

# JUNE

## open-source individual-based epidemiology simulation

Joseph Bullock<sup>1,2,\*</sup>, Carolina Cuesta-Lazaro<sup>1,3,\*</sup>, Arnau Quera-Bofarull<sup>1,3,\*</sup>, Miguel Icaza-Lizaola<sup>1,3,\*\*</sup>, Aidan Sedgewick<sup>1,4,\*\*</sup>, Henry Truong<sup>1,2,\*\*</sup>, Aoife Curran<sup>1,3</sup>, Edward Elliott<sup>1,3</sup>, Tristan Caulfield<sup>7</sup>, Kevin Fong<sup>8,9</sup>, Ian Vernon<sup>1,6</sup>, Julian Williams<sup>5</sup>, Richard Bower<sup>1,3</sup>, and Frank Krauss<sup>1,2,+</sup>

<sup>1</sup>Institute for Data Science, Durham University, Durham DH1 3LE, UK

<sup>2</sup>Institute for Particle Physics Phenomenology, Durham University, Durham DH1 3LE, UK

<sup>3</sup>Institute for Computational Cosmology, Durham University, Durham DH1 3LE, UK

<sup>4</sup>Centre for Extragalactic Astronomy, Durham University, Durham DH1 3LE, UK

<sup>5</sup>Institute for Hazard, Risk & Resilience, Durham University, Durham DH1 3LE, UK

<sup>6</sup>Department of Mathematical Sciences, Durham University, Durham DH1 3LE, UK

<sup>7</sup>Department of Computer Science, University College of London, London WC1E 6BT, UK

<sup>8</sup>Department of Science, Technology, Engineering and Public Policy, University College London, London WC1E 6BT, UK

<sup>9</sup>Department of Anaesthesia, University College London Hospital, London NW1 2BU, UK

\*Equal contribution

\*\*Equal contribution

+Corresponding author: [frank.krauss@durham.ac.uk](mailto:frank.krauss@durham.ac.uk)

**Abstract:** We introduce JUNE, an open-source framework for the detailed simulation of epidemics on the basis of social interactions in a virtual population constructed from geographically granular census data, reflecting age, sex, ethnicity, and socio-economic indicators. Interactions between individuals are modelled in groups of various sizes and properties, such as households, schools and workplaces, and other social activities using social mixing matrices. JUNE provides a suite of flexible parameterisations that describe infectious diseases, how they are transmitted and affect contaminated individuals. In this paper we apply JUNE to the specific case of modelling the spread of COVID-19 in England. We discuss the quality of initial model outputs which reproduce reported hospital admission and mortality statistics at national and regional levels as well as by age strata.

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>The Structure of the JUNE Modelling Framework</b>	<b>4</b>
<b>3</b>	<b>Population and its Static Properties</b>	<b>5</b>
<b>4</b>	<b>Simulating Social Interactions</b>	<b>10</b>
<b>5</b>	<b>Infection Modelling: Spreading and Health Impact</b>	<b>15</b>
<b>6</b>	<b>Mitigation Policies and Strategies</b>	<b>22</b>
<b>7</b>	<b>Discussion of model outputs</b>	<b>27</b>
<b>8</b>	<b>Fitting via Bayesian Emulation</b>	<b>30</b>
<b>9</b>	<b>Summary</b>	<b>32</b>
<b>A</b>	<b>Algorithms</b>	<b>33</b>
<b>B</b>	<b>Contact matrices</b>	<b>36</b>
<b>C</b>	<b>Details on modelling health trajectories</b>	<b>41</b>
<b>D</b>	<b>Calibration via Bayes Linear Emulation and History Matching</b>	<b>41</b>

# 1 Introduction

The spread of SARS-CoV-2 in populations with largely no immunological resistance, and the associated COVID-19 disease, have caused considerable disruption to health care systems and a large number of fatalities around the globe. The assessment of policy options to mitigate the impact of this and other epidemics on the health of individuals, and the efficiency of healthcare systems, relies on a detailed understanding of the spread of the disease, and requires both short-term operational forecasts and longer-term strategic resource planning.

There are various modelling approaches which aim to provide insights into the spread of an epidemic. They range from analytic models, formulated through differential or difference equations, which reduce numerous aspects of the society–virus–disease interaction onto a small set of parameters, to purely data-driven parametrizations, often based on machine learning, which inherently rely on a probability density that has been fitted to the current and past state of the system in an often untraceable way. As another class of approaches, agent-based models (ABMs) are “*particularly useful when it is necessary to model the disease system in a spatially-explicit fashion or when host behavior is complex[.]*” [1] p2:5<sup>1</sup>. Being the traditional tool of choice to analyse behavioural patterns in society, they find ample use in understanding and modelling the observed spread of infections and in leveraging this for intermediate and long-term forecasting [3, 4, 5]. Such models also provide the flexibility to experiment with different policies and practices, founded in realistic changes to the model structure, such as the inclusion of new treatments, changes in social behaviour, and restrictions on movement.

To simulate pandemics, specific realisations of ABMs, individual-based models (IBMs), have been developed in the past two decades, for example [6, 7]. In these models, the agents represent individuals constituting a population, usually distributed spatially according to the population density and with the demographics - age and sex - taken from census data. Within the existing taxonomy of agent-based models in epidemiology, see for instance [8, 9], these models often use a disease-specific modelling framework. Interactions between individuals in predefined social settings, systematically studied for the first time in [10], provide the background for disease spread, formulated in probabilistic language and dependent on the properties of the individuals and the social setting. The sociology of the population and the transmission dynamics are constrained separately using external datasets and available literature, and connected in the description of the spread of the disease. Calibration of such models to observed disease outcomes, such as hospital admission and mortality rates, is therefore reduced to the specific interface between the disease and the varying physiology across the broad population. Policy interactions and mitigation strategies can be flexibly encoded in detail as modifications of the social setting, and allow precise analysis of their efficacy that is not readily available in other approaches.

Evidence from disease data such as COVID-19 fatality statistics suggests that case and infection fatality rates are correlated, amongst other factors, to the age and socio-economics status of the population exposed to the etiological agent [11]. This necessitates the construction of a model with exceptional social and geographic granularity to exploit highly local heterogeneities in the demographic structure. In this publication we introduce a new individual-based model, JUNE<sup>2</sup>, a generalisable modular framework for simulating the spread of infectious diseases with a fine-grained geographic and demographic resolution and a strong focus on the detailed simulation of policy interventions. JUNE reaches a geographic resolution of societal factors similar to models that focus on single-site infection models, such as [12], where space, location and distance are carefully modelled. In addition, similar to approaches such as the STHAM model [13], the individuals in JUNE follow detailed spatio-temporal activity profiles that are informed by available data including time surveys, geographic and movement data. In contrast to such models that are usually constrained to a few tens of thousands of agents, JUNE simulates, simultaneously, the full population of a country in its spatio-temporal setting, and how a disease spreads through its population mediated by contacts between individuals. JUNE allows for flexible and precise parameterisations of policies that affect groups of individuals selected according to any of their characteristics. This allows modelling of policies to mitigate the further spread of a disease, realised as changes and restrictions on movement, to which we add the effectiveness of changes in social behaviour such as social distancing.

<sup>1</sup>Indeed, many models also feature some optimising behaviour of individuals as artificial intelligence-type actors against randomly drawn welfare functions, see for example [2].

<sup>2</sup>A full open source code base and implementation examples are linked here:  
github: <https://github.com/IDAS-Durham/JUNE>  
and pypi: <https://pypi.org/project/june/>.

The major cost for this level of detail in the model is in computational load; indeed, models such as JUNE would likely not have been possible prior to 2010 without using a prohibitive amount of computing power, see for instance [14].

As a first application of JUNE, we model the spread of COVID-19 in England. In this context, JUNE uses census, household composition, and workplace data to ensure that each of the 53 million people in England are assigned a specific, identifiable location at any point in time. Their activities, health, age and other demographic attributes are then modelled at a fine-grained geographical level, which helps to ensure that local heterogeneity in population and movement characteristics are well recovered. This societal structure, generated by the model, is validated against a series of datasets (among others this includes: surveys of household size and composition, location and size of businesses, size and type of schools by region). The calibration to observed data from the actual spread of SARS-CoV-2 is then limited to how the virus is transmitted in the community through person-to-person ‘contacts’ (in the sense of sufficient proximity and timing to transmit). This component of the infection is calibrated to the spatio-temporal development of hospitalisations and casualties during the COVID-19 outbreak in England, starting in early March 2020. Preliminary observations demonstrate that a detailed large-scale model of this type has important implications for intermediate- to long-term modelling of the SARS-CoV-2 spread in the UK and elsewhere.

