

Optimal non-pharmaceutical intervention policy for Covid-19 epidemic via neuroevolution algorithm

Arash Saeidpour^{a,b} and Pejman Rohani^{a,b,c,d,*}

^aOdum School of Ecology, University of Georgia, Athens, GA, USA, 30602

^bCenter for the Ecology of Infectious Diseases, University of Georgia, Athens, GA, USA, 30602

^cDepartment of Infectious Diseases, University of Georgia, Athens, GA, USA, 30602

^dCenter for Influenza Disease & Emergence Research (CIDER), Athens, Georgia, USA

*rohani@uga.edu

January 4, 2022

Abstract

National responses to the Covid-19 pandemic varied markedly across countries, from business-as-usual to complete shutdowns. Policies aimed at disrupting the viral transmission cycle and preventing the overwhelming of healthcare systems inevitably exact an economic toll. We developed an intervention policy model that comprised the relative human, implementation and healthcare costs of non-pharmaceutical epidemic interventions and identified the optimal strategy using a neuroevolution algorithm. The proposed model finds the minimum required reduction in transmission rates to maintain the burden on the healthcare system below the maximum capacity. We find that such a policy renders a sharp increase in the control strength during the early stages of the epidemic, followed by a steady increase in the subsequent ten weeks as the epidemic approaches its peak, and finally the control strength is gradually decreased as the population moves towards herd immunity. We have also shown how such a model can provide an efficient adaptive intervention policy at different stages of the epidemic without having access to the entire history of its progression in the population. This work emphasizes the importance of imposing intervention measures early and provides insights into adaptive intervention policies to minimize the economic impacts of the epidemic without putting an extra burden on the healthcare system.

1 Introduction

- 2 On March 11, 2020 the World Health Organization (WHO) announced that Covid-19,
- 3 caused by severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) [1], "can be
- 4 characterized as a pandemic" [2]. Within a month, most countries around the world had
- 5 taken public health measures to contain the spread of the novel virus [3]. However, the

6 type and severity of implemented measures and their subsequent success in minimizing
7 the public health impacts of the outbreak varied greatly by country [4]. This variation
8 in policies and their effectiveness reflects the complexity of finding the balance between
9 two often competing policy objectives: protecting the public's health versus minimizing
10 the economic impact of intervention measures [5].

11 Initially, without access to pharmaceuticals, studies focused on two distinct control
12 approaches: mitigation and suppression [6–8]. The mitigation strategy aims to reduce
13 transmission such that healthcare systems are not overwhelmed, while aiming to
14 maintain the chain of transmission in order to achieve herd immunity. In contrast, the
15 suppression strategy is aimed at virus elimination. In hindsight, countries that acted
16 early to suppress the disease have excelled at minimizing both the public health and
17 economic impact of the epidemic [9–11]. While early suppression measures appear to
18 outperform the mitigation strategy both in terms of public health goals and economic
19 costs, such policies would not necessarily be successful in countries where citizens are
20 more averse to government-enforced control and surveillance measures [12]. Moreover,
21 suppression measures would only be successful if implemented in the early stages of the
22 epidemic and sufficiently strictly as to curtail transmission effectively. In a number of
23 settings, however, suppression has been implemented in a piece-meal manner, leading to
24 periods of drastic interventions including lockdowns punctuated by relaxation of social
25 distancing measures and subsequent uptick in transmission [13, 14]. This prompted us
26 to examine the optimal mitigation strategy, which aims to manage or mitigate the
27 healthcare impacts of the epidemic while population approaches herd immunity.

28 Characterizing immediate and long-term economic, social and human burden of
29 Covid-19 epidemic is challenging and has led to several research efforts to examine the
30 optimal intervention policy from various perspectives. It is unfeasible to review
31 comprehensively this body of work, so we confine ourselves to a number of the key
32 studies. Rowthorn and Maciejowski [15] investigated the optimal uniform lockdown in
33 an *SIR* model assuming a variety of parameterizations [15]. Their objective function
34 assigned monetary values to costs arising from infection, lockdown, and value of life.
35 Their main finding was that in the medium term, a policy that maintains effective
36 reproduction number value close to 1 provides the best path. Bethune and Korinek [16]
37 contrasted the decisions made by rational, individual agents with the choices made by a
38 social planner who is able to coordinate the choices of individuals [16]. They found that
39 rational agents generate large externalities because they fail to internalize the effects of
40 their economic and social activities on others' risk of infection. Alvarez *et al.* formalized
41 the social planner's dynamic control using an *SIR* epidemiological model and a linear
42 economy. The best strategy starts with a severe lockdown two weeks after the epidemic,
43 covers 60% of the population after a month, and progressively decreases to 20% of the
44 population after three months. More recently, a number of studies have broadened this
45 exploration to identify age-specific optimal control strategies [17, 18].

46 Inspired by [19–21], we sought to use a neuroevolution strategy to finding the
47 optimal policy function which would dynamically determine the minimal required
48 reduction in transmission rates at each time instant, deemed as *control strength*
49 hereafter. Reductions in transmission may result from lower contacts (due to
50 isolation-in-place ordinances, movement restrictions, or lockdown policies), or the
51 adoption of personal protective measures that serve to curtail transmission upon contact
52 (such as the use of face masks and PPE), with varying societal impact. The fitness
53 function is expressed such that a strategy is rewarded for allowing the epidemic to
54 remove individuals from the susceptible pool without overwhelming the healthcare
55 capacity. The proposed neuroevolution strategy begins by initializing a population of
56 random policy functions. The generated policy functions are then used to simulate the
57 trajectory of the epidemic. The fitness of each function is then evaluated based on the

58 specified reward function. The most elite policy functions are then perturbed (mutated)
 59 to generate the next generation (offspring). The new population is then evaluated and
 60 this process is repeated for a pre-defined number of iterations. We also derived the
 61 optimal control solution via Pontryagin's maximum principle (PMP) [22] and compared
 62 the results with the optimal neuroevolution policy.

63 We have chosen the United Kingdom as our target population to implement the
 64 proposed approach. The choice of the UK as our target population was largely
 65 motivated by the frequent changes in the government's strategy to contain the
 66 epidemic [23], as summarized in Figure 1. The UK's initial response was a mitigation
 67 policy, majorly inspired by the response to the flu pandemic, with an emphasis on
 68 protecting the most vulnerable to avoid overburdening the healthcare system in an
 69 effort to achieve herd immunity [9]. This initial policy later changed to a suppression
 70 policy by implementing lock-downs and imposing face mask-wearing requirements.
 71 Looking back at the early days of the epidemic, this study aims to understand how an
 72 effective mitigation policy could have been implemented (see [9] for a comparison of
 73 initial responses to Covid-19 by different countries including United Kingdom).

74 Our study explores mechanisms for "flattening the curve" – it is motivated by the
 75 COVID-19 pandemic but need not be restricted to precise courses of action undertaken
 76 in the response to this pandemic event. Our findings are intended to be informative for
 77 future epidemic control, particularly at the early stages of an epidemic when there may
 78 be no effective pharmaceuticals in sight.

