The highest correlation  
72 appears on 5<sup>th</sup> January, inferring a long incubation period up to two weeks (Supp.  
73 Table 1).

#### 74 **Discussion**

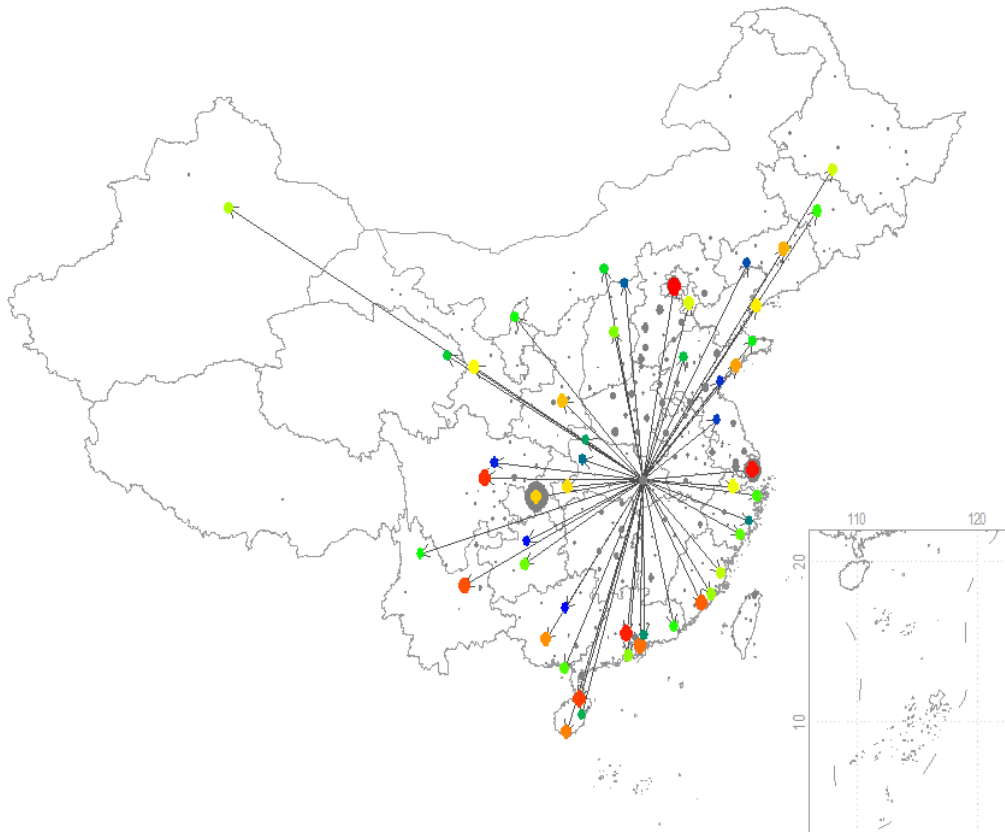
75 In the beginning of the spread of the pneumonia, there is a high correlation (0.71)  
76 between the number of confirmed infection cases and air passengers from Wuhan,  
77 proofed Wuhan the source of the pneumonia (7). As time going, local population  
78 played a more dominant role, because the local spread of 2019-nCoV is likely to  
79 happen. The basic reproductive number of the infection ( $R_0$ ) to be estimated as 3.8  
80 means 72-75% of transmissions must be prevented to stop the outbreak (4).  
81 Fortunately, the transmission of the virus was really controlled due to strict prevention  
82 measures carried out by Chinese government. Restricted population movements ban  
83 was enforced upon 16 cities in Hubei Province since 23<sup>rd</sup> January 2020 (8), resulting  
84 in significant decrease in passengers from Wuhan and adjacent cities, which  
85 effectively reduces the spread of the pneumonia. However, 3-5 million people had left  
86 Wuhan for numerous cities in China before the province lockdown (Supp. Fig.2), and  
87 among them a number of infected people have no clinical symptom yet infectious to  
88 others. We believe this is the highest challenge against the current national level  
89 antiviral campaign.

