

1 **A data-driven model for predicting the course of COVID-19 epidemic with**
2 **applications for China, Korea, Italy, Germany, Spain, UK and USA**

3
4 (Revised: April 5, 2020)

5
6 **Authors:**

7 Norden E. Huang#, Fangli Qiao# and Ka-Kit Tung*

8
9 **Affiliations:**

10 # *Data Analysis Laboratory, FIO, Qingdao 266061, China*

11 * *Department of Applied Mathematics, University of Washington, Seattle, WA 98195*

12 # Co-first authors: These authors contribute equally to this work.

13 FQ: qiaofl@fio.org.cn. NEH: norden@ncu.edu.tw

14 *Corresponding author: ktung@uw.edu

15
16 **KEYWORDS**

17 Covid-19; Epidemiology; Covid-19 predictions for Italy, Germany, Spain, UK and
18 USA; Data-driven approach; prediction of turning points; peak active infected
19 number; Covid-19 recovery period; end of Covid-19 epidemic; local-in-time metric
20 for epidemic management.

23 **ABSTRACT**

24

25 **For an emergent disease, such as Covid-19, with no past epidemiological data**
26 **to guide models, modelers struggle to make predictions of the course of the**
27 **epidemic (Cyranoski, *Nature News* 18 February 2020). The wildly varying**
28 **predictions make it difficult to base policy decisions on. On the other hand**
29 **much empirical information is already contained in data of evolving**
30 **epidemiological profiles. We offer an additional tool, based on general**
31 **theoretical principles and validated with data, for tracking the turning points,**
32 **peak and accumulated case numbers of infected and recovered for an**
33 **epidemic, and to predict its course. Ability to predict the turning points and**
34 **the epidemic's end is of crucial importance for fighting the epidemic and**
35 **planning for a return to normalcy. The accuracy of the prediction of the peaks**
36 **of the epidemic is validated using data in different regions in China showing**
37 **the effects of different levels of quarantine. The validated tool can be applied**
38 **to other countries where Covid-19 has spread, and generally to future**
39 **epidemics. US is found to have the largest net infection rate, and is predicted**
40 **to have the largest total infected cases (708K) and will take two weeks longer**
41 **than Wuhan to reach its turning point, and one week longer than Italy and**
42 **Germany.**

43

44 **SIGNIFICANCE:** We offer a practical tool for tracking and predicting the course of an
45 epidemic using the daily data on the infection and recovery. This data-driven tool
46 can predict the turning points two weeks in advance, with an accuracy of 2-3 days,
47 validated using data from various regions in China selected to show the effects of
48 quarantine. It also gives information on how rapid the rise and fall of the case
49 numbers are, and what the peak and total number of infected are. Although
50 empirical, this approach has a sound theoretical foundation; the main components
51 of the results are validated after the epidemic is near an end, as is the case for China,
52 and therefore is generally applicable to future epidemics.

53

54 **1. Introduction.**

55 The current COVID-19 epidemic is caused by a novel corona virus, designated
56 officially as SARS-CoV-2, spreading from Wuhan, the capital city of Hubei province in
57 China (2-4). The new virus seems to have characteristics different from SARS
58 (severe acute respiratory syndrome) (5, 6): it is less deadly but spreads more widely
59 (7-10). Modeling the epidemic as it develops has been difficult (1). Depending on
60 the model assumptions, predictions of when it “turns a corner” for China varies
61 greatly (11-21), up to after 650 million people have been infected before peaking;
62 many have now been shown to be inaccurate (22). Now as the epidemic has
63 subsided in China and become a global pandemic (23, 24), a reliable forecast of the
64 course of the outbreak in each region is critical for the management and
65 containment of the epidemic, and for balancing the impact from the public health
66 crisis vs the economic crisis. China has instituted some of the strictest quarantine
67 measures around Wuhan and Hubei, which may or may not be adoptable in other
68 countries (25-27). It would be useful to extract the dependence of the epidemic’s
69 evolution on the degree of quarantine to guide policy decisions, while also to
70 characterize properties of Covid-19 that are applicable to other countries.

71

72 Mainstream epidemiological models have their origin in the SIR (Susceptible,
73 Infected, and Recovered or Removed) model (28) and its many variations. We
74 explain in the Supplementary Information why existing model predictions vary so
75 widely by commenting on the assumptions underlying these models.

76

77 These SIR-type models, however, serves a critical purpose for long-range policy
78 planning, such as warning policy-decision makers of the gravity of the potential
79 impact and prompting them to take proper actions before it is too late. After the
80 breakout, more information is needed for more detailed planning, such as the
81 arrival of the critical turning points, the number of hospital bed we might need at
82 the peak, and the estimate for when to lift the quarantine, and when to return to
83 normalcy. We offer here an additional tool that has the advantage that it has does
84 not depend on the elusive infection rate or the susceptible population, information
85 needed for most models, but has the disadvantage that it cannot be used when the
86 epidemic first started and the data are inaccurate or incomplete. It is based on daily
87 case numbers (i.e. newly confirmed cases), $N(t)$, and recovered cases, $R(t)$.

88

89 Without universal testing, the confirmed case number might be only a subset of the
90 true total infected number, which may never be known unless frequent universal
91 testing is instituted. The asymptomatic infected who are not tested and then
92 recover on their own do not get counted but they also do not tax hospital resources.
93 Nevertheless, they can infect others and some of the latter may develop more
94 serious symptoms that require hospitalization. Then these secondary infections are
95 included in our case data. Our aim is to provide a tool that can be used for the
96 management of medical resources. Since we do not use a model to calculate how the
97 asymptomatic infectives infect others, we do not need to know either the infection
98 rate or the asymptomatic infective numbers. Since those who are admitted to the

99 hospital either recover after a hospital stay of T days, or dead after a similar number
100 of days, there should be a delayed relationship between $N(t)$ and $R(t)$, which we will
101 explore in the Theory section. Now that the epidemic in China appears to have come
102 to an end, the data from various regions in China can be used to validate the model.
103 After validation we then apply it to other regions in the world.

104
105 Our estimate of the end date of the epidemic is not based on the number of
106 susceptibles, S , approaching zero as in most models (i.e. most of the population is
107 infected, hence acquiring immunity), but $N(t)$ approaching zero and remaining so
108 for two incubation periods. The first incubation period is to allow the asymptomatic
109 infected to show symptoms and the second period to allow those that are infected
110 by the asymptomatic infected to show symptoms. For prediction purpose, the date
111 when the $N(t)$ is zero is estimated by 3 standard deviations from its peak. These
112 two quantities can be extracted from the data as the epidemic is developing. Our
113 estimate of the end of the epidemic is earlier than most model predictions, usually
114 significantly so, because it does not depend on the herd immunity concept.

115
116 As is true for all data-driven approaches, our result inevitably depends on the
117 quality of the data used, and some of the early data of the epidemic are not as good
118 as the later data, when better diagnostic methods and more complete reporting are
119 established. However, many of the metrics commonly in use require accumulation
120 of data from the beginning of the epidemic, and consequently are affected by poor
121 data or change of diagnostic methods along the way. We try to avoid accumulation
122 and use local-in-time metrics. Nevertheless, data problems cannot be avoided.
123 Sensitivity of our conclusion on data problems is extensively discussed in this work.
124 Figure S1 displays examples of data used in this study. One problem immediately
125 becomes obvious for the Chinese data: On 12 February, when Hubei changed its
126 definition of confirmed infection from the gold standard of nucleic acid gene-
127 sequencing tests to clinical observations and radiological chest scans, over 14,000
128 newly infected cases were added that day, creating a peak that has not been
129 exceeded since. Overwhelmed doctors in Wuhan pleaded for the change so that
130 they did not have to wait for the returned tests to confirm the infection. Outside
131 Hubei, there was no change in definition for the “infected”. How this artifact affects
132 our conclusion will be discussed.

133

134 **2. Model and its validation using data**

135 **Definition:** Let $I(t)$ be the number of active infected at time t . Its change is given
136 by;

$$137 \quad \frac{d}{dt}I = N(t) - R(t),$$

138 where $N(t)$ is the number of newly infected, and $R(t)$, designated as removed, is
139 the sum of the daily recovered and dead. For a disease such as Covid-19 with low
140 fatality rate, $R(t)$ consists of mainly recovered. However even for this disease, the
141 fatality rate in some regions, such as in Northern Italy, approaches 10%. For these
142 regions $R(t)$ should include the dead as well. Note that for the theory part, $N(t)$

143 includes both confirmed and unconfirmed cases. The term: Existing Infected Case
144 (*EIC*) number is used to denote the confirmed $I(t)$ when we deal with data.

145

146 Let t_p , the turning point defined as the peak of the active infected number. At this
147 point maximum medical resource is needed. This maximum occurs when

148 $\frac{d}{dt}I = 0$, implying $N(t_p) = R(t_p)$.

149 This is a local-in-time metric. There is therefore no need to first find $I(t)$ to locate
150 this peak. After the turning point, the newly recovered starts to exceed the newly
151 infected. The demand for medical resources, such as hospital beds, isolation wards
152 and respirators, starts to decrease.

153

154 The theoretical foundation for our model is given in Supplementary Information.
155 Here we discuss the main results and offer validation of these results using data
156 from China.

157

158 **Main Result:** The daily newly recovered/removed number $R(t)$, is related to the
159 daily newly infected number $N(t)$ as:

160

$$161 \quad R(t) = N(t - T),$$

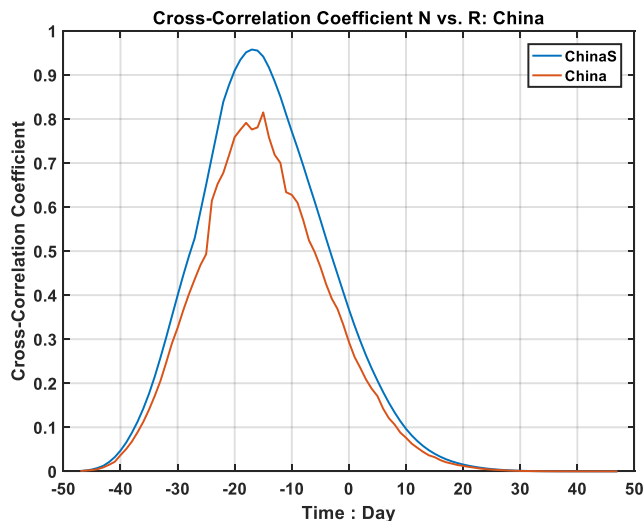
162 for $t > T$, where T is the mean recovery period.

163 This result can be rigorously derived using an age-structured population model (see
164 reference (36)). It is also common sense: the infected eventually recovered after a
165 number of days, or dead after a similar number of days. Of course, the number of
166 days a patient stays in the hospital before discharge depends on the efficacy of
167 treatment and so varies somewhat, and the time it takes for a patient to die may also
168 depend on the age and underlying conditions. T is therefore a statistical quantity.

169 **Validation:** This fundamental relationship can be validated statistically with data.

170 Figures 1, obtained using data from China during the Covid-19 epidemic, shows that
171 $N(t)$ and $R(t)$ are highly correlated: with correlation coefficient of 0.95 when both
172 distributions are smoothed with 5-point boxcar. The unsmoothed daily data also
173 yield a high correlation coefficient of 0.80, with $R(t)$ lagging $N(t)$ by $T \sim 15$ days. Both
174 correlation coefficients are statistically significant. A similar result is found for
175 Hubei (Figure S2) and other regions (not shown). This is one of the ways the *mean*
176 *recovery period* is determined statistically from data, but it is not practical in the
177 early phase of the epidemic. We will give different methods for the latter purpose.
178 The result on T is consistent with that estimated or predicted later using the slope of
179 the distribution in Figure 4. The latter, obtained by the intercept of the straight line,
180 is less accurate because of the slope is rather shallow.