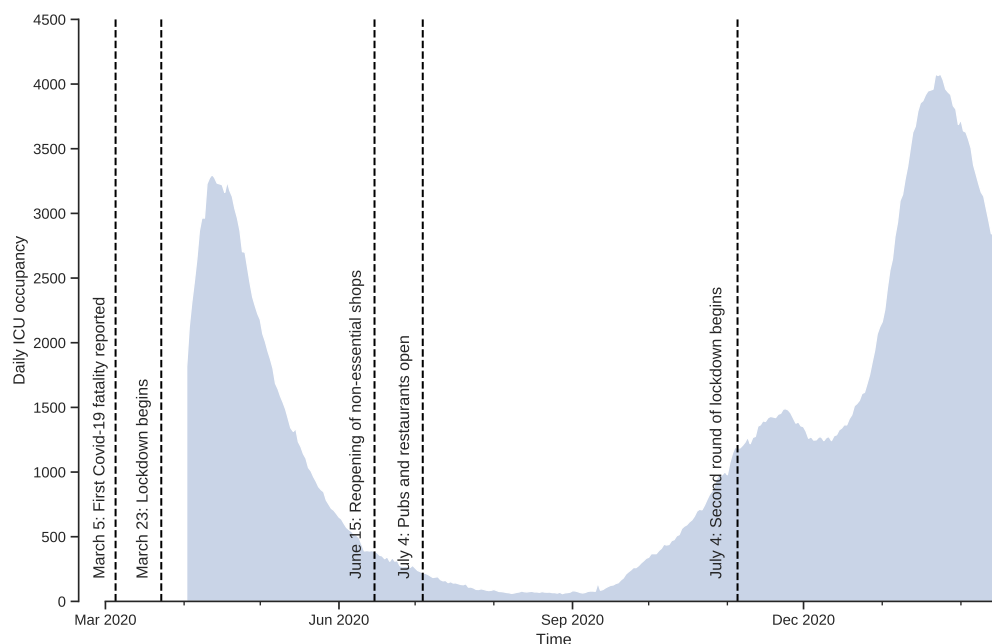


Fig 1. Number of Covid-19 patients in intensive care (ICU) and timeline of lockdowns in the UK.

79 We find that the ideal intervention policy results in a rapid increase in control
 80 strength early in the epidemic, followed by a sustained increase over the next ten weeks
 81 as the epidemic reaches its peak, and ultimately a progressive drop in control strength
 82 as the population achieves herd immunity. We have also shown how, without having
 83 access to the complete history of the epidemic's growth in the population, such a model
 84 may give an effective adaptive intervention policy at various stages of the epidemic.
 85 This study highlights the significance of implementing control measures as promptly as

86 possible and offers insights into adaptive intervention strategies aimed at reducing the
 87 economic effect of epidemics while avoiding undue strain on the healthcare system.

88 Materials and methods

89 Model structure

90 We used a deterministic, time-varying
 91 Susceptible-Exposed-Infectious-Recovered-Hospitalized in ICU (*SEIRH*) model [24] to
 92 characterize the transmission dynamics in the UK as described in Eqs. 1–5:

$$\dot{S} = \frac{dS}{dt} = -(1 - c(t)) \frac{\beta SI}{N} \quad (1)$$

$$\dot{E} = \frac{dE}{dt} = (1 - c(t)) \frac{\beta SI}{N} - \rho E \quad (2)$$

$$\dot{I} = \frac{dI}{dt} = \rho E - \gamma I - P_{Detection} \sigma_{ICU} \gamma_{ICU_{Delay}} I \quad (3)$$

$$\dot{R} = \frac{dR}{dt} = \gamma I + \gamma_{ICU_{Stay}} H \quad (4)$$

$$\dot{H} = \frac{dH}{dt} = P_{Detection} \sigma_{ICU} \gamma_{ICU_{Delay}} I - \gamma_{ICU_{Stay}} H \quad (5)$$

93 where β is the transmission rate, $1/\rho$ and $1/\gamma$ give the mean latent and infectious
 94 periods, respectively and $c(t) \in [0, 1]$ is the reduction in transmission (such that $c(t) = 1$
 95 signifies complete cessation of transmission). The state variable $H(t)$ denotes the
 96 number of occupied ICU beds and is determined by the probability that an infection is
 97 detected ($P_{Detection}$), the fraction of cases that require ICU treatment (σ_{ICU}) and the rate
 98 of admission to the ICU ($\gamma_{ICU_{Delay}}$). The mean duration of stay in the ICU is determined
 99 by $1/\gamma_{ICU_{Stay}}$. Model parameters and chosen values are presented in Table 1.

100 In our analyses, we examine changes in optimal intervention policy assuming policies
 101 are implemented starting at different points during the epidemic, T_0 . To identify the
 102 appropriate initial conditions at these different starting points, we used a particle
 103 filter [25] to estimate the effective retrospective daily $c(t)$ (where $t = 0, \dots, T_0$), thus
 104 yield the epidemiological state of the population at different stages of the epidemic. The
 105 agreement between our fitted *SEIRH* model and data is shown in Figure S2.

106 The Reward function

107 As discussed by [35], there is precedent for integrating modeling methodologies and
 108 health economic analyses to inform public health intervention decisions based on a
 109 willingness to pay for each Quality-Adjusted Life Year (QALY) saved [21, 36, 37]. Such
 110 an approach allows for allocating explicit monetary values to each term in the reward
 111 function [21]. While some cost-benefit analysis via this approach has been carried out in
 112 relation to Covid-19 [35], the pandemic's enormous scope renders traditional economic
 113 measurements largely impractical. As a result, a health-economic approach is not the
 114 main emphasis of this study. Instead, in order to capture the general societal impacts of
 115 pandemic mitigation efforts, we have employed a simple *relative* economic cost to
 116 formulate the reward function.

117 We first introduce the following multi-objective reward function to account for three
 118 opposing goals: i) Sustain viral transmission to achieve herd immunity, ii) Keep the ICU
 119 occupancy below the maximum capacity, and iii) Impose the minimum possible control:

Table 1. Parameters of SEIRH model

Parameter	Definition	Value	Source
N	Total population size	66,436,000	[26]
R_0	Basic reproduction number	2.3	[27, 28]
$1/\gamma$	Mean infectious period (days)	2.9	[27, 28]
$1/\rho$	Mean latent period (days)	3.4	[29]
β	Mean transmission rate (1/day)	0.793	Estimated
$P_{Detection}$	Ratio of confirmed cases to total infections	0.3	[30]
σ_{ICU}	Proportion of confirmed cases that end up in ICU	0.05	[31]
$1/\gamma_{ICU_{Delay}}$	Median time from symptoms onset to ICU admission (days)	10	[32]
$1/\gamma_{ICU_{Stay}}$	Mean ICU stay period (days)	9	[33]
H_{max}	Number of ICU beds	4074	[34]

The table presents the parameters of SEIRH model used to model the dynamics of Covid-19 transmission in the population of UK.

$$\begin{aligned}
 r_1(t) &= \alpha_1 r_1(t)_{Herd\ Immunity} - \alpha_2 r_1(t)_{Exceedance} - \alpha_3 c(t)^2 \\
 &= \alpha_1 E(t)/N - \alpha_2 (H(t) - H_{max})/H_{max} - \alpha_3 * c(t)^2.
 \end{aligned} \tag{6}$$