The remainder of this paper is as follows. Section 2 provides an overview of the structure of the JUNE framework. In Section 3, we detail the construction of a virtual population including a variety of demographic attributes. For the example case of England, we demonstrate that the constructed population reproduces the distributions of age, gender, ethnicity, socioeconomic indices, and the composition of the households they live in, all with a granularity of a few hundred people. The static properties of the population also include the assignment of students to schools and Universities and of employment in companies dis-aggregated by 21 industry sectors. In Section 4 we discuss the dynamics of the population model. We demonstrate how JUNE correctly reproduces the average time-profile of daily activities of individuals in England. We also describe in detail how we reconstruct movement and daily commute patterns based on publicly available data. Social interactions in various settings are modelled through parameters informed by social mixing matrices derived from surveys such as *PolyMod* [10] and the BBC Pandemic project [15]. In contrast to other models, JUNE also incorporates interactions in various social venues such as pubs, restaurants, cinemas and shopping, outside the more structured settings of households, work places and schools. Section 5 introduces the generalisable disease model with specific applications to COVID-19 — its transmission properties and the impact it has on infected individuals. We employ a probabilistic model for the former, while for the latter we incorporate data from the UK and other countries to characterise the journey of infected people through the healthcare system. In Section 6 we describe how JUNE models the impact of various policy interventions and other mitigation strategies. In Section 7 we show some first indicative results of JUNE highlighting its potential for future, more detailed studies. Section 8 introduces our approach to fitting the model using Bayesian Emulation. We summarize our work in Section 9, and conclude the paper with discussion of future work and improvements to the model.

## 2 The Structure of the JUNE Modelling Framework

The JUNE framework is built on four interconnected layers: **population**, **interactions**, **disease** and **policy**, the layers and their interfaces are illustrated in Fig. 1. In the context of this publication, we focus on the application of JUNE to England’s population, the spread of the COVID-19 disease, and policies that have been enacted by the UK Government in 2020. Clearly, a different population with different behavioural patterns will not only affect the distribution of individuals according to their personal characteristics, but it will also necessitate the adaptation of, e.g., social venues to these patterns and corresponding changes to the **population** and **interactions** layers. Similarly, modifications to the **disease** layer will allow application of the JUNE framework for a different disease or, possibly, even a range of competing diseases. This flexibility and adaptability is even more pronounced in the **policy** layer where the introduction of new policies in reaction to an epidemic depends on behavioural patterns or societal norms.

The **population** layer encodes the individuals in the model and constructs static social environments such as the households they live in, the schools and universities they study in, and the workplaces where

they work. The construction of the virtual population is informed by demographic data such as age, sex, and ethnicity distributions, the geographic location of their residence, and its composition. Depending on their age, individuals will attend school or university, work, or be retired.

The **interaction** layer models the social interactions of individuals, based on data about the frequency and intensity of contacts with other people in social settings. In addition to daily patterns of regular interactions with fixed groups of individuals such as household members, students and teachers in schools, and work colleagues, the **interaction** layer also models more randomised interactions. These include daily commute patterns to and from work, and more dynamic activities such as visits to restaurants, pubs, cinemas, and visits to other households.

The **disease** layer, which sits on top of the **population** and **interaction** layers, models the characteristics of disease transmission and the effects it has on those infected. In terms of disease transmission, the model incorporates the varying susceptibility of individuals, how likely individuals are to become infected when they mix with others in various locations, and how infectious they are over the course of their infection. In terms of disease progression, the model captures how likely individuals are to experience symptoms with varying severity, to be hospitalised, to be admitted to intensive care, or to die, as well as the timings associated with these events.

In response to the spread of a disease through its population, a government might introduce policy measures designed to control and reduce the impact of the disease. In the case of COVID-19 in England and many other countries, policies have included social distancing measures, the closure of schools, shops, restaurants, and other leisure venues, and restrictions on movement. In **JUNE**, these are modelled in the **policy** layer. The high level of detail present in the **population** and **interaction** layers allows policies to be modelled at a corresponding granularity. This enables **JUNE** to describe the impact of policies that can be applied to specific geographical regions, to specific venues or sectors, or to individuals with specific characteristics. Examples include, but are not restricted to, the closure of targeted different types of (or even singular) venues, the inclusion or exclusion of specific age groups when going to school, shielding of the older parts of the population, modification to inter-household visits, and self-isolation measures for infected individuals and their contacts, including variations of compliance with these measures.

### 3 Population and its Static Properties

**JUNE** creates a detailed virtual population at the individual level through its **population** layer, utilising a cross-section of demographic and geographic information. Since **JUNE** relies on multiple datasets, and is built to dynamically adapt to varying types of input, the approaches described in this section are generalisable to other settings with similar or complementary data availability. Given that different settings, e.g. countries, may have different methods and types of data collection, many of the input parameters described here are optional, allowing **JUNE** to be more easily adapted to differences in reporting.

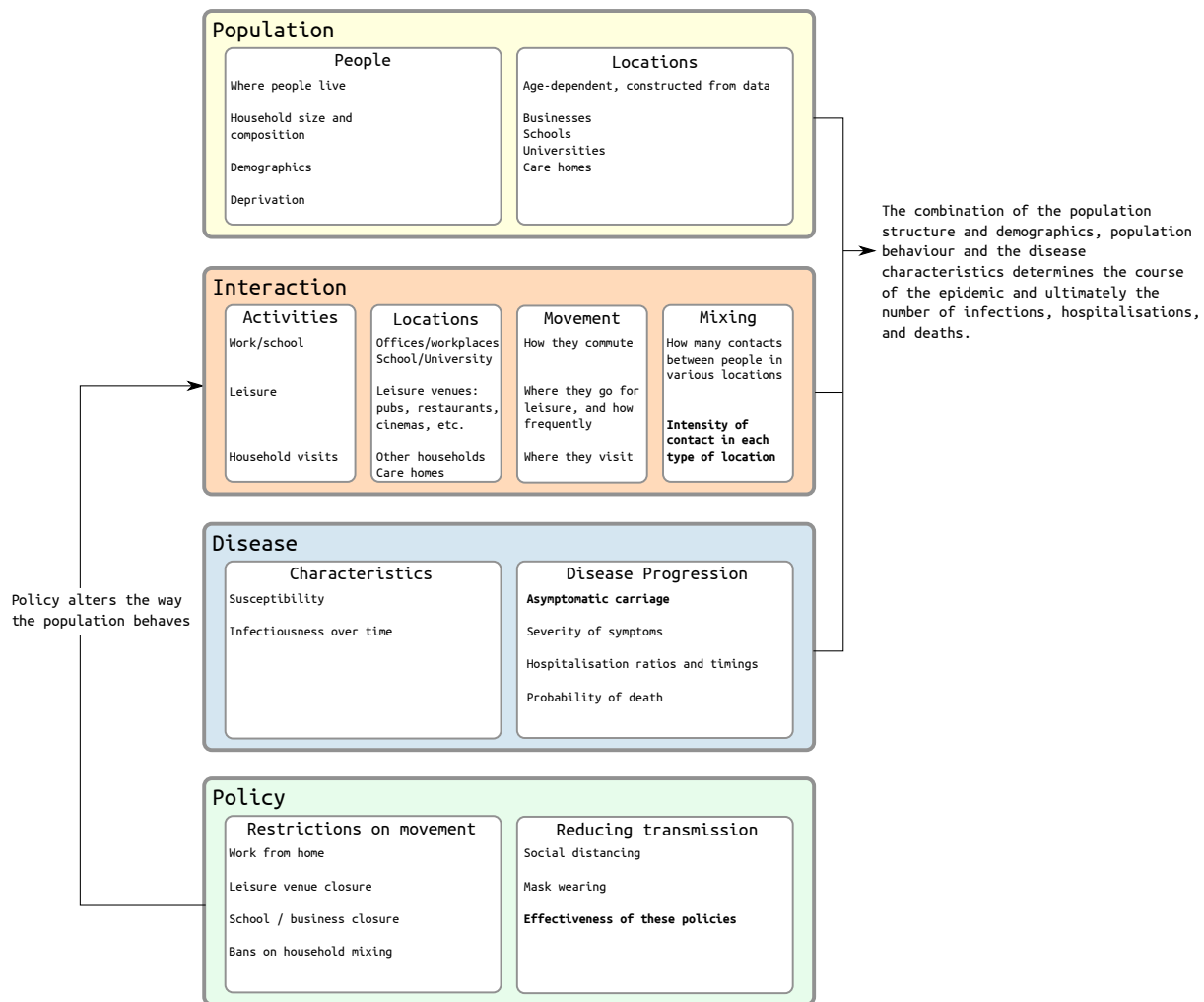
#### 3.1 Geography and demography

To facilitate generalisability across multiple settings, **JUNE** models the geographic distribution of a population using a hierarchy of three layers – regions, super areas and areas. Layering these geographies allows the use of data at different levels of aggregation and enables simple statistical projections of data between these levels.