90 Currently the first-level response to major public health emergencies has been  
91 initiated in 30 provinces, municipalities and autonomous regions in China on 25<sup>th</sup>  
92 January 2020 (9), so that strict control procedures are carried out to prevent the spread  
93 of the virus. According to an infectious disease model, it was estimated that the actual

94 number of infected people in Wuhan on 4<sup>th</sup> February 2020 may reach 250 thousand  
95 (prediction interval, 164,602 to 351,396) without any control (3). Whereas only 5,142  
96 confirmed infections were reported in China on 2<sup>nd</sup> February 2020. Starting from 3<sup>rd</sup>  
97 February, all suspected case in Wuhan will be taken in to medical care, and the virus  
98 spread in Wuhan can be controlled gradually.

99       The correlation coefficients between the number of confirmed cases from 26<sup>th</sup>  
100 January to 30<sup>th</sup> January and number of passengers from Wuhan were higher than that  
101 on 25<sup>th</sup> January (Supp. Table 1). We think the reason is that the confirmed cases on  
102 25<sup>th</sup> is too low due to lack of virus detection kit. After 25<sup>th</sup> January the supply of virus  
103 detection kit was enough and the number of confirmed cases reflect the real situation,  
104 which have very high correlation with the number of passengers from Wuhan to these  
105 cities during 1<sup>st</sup> January to 26<sup>th</sup> January. We notice the highest correlation appear on  
106 5<sup>th</sup> January, two to three weeks ahead of the confirmed cases in those cities  
107 (confirmation also needs several days to complete at that time), which infers the long  
108 latent period of 2019-nCoV.

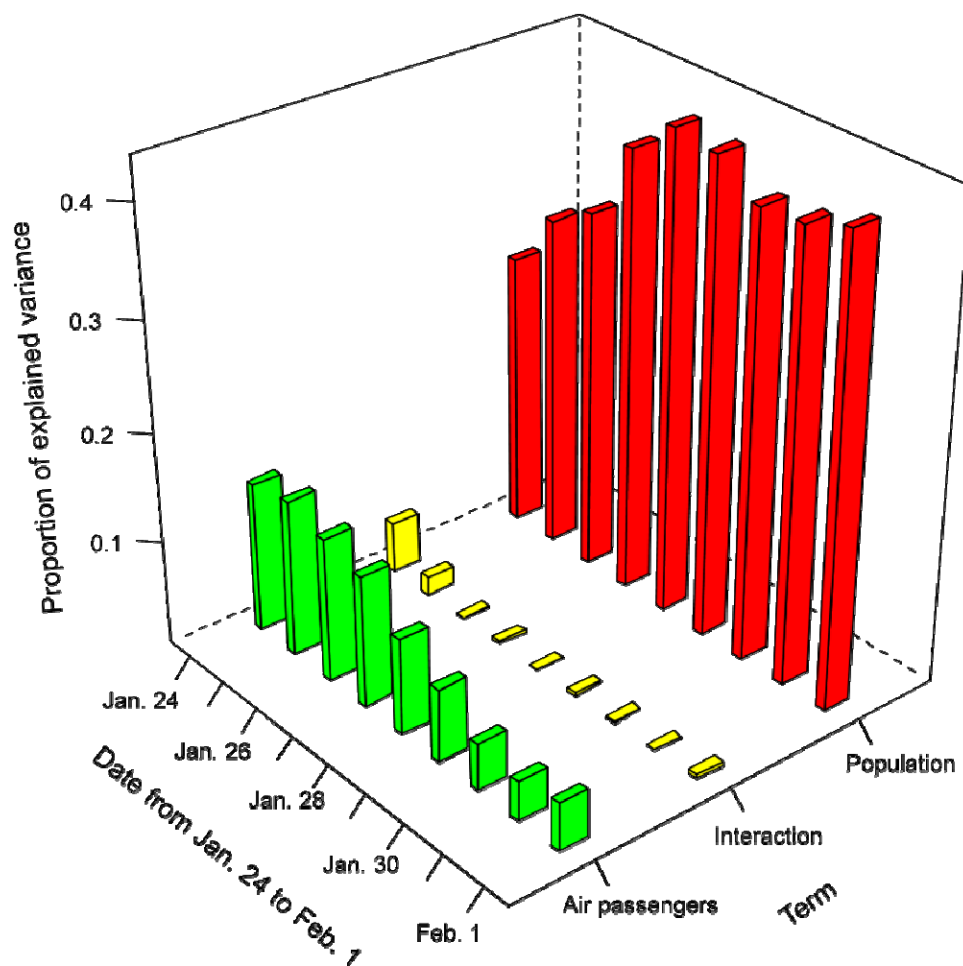
109       Combining the two variables (local population and passengers from Wuhan) to  
110 interpret the virus outbreak risk, we recommend that the preventing and controlling  
111 measures can be divided into two different stages. In the early stage, checking the  
112 passengers from Wuhan is more important. Some cities (e.g., Wenzhou in Zhejiang  
113 Province) with a large number of returnees from Wuhan have many cases even when  
114 they are far from Wuhan in space. The database of travel routes of confirmed patients  
115 has been developed and published for free use (10), in order to tracking and warning  
116 close contactors. In the late stage, local population become the most important factor  
117 in predicting the number of confirmed infection cases. This suggests us that the  
118 densely populated metropolitan areas, such as Shanghai, Beijing, Guangzhou, and  
119 Shenzhen should pay special attentions to preventing the second-generation infections.  
120 Moreover, the densely populated rural areas around Wuhan may face double threats of  
121 the spread of infectious people from Wuhan and a large number of local susceptible  
122 people. These areas, including areas of Hubei Province except for Wuhan, and  
123 surrounding area of Henan (Nanyang, Xinyang) and Chongqing, need to prepare for a  
124 surge in infection. The rural medical facilities in these areas are scarcer than in cities,  
125 clustering cases are more likely to happen in these places.