181



182

183

Figure 1. Lagged correlation of case numbers $R(t)$ and $N(t)$ for China as a whole.

184

185

Main Result: The natural logarithm of the ratio of N and R is a linear function of time for $t > T$. This relationship is important for the purpose of forecast because it is easy to extrapolate from a straight line into the future.

186

187

188

189

Validation of log NR as a straight line:

190

From data we use the report newly confirmed case number and the recovered case number to define NR ratio as

191

$$NR(t) = N(t)/R(t).$$

192

At t_p , $NR=1$.

193

194

195

We show in Figure 2, using the data of the epidemic for COVID-19, that the logarithm of $NR(t)$ lies on a straight line, with small scatter, passing through the turning point t_p . And data for various stages of the epidemic, from the initial exponential growth stage, to near the peak of EIC , and then past the peak, all lie on the same straight line. The intercept with $\log NR=0$ yields the turning point. This line, obtained by linear-least-square fit in the semi-log plot, is little affected by the rather large artificial spike in the data on 12 February because of its short duration and the logarithmic value. That reporting problem is necessarily of short duration because, on the date of definition change, previous week's cases of infected according to the new criteria were reported in one day. After that, the book is cleared, and $N(t)$ returned to its normal range.

200

201

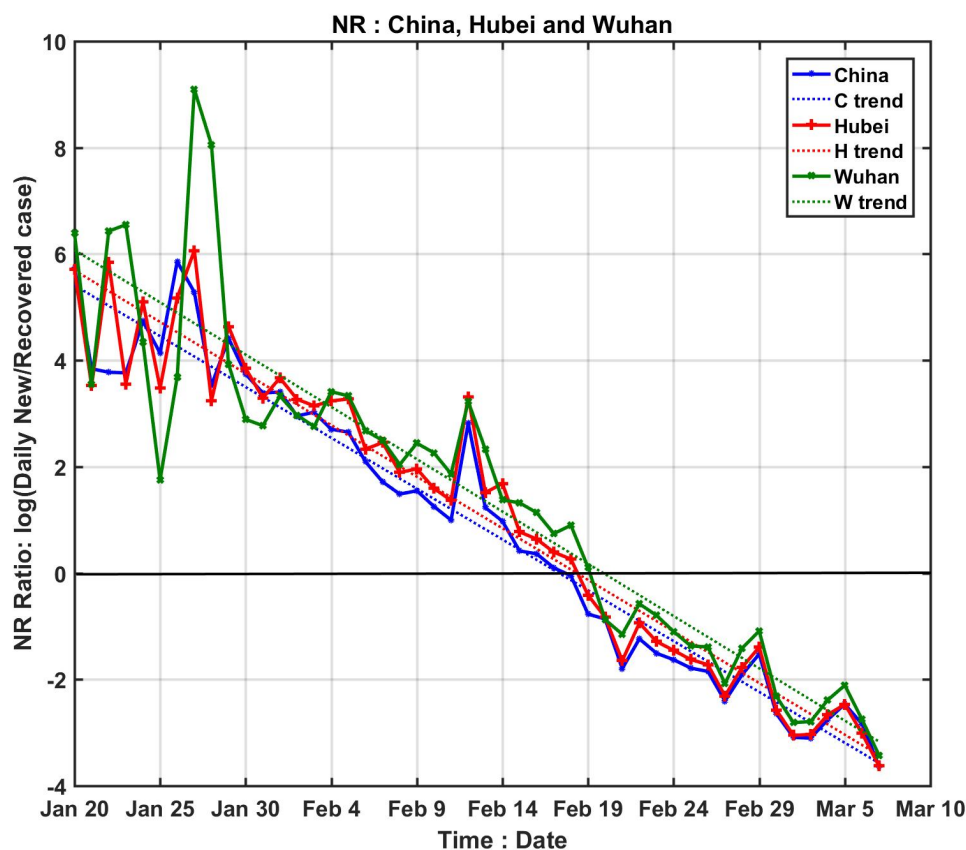
202

203

204

205

206



207
208

209 **Figure 2.** Logarithm of the ratio of daily newly infected to newly recovered. They lie
210 on straight lines with some small scatter. The dotted straight lines are obtained by
211 linear-least squares fit is. The slopes of the lines are almost the same but with
212 different intercept; the trend lines cross zero (the black solid line) at different time
213 for different regions indicating different peaking time for *EIC*. The epicenter Wuhan
214 (green) has latest turning point than its province Hubei (pink), which has a later
215 turning point than China as a whole (cyan).

216

217 The theoretical result in SI suggests that the slope of the linear line is $-T/\sigma_R^2$, where
218 σ_R is the standard deviation of the $R(t)$ profile. In general, the slope can be different for
219 different regions with different levels of quarantine and epidemic characteristics.
220 The hospital treatment efficacy would influence T directly. The effect of quarantine
221 would influence the value of σ_N , the standard deviation of the newly infected, and so
222 indirectly $R(t)$ and σ_R . Our empirical result from Fig. 2 however shows that the slope
223 is the almost the same for different regions in China, implying that efficacy of
224 treatment and level of quarantine affect T and σ^2 proportionally.

225

226 **Result:** The derivative of $\log N(t)$ and of $\log R(t)$ is each a linear function of time,
227 with known slope. Their intercept with the zero derivative line yields the time for
228 their respective peak.

229

230 **Validation**

231 Empirically, the derivative of $\log N(t)$ or $\log R(t)$ lies on a straight line, as shown in
232 Fig. 3 (although the scatter is larger as to be expected for any differentiation of
233 empirical data). The positive and negative outliers one day before and after 12 Feb
234 are caused by the spike up and then down, with little effect on the fitted linear trend
235 (but increases its variance and therefore uncertainty). Moreover, the straight line
236 extends without appreciable change in slope beyond the peak of $N(t)$, suggesting
237 that the distribution of the newly infected number is approximately Gaussian. The
238 mean recovery time T can be predicted as $t_R - t_N$, where t_R is the peak of $R(t)$ and t_N
239 is the peak of $N(t)$. These two peak times can be obtained by extending the straight
240 line in Fig. 3 to intersect the zero line. This predicted result can be verified
241 statistically after the fact by the lagged correlation of $R(t)$ and $N(t)$. If the
242 distribution is indeed Gaussian or even approximately so, the slope in Fig. 3 would
243 be proportional to the reciprocal of the square of its standard deviation, σ , as (See
244 SI):

$$245 \quad \frac{d \log N(t)}{dt} = \frac{-(t - t_N)}{\sigma_N^2}.$$

246

247 Similarly result holds for the daily number of recovered, $R(t)$.

248

249 After the epidemic is nearing the end as is the case in China, fitting the data to a
250 Gaussian can be done after the fact (see Figures S3 and S4). The fit is satisfactory
251 even without using any disposal parameters. The parameters used are determined
252 using slopes of $\log N$ and $\log R$ (see Table S1)

253

254 The inferred statistical characteristics of the Covid-19 epidemic are summarized in
255 Table S1 for various regions. The mean recovery time T , is about 13 days for China
256 as a whole. For Wuhan, the city at the epicenter whose hospitals were more
257 overwhelmed and the patients admitted into hospitals more seriously ill than those
258 in other provinces, $T \sim 16$ days, while that for Hubei is 14 days. The standard
259 deviation, σ , is found to be around 8 days, with slight difference between that for
260 $N(t)$ and for $R(t)$, with one exception for Hubei outside Wuhan. Such a fine
261 subdivision may not be practical for the data quality we have. The σ tends to be
262 smaller for China as a whole than Wuhan. One can see that T and σ^2 indeed varying
263 approximately in proportion.

264

265 **The peak infected cases:**

266 Writing $N(t_N) = N(t_B) \exp\left\{\int_{t_B}^{t_N} n(t) dt\right\}$, and noting $n(t) = \frac{d}{dt} \log N(t) = -\frac{(t - t_N)}{\sigma_N^2}$, the

267 exponent is $\int_{t_B}^{t_N} n(t) dt = \frac{(t_N - t_B)^2}{2\sigma_N^2} = \frac{1}{2}n(t_B)(t_N - t_B)$. Hence the peak infected case number

268 can be predicted, using the predicted value for t_N starting from a conveniently chosen

269 time t_B , such as the latest time with data available, as:

$$270 \quad N(t_N) = N(t_B) \exp\left\{\frac{n(t_B)(t_N - t_B)}{2}\right\}.$$

271

272 **Accumulated quantities.**

273 To calculate $I(t)$ using reported data, only confirmed cases are used (We call it EIC).

274 It is given by the accumulated newly confirmed cases minus the accumulated

275 confirmed recovered. Since the accumulation of early poor data can introduce

276 errors a more local-in-time formula is given as:

277

$$278 \quad I(t) = \int_{-\infty}^t N(t) dt - \int_{-\infty}^t R(t) dt = \int_{-\infty}^t N(t) dt - \int_{-\infty}^t N(t-T) dt \\ = \int_{t-T}^t N(t) dt.$$

279 That is, to find I at time t , one only needs to add up the daily newly infected case

280 numbers for a period of T preceding t . This is an almost local-in-time property even

281 for this accumulated quantity. For validation, we estimate the peak of the I case

282 number on 18 February by computing the sum of daily newly infected case numbers

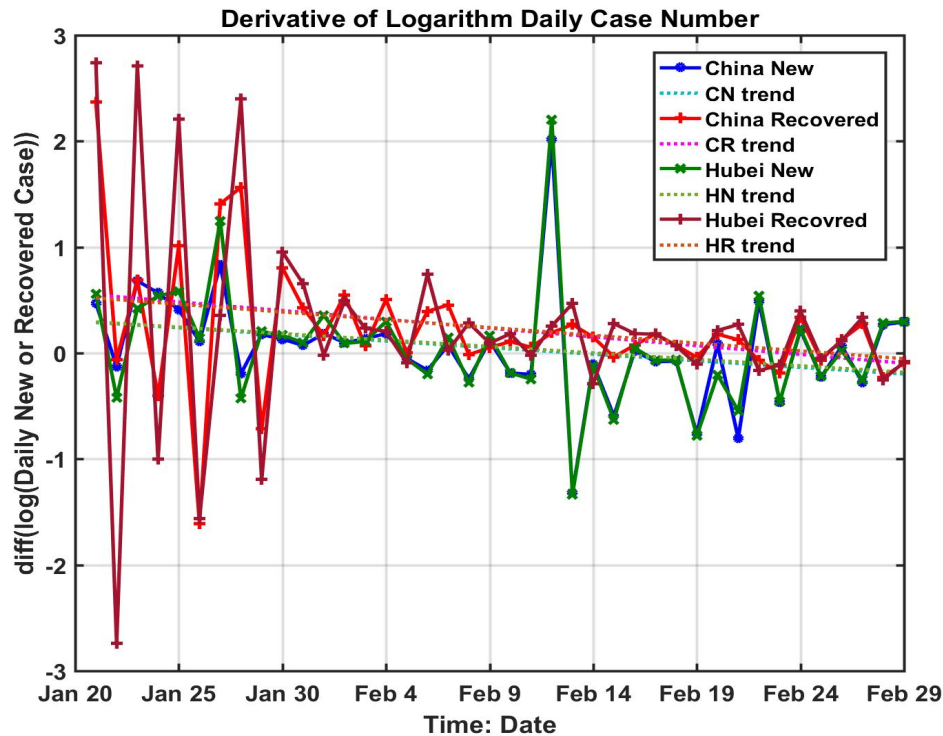
283 for 15 days, from February 4 to February 18, which yields a peak value for the total

284 infected cases on 18 February of 54,747. This is within 10% of the reported number

285 of 57,805, even after taking into account the deaths (by subtracting the

286 accumulated deaths of 2,004 from our estimate).