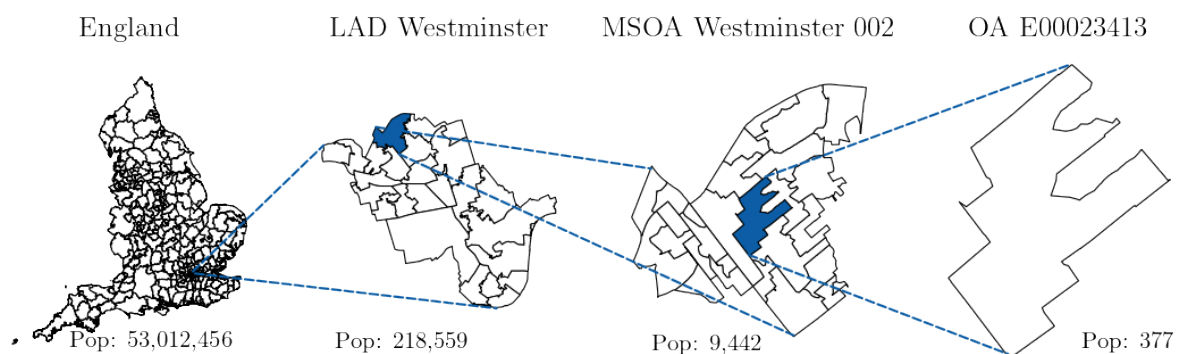
For the case of England, the construction of the virtual population in **JUNE** is largely based on data from the latest UK census, which was carried out in 2011. This data is accessible through **NOMIS**, an open-access database provided by the Office for National Statistics (ONS), and each dataset varies in its degree of aggregation. The three hierarchical geographical layers represented in Fig. 2 are:

1. regions – London, East Midlands, West Midlands, the North West, the North East, etc.;
2. super areas – approximately 7,200 middle layer super output areas (MSOAs);
3. areas – approximately 180,000 output areas (OAs).

The individuals in **JUNE**'s virtual population are constructed according to age and sex dis-aggregated information, the minimal information required by **JUNE**. In the case of England, the ONS census data provides this information at the OA (area) level [16, 17] such that **JUNE** naturally captures the population

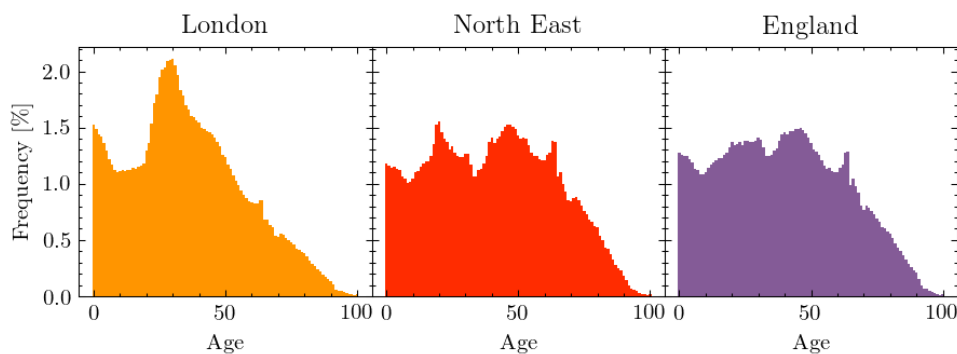


**Figure 1:** Overview of the structure of JUNE. Fitted parameters are shown in bold.



**Figure 2:** Graphical representation of how the census data for England are structured, from the level of local authority districts (LAD), down to the level of output areas (OA), with middle layer super output area (MSOA) in between.

density at the most fine-grained level. In Fig. 3 we show age distributions in different regions. We use data derived from the ONS to additionally assign one of five broad ethnic categories to individuals based on their age, sex and location of residence [18] and follow a similar procedure for the socio-economic



**Figure 3:** Age profiles in different regions of England, taken from the ONS database and implemented in JUNE

index, which we divide into centiles, according to the ranked English Index of Multiple Deprivation (IMD) [19].

### 3.2 Household construction

Once virtual individuals living in an area have been initialised, we cluster them into households. Given that the household setting can serve as a key location for disease transmission, it is important that their population construction is well represented. Depending on available data, households in JUNE can be constructed with a high degree of granularity, taking account of multiple demographic attributes.

The ONS census provides detailed records of both household type and composition in England at the OA (area) level. It decomposes households into the following broad categories: single, couple, family, student, communal, and other [20], and further specifies them by the number of old adults, aged over 65, adults, dependent adults (such as students), and children, to arrive at approximately 20 distinct classes. However, it is impossible to recover the exact composition for each household type. For example, the number of non-dependent children (people over the age of 18 living with their parents), the number of multi-generational families, and the exact distribution of adult groups sharing a household are not specified in these data-sets<sup>3</sup>. Given this lack of complete data, we populate the households so as to iteratively build up realistic household compositions, giving preference to those types for which we have the most precise data. Appendix A.1 details the procedures followed to seed and construct households in JUNE according to ONS data.

Since a large fraction of the vulnerable population resides in care homes, it is also important to model their composition correctly. We use census data to determine the number of communal establishment residents at the OA (area) level, for which care homes are considered a specific type [21]. The ONS also collects information on the age and sex distribution of residents of communal establishments at the MSOA (super area) level [22]. By combining these datasets, we infer the age and sex distribution of the care home population.

We also include all the other communal establishments specified in the census, including student accommodations and prisons. We currently do not model explicitly the age and sex distribution of these other communal establishments, however, since the age and sex distribution of the area’s population is biased towards the communal residents, their resident characteristics are realistic.

### 3.3 Construction of virtual schools and universities

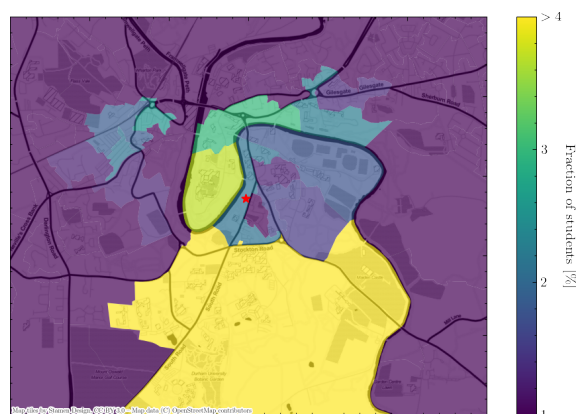
Having constructed our virtual population, including their residences, we address schools and universities by first constructing them on their spatial coordinates, and then allocating nearby children and students accordingly.

JUNE initialises schools according to their geo-coordinates and the age ranges in attendance at each school. Students are sent to one of the  $n$  nearest schools to their place of residence, according to which

<sup>3</sup>We will discuss how we use secondary data to constrain these, and their impact on the spread of COVID-19, in a forthcoming publication.

schools cater for their age. We form year groups which include all students of the same age. The formation of year groups, and classes within them, allows JUNE to control mixing within and between children of different ages in the school environment.

When modelling schools in England, we use data provided by [23] to determine the location of schools and their age brackets. We assume that children between the ages 0-19 can attend school, with mandatory attendance between 5-18. Since 19 year-olds can attend school, university, work or none of these, the institution they attend is determined by the number of vacancies in schools accepting students of that age group. We send children to one of the  $n = 10$  nearest schools where classes sizes are limited to 40. One way in which we validate our assumptions is by comparing average travel distance to schools of different types. In JUNE, we find 1.7 km and 5.0 km for primary and secondary school students, compared with 2.6 km and 5.5 km respectively from the 2014 national travel survey [24]. Teachers are chosen from the population in a manner similar to the case of medics working in hospitals, see below. The number of teachers assigned to a particular school, and therefore the number of classes, is determined by sampling the ratio of students to teachers from a Poisson distribution with mean equal to the UK national average, separately for primary (mean of 21) and secondary schools (mean of 16), or a random choice of the two for mixed schools [25]. JUNE's recovered student-teacher ratios are 22.0 and 17.8 for primary and secondary schools respectively.



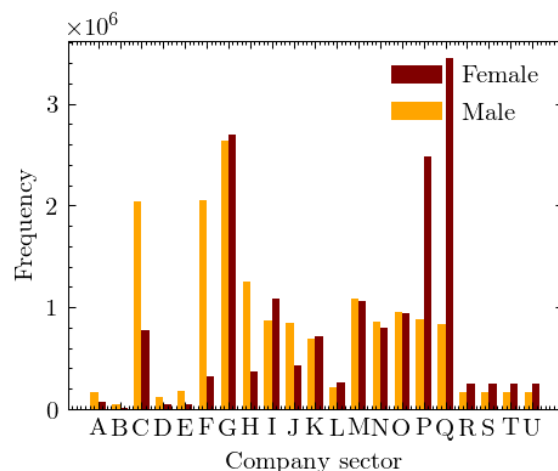
**Figure 4:** A geographical visualisation of the location of student residences in Durham in JUNE, with the university location represented as a red star in the middle. Output areas are colour-coded according to the fraction of students they host. Note that the large southern area is where most of the university accommodation blocks are located.