120 We defined $r_1(t)$ for the sake of mathematical simplicity in deriving PMP solution
 121 and it is only used to compare the optimal NPI policies obtained from neuroevolution
 122 and PMP methods. For the remainder of this study, we use a slightly different objective
 123 function, $r_2(t)$, defined as follows:

$$\begin{aligned}
 r(t) &= \alpha_1 r_{Herd\ Immunity}(t) + \alpha_2 r_{Exc}(t) + \alpha_3 r_{Control}(t), \\
 &= \alpha_1 (R(t)/N) - \alpha_2 Relu((H(t) - H_{max})/N) - \alpha_3 * c(t).
 \end{aligned} \tag{7}$$

124 In both reward functions (equations (6) & (7)), the terms α_1 , α_2 and α_3 modulate
 125 the relative importance of herd immunity, healthcare burden and societal costs,
 126 respectively. The goal, therefore, is to identify the optimal intervention function $c(t)$
 127 that maximizes the sum of rewards, J , during the course of the epidemic:

$$\max_{c(t)} J = \int r_i(t) dt, i \in 1, 2. \tag{8}$$

128 Pontryagin's maximum principle (PMP)

129 In this section we first derive the necessary conditions for optimal control via
 130 Pontryagin's maximum principle, and describe the iterative numerical algorithm (the
 131 forward-backward sweep method) used to find the optimal solution. First, we form the
 132 following Hamiltonian function:

$$\mathcal{H}(t, \mathbf{s}(t), c(t), \lambda_{\mathbf{s}}(t)) = r(t) + \lambda_S(t)\dot{S} + \lambda_E(t)\dot{E} + \lambda_I(t)\dot{I} + \lambda_R(t)\dot{R} + \lambda_H(t)\dot{H}. \tag{9}$$

133 Here, $\lambda_{\mathbf{s}}(t)$ are adjoint functions satisfying the adjoint system:

$$\dot{\lambda}_{\mathfrak{s}}(t) = -\frac{\partial \mathcal{H}(t, \mathfrak{s}^*(t), c^*(t), \lambda_{\mathfrak{s}}^*(t))}{\partial \mathfrak{s}}, \mathfrak{s} \in \{S, E, I, R, H\}, \quad (10)$$

$$\lambda_{\mathfrak{s}}(T) = 0 \text{ (Transversality condition)}. \quad (11)$$

Expanding equation 10 yields:

$$\dot{\lambda}_S(t) = -\partial \mathcal{H} / \partial S(t) = (\lambda_S - \lambda_E) \frac{(1-c)\beta I}{N} \quad (12)$$

$$\dot{\lambda}_E(t) = -\partial \mathcal{H} / \partial E(t) = (\lambda_E - \lambda_I)\rho - \frac{\alpha_1}{N} \quad (13)$$

$$\begin{aligned} \dot{\lambda}_I(t) = -\partial \mathcal{H} / \partial I(t) = & (\lambda_E - \lambda_S) \frac{(1-c)\beta SI}{N} + (\lambda_I - \lambda_R)\gamma + \\ & (\lambda_I - \lambda_H)\gamma_{ICU\text{Delay}} P_{\text{Detection}} \sigma_{ICU} \end{aligned} \quad (14)$$

$$\dot{\lambda}_R(t) = -\partial \mathcal{H} / \partial R(t) = 0 \quad (15)$$

$$\dot{\lambda}_H(t) = -\partial \mathcal{H} / \partial H(t) = (\lambda_H - \lambda_R)\gamma_{ICU\text{Stay}} + \frac{\alpha_2}{H_{max}}. \quad (16)$$

134 The necessary conditions for the optimal control is obtained by maximizing the
135 Hamiltonian (equation 9) with respect to $c(t)$:

$$\frac{\partial \mathcal{H}}{\partial c} = 0 \text{ at } c^*t \rightarrow c^*(t) = (\lambda_S - \lambda_E) \frac{\beta I}{2\alpha_3 N}, c^*(t) \in [0, 1] \quad (17)$$

136 The state equations (equations 1-5) and adjoint equations (equations 10-16) together
137 with state initial conditions and transversality conditions (equation 11) form the
138 *Optimality system*. The explicit solution can not be analytically derived. Thus we
139 turned to an iterative numerical method, *Forward-backward Sweep*, to solve the
140 *Optimality system*.

141 Neuroevolution algorithm

142 The optimal policy function, π_{θ} , is a feed forward neural network, parameterized by θ
143 which takes the state of the system at current time t , $\{S(t), E(t), I(t), R(t)\}$ as input
144 and returns the control strength, $c(t)$. The neuroevolution strategy aims to find the
145 optimal policy function, $\mathcal{P}_{\text{Most elite}}^G$, with highest fitness score. Fitness score of policy
146 function j in generation i , f_j^i , is equal to the sum of rewards, J (equation 8) and is
147 obtained by running the *SEIRH* model with the corresponding policy function. First,
148 M policy functions (\mathcal{P}_j^1) are randomly initialized. For each policy function, a trajectory
149 is rolled out and fitness score is calculated at the end of simulation, as shown in figure 2.
150 The L policy functions with the highest fitness scores are mutated to generate the next
151 generation of policy functions. Mutation is implemented by adding a random Gaussian
152 noise, scaled by the mutation rate, σ , to θ parameters of elite policy functions. The new
153 offspring policy functions served as the parents of next generation. This process
154 continues to find a policy function with a sufficiently high fitness score, $\mathcal{P}_{\text{Most elite}}^G$. We
155 used a fully-connected feed-forward network with 3 16-unit hidden layers and one tanh
156 output layer to model the policy function. Pseudocode for the neuroevolution algorithm
157 used in this study is provided in Algorithm 1.

158 Results

159 Which optimization algorithm?

160 We compared the optimal intervention policies obtained from PMP and neuroevolution
161 policies (Fig S1). The policies are obtained using the r_1 reward function (equation 6)

Algorithm 1 Neuroevolution algorithm

Require: Population size M , Number of generations G , Elite population size L , Mutation rate σ
Initialize M policy functions, \mathcal{P}_j^1 , with random initial weights θ_j^1
for i **do**=1 to G . # Iterate G generations
 for j **do**=1 to M
 $f_j \leftarrow$ Roll out a trajectory by running the model using \mathcal{P}_j^i # Fitness score
 end for
 Sort θ_j^i by f_j in descending order
 $\theta_{Elite}^i = \{\theta_j^i | j < L\} \cup \theta_{Most\ elite}^{i-1}$
 for j **do**=1 to M
 Draw sample $t \sim U(1, L)$ # Select a parent
 Draw sample $\epsilon \sim \mathcal{N}(0, 1)$ # Gaussian noise
 $\theta_j^{i+1} = \theta_t^i + \sigma\epsilon$ # Mutate
 end for
end for
return $\mathcal{P}_{Most\ elite}^G$