287



288
289
290
291
292
293
294
295

Figure 3 The derivative of the logarithm of daily newly infected or recovered. Notice the clear separation of the new and recovered cases and also the subtle difference of their slopes. The zero crossings of the trend line give the peak dates of the new and recovered case respectively. And the slopes give an estimate of σ values. In this Figure, the following abbreviations are used: C=China; H=Hubei; N=New Case; R=Recovered.

296
297
298
299
300
301
302
303
304
305
306

3. Predictability

Prediction of the turning point using *NR* ratio.

We first discuss how the true turning point can be determined from data after it has occurred. Then we give a method for predicting this true value in advance and assess the accuracy of the forecasts as function of days in advance when the prediction is made. A note: after this manuscript was submitted for review the predictions that we made previously have come to pass. Although consequently the value of our predictions has greatly diminished, it gives us a chance to compare our predictions against the truths. This model validation process is important if we are to apply the same method for prediction to other regions.

307
308
309
310
311
312
313

The turning point and the end of the epidemic are the two most watched markers on its development (28, 29), along with the number of infected at each stage of the epidemic. There are various definitions of the turning point. A common one defines the turning point of the epidemic as the reported daily number of newly infected reaching a peak and then declining. This is the one touted in the various news announcements, and also used by some research groups (22). The fact that the number of newly infected reaching a peak and then declining does not necessarily

314 imply that the epidemic has “turned a corner”, because the total number of active
315 infected can still be rising with the associated urgent need for additional medical
316 resources, such as hospital beds, isolation wards and ventilators. Furthermore,
317 locating this peak is highly susceptible to data glitches and change in diagnostic
318 definition. A more meaningful turning point should be based on the number of
319 confirmed infected individuals, designated as EIC (15), reaching a peak and then
320 starting to decline. EIC is in theory obtainable from data of the daily number of new
321 confirmed cases, $N(t)$, and the daily number of newly recovered, $R(t)$, by subtracting
322 the accumulated sum of $R(t)$ from the accumulated sum of $N(t)$. Analysis of this
323 accumulated quantity is sensitively affected by accumulation of poorer early data of
324 reported cases, including under-reporting and under-detection of the number of
325 infected caused by insufficient test kits, in addition to the history of changing
326 diagnostic criteria. Moreover in practice its peak is often not detected until several
327 weeks after it has occurred.

328
329 Since the maximum of EIC , can be located by the zero of its derivative, we propose
330 using a local-in-time metric of $N(t_p)=R(t_p)$ at the peak of EIC , t_p .

331
332 Referring to Figure S1, for China as whole, t_p is found to be February 18; for Hubei,
333 the province of the epicenter Wuhan, t_p is found to be 19 February, and for China
334 outside Hubei (China exHubei), 12 February, coincidentally on the same day as the
335 Hubei data spike. However there is no such bump in the data outside Hubei, and so
336 is not likely the result of the data artifact. These results, even including that for
337 Hubei, are not affected by the historical data problems because of our local-in-time
338 method for determining the turning point.

339
340 Can such a turning point be predicted before it happened, and if so by how many
341 days in advance? Since the logarithm of NR lies on a straight line passing through
342 the turning point of EIC , it would be interesting to explore if the turning point can be
343 predicted by extrapolation using data weeks before it happened by extrapolation
344 along the straight line (see Figure S5). How far in advance this can be done appears
345 to be limited by the poor quality of the initial data. Fig. 4 shows the results of such
346 predictions. The horizontal axis indicates the last date of the data used in the
347 prediction. The beginning date of the data used is 24 January for all experiments.
348 Prior to that day, data quality was poor and the newly recovered number was zero
349 in some days, giving an infinite NR ratio.

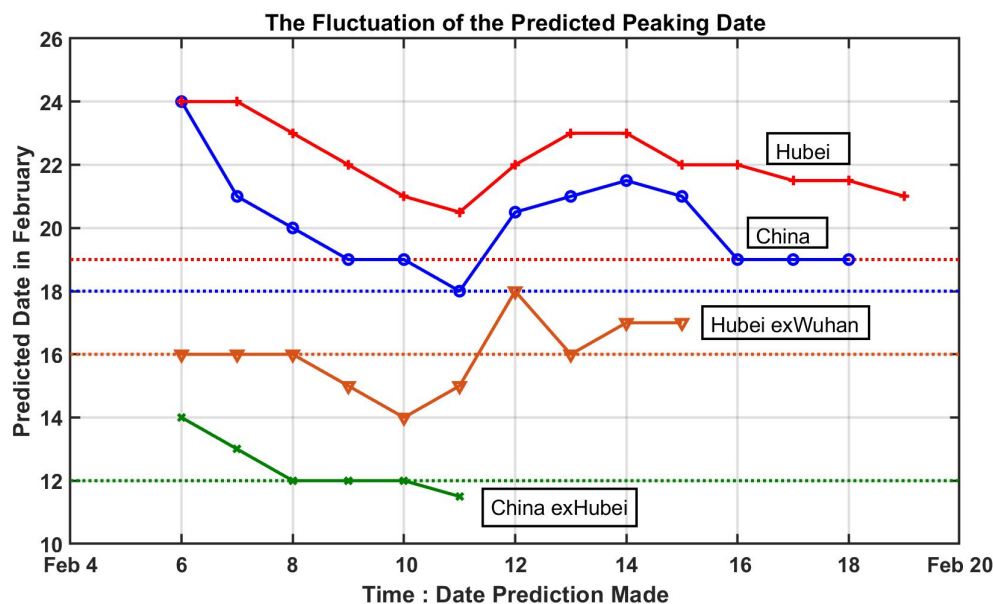
350
351 For China outside Hubei, the prediction made on 6 February gives the turning point
352 as 14 February, two days later than the truth. A prediction made on 8 February
353 already converged to the truth of 12 February, and stays near the truth, differing by
354 no more than fractions of a day with more data.

355
356 The huge data glitch on 12 February in Hubei affected the prediction for Hubei, for
357 China as whole, and for Hubei-exWuhan. These three curves all show a bump up
358 starting 12 February, as the slope of $N(t)$ is artificially lifted. Ironically, predictions
359 made earlier than 12 February are actually better. For example, for China as a whole,

360 predictions made on 9 February and 10 February both give 19 February as the
 361 turning point, only one day off the truth of 18 February. A prediction made on 11
 362 February actually gives the correct turning point that would occur one week later.
 363 At the time these predictions are made, the newly infected cases were rising rapidly,
 364 by over 2,000 each day, and later by over 14,000. It would have been incredulous if
 365 one were to announce at that time that the epidemic would turn the corner a week
 366 later.

367
 368 Even with the huge spike for the regions affected by the Hubei's changing of
 369 diagnosis criteria, because of its short duration the artifact affects the predicted
 370 value by no more than 3 days, and the prediction accuracy soon recovers for China
 371 as a whole. For Hubei, the prediction never converges to the true value, but the
 372 over-prediction is only 2 days. This smallness of the error is remarkable given that
 373 other model predictions differ by weeks or months.

374
 375 Table S2 lists the mean and standard deviation of the predictions. For applications
 376 to other countries and to future epidemics without a change in the definition of the
 377 "infection" to such a large extent, we expect even better prediction accuracy.
 378
 379



380
 381
 382 **Figure 4** Prediction of the turning point in *EIC* by extrapolating the trend in
 383 logarithm of *NR*. The horizontal axis indicates the date the prediction is made using
 384 data prior to that date. The vertical axis gives the dates of the predicted turning
 385 point. Dashed horizontal lines indicated the true dates for the turning point, as
 386 determined from Fig. S1.
 387

388 **Estimate of “all clear” declaration**

389 We can now estimate a time for a declaration of “all clear”. No verification is yet
390 possible as the predicted date has not occurred. At the turning point, the *EIC* is still
391 at its peak. For the disease to have run its course, and an “all clear” declaration can
392 be announced, we require that the newly infected case number to drop to zero. For
393 prediction practice this “zero” is measured by three standard deviations from the
394 peak of $N(t)$. Then we wait for two incubation periods, each 14 days, to pass, before
395 we declare “all clear”. Using the inferred disease characteristics in Table S1, our
396 prediction is, for China outside Hubei: the last week of March. For China as a whole:
397 the first week of April, barring “imports” of infected from abroad. At this point there
398 may still be some patients in the hospital who are infected with the virus. The “all
399 clear” call assumes that these patients are not roaming freely to cause new
400 infections.

401

402 **Prediction for South Korea**

403 Figure S6 summarized the available data for Korea at the present. The recovered
404 case numbers hovered around 1 and 2 daily up to March 1st. It only picked up
405 toward the end. Starting from 19 February, there seems to be enough new daily
406 infected cases. The South Korea Government has identified that the epic center of
407 the epidemic was at church gatherings in the city of Daegu and North Gyeongsang
408 province, where 90% of the cases are found. Specifically, a confirmed COVID-19
409 patient was reported to have attend the Shincheonji Church of Jesus services twice
410 on February 9th and 16th. Given the incubation period of 7 to 14 days, the initial
411 explosion at February 19th and the first peak value around February 24th are not
412 accidents.

413

414 If we use the available daily new cases data, we can get the statistical characteristics
415 of the distribution of the daily new cases from Figure S7, which gives the t_N as March
416 3rd and a σ_N value of 4.5 days. If we further use the turning point as approximately
417 $t_N + T/2$, then the turning point should fall on March 10, assuming T as 14 days based
418 on the over all mean from different regions in China.

419

420 For the NR ratio, it is limited by the availability of recovered case number. If we use
421 the limited recovered cases starting from March 1st, we have 7 days of data. The
422 computed the NR ratio together with the trend is given in Figure S8. The turning
423 point, at the zero-crossing of the extended trend line, would occur between March
424 11th and 12th. This approach does not need to use a value for T .

425

426 An estimate of the end of the epidemic can be given as the second week of April,
427 using the estimated value for $t_N = 3$ March, $\sigma = 4.5$ days. Remarkably, this date is
428 around the same time as for Wuhan, China. South Korea owes its quick turning point
429 and end of the epidemic date to its ability to identify the first infection and the
430 secondary infections at Shincheonji Church (31), where most of the infected were
431 concentrated. This is reflected in the data: σ for South Korea is only half that of
432 China, with a more rapid rise and fall of the newly infected. Its data for the newly

433 infected are probably more accurate compared to other countries in similar stage of
434 the epidemic, due to its massive and speedy (within 6 hours) testing of the
435 population in its “trace, test and treat” policy.

436

437 **4. Effects of quarantine judging from data on the net infection rate.**

438 Following (15), we define a time-dependent *net infection rate* as:

439
$$\alpha(t) = \frac{dI/dt}{I} = \frac{d}{dt} \log I(t).$$

440 In traditional models, such as the SIR model, there is also a time-dependent
441 infection rate, which at $t=0$ is related to the *Basic Reproductive Number* R_0 . If this
442 number is greater than 1 then an epidemic will ensue, i.e. the infected population
443 will increase exponentially after the introduction to a susceptible population S at $t=0$
444 some initial infected. That is, from the SIR model equation:

445
$$dI/dt = aSI - bI = bI(aS/b - 1),$$

446 where $aS(t)$ is the infection rate and $b = 1/T$ is the mean recovery rate.

447 Therefore $\alpha(t) = \frac{dI/dt}{I} = b\left(\frac{aS}{b} - 1\right)$. $\alpha(0) = b(R_0 - 1)$, with $R_0 = \frac{aS(0)}{b}$.