We construct universities according to their address as recorded in the UK Register of Learning Providers (UKRLP) [23]. This dataset also provides data on the number of students attending each university, which we use to assign a subset of the local population to the university, reflecting the fact that the ONS census uses the term-time address of students. Students are sampled from adults between the ages 18-25 with a preference given to those previously assigned to living in student or communal households in a given radius around the university. Fig. 4 shows an example of a university city, Durham, in which we highlight the modelled regions inhabited by students. To date, we have not explicitly constructed the employees at Universities and their interactions with the student body.

### 3.4 Construction of work places

After schools and universities, we construct work places for the employed subset of the population. We divide employment structures into three categories: work in companies with employees; work outside fixed company structures; work in hospitals and schools. To distribute the workforce over workplaces, JUNE first initialises companies based on data containing their locations, sizes and the sectors in which they operate. In a next step, individuals who are eligible to work (i.e. between the ages of 18-65) are assigned the industry sector based on the geographic distribution of the workforce per sector. This results in origin-destination matrices which are used to match workers to their workplace and to optimize the distribution of individual company to reproduce sector-dependent distributions.





**Figure 5:** Number of workers by sex and company sector (denoted by SIC code identifiers, see Table 1) in JUNE.

In England, the ONS database contains highly detailed information on companies and workforce by industry type. Industries and companies are categorised according to 21 sectors following the Standard Industrial Classification (SIC) code convention [26] (see Table 1) and information about company numbers per sector, and company sizes is available at the MSOA (super area) level [27]. Similarly, the ONS data also contain the size and sex distribution of the workforce by sector at the MSOA level, as well as the location of their employment [28, 29]. This enables the construction of an origin-destination matrix and allows us to distribute the workforce accordingly. More details on this specific procedure for initialising companies in JUNE and matching working individuals to these companies can be found in Appendix A.3.

The resulting distribution of our procedure assigning individuals an industry sector can be seen in Fig. 5. JUNE captures many of the sex-dependent features of the job market such as females dominating the healthcare profession and males the manufacturing sector. Recovering these sector-level sex imbalances can be crucial to reproducing and predicting potential sex imbalances in disease spread.

JUNE fixes the location of employment through data specifying the physical position of, e.g., company buildings. This, however, does not capture other modes of employment. We model people working from home through the specification of single-person companies in the same location as their place of residence. It should be noted that we do not currently explicitly model those workers who may not work in formal company buildings but also do not work from home, such as contractors who may interact with a household of the people they are visiting for maintenance work.

Hospitals play a dual role in JUNE, both as an essential part of patient’s possible medical journey and as workplaces. We will discuss the role of hospitals for the former case in Section 5. For both purposes, hospitals are initialised like many other locations in JUNE, based on available data regarding their location and capacity. Hospitals can be modelled individually or as clusters; in the latter case we represent the full cluster by one hospital. For our simulation of COVID-19 in England, we define the relevant National Health Service (NHS) trusts as those that reported disease-related casualties – this amounts to a total of 129 trusts – and we cluster them into single hospitals<sup>4</sup>. The clustering of hospitals is in fact a better representation of the situation in England. The aggregation of data by NHS trust allows for a more detailed comparison of the number and geographical spread of hospital admissions with available data. We assign medical workers to hospitals based on the same origin-destination matrix at the MSOA (super area) level as derived above, by choosing from those who work in the healthcare sector (“Q”), with the additional constraint of assuming a fixed ratio of 10 hospital beds per medic – nurse or doctor. Teachers are chosen from the population in a similar matter by using the origin-destination matrix and choosing from those in the education sector (“P”).

<sup>4</sup>Some NHS trusts share resources and exchange patients across regions in an ad-hoc manner, however, this is not modelled explicitly.

SIC code identifier	Description
A	Agriculture, forestry and fishing
B	Mining and quarrying
C	Manufacturing
D	Electricity, gas, steam and air conditioning supply
E	Water supply; sewerage, waste management and remediation activities
F	Construction
G	Wholesale and retail trade; repair of motor vehicles and motorcycles
H	Transportation and storage
I	Accommodation and food service activities
J	Information and communication
K	Financial and insurance activities
L	Real estate activities
M	Professional, scientific and technical activities
N	Administrative and support service activities
O	Public administration and defence; compulsory social security
P	Education
Q	Human health and social work activities
R	Arts, entertainment and recreation
S	Other service activities
T	Activities of households as employers; undifferentiated goods-and services-producing activities of households for own use
U	Activities of extraterritorial organisations and bodies

**Table 1:** Standard Industrial Classification (SIC) code identifiers for the 21 work place sectors modelled in JUNE and used by the ONS to categorise companies [26].

## 4 Simulating Social Interactions

The `interaction` layer allows for careful mapping of spatial movement, and location and intensity of social interactions, of the virtual population. This is used to create as close a match as possible to observations and available data. Comparable IBMs, such as [6, 7], simulate social interactions in either static environments, such as households, schools, or work places, in a similar manner to that described in the previous section, or in an less specific way determined by gravity models. In contrast, JUNE allows for the specification of additional social settings, and directly connects them to geographical locations, such as shops or restaurants. We can also model transport routes of different types between specified geographical start and end-points. This granularity is further increased through the addition of social mixing matrices which parametrise differences in frequency and intensity of contacts between individuals in various settings [10, 30].

### 4.1 A day in a virtual person’s life

Calendar days, decomposed into time-steps of varying length given in units of hours, are the background for our simulation of the social interactions of our virtual population. The use of a calendar allows JUNE to distinguish between week-day and weekend activity profiles, which is relevant for time spent at work or in school. Each day can have a number of fixed, static, activities, such as 8 hours of work at the workplace or 10 hours at home overnight, supplemented with other activities, denoted as “other”, that are distributed dynamically. During each time-step during which an “other” activity is allowed, each person who is not otherwise occupied, e.g., working, or ill at the hospital, is assigned a set of probabilities for undertaking other activities in the model. These probabilities are part of our social interaction model, and depend on the age and sex of the person<sup>5</sup>. Given  $N$  possible activities with associated probabilities per hour given by  $\lambda_1, \dots, \lambda_N$ , for a person with characteristic properties  $\{p\}$ , the overall probability  $\mathcal{P}$

<sup>5</sup>These probabilities can be generalised to depend on any attributes of the individual given reliable data.

of being involved with any activity in a given time interval  $\Delta t$  is modelled through a Poisson process,

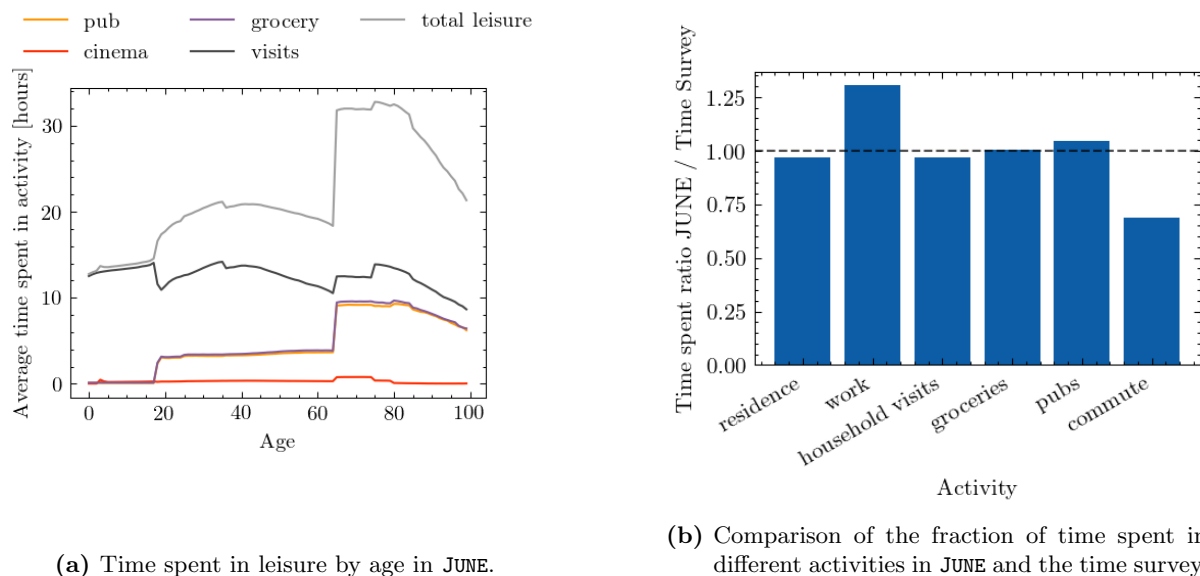
$$\mathcal{P} = 1 - \exp\left(-\sum_{i=1}^N \lambda_i(\{p\})\Delta t\right). \quad (4.1)$$

If the individual is selected to participate in one of these activities, the chosen activity,  $i$ , is then selected according to its probability

$$\mathcal{P}_i = \frac{\lambda_i(\{p\})}{\sum_{j=1}^N \lambda_j(\{p\})}. \quad (4.2)$$

The person is then moved to the relevant location corresponding to this activity. If no activity is selected, the individual will stay at home.