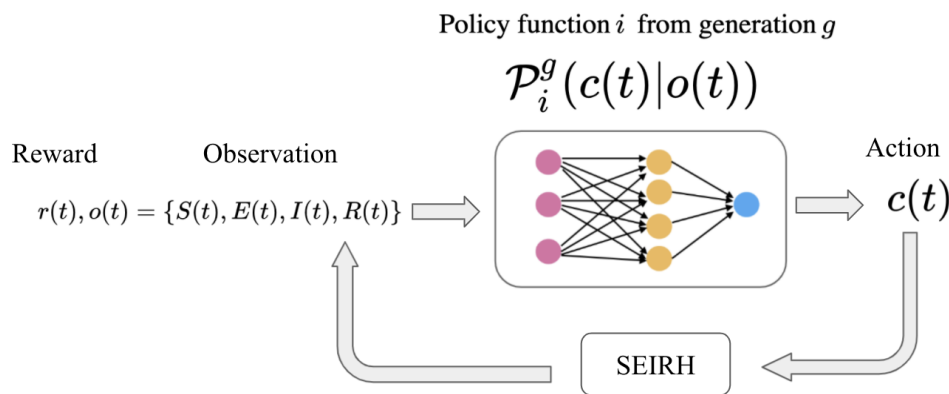


Fig 2. Schematic representation of policy function \mathcal{P}_i^g , represents the policy function i of generation g . The L most elite policy functions of each generation are mutated to generate the M policy functions of next generation.

162 with $\alpha_2 = 1e - 1, \alpha_3 = 5e - 3$ and same initial conditions. We found the optimal
163 policies obtained from both methods to be very similar. In simpler problems where an
164 analytic solution can be obtained for the optimality system, the PMP method can
165 provide more insights about the optimal control solution and the dynamics of the
166 system. Otherwise, a neuroevolutionary approach is computationally advantageous since
167 the resulting policy function provides an optimal strategy for a broad range of initial
168 conditions at a substantially smaller computations cost. That is, the PMP optimal
169 intervention for a given initial condition is obtained by solving the boundary-value
170 problem formulated in equations (1-5) and (10-16). For a new boundary condition, the
171 numerical solution must be repeated to solve the new boundary-value problem. In the
172 remainder of the paper, our optimal solutions are obtained via the neuroevolutionary
173 approach.

174 Reward function exploration

175 The relative economic burden of different objectives in the reward function is
 176 determined by the weights, $\{\alpha_1, \alpha_2, \alpha_3\}$. Thus, we examined the effects of variation in
 177 these parameters on the resulting optimal policy (see Figure S3). We constrained α_1 to
 178 be 1 and changed the values of α_2 and α_3 over a logarithmic grid. For each parameter
 179 set, we trained the neuroevolution algorithm for 2000 generations with a population size
 180 of 256. The resulting policy functions (purple lines) and corresponding ICU occupancy
 181 trajectories of the 10 best-performing agents for each parameter set are depicted in
 182 Figure S3. We found the reward function to be consistently robust to variation in the
 183 values of α_2 . That is, the tested range of α_2 values makes the cost of ICU overflow
 184 sufficiently prohibitive, leading to high-fitness strategies ensuring ICU maximum
 185 capacity is not exceeded (note that the ICU overflow reward is equal to 0 while the ICU
 186 occupancy is below the maximum capacity and negative otherwise). Evidently, making
 187 α_2 smaller would eventually deprioritize the goal of maintaining the ICU occupancy
 188 below the limit. Without loss of generality, we will use $\alpha_2 = 1e9$ in the remainder of
 189 this paper. In contrast, we found the reward function to be highly sensitive to variation
 190 in α_3 . For $\alpha_3 > 10^{-4}$, the relative cost (negative reward) of imposing control becomes
 191 prohibitive and leads to one of the extreme intervention strategies: Suppression policy
 192 to end the endogenous transmission at the earliest possible time and avoid imposing
 193 lengthy control measures; or a no-intervention policy which plainly leads to the
 194 minimum relative control cost. In practice, the inclination for a specific intervention
 195 strategy depends on the policy maker's priorities. We observed pronounced variation in
 196 the optimal policies and resulting ICU occupancy trajectories for smaller values of α_3
 197 (compare the first and third columns, Fig. S3). In Figure S4, we demonstrate this
 198 variation for each parameter set and across the values of α_3 . As shown in Fig. S4A,
 199 values of α_3 smaller than 10^{-4} result in greater *Cumulative herd immunity reward*.
 200 Thus, when the relative cost of control is modest, the optimal policy function will tend
 201 to maximize the reward by increasing the number of individuals removed from the
 202 susceptible pool, which in turn leads to greater *Cumulative control reward* (Fig. S4B)
 203 and longer epidemic duration (Fig. S4C). Therefore, among the tested values,
 204 $\alpha_3 = 1e - 4$ represents the middle ground between prolonged intervention and
 205 suppression policies, and is the value that we have used in the rest of this paper.

206 No-intervention policy, uniform intervention policy and optimal 207 policy

208 Figure 3 presents a comparison between the optimal intervention policy identified via
 209 our neuroevolution algorithm, a uniform intervention policy and no-intervention policy.
 210 The uniform intervention policy is implemented by imposing a constant reduction in
 211 transmission throughout the epidemic, $c(t) = c_u$. The value of control strength, c_u , is
 212 estimated such that the peak ICU occupancy tangents the maximum capacity. Figure
 213 3A depicts the ICU occupancy trajectories of these three policies. As expected, the
 214 no-intervention policy leads to ICU burdens well beyond the threshold capacity for more
 215 than two months (67 days). The other notable observation is the difference between the
 216 optimal and uniform policies in managing the ICU burden: the optimal policy
 217 maintains the ICU occupancy near the maximum capacity throughout the epidemic, but
 218 not beyond it. Figure 3B depicts the implemented control strength in time for optimal
 219 and uniform policies. Except for a period of time less than 10 weeks at the onset of the
 220 epidemic, the control strength of the optimal policy is below the uniform intervention
 221 policy. The difference in the imposed control between two policies is better illustrated
 222 by Figure 3C, where a widening gap between the cumulative imposed control of the
 223 two policies emerges after day 200. In Figure 3D, we present the recovered individuals

224 for each policy. Unlike the optimal policy, the final fraction of recovered individuals in
 225 the uniform intervention policy case is well below the theoretical herd immunity
 226 threshold. This suggests that the any reduction in the control strength, could lead to
 227 another epidemic wave given the large fraction of susceptible individuals.