126 Fig. 1 The distribution of 854,424 air passengers from Wuhan Tianhe Airport to 49 cities in China  
127 during the period from 30<sup>th</sup> December 2019 to 20<sup>th</sup> January 2020. The sizes of the points indicate  
128 the number of passengers, which are also represented in color series: blue, green, yellow, brown,  
129 and red.



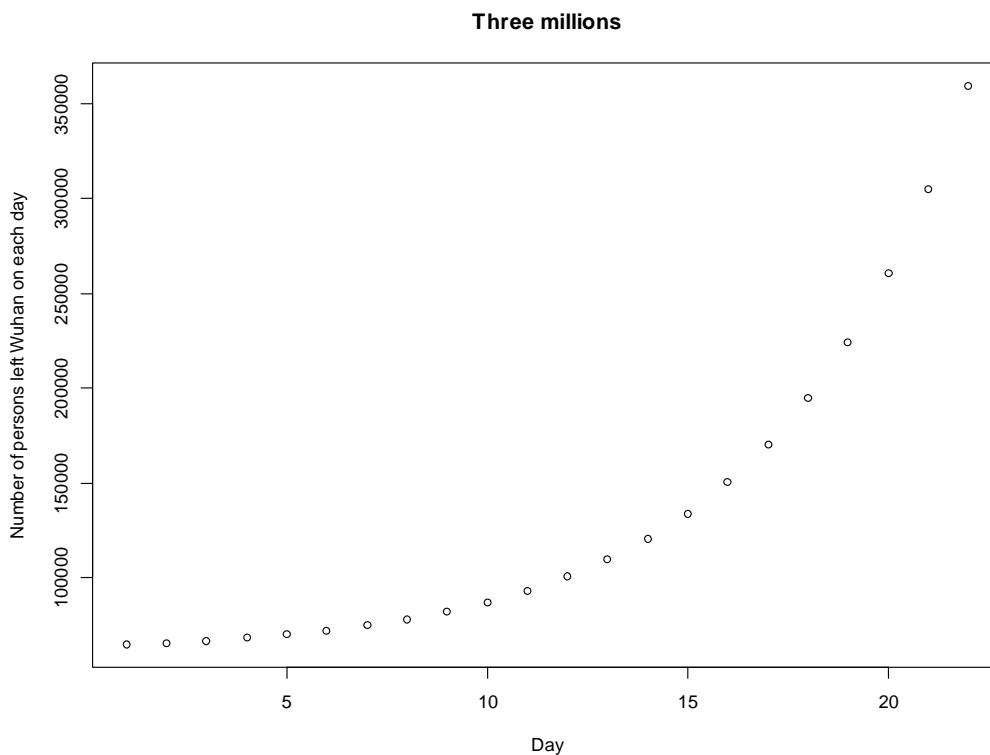
130 Fig. 2 The spatial distribution of 24,281 confirmed cases (incomplete data) of 2019-nCoV on 5<sup>th</sup>  
131 February in 299 cities in China. The sizes of red points represent the number of confirmed cases.  
132 The sizes of grey points show the population of the cities.