448 Our time-dependent net infection rate generalizes this concept to be independent of
449 the SIR or other models and be applicable at later times as well: If in the course of an
450 epidemic, $\alpha(t)$ is positive, the number of infected will grow exponentially, reaching
451 a peak number of infected when $\alpha(t) = 0$ at $t = t_p$. Then the total number of active

452 infected will decrease exponentially. One could in analogy to R_0 , define a time-

453 dependent Reproductive Number $R_t = \alpha(t)T - 1$, so that if this number is greater
454 (less) than 1 the number of infected will grow (decrease) at time t . We will here
455 use $\alpha(t)$ directly.

456

457 Since $\alpha(t)$ has a zero and its first derivative near the zero is nonzero (viz negative)
458 because t_p is the maximum of $I(t)$, it is a linear function of t with negative slope in
459 that neighborhood. This expectation is verified empirically, using data for the total
460 existing case numbers. We find that this time dependent infection rate is
461 approximately a linear function of time in the neighborhood of its zero over a period
462 of a few weeks.

463

464 This gives another way to predict the turning point t_p which can be used instead of
465 (but is less accurate than) the NR ratio, during the early stage of the epidemic when
466 not enough $R(t)$ data is available, as is the case currently for US.

467

468 The peak *EIC* number can be predicted as

469
$$I(t_p) = I(t_b) \exp\left\{\frac{1}{2}\alpha(t_b)(t_b - t_p)\right\}$$
, where t_b is the last available data before the

470 turning point. It is assumed that $\alpha(t)$ lies on a straight line between t_b and t_p .

471 **Predicting the peak active infected cases**

472 Since the turning point can be predicted two weeks in advance, the above formula
473 can be used to predict the peak EIC numbers.

474

475 **Predicting the total infected cases (TIC):**

476
$$TIC(t) = \int_0^t N(t) dt.$$

477 To predict the total infected cases for the epidemic for a region, we need to do the
478 above accumulation of $N(t)$ into the future, to the end date of the epidemic. That
479 total is approximately:

480
$$TIC_{\infty} = 2 \cdot TIC(t_N),$$

481 assuming that $N(t)$ is approximately symmetric about its peak at t_N . In reality $N(t)$
482 may not be symmetric and likely has a long tail. However, since the number of cases
483 along the tail is small, the above approximation for the total is still good. If the
484 present time t_B is before t_N , we need a way to predict $TIC(t_N)$. Let the *total infection*
485 *rate* be defined as:

486
$$\beta(t) = \frac{d}{dt} \log TIC(t) = \frac{\frac{d}{dt} TIC(t)}{TIC(t)} = \frac{N(t)}{TIC(t)}.$$

487 By extrapolate the total infection rate forward in time we can predict:

488
$$TIC(t_N) = TIC(t_B) \exp\left\{\int_{t_B}^{t_N} \beta(t) dt\right\}.$$

489 **USA:** t_N is predicted to be 2.2 days from today, April 5. Today's TIC is 308,850. So the
490 predicted total infected cases for the epidemic when it is over is predicted to be
491 708,750. This is a very large number, nine times larger than that of China, which has
492 a much larger population, but is nevertheless much lower than some other
493 predictions of a few million infected. Its current EIC is 285,000, and is predicted to
494 peak 12 days from now at 547,000. This is the peak demand for hospital beds.

495 **Germany:** Germany has good data. Its t_N was 2 days ago. On that day its TIC was
496 85,000. So the total infected cases for the epidemic when it is over is predicted to be
497 2 times that: 170,000. Its current EIC is 68,248. The peak EIC is 69,627, which will
498 occur two days from now.

499 **Spain:** t_N occurred on March 31. On that day its TIC was 94,417. The total TIC when
500 the epidemic is over is predicted to be twice that, 188,800. Its peak EIC is predicted
501 to be 84,250, to occur 4 days from now.

502 **UK:** Its current TIC is 41,000. Its t_N is 7.5 days from now. Calculating its TIC at that
503 time and then doubling it yields a total TIC of 134,000 when the epidemic is over.
504 The peak EIC is predicted to be 83,200, 16 days from now. This is the peak demand
505 for hospital beds. (It should be noted that the recovered case number for UK is
506 unusually low, currently at 209, while the dead is much higher, at 4,900. The data
507 may be doubtful.)

508

509 **Comments on the effects of quarantine**

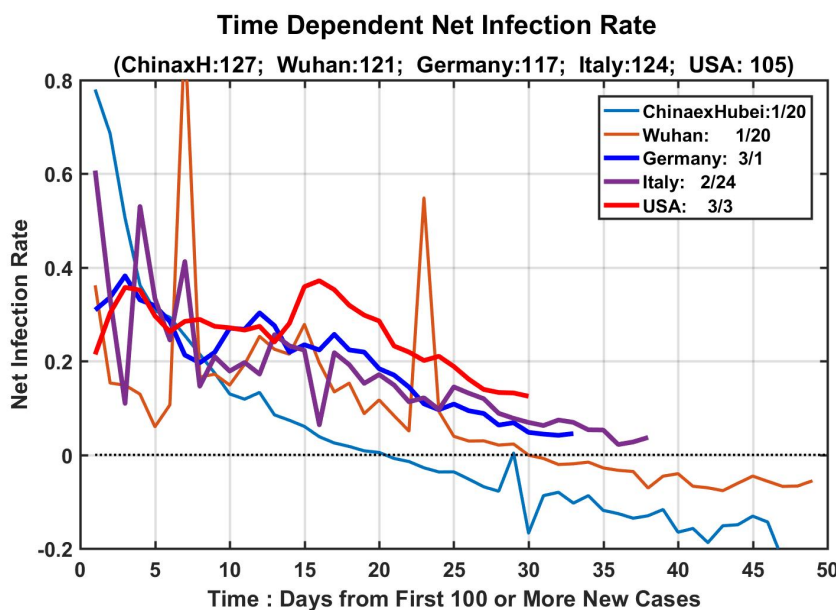
510 Additionally the net infection rate reveals the effect of measures taken with social
511 distancing and quarantine. Figure 5 shows the time-dependent net infection rate for
512 each region starting when the newly confirmed cases exceed 100. This way of
513 plotting facilitates comparison of different regions at the same stage of the epidemic.
514 First, China outside Hubei has the lowest time-dependent infection rate after Hubei
515 was lockdown. Germany and Italy have similar exponential growth rate of the net
516 infected case numbers, both higher than even Wuhan, the epicenter in China. More
517 surprisingly, US has the highest exponential net infection rate, higher than Germany
518 and Italy and China. This can be attributed to the fact that US so far does not have a
519 nation-wide shutdown, unlike these other countries. Secondly, China outside Hubei
520 reached its turning point early, in fact 20 days earlier than the epicenter, Wuhan.
521 We had previously predicted this, which is qualitatively different than many model
522 predictions, which had the epicenter achieving its turning point 1-2 weeks earlier
523 than China outside Hubei (13). Italy and Germany are predicted to take a week
524 longer to reach their turning point, while US will take another week longer than it
525 will take Germany and Italy.

526

527 That Italy would take the same amount of time to reach its turning point as Germany
528 and has approximately the same net infection rate may be due to two reasons: Italy
529 does not test as widely as Germany and so the case numbers represent a smaller
530 portion of the infected. Secondly Germany has the lowest fatality rate while Italy one
531 of the highest. The number of the dead is about the same as the number of cured for
532 Italy. The dead is included in recovered/removed. If it were not included, *EIC*
533 would have been higher for Italy. Nevertheless, whether dead or cured, the hospital
534 bed is vacated.

535

536



537

538

539

540 **Figure 5.** The time-dependent net infection rate (in units of 1/day) as a function of
541 time starting on the date (listed in the inset) when the newly confirmed case
542 number exceeds 100 for each region. To obtain the actual calendar date, add the
543 dates on the horizontal axis to the starting date indicated in the inset. The number of
544 confirmed cases on the starting date is listed at the top.

545

546 **5. Conclusion.**

547 We offer an additional data-driven approach to track and predict the course of the
548 epidemic. Many parameters characterizing an epidemic can be determined from
549 local-in-time data. Validated by real data, we suggest that our approach could be
550 applied not just to the current Covid-19 epidemic, but also generally to future
551 epidemics. It could also be used as a practical tool for epidemic management
552 decisions such as quarantine institution and medical resource planning and
553 allocations (32-35).

554

555 Two results are of special significance for future policy makers. First the turning
556 point for the epidemic in China exHubei occurred a week earlier than that for
557 Wuhan. Second the US will take 2 weeks longer to reach its turning point than even
558 Wuhan. After its lockdown, Wuhan, with a large susceptible population of 11 million,
559 enforced straight social distancing, which is more strict than that adopted in the US..
560 As a consequence, even with a large pool of potential susceptible, the outbreak could
561 end sooner, as compared to the time it will take US and Europe. For China, the
562 lockdown of Wuhan and Hubei was the reason why the epidemic outside Hubei was
563 under control, and the turning point occurred earlier. In Wuhan, with hospitals
564 facing the number of infected patients far exceeding available hospital beds in the
565 initial period, some infected patients were not adequately isolated. The infected
566 were sent home and caused secondary infections among family members. This
567 might have played a role in delaying the turning point. On the other hand, outside
568 Hubei, hospitals were not as overwhelmed because of the strict quarantine placed
569 on Hubei, which drastically reduced the import of the disease originating from
570 Hubei. The infected were better isolated, reducing further spread, and treated in
571 hospitals, resulting in shorter time to recovery (see Table S1). This is evidence of the
572 effectiveness of the city and province-wide lockdown in “flattening the curve”
573 outside.

574

575 The additional and surprising finding that the net infection rates in Italy, Germany
576 and the US are higher than even Wuhan also manifests the effect of the enforcement
577 of lockdown, stay-at-home and strict social distancing policy in Wuhan, which was
578 much stricter than those adopted currently in Europe and US. The more lax latitude
579 and lack of enforcement in the US will lead to a longer period of the epidemic, longer
580 than even Italy, and the largest number of total infected cases of 710,000 before it is
581 over.