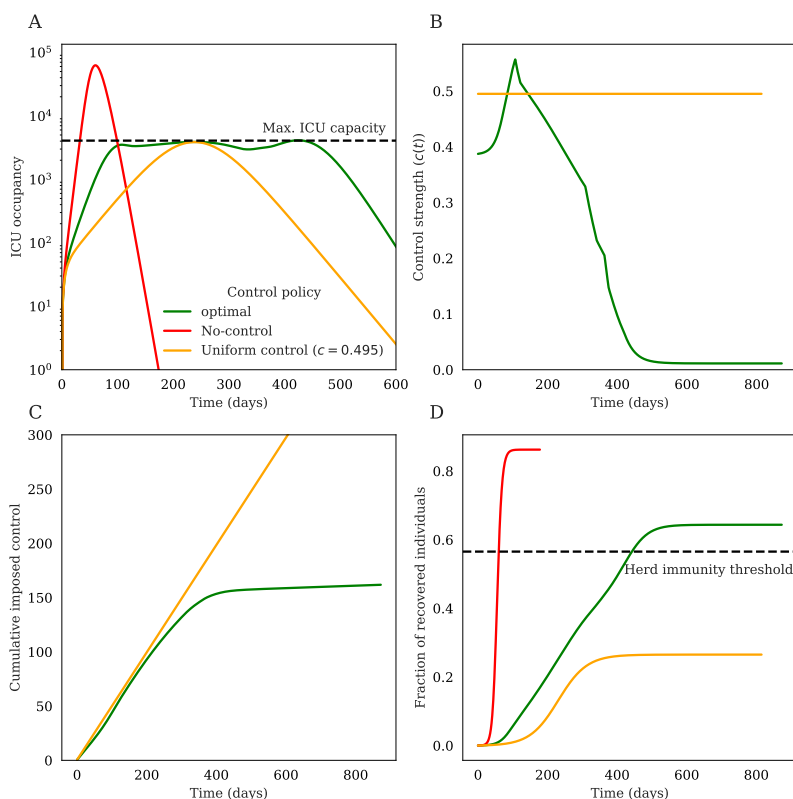


Fig 3. No-intervention policy, Uniform intervention policy and optimal policy The figure presents the (A) ICU occupancy (B) Control strength (C) Cumulative imposed control and (D) recovered individuals for three different policies: No-intervention policy, Uniform intervention policy and optimal policy.

228 The sooner the better

229 We have estimated the optimal intervention policy initiated at different stages of the
 230 epidemic, as shown in Figure 4. Each scenario corresponds to a particular start date for
 231 the roll out of the optimal intervention policy. Figure 4A depicts the scenario in which
 232 optimal intervention policy starts on March 1st, which coincides with a surge in cases in
 233 the UK. The optimal intervention policy starts with $c(t) = 0.33$ (a 33% reduction in
 234 transmission rates) and is gradually increased to $c(t) = 0.54$ by mid-May. The control
 235 strength tapers off to 0 by June 2021. This scenario leads to two peaks in ICU
 236 occupancy, in November 2020 and June 2021. Figures 4B-E depict the optimal
 237 intervention policy starting at intermediate stages of the epidemic. As mentioned above,
 238 we estimated the initial conditions for each scenario by fitting our *SIER* model to
 239 fatality data using particle filtering, a Monte Carlo likelihood estimation algorithm for
 240 hidden state-space dynamical systems [38]. Comparing the optimal intervention policy
 241 curves in different scenarios depicts how implementing transmission reduction measures
 242 at earlier stages of the epidemic will eventually shorten the epidemic: The termination

243 of optimal intervention policy is delayed from June 2021 (in Figure 4A) to February
 244 2022 (in Figure 4D). The only exception is Figure 4E, in which the optimal intervention
 245 policy terminates slightly sooner than in Figure 4D. This is most likely due to the
 246 emergence of new variants with higher transmissibility [39] which gave rise to a faster
 247 depletion of the susceptible pool than accounted for in our model.

248 To better illustrate the importance of implementing early control measures, we have
 249 demonstrated the *Total duration of intervention policy implementation* and *Cumulative*
 250 *imposed control* for different scenarios in Figure 5. The *Total duration of intervention*
 251 *policy implementation* represents the time period between March 1st 2020 and the
 252 termination date of intervention policy for each scenario. The *Cumulative imposed*
 253 *control* is obtained by summing the daily implemented control strength ($c(t)$), divided
 254 by total number of days with $c(t) > 0$ for each scenario. As shown in Figure 5A, the
 255 *Total duration of intervention policy implementation* increases from 442 days in the first
 256 columns to 700 days in the last one. Figure 5B also confirms the fact that implementing
 257 the optimal intervention policy from earlier stages of epidemic would reduce the overall
 258 required control measures. Note that depicted *Cumulative imposed control* values do not
 259 include the actual imposed control strength ($c(t)$) before the start of optimal
 260 intervention policy and adding those values would only widen their differences. Also,
 261 the *Cumulative imposed control* is a linear measure of overall imposed control, however,
 262 the actual economic cost would not necessarily change linearly with duration and
 263 strength of imposed intervention policy.

264 Finding the balance

265 Figure 6 paints an overall picture of how the optimal policy fine tunes the transmission
 266 rates to sustain endogenous transmission in the population without overburdening the
 267 ICU capacity. Figure 6A demonstrates the variation of effective reproductive ratio
 268 (R_{eff}) throughout the epidemic (black line), the control strength is also shown (blue
 269 dashed line). At the onset of the epidemic, R_{eff} is instantly reduced to 1.52 from 2.3 by
 270 imposing a 0.33 reduction in contact rates ($c(t) = 0.33$) and further decreased to
 271 $R_{eff} \approx 1$ by mid-may (point i) to stall the epidemic growth. From point i to point ii,
 272 The R_{eff} is maintained close to 1 to maintain the ICU occupancy close to the maximum
 273 capacity. At this point, $c(t)$ is slightly increased which leads to a sharp decrease of R_{eff}
 274 to 0.89 in point iii. This is followed by a steep decrease in $c(t)$ to bring the R_{eff} above 1
 275 to sustain the transmission. To summarize, the optimal mitigation policy is achieved by
 276 finding the balance between two extreme scenarios: Suppression policy which aims to
 277 stall the endogenous transmission in the population, and "No-intervention" which leads
 278 to exponential epidemic growth and the overburdening of healthcare capacity.

279 Discussion

280 More than eighteen months into the SARS-CoV-2 pandemic, it is becoming increasingly
 281 clear that countries that implemented suppression strategies early on experienced
 282 greater success in managing both the public health and economic burden of the
 283 epidemic [9–11]. However, such strategies work best when employed early in the
 284 epidemic, when number of cases is relatively small. Moreover, in countries where
 285 government-imposed restrictions are not well received by the public, implementation of
 286 such policies will be challenging. Looking back at the early stages of the epidemic, our
 287 work provides a dynamic mitigation strategy that sustains the community transmission
 288 without overwhelming the healthcare capacity.

289 A number of previous studies on optimal non-pharmaceutical interventions have
 290 used quadratic cost expressions for the control term in the cost function [18, 40, 41].