133 Fig. 3 The proportion of explained variance in linear regression: Number of confirmed cases in 34  
134 provinces ~ Air passengers from Wuhan \* Population of the province for the date from 24<sup>th</sup>  
135 January to 1<sup>st</sup> February.

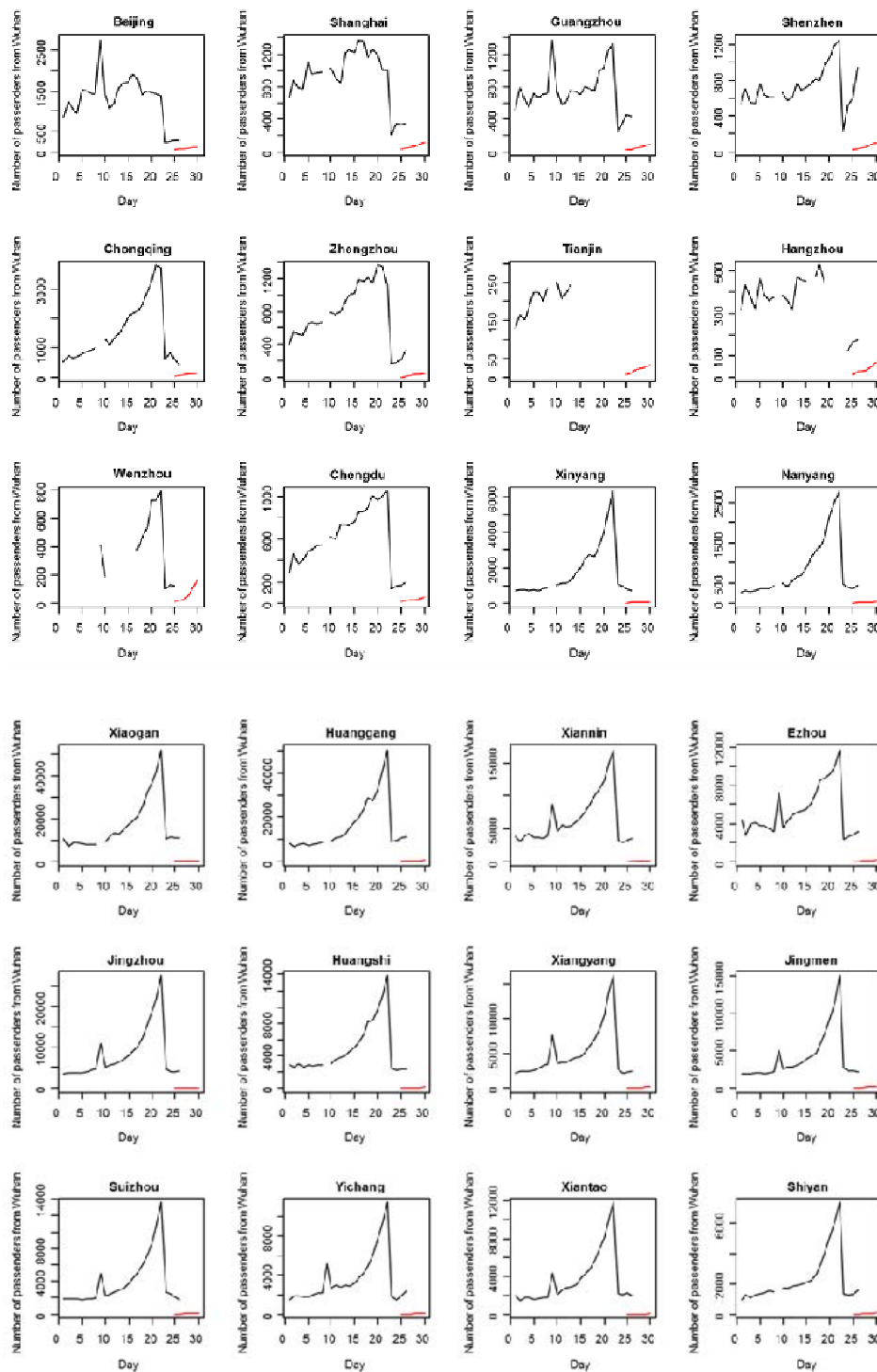
136 **Supplementary documents**

137 Supp. Fig. 1 The daily passenger (include all vehicles) numbers from Wuhan to top 50 cities in  
138 China from 1st January to 26th January are showed by an animated GIF file  
139 2019-nCoV\_spread\_1-26 Jan.gif (Supp. Fig. 1).



140 Supp. Fig. 2 Estimated daily number of passengers from Wuhan using a logistic curve  $y = 60000 +$   
141  $10000/\exp(1-0.2*x)$  from 1<sup>st</sup> January to 22<sup>nd</sup> January (the day before city lockdown). The  
142 assumption is that people tend to leave Wuhan just before the Chinese New Year on 25<sup>th</sup> January.





143 Supp. Fig. 3 Estimated daily number of passengers from Wuhan to 12 major cities in China (upper  
 144 panel) and 12 major cities in Hubei Province (lower panel) base on Baidu Bigdata server. The red  
 145 lines show the daily accumulated confirmed case in each city from 25<sup>th</sup> January to 30<sup>th</sup> January.

146 **References**

- 147 (1) Dingxiangyuan:[https://ncov.dxy.cn/ncovh5/view/pneumonia\\_peopleapp?scene=126&clicktime=1579832412&from=timeline&isappinstalled=0](https://ncov.dxy.cn/ncovh5/view/pneumonia_peopleapp?scene=126&clicktime=1579832412&from=timeline&isappinstalled=0)