582

583

584

585

586 **References**

587

- 588 1. David Adam, Modelers Struggle to Predict the Future of the COVID-19 Pandemic.
589 The Scientist, [https://www.the-scientist.com/news-opinion/modelers-struggle-](https://www.the-scientist.com/news-opinion/modelers-struggle-to-predict-the-future-of-the-covid-19-pandemic-67261)
590 [to-predict-the-future-of-the-covid-19-pandemic-67261](https://www.the-scientist.com/news-opinion/modelers-struggle-to-predict-the-future-of-the-covid-19-pandemic-67261) (2020).
- 591 2. WHO, Laboratory testing of human suspected cases of novel coronavirus (nCoV)
592 infection: interim guidance, World Health Organization, Geneva (2020).
- 593 3. N. Zhu, D. Zhang, W. Wang, X. Li, B. Yang, J. Song, X. Zhao, B. Huang, W. Shi, R. Lu,
594 P. Niu, F. Zhan, A novel coronavirus from patients with pneumonia in China,
595 2019. *N. Engl. J. Med.* **382**, 727-733 (2020).
- 596 4. R. Lu, X. Zhao, J. Li, P. Niu, B. Yang, H. Wu, W. Wang, H. Song, B. Huang, N. Zhu, Y.
597 Bi, X. Ma, F. Zhan, L. Wang, T. Hu, H. Zhou, Z. Hu, W. Zhou, L. Zhao, ..., W. Tan,
598 Genomic characterisation and epidemiology of 2019 novel coronavirus:
599 implications for virus origins and receptor binding. *The Lancet* **395**(10224),
600 565-574 (2020).
- 601 5. Y. Liu, A. A. Gayle, A. Wilder-Smith, J. Rocklöv, The reproductive number of
602 COVID-19 is higher compared to SARS coronavirus. *J. Travel Med.* taaa021
603 (2020).
- 604 6. J. W. Glasser, N. Hupert, M. M. McCauley, R. Hatchett, Modeling and public health
605 emergency responses: Lessons from SARS. *Epidemics* 3: 32-37 (2011),
606 doi:10.1016/j.epidem.2011.01.001.
- 607 7. P. Zhou, X. Yang, X. Wang, B. Hu, L. Zhang, W. Zhang, H. Si, Y. Zhu, B. Li, C. Huang,
608 H. Chen, J. Chen, ..., Z. Shi, A pneumonia outbreak associated with a new
609 coronavirus of probable bat origin. *Nature* **579**, 270-273 (2020).
- 610 8. C. Huang, Y. Wang, X. Li, L. Ren, J. Zhao, Y. Hu, L. Zhang, G. Fan, J. Xu, X. Gu, Z.
611 Cheng, T. Yu, J. Xia, Y. Wei, W. Wu, X. Xie, W. Yin, H. Li, M. Liu, Y. Xiao, H. Gao, L.
612 Guo, J. Xie, G. Wang, R. Jiang, Z. Gao, Q. Jin, J. Wang, B. Cao, Clinical features of
613 patients infected with 2019 novel coronavirus in Wuhan, China. *The Lancet*
614 **395**(10223), 497-506 (2020).
- 615 9. J. F.-K. Chan, S. Yuan, K.-H. Kok, K. K.-W. To, H. Chu, J. Yang, F. Xing, J. L. BNurs, C.
616 C.-Y. Yip, R. W.-S. Poon, H.-W. Tsoi, S. S.-F. Lo, K.-H. Chan, V. K.-M. Poon, W.-M.
617 Chan, J. D. Lp, J.-P. Cai, V. C.-C. Cheng, H. Chen, C. K.-M. Hui, K.-Y. Yuen, A familial
618 cluster of pneumonia associated with the 2019 novel coronavirus indicating
619 person-to-person transmission: a study of a family cluster. *The Lancet*
620 **395**(10223), 514-523 (2020).
- 621 10. X. Xu, P. Chen, J. Wang, J. Feng, H. Zhou, X. Li, W. Zhong, P. Hao, Evolution of the
622 novel coronavirus from the ongoing Wuhan outbreak and modeling of its spike
623 protein for risk of human transmission. *Sci. China Life Sci.* **63**, 457-460 (2020).
- 624 11. Z. Chen, W. Zhang, Y. Lu. C. Guo, Z. Guo, C. Liao, X. Zhang, Y. Zhang, X. Han, Q. Li, W.
625 lan Lipkin, J. Lu, From SARS-CoV to Wuhan 2019-nCoV Outbreak: Similarity of
626 Early Epidemic and Prediction of Future Trends. *Biorxiv* preprint (2020),
627 doi: <https://doi.org/10.1101/2020.01.24.919241>.
- 628 12. J. M. Read, J. R. E. Bridgen, D. A. T. Cummings, A. Ho, C. P. Jewell, Novel
629 coronavirus 2019-nCoV: early estimation of epidemiological parameters and
630 epidemic predictions. *medRxiv* preprint (2020), doi:
631 <https://doi.org/10.1101/2020.01.23.20018549>.

- 632 13. J. T. Wu, K. Leung, G. M. Leung, Nowcasting and forecasting the potential
633 domestic and international spread of the 2019-nCoV outbreak originating in
634 Wuhan, China: a modelling study. *The Lancet* **395**(10225), 689-697 (2020).
- 635 14. S. Zhao, S. S. Musa, Q. Lin, J. Ran, G. Yang, W. Wang, Y. Lou, L. Yang, D. Gao, D. He,
636 M. S. Wang, Estimating the Unreported Number of Novel Coronavirus (2019-
637 nCoV) Cases in China in the First Half of January 2020: A Data-Driven Modelling
638 Analysis of the Early Outbreak. *J. Clin. Med.* **9**, 388 (2020).
- 639 15. N. E. Huang, F. Qiao, A data driven time-dependent transmission rate for tracking
640 an epidemic: a case study of 2019-nCoV. *Sci. Bull.* **65**, 425-427(2020) ,
641 <https://doi.org/10.1016/j.scib.2020.02.005>.
- 642 16. Q. Li, W. Feng, Trend and forecasting of the COVID-19 outbreak in China. *J.*
643 *Infection* arXiv:2002.05866v1, (2020).
- 644 17. H. Xiong, H. Yan, Simulating the infected population and spread trend of 2019-
645 nCov under different policy by EIR model. *medRxiv* preprint (2020), doi:
646 <https://doi.org/10.1101/2020.02.10.20021519>.
- 647 18. L. Damon, E. Brooks-Pollock, M. Bailey, M. J. Keeling, A spatial model of CoVID-19
648 transmission in England and Wales: early spread and peak timing. *medRxiv*
649 preprint (2020), doi: <https://doi.org/10.1101/2020.02.12.20022566>.
- 650 19. H. Sun, Y. Qiu, H. Yan, Y. Huang, Y. Zhu, S. Chen, Tracking and Predicting COVID-
651 19 Epidemic in China Mainland. *Medrxiv* preprint (2020),
652 doi: <https://doi.org/10.1101/2020.02.17.20024257>.
- 653 20. Q. Liu, Z. Liu, D. Li, Z. Gao, J. Zhu, J. Yang, Q. Wang, Assessing the Tendency of
654 2019-nCoV (COVID-19) Outbreak in China. *medRxiv* preprint (2020), doi:
655 <https://doi.org/10.1101/2020.02.09.20021444>.
- 656 21. L. Peng, W. Yang, D. Zhang, C. Zhuge, L. Hong, Epidemic analysis of COVID-19 in
657 China by dynamical modeling. *arXiv* 2002.06563, (2020).
- 658 22. D. Cyranoski, When will the coronavirus outbreak peak? *Nature news* (2020).
- 659 23. C. R. MacIntyre, Global spread of COVID-19 and pandemic potential. *Global*
660 *Biosecurity* **1**(3), (2020).
- 661 24. WHO, Coronavirus latest: WHO describes outbreak as pandemic, *Nature news*
662 (2020), <https://www.nature.com/articles/d41586-020-00154-w>.
- 663 25. K. Kupferschmidt, J. Cohen, Can China's COVID-19 strategy work elsewhere?
664 *Science* **367**(6482), 1061-1062 (2020).
- 665 26. J. M. Read, J. R. Bridgen, D. A. Cummings, A. Ho, C. P. Jewell, Novel coronavirus
666 2019-nCoV: early estimation of epidemiological parameters and epidemic
667 predictions. *medRxiv* (2020), doi:10.1101/2020.01.23.20018549.
- 668 27. S. Zhao, Q. Lin, J. Ran, S. S. Musa, G. Yang, W. Wang, Y. Lou, D. Gao, L. Yang, D. He,
669 M. H. Wang, Preliminary estimation of the basic reproduction number of novel
670 coronavirus (2019-nCoV) in China, from 2019 to 2020: A data-driven analysis in
671 the early phase of the outbreak. *Int. J. Infect. Dis.*, 92: 214-217 (2020),
672 <https://doi.org/10.1016/j.ijid.2020.01.050>.
- 673 28. W. Kermack, A. McKendrick, A contribution to the mathematical theory of
674 epidemics. *Proc. Roy. Soc. London A* **115**, 700-721 (1927).
- 675 29. D. L. Heymann, N. Shindo, COVID-19: what is next for public health?
676 *Lancet*, 395(10224): 542-545 (2020), [https://doi.org/10.1016/S0140-](https://doi.org/10.1016/S0140-6736(20)30374-3)
677 [6736\(20\)30374-3](https://doi.org/10.1016/S0140-6736(20)30374-3).

- 678 30. The Novel Coronavirus Pneumonia Emergency Response Epidemiology Team,
679 The Epidemiological characteristics of an outbreak of 2019 novel coronavirus
680 disease (COVID-19)-China, 2020, *China CDC Weekly* (2020).
- 681 31. E. Shim, A. Tariq, W. Choi, Y. Lee, G. Chowell, Transmission potential of COVID-19
682 in South Korea. *medRxiv* preprint, (2020), doi:
683 <https://doi.org/10.1101/2020.02.27.20028829>.
- 684 32. C. M. Peak, L. M. Childs, Y. H. Grad, C. O. Buckee, Comparing nonpharmaceutical
685 interventions for containing emerging epidemics. *Proc. Natl. Acad. Sci.* 114(15):
686 4023-4028 (2017), doi:10.1073/pnas.1616438114.
- 687 33. R. S. Dhillon, D. Srikrishna, When is contact tracing not enough to stop an
688 outbreak? *Lancet Infect. Dis.*, 18: 1302-1304 (2018),
689 [https://doi.org/10.1016/S1473-3099\(18\)30656-X](https://doi.org/10.1016/S1473-3099(18)30656-X).
- 690 34. X. Pang, Z. Zhu, F. Xu, J. Guo, X. Gong, D. Liu, Z. Liu, D. P. Chin, D. R. Feikin,
691 Evaluation of control measures implemented in the severe acute respiratory
692 syndrome outbreak in Beijing, 2003. *JAMA*, 290(24): 3215-3221 (2003).
- 693 35. G. Wang, N. E. Huang, F. Qiao, Quantitative evaluation on control measures for an
694 epidemic: A case study of COVID-19. *Sci. Bull.* **65** (2020), doi: 10.1360/TB-2020-
695 0159.
- 696 36. J.D. Murray: *Mathematical Biology I*. Third Edition. Springer, 551pp.
697
698
699
700
701

702 **Acknowledgements:** NEH and FQ are supported by the National Natural Science
703 Foundation of China under Grant 41821004. KKT's research is supported by the
704 Frederic and Julia Wan Endowed Professorship.

705 **Competing Interests:** The authors declare no competing interests.

706 **Data Availability:** All data in this study are publicly available from World Health
707 Organization (WHO) at [https://www.who.int/emergencies/diseases/novel-](https://www.who.int/emergencies/diseases/novel-coronavirus-2019/situation-reports/)
708 [coronavirus-2019/situation-reports/](https://www.who.int/emergencies/diseases/novel-coronavirus-2019/situation-reports/)

709 and on the Daily Brief site of the China's National Health Commission at
710 <http://en.nhc.gov.cn/>

711 The Korean data is available at

712 <https://sa.sogou.com/new-weball/page/sgs/epidemic>

713 Coronavirus COVID-19 Global Cases by Johns Hopkins CSSE

714 [https://gisanddata.maps.arcgis.com/apps/opsdashboard/index.html#/bda759474](https://gisanddata.maps.arcgis.com/apps/opsdashboard/index.html#/bda7594740fd40299423467b48e9ecf6)
715 [0fd40299423467b48e9ecf6](https://gisanddata.maps.arcgis.com/apps/opsdashboard/index.html#/bda7594740fd40299423467b48e9ecf6)

716

717

718 **Supplementary Information:**

719

720 **Comments on SIR-type of models**

721 Mainstream epidemiological models have their origin in the SIR (Susceptible,
722 Infected, and Recovered or Removed) model (28) and its many variations. We
723 discuss here why existing model predictions vary so widely by commenting on the
724 assumptions underlying these models, using the classical SIR model as an example.