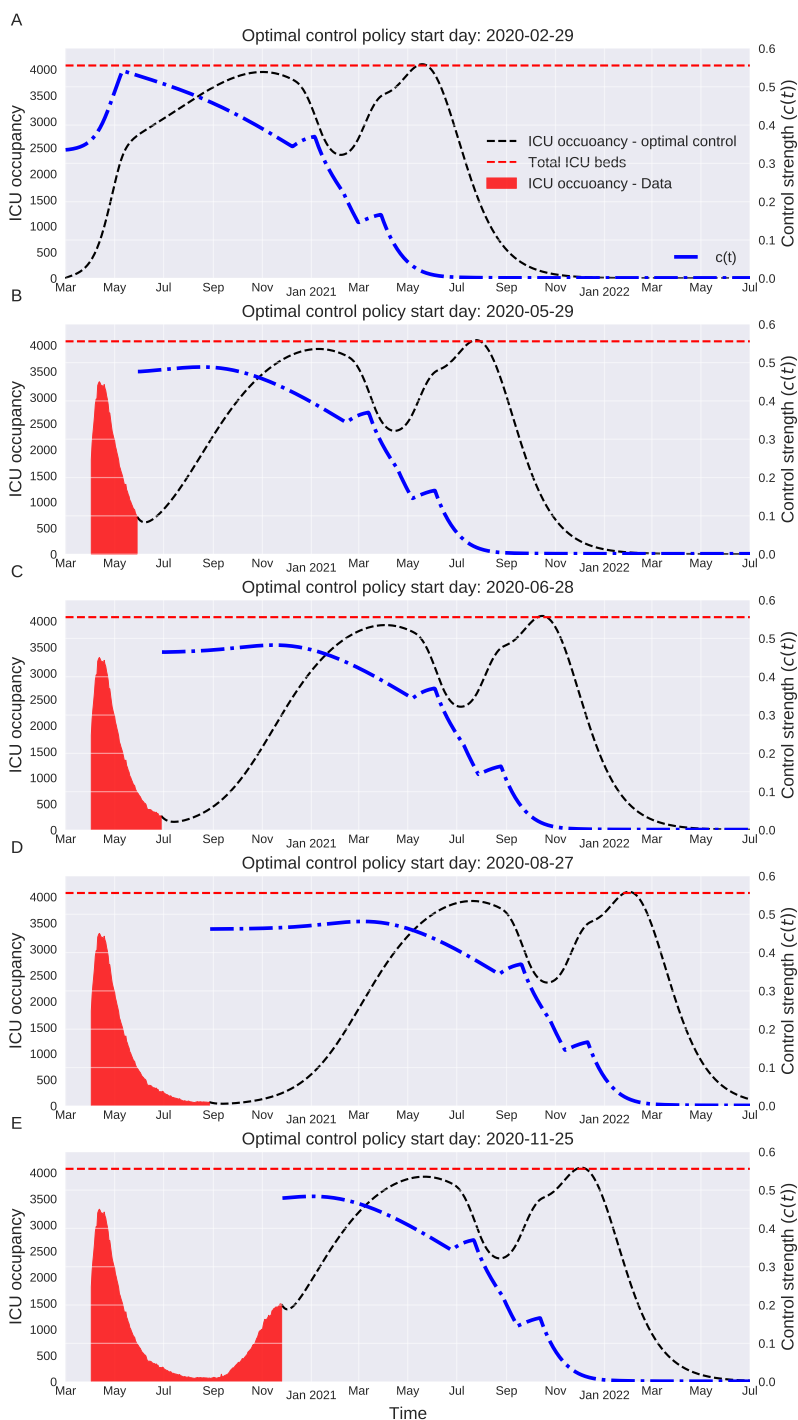


Fig 4. Optimal intervention policy at different stages of epidemic The figure depicts the optimal intervention policy starting at different stages of epidemic. For each scenario, the number of susceptible, exposed, infectious and recovered individuals is estimated from a *SEIRH* model fitted to the UK fatality data and used as initial condition to derive the optimal intervention policy.

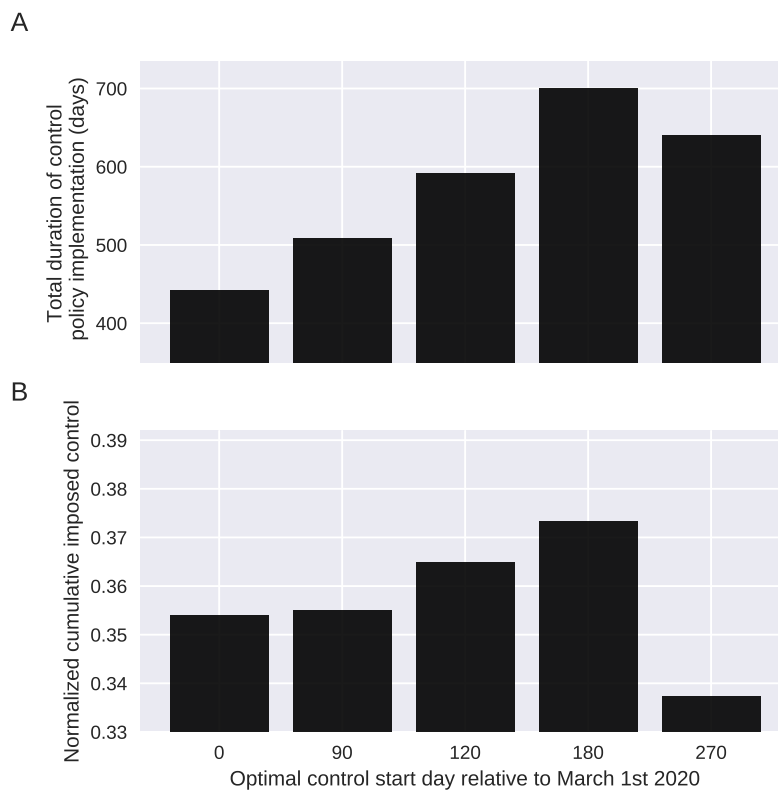


Fig 5. Implementing the optimal intervention policy will reduce the overall impact of control measures the *Total duration of intervention policy implementation* and *Cumulative imposed control* for different scenarios. The *Total duration of intervention policy implementation* represents the time period between March 1st 2020 and termination date of intervention policy for each scenario. The *Cumulative imposed control* is obtained by adding up the implemented control strength ($c(t)$) in each day, divided by total number of days with $c(t) > 0$ for each scenario.

291 This is mainly because when the cost function is quadratic with respect to the control,
 292 the differential equations arising from the necessary conditions for an optimal control
 293 have a known solution. Other functional forms frequently provide difficult-to-solve
 294 systems of differential equations. To circumvent this, we employed a neuroevolution
 295 algorithm which enabled us also to explore non-quadratic functions. The neuroevolution
 296 algorithm was used to train a policy function that takes the epidemiological state of
 297 population (the numbers of susceptible, exposed, infectious and recovered individuals)
 298 on each time day and provides the corresponding control strength. We defined a
 299 multi-objective reward function to account for three conflicting goals: Sustain the
 300 transmission to achieve herd immunity when suppression is not feasible, maintaining the
 301 ICU occupancy below the maximum capacity and imposing minimum possible control
 302 measures to reduce the contact rates. A relative weighting parameter was assigned to
 303 corresponding terms of each of these objectives in the reward function. The sensitivity
 304 analysis indicated that the resulting policy function is highly sensitive relative weighting
 305 of the control term and found a optimal range of of values for it. We chose United
 306 Kingdom as our target population and fitted an *SEIRH* model to fatality data to
 307 estimate the initial conditions at different stages of the epidemic.