- 148
- 149 (2) Guo Q, Li M, Wang C, Wang P, Fang Z, Tian J, Wu S, Xiao Y, Zhu H. Host and infectivity
- 150 prediction of Wuhan 2019 novel coronavirus using deep learning algorithm. 2020, bioRxiv.
- 151 <http://dx.doi.org/10.1101/2020.01.21.914044>
- 152 (3) Riou J, Althaus CL. Pattern of early human-to-human transmission of Wuhan 2019-nCoV.
- 153 2020, bioRxiv. <http://dx.doi.org/10.1101/2020.01.23.917351>
- 154 (4) Read JM, Bridgen JRE, Cummings DAT, Ho A, Jewell C. Novel coronavirus 2019-nCoV:
- 155 early estimation of epidemiological parameters and epidemic predictions. 2020, medRxiv.
- 156 (5) Baidu Map Huiyan platform: <https://qianxi.baidu.com/?city=420100> (in Chinese)
- 157 (6) Aviationtalk : <https://mp.weixin.qq.com/s/7ynWYxB-s7nfz7rmjBLSpQ> (in Chinese)
- 158 (7) Zhu N, Zhang D, Wang W, Li X, Yang B, Song J, Zhao X, Huang B, Shi W, Lu R, Niu P, Zhan
- 159 F. A novel coronavirus from patients with pneumonia in China, 2019. 2020, New England
- 160 Journal of Medicine. DOI: 10.1056/NEJMoa2001017.
- 161 (8) China news : <http://www.chinanews.com/gn/2020/01-25/9069668.shtml> (in Chinese)
- 162 (9) Xinhua net: [http://www.xinhuanet.com/2020-01/26/c\\_1125503530.htm](http://www.xinhuanet.com/2020-01/26/c_1125503530.htm) (in Chinese)
- 163 (10) Close contactor detection tool:
- 164 <http://2019ncov.nosugartech.com/?from=timeline&isappinstalled=0> (in Chinese)

Jan. 1