725

726 The number of the infected is increased by the rate of infection, which is assumed to
727 be aIS , and decreased by the rate of recovery/removed, who are assumed to be
728 immune to further infection, bI . ($dI/dt = aIS - bI$). This latter category, R , increases
729 at the rate of bI . ($dR/dt = bI$). Both of the parameters, especially the infection rate
730 a , is largely unknown for an emergent disease (1), and have to be either estimated
731 based on limited statistics the number of contacts a single infected person may have
732 and how many of the contacted people will be infected, or obtained by curve fitting
733 of reported $I(t)$. This approach has problems. First, for Covid-19, there is a
734 population of asymptomatic infected, which may be larger than the
735 reported/confirmed I during the early stages of the epidemic, when wide-scale
736 testing is not available. Since this asymptomatic infected is also infectious, and
737 some of those infected by them could later become confirmed infected, fitting the
738 rate of increase of reported I to estimate the infection rate would inevitably yield a
739 much larger a . Consequently the model predictions using this estimated parameter
740 value may yield a very large peak infection number. Secondly, since
741 $dI/dt = aI(S - b/a)$, the number of infected will increase (approximately
742 exponentially) for $S > b/a$, but for a new virus to which the population has no
743 immunity, the susceptible population could be very large. For China it could be as
744 large as 1.4 billion. For Hubei it is 65 million. The susceptible population should be
745 lower if the quarantine of Wuhan were tight, but it is difficult to estimate what it is
746 under realistic conditions. Thirdly, predictions of the end of the epidemic vary
747 widely because of the basic assumption in some SIR type of models that if most of
748 the population is infected and recovered (and hence acquired immunity) or dead
749 (and hence no longer infectious). The concept of "herd immunity" is rooted in the
750 idea that with enough of the population acquiring immunity this way (or from a
751 vaccine if one is developed in time), the susceptible population is reduced to
752 $S < b/a$, and the rate of increase of I will turn negative. This will require as much as
753 70%-90% of the population be infected, a huge number. Modern models take into
754 account of the effect of quarantine and isolation in reducing the number of people
755 each infected can contact, thus reducing the infection rate a , so that S does not have
756 to be reduced to such an extent for the epidemic to peak and I starts to decrease.
757 However, such estimates are very dependent on model assumptions.

758

759 Of course the modern models of epidemiology are more sophisticated than the SIR
760 model (12-14, 17, 21).

761

762 We offer here an additional tool that has the advantage that it does not depend
763 on the elusive infection rate or the susceptible population, information needed for
764 most models, but has the disadvantage that it cannot be used when the epidemic
765 first started and the data are inaccurate or incomplete. It is based on daily case
766 numbers (i.e. newly confirmed cases), $N(t)$, and recovered cases, $R(t)$.

767
768 Our estimate of the end date of the epidemic is not based on the number of
769 susceptibles, S , approaching zero as in most models (i.e. most of the population is
770 infected, hence acquiring immunity), but $N(t)$ approaching zero and remaining so
771 for two incubation periods. The first incubation period is to allow the asymptomatic
772 infected to show symptoms and the second period to allow those that are infected
773 by the asymptomatic infected to show symptoms. For prediction purpose, the date
774 when the $N(t)$ is zero is estimated by 3 standard deviations from its peak. These
775 two quantities can be extracted from the data as the epidemic is developing. Our
776 estimate of the end of the epidemic is earlier than most model predictions, usually
777 significantly so, because it does not depend on the herd immunity concept.

778
779

780 **THEORY:**

781 **Definition:** Let $I(t)$ be the number of active infected at time t . Its change is given
782 by;

$$783 \frac{d}{dt}I = N(t) - R(t),$$

784 where $N(t)$ is the number of newly infected, and $R(t)$ that of the newly recovered or
785 removed (dead). Note that for the theory part, $N(t)$ includes both confirmed and
786 unconfirmed cases. The term: Existing Infected Case (*EIC*) number is used to denote
787 the confirmed $I(t)$ when we deal with data.

788

789 Let t_p , the turning point defined as the peak of the active infected number. At this
790 point maximum medical resource is needed. This maximum occurs when

$$791 \frac{d}{dt}I = 0, \text{ implying } N(t_p) = R(t_p).$$

792 There is no need to first find $I(t)$ to locate this peak. After the turning point, the
793 newly recovered starts to exceed the newly infected. The demand for medical
794 resources, such as hospital beds, isolation wards and respirators, starts to decrease.

795

796 Let $t=0$ be when the first infection began. For Wuhan, China, this date is near the
797 end of 2019, perhaps even earlier. Let t_b be the beginning of the better quality data.

798 This time is beyond the initial incubation period of the disease and it can be
799 assumed that at that time there is already a population of infected, some of them
800 asymptomatic but nevertheless infectious.

801

802 Let $X(t,s)$ be the number of infected cases at time t , with s being the “age”
803 distribution, i.e. number of days sick.
804 The total number of infected is given by:

$$805 \quad I(t) = \int_0^T X(t,s) ds .$$

806 After being sick for T days, a patient either recovers or is removed (dead). T is called
807 the recovery period (or removal period). It is also called the infectious time if the
808 patient is infectious during this period. Of course its value varies by patient and by
809 the efficacy of treatment for each hospital. For the removed it also depends on the
810 age of the patient and whether there are underlying medical conditions. Only a
811 mean recovery period is obtainable from data, and so this is in reality a statistical
812 quantity. We will discuss later how this statistical quantity can be obtained from
813 data.

814

815 **Conservation law** (see Murray: Mathematical Biology Part I, Chapter 1 (36)):
816 After first infected and until removed or cured, we have:

$$dX(t,s) = \frac{\partial}{\partial t} X \cdot dt + \frac{\partial}{\partial s} X \cdot ds = 0, \quad 0 < s < T.$$

817 So, since $ds/dt = 1$,

$$\frac{\partial}{\partial t} X + \frac{\partial}{\partial s} X = 0.$$

818 This equation is to be solved using the method of characteristics as

$X(t,s) = \text{constant}$ along characteristics defined by $ds/dt = 1$.

Boundary condition: $X(t,0)$, specifies the “birth” process,

819 i.e. how the disease spawns newly infected (with “age” $s = 0$).

Initial condition: $X(0,s) = 0$ for $s > 0$, specifies the initial age distribution at $t = 0$

820 There are two types of characteristics:

821 (i) $s > t$, (ii) $s < t$.

822 The first type of characteristics intersects the $t = 0$ axis, and since the initial
823 condition is zero, we have the solution:
824

$$825 \quad X(t,s) \equiv 0 \text{ for } s > t.$$

826

827 That is, there is no infected population who is sick for more days than the lapsed
828 time since the first infection occurred.

829

830 For the second type of characteristics, $t > s$ the solution is

831

$$832 \quad X(t,s) = f(s-t)$$

833 with the form of f to be determined by the boundary condition. Even without
834 determining the form of f we have the following general results:

835 For $t > T$, and therefore $t > s$:

836
$$I(t) = \int_0^T X(t,s) ds = \int_0^T f(s-t) ds$$
$$= \int_{t-T}^t f(p) dp.$$

837
$$\frac{d}{dt}I = f(t) - f(t-T).$$

838

839 Since the rate of increase of confirmed $I(t)$ is by definition equal to the newly
840 confirmed infected number, $N(t)$, minus the newly recovered (or removed) number,
841 $R(t)$, we have:

842
$$N(t) - R(t) = \frac{d}{dt}I = f(t) - f(t-T).$$

843

844 For a fatal disease with low fatality rate, where almost all infected cases eventually
845 recover after a hospital stay of T days, we can identify
846 $f(t)$ with $N(t)$, and $f(t-T)$ with $R(t)$.

847 If the disease has a non-negligible fatality rate, we include the dead in $R(t)$.

848 **Main Result:** The daily newly recovered/removed number $R(t)$, is related to the
849 daily newly infected number $N(t)$ as, for $t > T$:

850

851
$$R(t) = N(t-T).$$

852

Validation: This fundamental relationship can be validated statistically with data.
853 Figures 1, obtained using data from China during the Covid-19 epidemic, shows that
854 $N(t)$ and $R(t)$ are highly correlated: with correlation coefficient of 0.95 when both
855 distributions are smoothed with 5-point boxcar. The unsmoothed daily data also
856 yield a high correlation coefficient of 0.80, with $R(t)$ lagging $N(t)$ by $T \sim 15$ days. Both
857 correlation coefficients are statistically significant. A similar result is found for
858 Hubei (Figure S2) and other regions (not shown). This is one of the ways the *mean*
859 *recovery period* is determined statistically from data, but it is not practical in the
860 early phase of the epidemic. We will give different methods for the latter purpose.
861 The result on T is consistent with that estimated or predicted later using the slope of
862 the distribution in Figure 4. The latter, obtained by the intercept of the straight line,
863 is less accurate because of the slope is rather shallow. For the regions considered in
864 Figure 1, the fatality rate is small and so the dead are not included in $R(t)$ for
865 convenience.
866

867 The second type of characteristics intersects the boundary $s = 0$. The boundary
868 condition itself needs to be solved as a function of t to describe how new infection
869 (at $s = 0$) occurs. This can be done using a birth model, such as Eq. (1.56) in (36).
870 For our purpose we assume that the solution of this model yields a distribution with
871 age that has a full spectrum $0 < s < T$ of infectives at a time t_b , long after a full
872 incubation period has passed.

873 $X(t_B, s) = f_0(s) = A \exp\left\{-\frac{(s-s_0)^2}{2b^2}\right\}$; A independent of s .
 874 $b = \frac{1}{2}T$.

875 Therefore the solution is, for $t > t_B > 0$:

876
$$X(t, s) = X(t_B, s-t) = f_0(s-t)$$

$$= A \exp\left\{-\frac{(s-s_0-t)^2}{2b^2}\right\}.$$

877
$$N(t) = A \exp\left\{-\frac{(t-s_0)^2}{2b^2}\right\} = A \exp\left\{-\frac{(t-t_N)^2}{2b^2}\right\}$$

878
$$R(t) = N(t-T) = A \exp\left\{-\frac{(t-t_R)^2}{2b^2}\right\},$$

879 where t_N is the peak of $N(t)$, and $t_R = t_N + T$ is the peak of $R(t)$.

880 Both distributions are Gaussians.

881

882 For $t_B < t < T$,

883
$$I(t) = \int_0^t X(t, s) ds + \int_t^T X(t, s) ds$$

$$= \int_0^t f(s-t) ds + \int_t^T 0 ds$$

$$= \int_{-t}^0 f(p) dp$$

884
$$\frac{d}{dt} I = f(-t) = A \exp\left[-\frac{(t-s_0)^2}{T^2}\right].$$

885
$$N(t) = A \exp\left[-\frac{(t-t_N)^2}{T^2}\right]$$

$$R(t) = 0.$$

886 Again, $N(t)$ is Gaussian, but there is no recovered or removed during this early stage.

887

888 **Main Result:** The natural logarithm of the ratio of N and R is a linear function of
 889 time for $t > T$:

890

891
$$NR(t) = \frac{N(t)}{R(t)} = \frac{N(t)}{N(t-T)}$$

$$= \exp\left\{-\frac{(t-t_N)^2}{2b^2} + \frac{(t-t_N-T)^2}{2b^2}\right\} = \exp\left\{\frac{T^2}{2b^2} - \frac{T(t-t_N)}{b^2}\right\}.$$

$$\log NR = -\frac{T(t-t_N) - \frac{1}{2}T^2}{b^2}$$

892 a linear function of t . This relationship is important for the purpose of forecast
 893 because it is easy to extrapolate from a straight line into the future.

894

895 It intersects 0 at $t - t_N = \frac{1}{2}T$. This yields the turning point, when the NR ratio is 1, and
 896 therefore its logarithm is zero.

897

898 **Result:** The turning point, defined as the maximum of $I(t)$, is given by $t_p = t_N + \frac{1}{2}T$.