308 The optimal intervention policy confirmed the importance of early interventions to

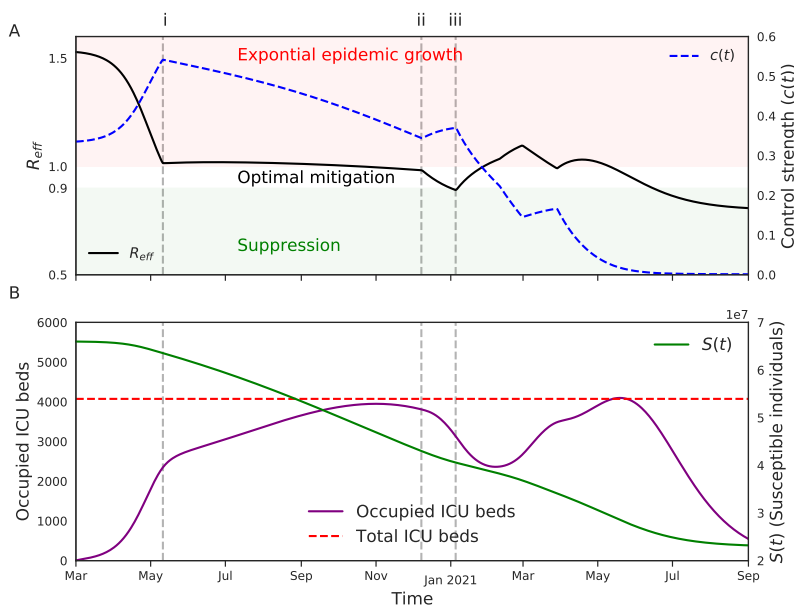


Fig 6. The optimal intervention policy maintains the effective reproductive ratio (R_{eff}) close to 1: The figure displays the changes in effective reproductive ratio when implementing the optimal intervention policy. The control strength ($c(t)$) is sharply increased at early stages of epidemic to stall the epidemic growth and keep healthcare capacity from being overwhelmed. The R_{eff} is maintained close to 1 by gradually reducing the $c(t)$ as the size of susceptible pool shrinks. Once the value of R_{eff} reaches below 0.9, $c(t)$ is increased to sustain the transmission in the population, while keeping the occupied ICU beds below the maximum capacity.

309 reduce the contact rates in the population, as highlighted in the previous studies [15,41].
 310 An initial 34% reduction in transmission at the onset of the epidemic, gradually
 311 increasing to 50% in the next 10 weeks is required to bring the R_{eff} near 1. After that,
 312 the restrictions are constantly decreased as the the size of susceptible pool diminishes.
 313 The association between the control strength and the size of the susceptible pool
 314 (except the first initial 10 weeks) highlights the importance of reliable and widespread
 315 serosurveys in order to inform policy decision making.

316 Our study highlights the neuroevolution algorithm, a gradient-free approach, as an
 317 efficient alternative to traditional PMP method for finding the optimal
 318 non-pharmaceutical intervention policy in dynamical disease transmission system. Past
 319 studies have demonstrated that in many challenging reinforcement learning tasks
 320 neuroevolution algorithm rivals (or even outperforms in some domains) state-of-the art
 321 gradient-based methods such as Q-learning and A3C [20]. Interestingly, the
 322 forward-backward sweep technique that we used obtain the optimal solution via PMP
 323 closely resembles the backpropagation, the algorithm used to train the gradient-based
 324 reinforcement learning methods [42]. Ultimately, we found the neuroevolution algorithm
 325 to be computationally advantageous to the PMP method as the former algorithm
 326 provides the optimal intervention policy for a broad range of initial values after initial
 327 training (as shown in Fig. 4) while the numerical solution to obtain the optimal control
 328 via PMP must be repeated for a new initial condition.

329 A key component of our neuroevolution algorithm is the assumption that the full
 330 epidemiological state of the population is observable at each time step. In reality,
 331 however, the observable data provide an incomplete and potentially biased picture of

332 epidemiology since they are based on reported incidence, hospitalization and fatality
 333 data in addition to seroprevalence surveys. Besides assuming complete epidemiological
 334 information, our approach also assumed that the optimal intervention policy is
 335 implemented in deterministically; that is, the output action is perfectly implemented at
 336 each time instant and the resulting new state given the corresponding action is always
 337 the same - something that is not practical. An important next step in this area would
 338 be to extend our novel framework to identify the optimal intervention strategies with
 339 hidden states in a stochastic setting. Furthermore, while this study addresses the
 340 optimal reduction in the contact rates over time, the economic cost and effectiveness of
 341 various non-pharmaceutical intervention mechanisms [43,44] to achieve the optimal
 342 policy reduction requirements must also be examined.

343 Acknowledgments

344 Research reported in this publication was supported by the National Institute Of
 345 General Medical Sciences of the National Institutes of Health under Award Number
 346 R01GM123007. The content is solely the responsibility of the authors and does not
 347 necessarily represent the official views of the National Institutes of Health.

References

1. Zhu N, Zhang D, Wang W, Li X, Yang B, Song J, et al. A novel coronavirus from patients with pneumonia in China, 2019. *New England journal of medicine*. 2020;382:727–733. doi:10.1056/NEJMoa2001017.
2. WHO Director-General’s opening remarks at the media briefing on COVID-19 - 11 March 2020;. Available from: <https://www.who.int/director-general/speeches/detail/who-director-general-s-opening-remarks-at-the-media-briefing-on-covid-19---11-march-2020>.
3. COVID-19 GOVERNMENT RESPONSE TRACKER;. Available from: <https://www.bsg.ox.ac.uk/research/research-projects/covid-19-government-response-tracker>.
4. Brauner JM, Mindermann S, Sharma M, Johnston D, Salvatier J, Gavenčiak T, et al. Inferring the effectiveness of government interventions against COVID-19. *Science*. 2021;371(6531).
5. Brett TS, Rohani P. Transmission dynamics reveal the impracticality of COVID-19 herd immunity strategies. *Proceedings of the National Academy of Sciences*. 2020;117(41):25897–25903.
6. Prem K, Liu Y, Russell TW, Kucharski AJ, Eggo RM, Davies N, et al. The effect of control strategies to reduce social mixing on outcomes of the COVID-19 epidemic in Wuhan, China: a modelling study. *The Lancet Public Health*. 2020;5(5):e261–e270.
7. Walker PG, Whittaker C, Watson OJ, Baguelin M, Winskill P, Hamlet A, et al. The impact of COVID-19 and strategies for mitigation and suppression in low-and middle-income countries. *Science*. 2020;369(6502):413–422.
8. Davies NG, Kucharski AJ, Eggo RM, Gimma A, Edmunds WJ, Jombart T, et al. Effects of non-pharmaceutical interventions on COVID-19 cases, deaths, and demand for hospital services in the UK: a modelling study. *The Lancet Public Health*. 2020;5(7):e375–e385.