A summary of how much time is spent each week on various activities as a function of age is reported in Fig. 6a. In Fig. 6b we show a comparison of the amount of time spent at home, work, grocery shopping, eating at restaurants/pubs, and commuting between JUNE and the UK Time Use Survey, 2014-2015 [31]. Care home and cinema visits are not accounted for in the time survey.

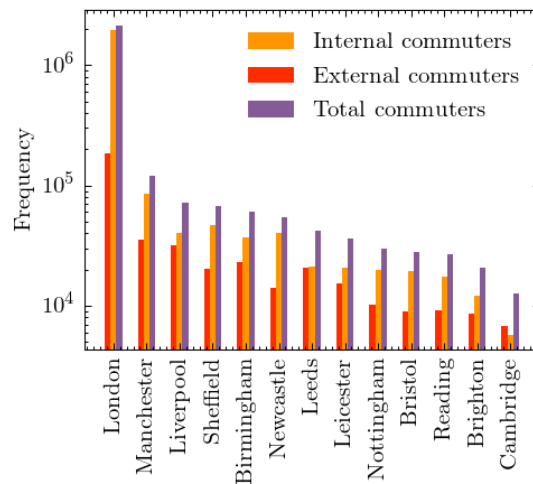


**Figure 6:** Leisure activities in JUNE

## 4.2 Localised activities

In the following we will detail JUNE’s model of social activities outside the static home, work and school settings. In the example of the spread of COVID–19 through the English population, it is broadly based on time surveys provided by the ONS [29], which identify a variety of activities, including time spent at home, work, or in school. In addition, we have identified five additional settings which we assume are similarly relevant for the spread of the disease and have a similar level of social mixing: visits to pubs or restaurants (“pubs”), cinema visits (“cinemas”), shopping (“groceries”), visiting friends or relatives in their homes (“household visits”), or visiting family members in care homes (“care home visits”).

In JUNE we locate 120,000 pubs and restaurants according to their geo-coordinates, as well as 32,000 stores and 650 cinemas. Each time a person is assigned to any of “pubs”, “groceries”, or “cinemas”, we pick a random venue from the  $n$  venues closest to their place of residence, or the closest venue if the distance to any of them is greater than 5km. We have chosen  $n = 7$  for pubs,  $n = 15$  for shopping stores, and  $n = 5$  for cinemas. Note that there are no permanent “workers” in these venues who return to a single venue daily; only “attendees” who choose their venue at random. Further locations such as gyms and places of worship can be easily added to the activity model, and, of course, it can easily be adjusted to other societies.



**Figure 7:** Number of internal and external commuters by city as modelled in JUNE.

In addition, we model interactions in naively constructed social networks, by linking each household to a list of up to  $\mathcal{N}$  other households in the same super area. One of the households in this list is selected, if “household visits” is chosen as activity during a time step. Residents will stay at home to receive the incoming visitor, who in turn may also bring their whole household with them according to a probability described by an external parameter. Comparison with national surveys suggests that setting the number of linked households  $\mathcal{N} = 3$  provides realistic movement profiles. While care home residents in JUNE cannot visit other people, each resident is connected to a household of them in the local super output area from whom they can receive visitors. JUNE also models the interactions that result from elderly people needing help in their daily activities. Each person older than 65 years old has a probability, increasing by age, of needing some kind of assistance in their daily activities. We therefore assign a member of the local super area to be the carer of an elderly person, following the data available in [32]. Every weekday, the carer spends their leisure time visiting the household of the person needing domestic care.

### 4.3 Modelling mobility: commute patterns

In addition to its high resolution in the type and location of social interactions, JUNE’s granularity also lends itself to the modelling of travel patterns of individuals which can be important for the geographical spread of an infection. In particular, this is true for commuter travel, which may have a significant impact in metropolitan areas given the high density of public transport users.

A necessary ingredient for any realistic mobility model is a directed network graph, or similar, of interconnected nodes with defined geographical locations, where the edges between the nodes correspond to possible transit routes. Travellers move between nodes in this network, and may share their means of transportation. This framework can easily be generalised to other settings, such as holiday travel, as long as data is available to facilitate the construction of an origin–destination matrix for inter-node movement.

In our example case of England, and for the spreading of an infectious disease, we assume that private means of transport relate to individual traffic with very little transmission, and focus solely on creating a model for public transport. Using commuting and rail travel data provided by the UK Department for Transport [33], we select the major transit cities to be nodes defining our simplified model of the national transport network. This is a reflection on the importance of commuting into and within major cities as a major potential transmission avenue, since it naturally induces social mixing between many people who may not usually come into contact otherwise.

To fill our origin–destination matrix we use information contained in the ONS database concerning the mode of commuting of individuals at the area (OA) level [34], to distribute commuting modes probabilistically. We define two modes of public transport, “external” which defines those commuting in and out of metropolitan areas, and “internal” which defines those commuting within these areas.



**Figure 8:** Commuting maps for London as derived from JUNE.

Metropolitan areas are defined using data obtained from the ONS [35]. For the sake of computational efficiency we model only the travel patterns of those working inside metropolitan areas, who in fact represent the overwhelming majority of public transport commuters. This includes commuters who live and work in the city, as well as those who are entering the metropolitan area from outside. The number of internal and external commuters by city in England is given in Fig. 7. The cities included are geographically spread across England thereby accounting for major commuting patterns in most regions modelled. In total, we explicitly model commuting into 13 out of a possible 109 cities in England, which accounts for 60% of all metropolitan commuters and 46% of all those using public transport to commute to work. Fig. 8 shows maps of the residences of internal and external commuters in two cities in our model, where the inner section in white denotes the respective metropolitan areas. Specifically, from Fig. 8b we can see that, given the large commute radius of cities like London (we observe a similarly large radius for Birmingham and several other cities), commuting can be a key driver for the inter-regional spread of infectious-diseases.

Travelling within a metropolitan area, i.e. the internal commuting mode, is modelled as a self-connected loop – practically speaking this means that internal commuters may in principle interact, irrespective of the actual movement inside the city. For external commuting, the travel into and out from the metropolitan area, we identify shared routes for commuters living in neighbouring areas and super areas. The number of possible routes into each city, and therefore the number of ways to divide regions around the cities, is informed by the approximate number of rail network lines into each city – currently this is set to 8 in London and 4 for each of the other 12 cities [36].

We randomly partition people sharing the same commuting route into subgroups, “carriages”, which define the environment in which social interactions take place. The commuting time-step is run twice a day and in each run the travellers are randomly distributed into carriages. The number of people per carriage is determined by city-dependent data obtained from the UK Department for Transport [33]. More details on the specific algorithm for modelling commuting in JUNE can be found in Appendix A.4.

#### 4.4 Social interaction frequencies and intensities

Social contact matrices [10, 15] provide information about the age-dependent frequency and intensity of in-person contact in different social settings, an important ingredient to many epidemiological simulations. They measure the average daily number of conversational and physical contacts between individuals of different ages. This means that they are normalised to the size of the population in the respective age bins, but do not account for whether they can take part in such contacts. To use them within JUNE we therefore have to account for the fact that social settings define the group of people coming into contact with each other. To exemplify this, consider the case of contacts between adults and students in schools. While the social contact matrices in the literature normalise the number of contacts of a 30-year old with children of a certain age to the number of 30-year old adults in the population, in

JUNE only a subset of 30-year old adults work as teachers and can therefore interact with the children. In the construction of matrices specific for JUNE, we therefore combine the results from [10, 15] with simple assumptions about possible participants in contacts.

Averaging over age ranges in different settings, we arrive at simplified social mixing matrices,  $\chi_{si}^{\mathcal{L}}$ , which will be comparable to the inputs from literature upon combination with the model results for the composition of social environments. Below we list our simplified social mixing matrices inferred from literature, with  $\mathcal{L} \in \{(H), (S), (W)\}$  (home, school, work place), as well as the relative proportions,  $\phi_{si}^{\mathcal{L}}$ , of physical contacts. The latter are relevant, since in line with standard approaches, closer physical contact in JUNE is proportional to a higher propensity for transmission for the etiological agent.