899 **Result:** The slope of $\log NR$ is equal to $4/T$.

900

901 When time is normalized by T , the derivative is given by:

902
$$\frac{d \log NR}{d(t/T)} = \frac{T^2}{b^2} = 4, \text{ a dimensionless constant.}$$

903

904 **Heterogeneous Data**

905 The above results are obtained for the case of a single introduction into a region of
 906 infected at $t=0$ and we solve for the subsequent development of the epidemic from
 907 that single source. Consider now a large region consisting of a number of small
 908 regions, and the “seeding” of the infected occurs at different times for different
 909 regions. The large region could be China, and the first infection could be Wuhan,
 910 Hubei and then the regions outside Hubei. Then we may have for the China as a
 911 whole data for the newly infected a sum of several Gaussians staggered in time. As
 912 long as the Gaussians are not separated so much that there are different peaks in the
 913 combined data, the combined data can still be considered as Gaussian, as is the case
 914 in the real data. However, the standard deviation σ of the combined Gaussian is
 915 inevitably larger and is no longer given by b :

916
$$N(t) = \frac{B}{\sqrt{2\pi}\sigma_N} \exp\left\{-\frac{(t-t_N)^2}{\sigma_N^2}\right\}.$$

917 We still have $R(t) = N(t-T)$ since this result holds for each sub-region. The result
 918 that $\log NR$ is a linear function of time still holds:

919
$$\log NR(t) = \log \frac{N(t)}{N(t-T)} = -\frac{T(t-t_N) - \frac{1}{2}T^2}{\sigma_N^2}.$$

920 The slope of the straight line is T/σ_N^2 .

921

922 Since the hospital state can add as a smoothing filter on $N(t)$ to yield $R(t)$, the
 923 standard deviation for $R(t)$ could be slightly wider than that for $N(t)$. So we could
 924 have two different Gaussians (but their integral over all time should be the same):

925
$$N(t) = \frac{B}{\sqrt{2\pi}\sigma_N} \exp\left\{-\frac{(t-t_N)^2}{2\sigma_N^2}\right\}; \quad R(t) = \frac{B}{\sqrt{2\pi}\sigma_R} \exp\left\{-\frac{(t-t_R)^2}{2\sigma_R^2}\right\}.$$

926 Taking this into account, we have, denoting $T = t_R - t_N$:

927
$$\log NR - \log \frac{\sigma_R}{\sigma_N} = - \left\{ \frac{(t-t_N)^2}{2\sigma_N^2} - \frac{(t-t_N-T)^2}{2\sigma_R^2} \right\} = - \frac{(\sigma_R^2 - \sigma_N^2)(t-t_N)^2 + 2\sigma_N^2 T(t-t_N) - \sigma_N^2 T^2}{2\sigma_N^2 \sigma_R^2}$$

928

929 As the values of σ_N and σ_R are very close based on the empirical data, the quadratic
 930 term is always small comparing to the other terms for the length of time we are
 931 considering here. Hence.

932
$$\log NR(t) = \frac{1}{\sigma_R^2} \left\{ -T(t-t_N) + \frac{1}{2} T^2 \right\},$$

933 a linear function of time. Its slope is $-\frac{T}{\sigma_R^2}$.

934 **Result:** The natural logarithm of the ratio of two Gaussians of slightly different
 935 standard deviation is approximately a straight line.

936

937 **Validation of log NR as a straight line:**

938 From data we use the report newly confirmed case number and the recovered case
 939 number to define NR ratio as

940
$$NR(t) = N(t)/R(t).$$

941 At t_p , $NR=1$.

942

943 We show in Figure 2, using the data of the epidemic for COVID-19, that the
 944 logarithm of $NR(t)$ lies on a straight line, with small scatter, passing through the
 945 turning point t_p . And data for various stages of the epidemic, from the initial
 946 exponential growth stage, to near the peak of *EIC*, and then past the peak, all lie on
 947 the same straight line. The intercept with $\log NR=0$ yields the turning point. This
 948 line, obtained by linear-least-square fit in the semi-log plot, is little affected by the
 949 rather large artificial spike in the data on 12 February because of its short duration
 950 and the logarithmic value. That reporting problem is necessarily of short duration
 951 because, on the date of definition change, previous week's cases of infected
 952 according to the new criteria were reported in one day. After that, the book is
 953 cleared, and $N(t)$ returned to its normal range.

954

955

956 The theoretical result suggests that the slope of the linear line is $-T/\sigma_R^2$, where σ_R is
 957 the standard deviation of the $R(t)$ profile. In general, the slope can be different for
 958 different regions with different levels of quarantine and epidemic characteristics.
 959 The hospital treatment efficacy would influence T directly. The effect of quarantine
 960 would influence the value of σ_N , the standard deviation of the newly infected, and so
 961 indirectly $R(t)$ and σ_R . Our empirical result from Fig. 2 however shows that the slope
 962 is the almost the same for different regions in China, implying that efficacy of
 963 treatment and level of quarantine affect T and σ^2 proportionally.

964

965 **Validation of the slope of log N(t) and log R(t)**

966 Interestingly, the derivative of $\log N(t)$ or $\log R(t)$ also lies on a straight line, as
 967 shown in Fig. 3 (although the scatter is larger as to be expected for any

968 differentiation of empirical data). The positive and negative outliers one day before
969 and after 12 Feb are caused by the spike up and then down, with little effect on the
970 fitted linear trend (but increases its variance and therefore uncertainty). Moreover,
971 the straight line extends without appreciable change in slope beyond the peak of
972 $N(t)$, suggesting that the distribution of the newly infected number is approximately
973 Gaussian. The mean recovery time T can be predicted as $t_R - t_N$, where t_R is the peak
974 of $R(t)$ and t_N is the peak of $N(t)$. These two peak times can be obtained by extending
975 the straight line in Fig. 3 to intersect the zero line. This predicted result can be
976 verified statistically after the fact by the lagged correlation of $R(t)$ and $N(t)$. If the
977 distribution is indeed Gaussian or even approximately so, the slope in Fig. 3 would
978 be proportional to the reciprocal of the square of its standard deviation, σ , as:

979
$$\frac{d \log N(t)}{dt} = \frac{-(t - t_N)}{\sigma_N^2}.$$

980

981 Similarly result holds for the daily number of recovered, $R(t)$.

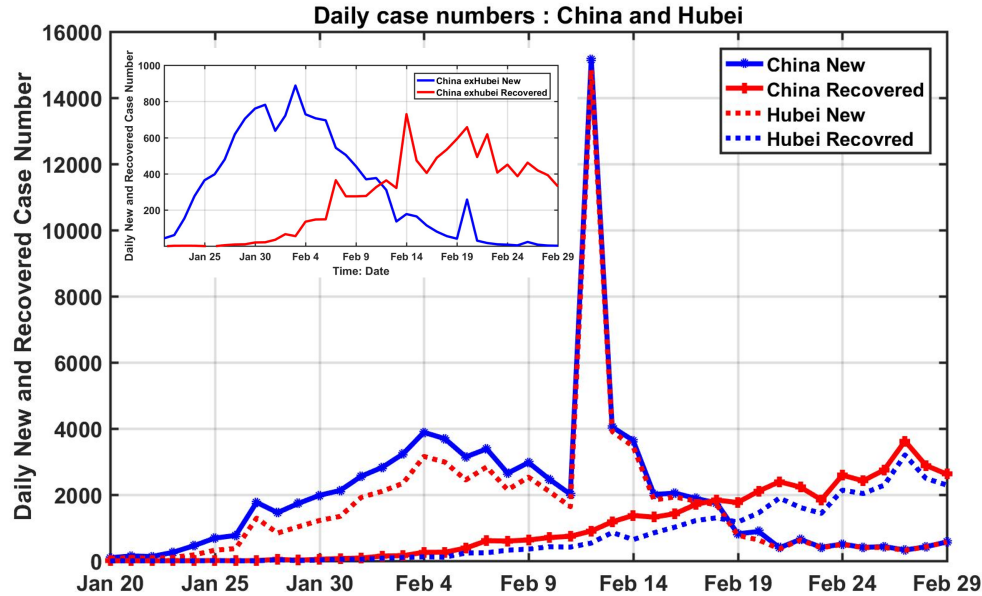
982

983 After the epidemic is nearing the end as is the case in China, fitting the data to a
984 Gaussian can be done after the fact (see Figures S3 and S4). The fit is satisfactory
985 even without using any disposal parameters. The parameters used are determined
986 using slopes of $\log N$ and $\log R$ (see Table S1)

987

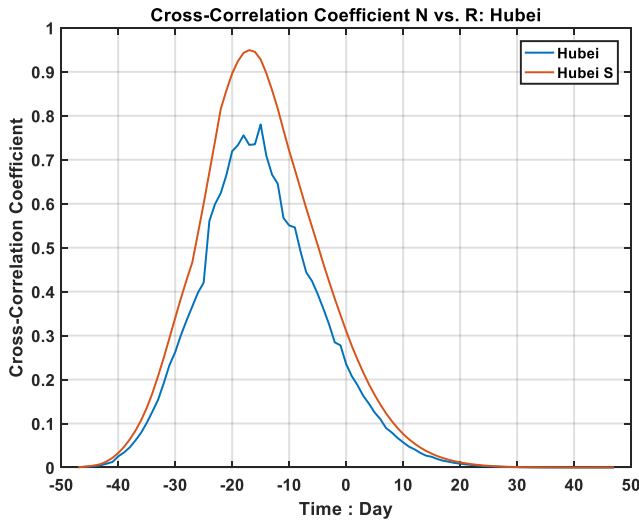
988 The inferred statistical characteristics of the Covid-19 epidemic are summarized in
989 Table S1 for various regions. The mean recovery time T , is about 13 days for China
990 as a whole. For Wuhan, the city at the epicenter whose hospitals were more
991 overwhelmed and the patients admitted into hospitals more seriously ill than those
992 in other provinces, $T \sim 16$ days, while that for Hubei is 14 days. The standard
993 deviation, σ , is found to be around 8 days, with slight difference between that for
994 $N(t)$ and for $R(t)$, with one exception for Hubei outside Wuhan. Such a fine
995 subdivision may not be practical for the data quality we have. The σ tends to be
996 smaller for China as a whole than Wuhan. One can see that T and σ^2 indeed varying
997 approximately in proportion.

998



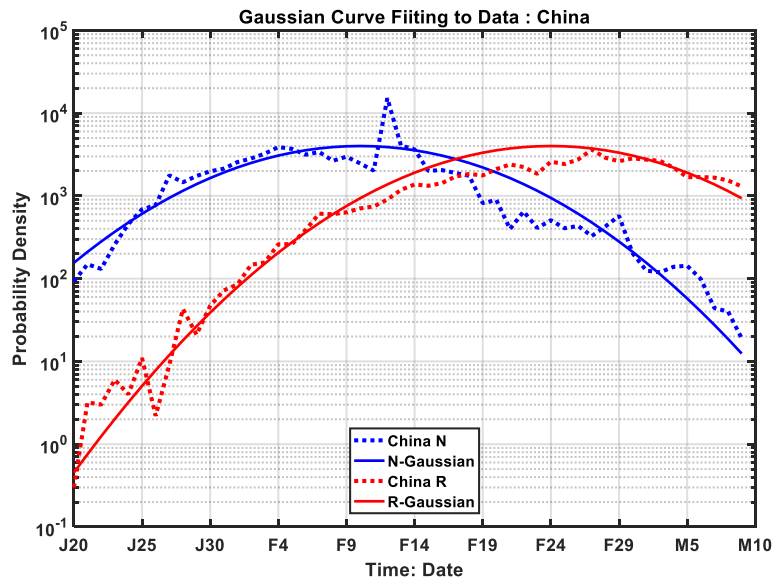
999
1000
1001
1002
1003
1004
1005

Figure S1. The daily newly infected (in blue) and the daily newly recovered (in red), as a function of time for China as a whole (in solid lines) and Hubei (in dotted lines). The turning point is determined by when the red and blue curves cross.
Inset: For China outside Hubei.