9. Hassan I, Mukaigawara M, King L, Fernandes G, Sridhar D. Hindsight is 2020? Lessons in global health governance one year into the pandemic. *Nature Medicine*. 2021;27(3):396–400.
10. Dong E, Du H, Gardner L. An interactive web-based dashboard to track COVID-19 in real time. *The Lancet infectious diseases*. 2020;20(5):533–534.
11. Kocharczyk M, Lipniacki T. Pareto-based evaluation of national responses to COVID-19 pandemic shows that saving lives and protecting economy are non-trade-off objectives. *Scientific reports*. 2021;11(1):1–9.
12. Anderson RM, Heesterbeek H, Klinkenberg D, Hollingsworth TD. How will country-based mitigation measures influence the course of the COVID-19 epidemic? *The lancet*. 2020;395(10228):931–934.
13. Ragonnet-Cronin M, Boyd O, Geidelberg L, Jorgensen D, Nascimento FF, Siveroni I, et al. Genetic evidence for the association between COVID-19 epidemic severity and timing of non-pharmaceutical interventions. *Nature Communications*. 2021;12(1):2188. doi:10.1038/s41467-021-22366-y.
14. Hollingsworth TD, Klinkenberg D, Heesterbeek H, Anderson RM. Mitigation Strategies for Pandemic Influenza A: Balancing Conflicting Policy Objectives. *PLoS Computational Biology*. 2011;7(2):e1001076. doi:10.1371/journal.pcbi.1001076.s001.
15. Rowthorn R, Maciejowski J. A cost–benefit analysis of the COVID-19 disease. *Oxford Review of Economic Policy*. 2020;36(Supplement_1):S38–S55.
16. Bethune ZA, Korinek A. Covid-19 infection externalities: Trading off lives vs. livelihoods. National Bureau of Economic Research; 2020.
17. Acemoglu D, Chernozhukov V, Werning I, Whinston MD. Optimal targeted lockdowns in a multi-group SIR model. National Bureau of Economic Research; 2020.
18. Richard Q, Alizon S, Choisy M, Sofonea MT, Djidjou-Demasse R. Age-structured non-pharmaceutical interventions for optimal control of COVID-19 epidemic. *PLoS computational biology*. 2021;17(3):e1008776.
19. Salimans T, Ho J, Chen X, Sidor S, Sutskever I. Evolution strategies as a scalable alternative to reinforcement learning. arXiv preprint arXiv:170303864. 2017;.
20. Such FP, Madhavan V, Conti E, Lehman J, Stanley KO, Clune J. Deep neuroevolution: Genetic algorithms are a competitive alternative for training deep neural networks for reinforcement learning. arXiv preprint arXiv:171206567. 2017;.
21. Riolo MA, Rohani P. Combating pertussis resurgence: One booster vaccination schedule does not fit all. *Proceedings of the National Academy of Sciences of the United States of America*. 2015;112(5):E472 – E477. doi:10.1073/pnas.1415573112.
22. Pontryagin LS. *Mathematical theory of optimal processes*. CRC press; 1987.
23. Cameron-Blake E, Tatlow H, Wood A, Hale T, Kira B, Petherick A, et al. Variation in the response to COVID-19 across the four nations of the United Kingdom. Blavatnik School of Government working paper series. 2020;.

24. Keeling MJ, Rohani P. *Modelling Infectious Diseases: In Humans and Animals*. Princeton University Press. Princeton University Press; 2008.
25. Shrestha S, King AA, Rohani P. Statistical inference for multi-pathogen systems. *PLoS Computational Biology*. 2011;7(8):e1002135. doi:10.1371/journal.pcbi.1002135.
26. Park N. Population estimates for the UK, England and Wales, Scotland and Northern Ireland, provisional: mid-2019; 2020. Available from: <https://www.ons.gov.uk/peoplepopulationandcommunity/populationandmigration/populationestimates/bulletins/annualmidyearpopulationestimates/mid2018#:~:text=The%20population%20of%20the%20UK,any%20year%20since%20mid%2D2004>.
27. Li Q, Guan X, Wu P, Wang X, Zhou L, Tong Y, et al. Early transmission dynamics in Wuhan, China, of novel coronavirus-infected pneumonia. *New England journal of medicine*. 2020;.
28. Zhang J, Litvinova M, Wang W, Wang Y, Deng X, Chen X, et al. Evolving epidemiology and transmission dynamics of coronavirus disease 2019 outside Hubei province, China: a descriptive and modelling study. *The Lancet Infectious Diseases*. 2020;20(7):793–802.
29. Li R, Pei S, Chen B, Song Y, Zhang T, Yang W, et al. Substantial undocumented infection facilitates the rapid dissemination of novel coronavirus (SARS-CoV-2). *Science*. 2020;368(6490):489–493.
30. Giattino C. How epidemiological models of COVID-19 help us estimate the true number of infections; 2020. Available from: <https://ourworldindata.org/covid-modelse>.
31. Coronavirus (Covid-19) in the UK; 2020. Available from: <https://coronavirus.data.gov.uk/details/healthcare>.
32. Wang D, Hu B, Hu C, Zhu F, Liu X, Zhang J, et al. Clinical characteristics of 138 hospitalized patients with 2019 novel coronavirus-infected pneumonia in Wuhan, China. *Jama*. 2020;323(11):1061–1069.
33. Grasselli G, Zangrillo A, Zanella A, Antonelli M, Cabrini L, Castelli A, et al. Baseline characteristics and outcomes of 1591 patients infected with SARS-CoV-2 admitted to ICUs of the Lombardy Region, Italy. *Jama*. 2020;323(16):1574–1581.
34. Critical Care Bed Capacity and Urgent Operations; 2020. Available from: <https://www.england.nhs.uk/statistics/statistical-work-areas/critical-care-capacity/>.
35. Moore S, Hill EM, Dyson L, Tildesley MJ, Keeling MJ. Modelling optimal vaccination strategy for SARS-CoV-2 in the UK. *PLoS computational biology*. 2021;17(5):e1008849.
36. Kim JJ, Goldie SJ. Cost effectiveness analysis of including boys in a human papillomavirus vaccination programme in the United States. *Bmj*. 2009;339.
37. Hill EM, Petrou S, de Lusignan S, Yonova I, Keeling MJ. Seasonal influenza: Modelling approaches to capture immunity propagation. *PLoS computational biology*. 2019;15(10):e1007096.

38. Doucet A, Johansen AM, et al.. A tutorial on particle filtering and smoothing: Fifteen years later; 2009.
39. Davies NG, Abbott S, Barnard RC, Jarvis CI, Kucharski AJ, Munday JD, et al. Estimated transmissibility and impact of SARS-CoV-2 lineage B.1.1.7 in England. *Science (New York, NY)*. 2021;doi:10.1126/science.abg3055.
40. Lee S, Chowell G, Castillo-Chávez C. Optimal control for pandemic influenza: the role of limited antiviral treatment and isolation. *Journal of Theoretical Biology*. 2010;265(2):136–150.
41. Djidjou-Demasse R, Michalakis Y, Choisy M, Sofonea MT, Alizon S. Optimal COVID-19 epidemic control until vaccine deployment. *MedRxiv*. 2020;.
42. LeCun Y, Touresky D, Hinton G, Sejnowski T. A theoretical framework for back-propagation. In: *Proceedings of the 1988 connectionist models summer school*. vol. 1; 1988. p. 21–28.
43. Liu Y, Morgenstern C, Kelly J, Lowe R, Jit M. The impact of non-pharmaceutical interventions on SARS-CoV-2 transmission across 130 countries and territories. *BMC medicine*. 2021;19(1):1–12.
44. Courtemanche C, Garuccio J, Le A, Pinkston J, Yelowitz A. Strong Social Distancing Measures In The United States Reduced The COVID-19 Growth Rate: Study evaluates the impact of social distancing measures on the growth rate of confirmed COVID-19 cases across the United States. *Health Affairs*. 2020;39(7):1237–1246.