For the households social mixing matrices, we define four categories, young children ( $K$ ), young dependent adults of age 18 or more ( $Y$ ) that still live with their parents, adults ( $A$ ), and older adults ( $O$ ) of age 65 and over. We use:

$$\chi_{ij}^{(H)} = \begin{pmatrix} 1.2 & 1.69 & 1.69 & 1.69 \\ 1.27 & 1.34 & 1.47 & 1.50 \\ 1.27 & 1.30 & 1.34 & 1.34 \\ 1.27 & 1.50 & 1.34 & 2.00 \end{pmatrix} \quad \text{and} \quad \phi_{ij}^{(H)} = \begin{pmatrix} 0.79 & 0.70 & 0.70 & 0.70 \\ 0.70 & 0.34 & 0.40 & 0.40 \\ 0.70 & 0.40 & 0.62 & 0.40 \\ 0.70 & 0.40 & 0.40 & 0.56 \end{pmatrix} \quad (4.3)$$

For household visits, we make the simplifying assumption that the same matrices also describe the contacts between visitors and residents. For visits to care homes, we believe that visitors come into contact only with residents and care home workers, and not with other visitors. We therefore hypothesize 6 conversational contacts with residents and 1.5 with care home workers.

Social contacts in schools identify teachers ( $T$ ) and students ( $S$ ), the latter are organised in year groups and further divided into classes of up to 40 students. In our age-averaging, we implicitly assume that the number and character of teacher-student contacts is independent of the age of the students. Student-student contacts are assumed to be most frequent within a class or year group, and fall off steeply with the age difference. This behaviour is captured by fitting a matrix with values for the age-diagonal elements and a fall-off per year age-difference by a factor of 3. Therefore we have

$$\chi_{ij \in \{T, S\}}^{(S)} = \begin{pmatrix} 4.8 & 0.75 \\ 15 & \chi_{SS}^{(S)} \end{pmatrix} \quad \text{and} \quad \phi_{ij \in \{T, S\}}^{(S)} = \begin{pmatrix} 0.05 & 0.08 \\ 0.1 & \phi_{SS}^{(S)} \end{pmatrix}, \quad (4.4)$$

with the student-student matrices taking the following form

$$\chi_{SS}^{(S)} = \begin{pmatrix} 2.5 & 0.75 & 0.25 & \dots \\ 0.75 & 2.5 & 0.75 & \dots \\ 0.25 & 0.75 & 2.5 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix} \quad \text{and} \quad \phi_{SS}^{(S)} = 0.15 \quad \forall i, j \in \{S\}. \quad (4.5)$$

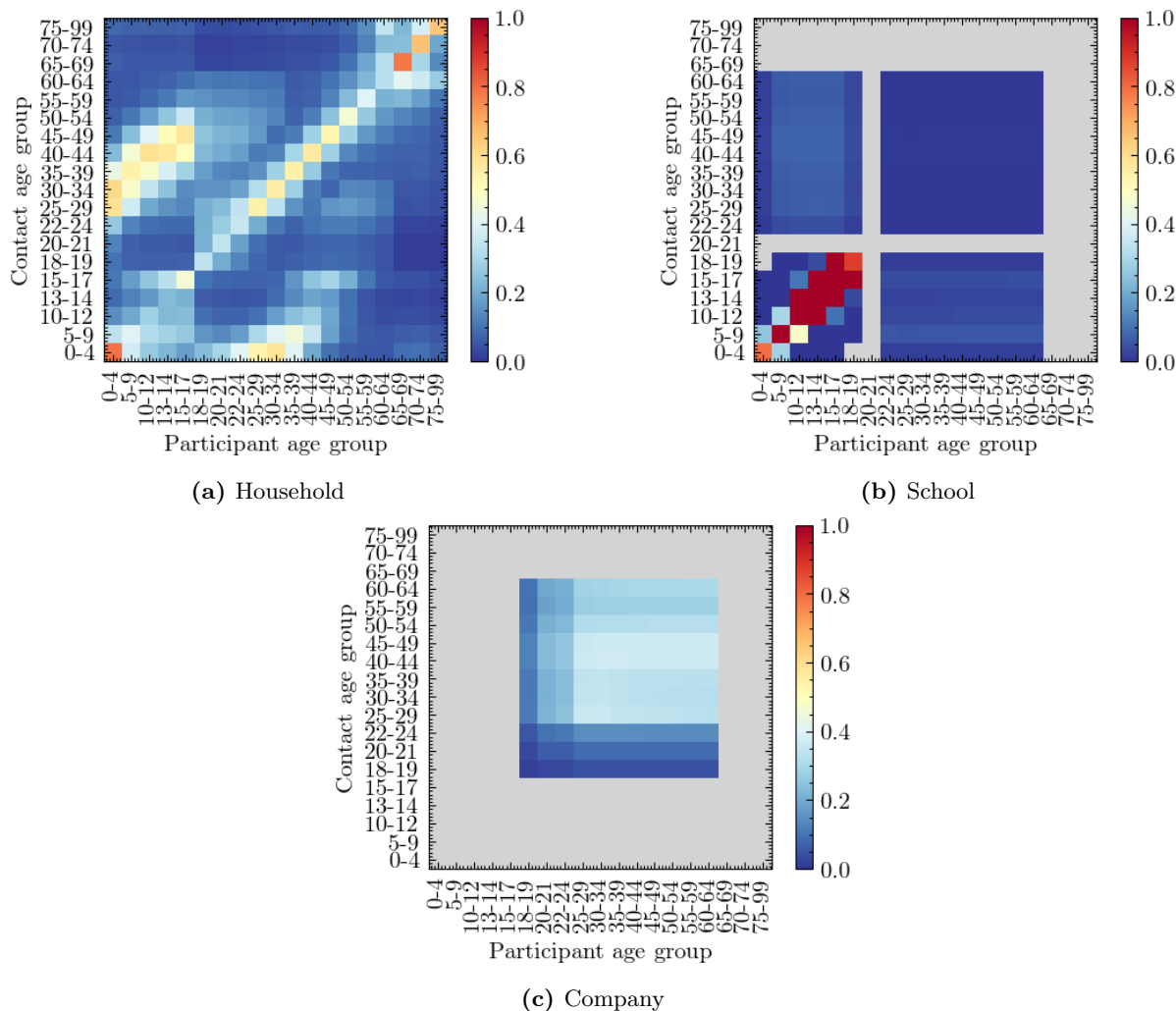
For the contacts at work we do not take into account of any age-dependence and, in the absence of data, do not model any sector-dependent variation of their number of intensity, thus

$$\chi_{si}^{(W)} = 4.8 \quad \text{and} \quad \phi_{si}^{(W)} = 0.07, \quad (4.6)$$

In Appendix B we detail the algorithms used to construct the social mixing matrices used in JUNE including the matrices for other locations not listed here.

These social mixing matrices in JUNE are defined for a setting-specific characteristic time  $t_{\text{char}}$ , so the total number of contacts in a time interval  $\Delta t$  in a given setting is then modified by a factor  $\Delta t/t_{\text{char}}$ .

To validate these simplified matrices, we include them within JUNE where they are combined with the composition of the specific social settings. In Fig. 9 we show the resulting contact matrices as ‘‘measured’’ from the JUNE simulation. The effect of the combination with the composition is most pronounced in the household matrices which exhibit textures that can be directly traced back to the age intervals of children, dependent children/young adults, adults and older adults that JUNE inherits from the ONS data. These matrices naturally recover much of the structure present in those recorded in [10, 15]. Further details on the methodology for extracting these matrices from JUNE can be found in Appendix B.5.



**Figure 9:** Social contact matrices for England derived from JUNE, before any mitigation strategies are implemented.

## 5 Infection Modelling: Spreading and Health Impact

The transmission of infection through social interactions described in the `interaction` layer, and the progression of the disease and its impact on the individual, are both modelled in the `disease` layer. Although we focus on the case of COVID-19 here, this layer is designed to be generalisable and can contain more than one circulating etiological agent.