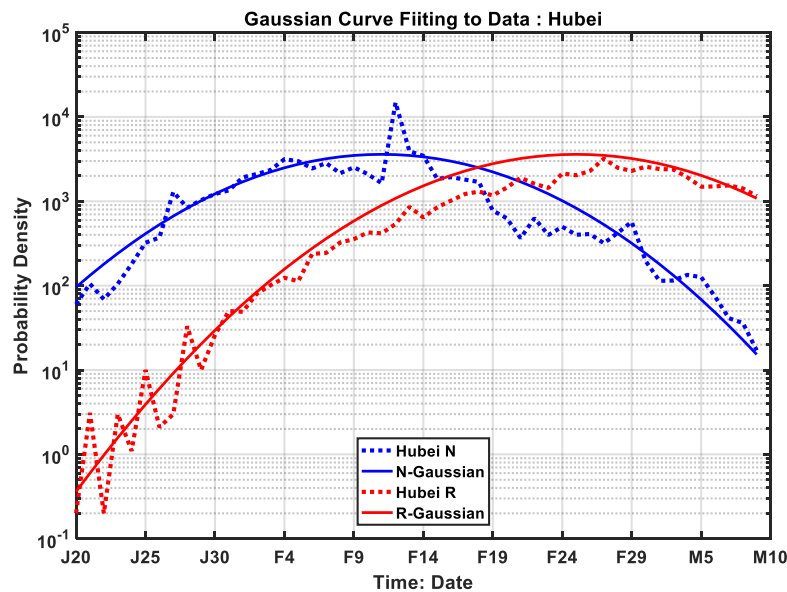
1006
1007
1008
1009

Figure S2. Lagged correlation of $R(t)$ with $N(t)$ for Hubei province.



1010
1011
1012

Figure S3. Gaussian fit of $N(t)$ and $R(t)$, for China as a whole.



1013
1014
1015

Figure S4. Gaussian fit of $N(t)$ and $R(t)$, for Hubei Province.

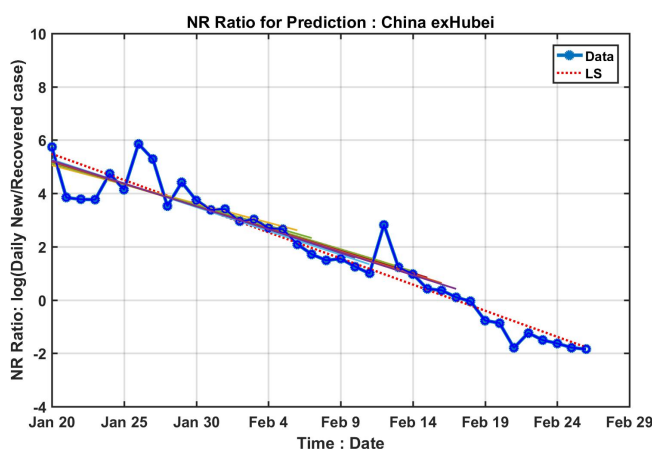
	Crossing	t_0	t_1	T	Sigma N	Sigma R
China	2/18	2/11	2/24	13	7.5	7.7
Hubei	2/20	2/12	2/26	14	7.9	8.5
Wuhan	2/21	2/14	3/01	16	8.8	8.8
C exHubei	2/12	2/08	2/21	13	7.5	7.1
H exWuhan	2/16	2/13	2/27	14	5.0	8.8

--	--	--	--	--	--	--

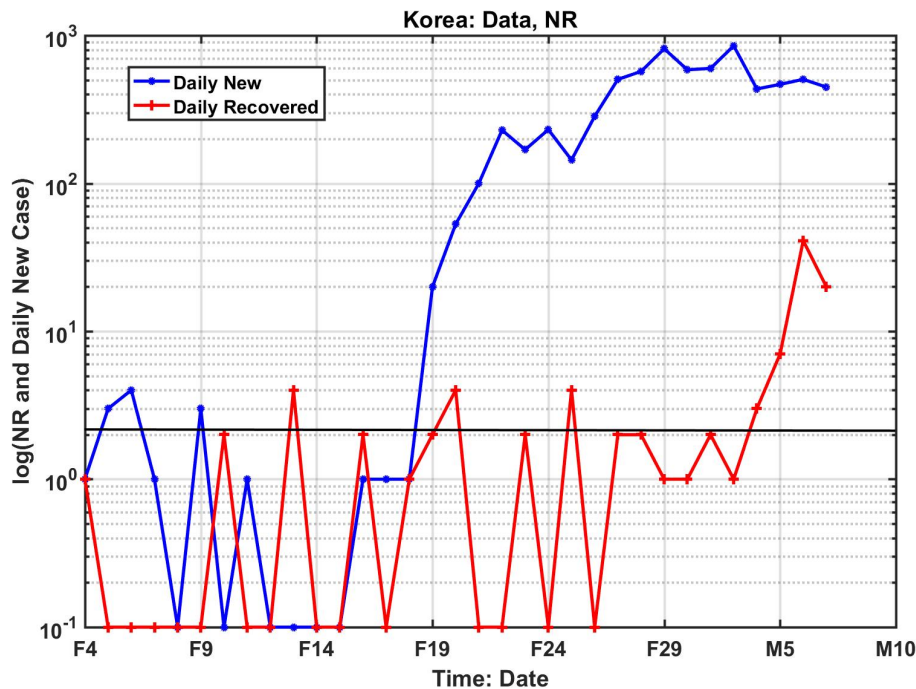
1016 **Table S1:** Statistical characteristics of the COVID-19 epidemic in different regions in
 1017 China inferred from data, for $N(t)$, the daily number of newly infected and for $R(t)$,
 1018 the daily number of recovered.
 1019

	China	Hubei	China-Hubei	Hubei-Wuhan
Truth (data)	18	19	12	15
NR Ratio	20.3±1.6 (Feb 20 nd)	22.3±1.0 (Feb 22 rd)	12.4±0.9 (Feb 12 th)	16.0±1.2 (Feb 16 th)

1020
 1021 **Table S2:** Predicted turning point dates. Shown are the mean and standard
 1022 deviation of the predictions over the prediction period, using the NR ratio method
 1023



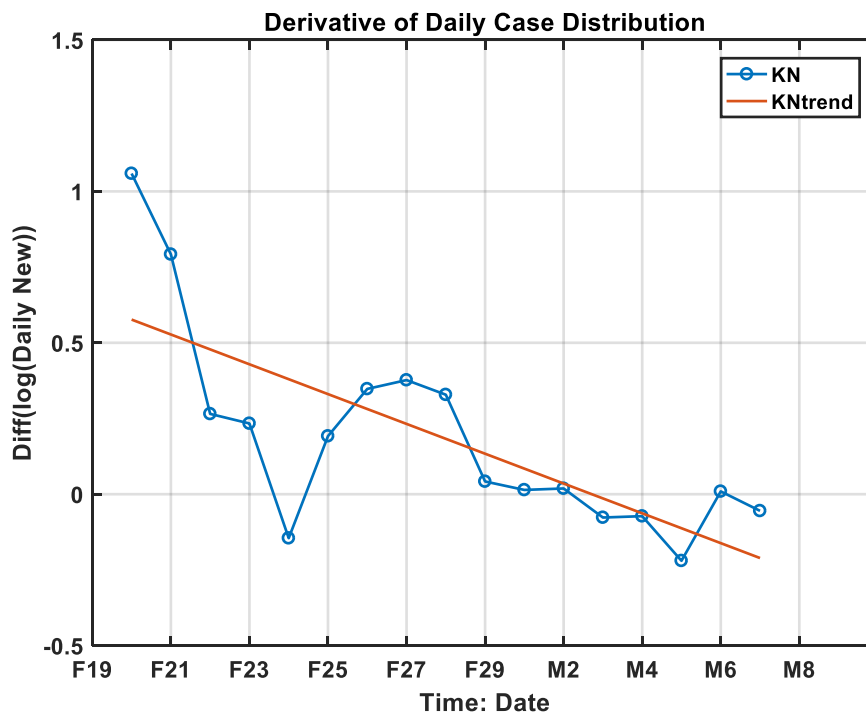
1024
 1025 **Figure S5.** Prediction of the turning point of EIC using linear least-squares trends
 1026 using various data lengths for China exHubei. All data used start from 24 January.
 1027 Different colored straight lines show the linear trend calculated from 24 January to a
 1028 particular date. The spread is over a very small range. Then these trends are
 1029 extrapolated (extrapolations not shown) to intersect the zero line to yield a
 1030 prediction for the turning point. The blue dots are the data.
 1031
 1032
 1033
 1034



1035

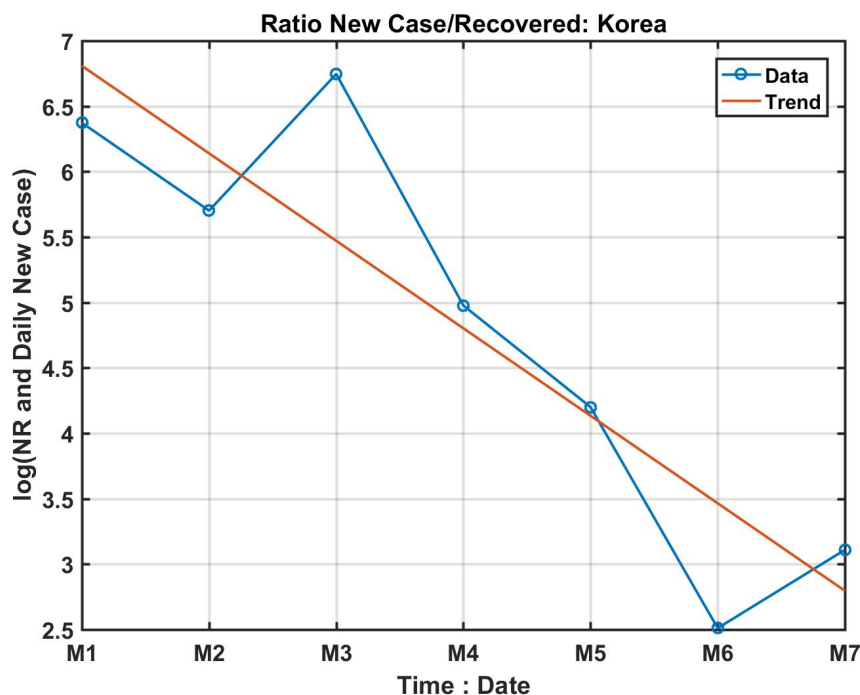
1036

1037 **Figure S6:** The available data from South Korea (as of March 7th). The sporadic
1038 recovered case numbers are mostly in the single digit. If we use the sudden increase
1039 of recovered case matching with the sudden explosive increase of new infected, the
1040 distance is approximately 14 days, a reasonable T value when compared to the
1041 mean value in China. For our data analysis, we used daily newly cases starting
1042 February 19th, for the derivative of individual distribution study; we used data case
1043 from March 1st, for the NR ratio study, in order to have enough recovered cases.



1044
1045
1046
1047

Figure S7: The derivative of the logarithmic value of daily new infected case distribution.



1048
1049
1050
1051

Figure S8: The *NR* ratio from 7 days of data from March 1st to 7th. The estimated zero-crossing time would occur between March 11th and 12th, a value consistent with the statistics from the daily new case distribution on March 10th.