### 5.1 Infection Transmission

JUNE models the transmission of an infection from infecting individual,  $i$ , to susceptible individual,  $s$ , in a probabilistic way. The probability of infection in a social setting within a group of people,  $g$ , at a location,  $L$ , depends on a number of factors:

- the number,  $N_i$ , of infectious people  $i \in g$  present;
- the infectiousness of the infectors,  $i$ , at time  $t$ ,  $I_i(t)$ ;
- the susceptibility,  $\psi_s$ , of the potential infectee,  $s$ ;
- the exposure time interval,  $[t, t + \Delta t]$ , during which the group,  $g$ , is at the same location;

- the number of possible contacts,  $\chi_{si}^{(L)}$ , and the proportion of physical contacts,  $\phi_{si}^{(L)}$ , in location  $L$ , both taken from Eqs. (4.3)-(4.6) in Section 4.4;
- and the overall intensity,  $\beta^{(L,g)}$ , of group contacts in location  $L$ .

Most of these ingredients depend on the time,  $t$ , of the contact. For example, the number of contacts,  $\chi_{si}^{(L)}$ , and the proportion of physical contacts,  $\phi_{si}^{(L)}$ , and the overall contact intensity,  $\beta^{(L,g)}$ , will change with the implementation of social distancing policies. To simplify notation, we introduce a combined contact intensity for a group  $g$  with size  $N_g$  at location  $L$ ,

$$\beta_{si}^{(L,g)}(t) = \beta^{(L,g)} \cdot \frac{\chi_{si}^{(L)}(t)}{N_g} \left\{ 1 + \phi_{si}^{(L)}(t) \left[ \alpha(t) - 1 \right] \right\}, \quad (5.1)$$

where the ratio  $\chi/N_g$  provides a simple parametrization of the probability of  $s$  being in contact with another individual in the group, and  $\alpha(t) > 0$  describes the relative impact of close physical contacts. Both the factor  $\alpha(t)$ , which we assume to be the same for all locations, and the location- and group-specific contact intensities,  $\beta^{(L,g)}$ , are taken from fits to data.

In the construction of an infection probability for a susceptible individual,  $s$ , we make a number of assumptions. First of all, we model the probability of being infected as a Poisson process. In keeping with the probabilistic process, the argument of the Poissonian is given by a sum over individual pairs of infectious individuals with the susceptible person, implying a simple superposition of individual infectiousness. The underlying individual transmission probabilities are written as the product of the susceptibility of the susceptible individual, the infectiousness of the infected person, and the contact intensity, all integrated over the time interval in which the interaction occurs. The integration over time ensures that the transmission probability increases with the time of exposure. We therefore arrive at the transmission probability, i.e. a probability for  $s$  to be infected as:

$$\bar{P}_s(t, t + \Delta t) = 1 - \exp \left[ -\psi_s \sum_{i \in g} \int_t^{t+\Delta t} \beta_{si}^{(L,g)}(t') \mathcal{I}_i(t') dt' \right]. \quad (5.2)$$

Note that in the actual implementation, we approximate the integral over time with a simple product,

$$\int_t^{t+\Delta t} \beta_{si}^{(L,g)}(t') \mathcal{I}_i(t') dt' \rightarrow \beta_{si}^{(L,g)}(t) \mathcal{I}_i(t) \Delta t. \quad (5.3)$$

This leaves us to fix the last two ingredients in Eq. (5.2), the individual susceptibility,  $\psi_s$ , and the infectiousness,  $\mathcal{I}_i(t)$ . Contemporary peer reviewed academic research on susceptibility to infection by the etiological agent with or without the onset of disease symptoms is sparse and inconsistent. Following some evidence, for example in [37] and [38], on transmission and susceptibility of children (using the UN classification), we fix  $\psi_s = 0.5$  for children under the age of 12, and  $\psi_s = 1$  for everybody else. The infectiousness of individuals,  $\mathcal{I}_i$ , changes with time, and it is not directly measurable. To model its behaviour we use the temporal dependence of viral shedding as a proxy for infectiousness. Studies in the context of COVID-19 have shown that viral shedding peaks at or slightly before the onset of symptoms, and then begins to decrease [39]. In JUNE we use a globally defined temporal dependence of infectiousness,  $f_I(t)$ , and multiply it with a peak value,  $\mathcal{I}_{i,\max}$ , which depends on the infected individual,

$$\mathcal{I}_i(t') = \mathcal{I}_{i,\max} \cdot f_I(t'). \quad (5.4)$$

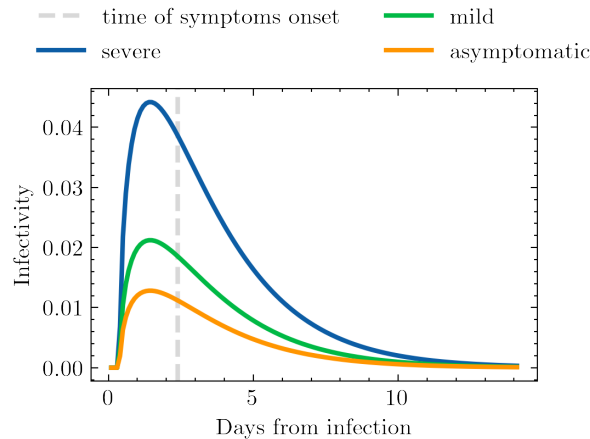
We choose the maximal infectiousness according to a log-normal distribution parameterised by its median  $\exp(\mu) = 1$  and shape  $\sigma = 0.25$ . The long right tail of the log-normal distribution allows for small numbers of highly infectious individuals more likely to precipitate superspreading events (SSEV). We also capture the conjectured reduced infectiousness of individuals with no or only mild symptoms. Following a similar parameterisation to that in [40], we multiply the maximal infectiousness of asymptomatic individuals by 0.33, and of individuals with mild individuals by 0.72. In Fig. 10, we show an example of the time



evolving profile for an infected individual in JUNE, comparing the resulting infectiousness for different symptoms. For the time-dependent profile, we use the gamma distribution as fitted in [39],

$$f_I(\tau = t' - t_0 - t_{inc}, a) = \frac{\tau^{a-1} e^{-\tau}}{\Gamma(a)}, \quad (5.5)$$

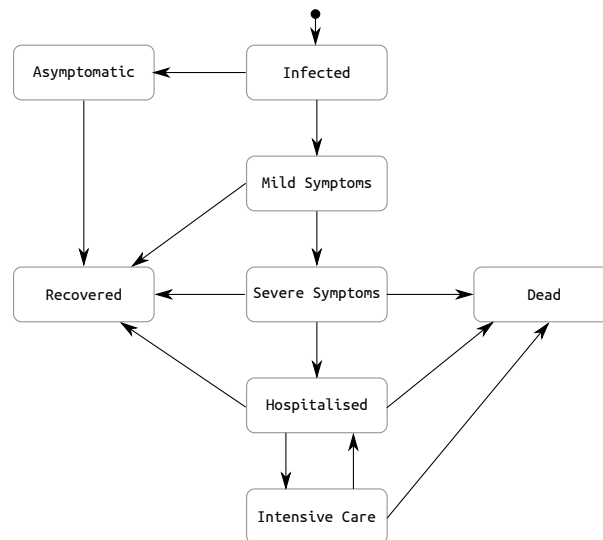
where  $t_0$  is the time of infection,  $t_{inc}$  is the incubation period, sampled from a normal distribution centered at two days prior onset of symptoms and with a width of half a day,  $a$  is the shape parameter of the gamma distribution, and  $\Gamma(a)$  is the the gamma function.



**Figure 10:** Time-dependent infectiousness profile,  $f_I(t')$ , shown for the same realisation of the infection but where the infected person attains severe symptoms, mild symptoms or becomes asymptomatic.

## 5.2 Infection Progression

Once individuals are infected, they will experience different impacts on their health. Fig. 11 presents



**Figure 11:** Pathways for the infection progression and possible outcomes.

the paths available in JUNE for the progression of the infection that aim to capture different symptom severities, outcomes, and their operational impact on the healthcare system, i.e. whether patients are hospitalised or admitted to intensive care or treatment units (ICU/ITU). Once an individual is infected,

JUNE selects their specific complete path according to these probabilities. These paths are codified as a sequence of possible different stages of the disease (“infected”, “asymptomatic”, “mild”, “severe”, “hospitalised”, “ICU/ITU”, “dead”, “recovered”) in addition to characteristic time intervals for each stage. The latter are chosen randomly according to probability functions informed by available data. The paths terminate with the individuals either dead and taken out of the simulation, or recovered, in which case their susceptibility is set to 0, making them immune to re-infection<sup>6</sup>

JUNE distinguishes the following different routes for the progression of the infection with rates depending on the characteristics of the infected individual (currently age, sex), summarily denoted by  $p$ :

1. asymptomatic individuals, rate  $R_{I \rightarrow A}(p)$ , continue their life normally;
2. individuals with mild symptoms, rate  $R_{I \rightarrow M}(p)$ , usually continue their lives as normal, except if certain policies are activated;
3. individuals with severe but not lethal symptoms, rate  $R_{I \rightarrow S}(p)$ , stay at home until recovery;
4. individuals with severe symptoms who will eventually die in their residences, with rate  $R_{I \rightarrow DR}(p)$ ;
5. individuals who are admitted to hospital but will recover, with rate  $R_{I \rightarrow H}(p)$ ;
6. individuals who are ultimately admitted to ICU/ITU before recovering, with rate  $R_{I \rightarrow ICU}(p)$ ;
7. individuals who are admitted to hospital and will die there, with rate  $R_{I \rightarrow DH}(p)$ .
8. individuals who are admitted to ICU/ITU and die there, with rate  $R_{I \rightarrow DICU}(p)$ .

The determination of probabilities for the different paths is based on COVID–19 data that are not entirely sufficient to develop a complete and detailed picture. As a consequence, we supplement them with assumptions by inferring some properties through cross-relating datasets. In the following, we will outline our procedure which is largely predicated by our choice of the example at hand - the spread of COVID–19 in England. We will use a notation where  $N_X(p)$  denotes the number of cases satisfying criterion  $X$  for people with characteristic properties  $p$ .

The construction of reasonable progression paths, and their probabilistic distribution, relies critically on the knowledge of how many people have been infected, as well as the dependence on attributes such as age and sex. COVID–19 tests between February and May 2020 in the UK were mostly administered to people presenting symptoms or people that have been in close contact with confirmed cases in hospital, thereby biasing the results. We therefore need to infer the number of infections from other controlled studies, such as antibody tests. In [41] the seroprevalence,  $r_{sp}(p)$ , of COVID–19 in the adult population in England was determined through a sample of more than 100,000 adults, showing a reduction in seroprevalence with increasing age. Because the seroprevalence is an estimate of all people that were infected up to the time of the test and – most importantly – survived, we need to correct for those who died of the disease until this point. This turns out to be an important correction, especially in older age bins due to the non-negligible probability of elderly who died. We therefore add the age- and sex-dependent number of deaths,  $N_D(p)$ , reported by the ONS [42], to the corresponding numbers inferred from the seroprevalence to arrive at the total number of cases,  $N_{tot}(p)$ :

$$N_{tot}(p) = r_{sp}(p)N(p) + [1 - r_{sp}(p)]N_D(p), \quad (5.6)$$

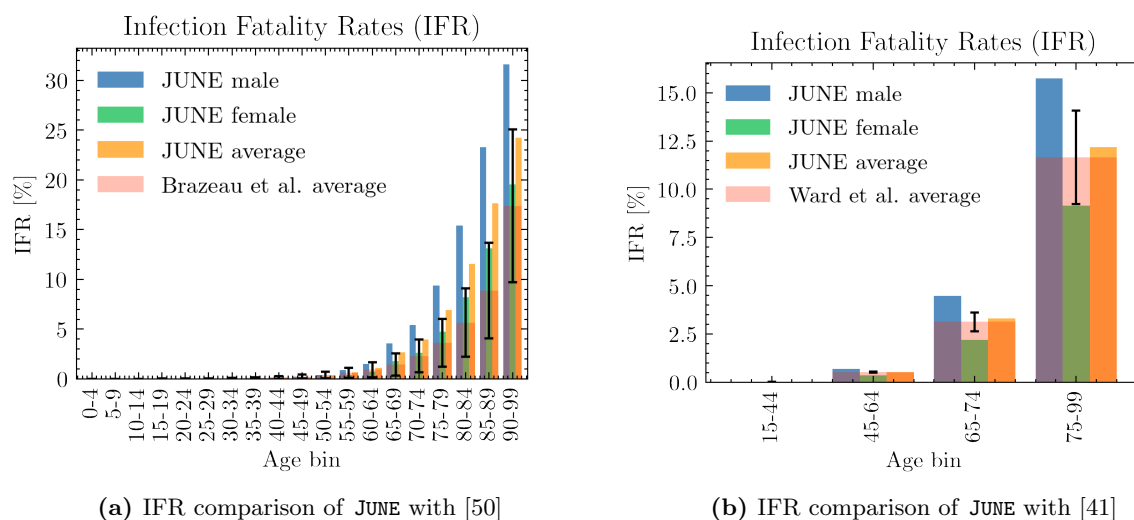
where  $N(p)$  is the total population number in England with characteristics  $p$ . We note that there were two population groups excluded from the serology survey: people under the age of 18, and care home residents. For the former, we assume that their seroprevalence by age is identical to the population group aged 18, while for the latter, we set a flat seroprevalence by age at 11% value as reported in the Vivaldi report of the UK Department of Health and Social Care [43] in the beginning of July 2020.

We now turn to the determination of the rates for different trajectories,  $R_{I \rightarrow X}$ , where  $X$  is one of the 8 trajectories listed previously. The asymptomatic rate,  $R_{I \rightarrow A}$ , and the mild case rate,  $R_{I \rightarrow M}$ , are taken from a calibration done in [40] from [51] and [52]. To calculate the different hospitalisation and fatality

<sup>6</sup>JUNE could also model reinfection of individuals but to date there are no data constraining this in the context of COVID–19.

Quantity	Source
Population by age, sex and residence type	[44], [45]
Seroprevalence in GP by age	[41]
Seroprevalence in CH by age	[43]
Deaths by place of occurrence and residence type	[42]
Deaths profile by age and sex	[42]
Deaths in CH profile by age and sex	[46]
Hospital deaths profile by age, sex	[47]
Hospital deaths in CH profile by age, sex	[48]
ICU/ITU deaths profile by age, sex	[48]
Total hospital admissions	[49]
Hospital admissions profile by age, sex	[47]
ICU/ITU admissions profile by age, sex	[48]
Hospital admissions in CH profile by age, sex	[48]

**Table 2:** Datasets used in the derivation of mortality and hospitalisation rates. GP stands for people living in a household, and CH stands for people living in care homes. If not specified, datasets involve people from both populations. All data is taken until the 13th of July 2020, consistently with the seroprevalence study [41].

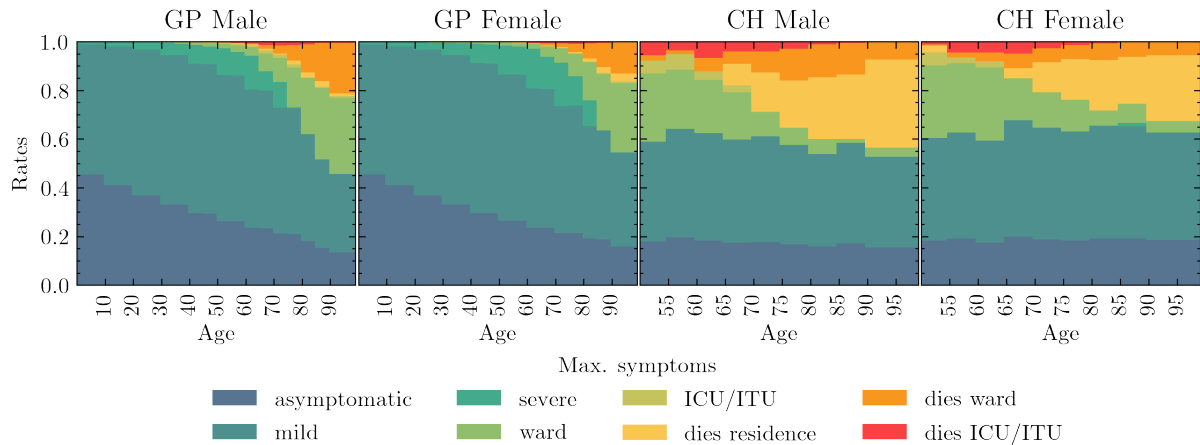


**Figure 12:** IFR comparison of JUNE with other works.

rates, we have used a series of datasets listed in Table 2, all of them containing data until the 13th of July 2020, to be consistent with the considered seroprevalence values. In order to avoid possible irregularities in our results derived from the use of different data sources, we normalise all our death data to the ONS reported numbers of total deaths (51,443), hospital deaths (32,164), and residence deaths (19,279), [42] and then use more granular data to distribute deaths by age and sex for each place of death occurrence [48], [47]. Likewise, the total number of hospital admissions is taken from [49], and distributed by age, sex, and residence type also using [48] and [47]. The number of deaths in care homes reported in [53] is only reported by age until late June, so we assume that the distribution does not change until the 13th of July. We also ensure that we correctly account for differences in reporting times. As a first step, we calculate the overall infection fatality rate (IFR) for the general population outside care homes (GP),

$$R_{I \rightarrow D}^{GP}(p) = \frac{N_D(p) - N_D^{ch}(p)}{N_{tot}(p) - N_{tot}^{ch}(p)}, \quad (5.7)$$

which can be directly compared to the results from the REACT2 study [41] (right panel of Figure 12), and the Imperial College London COVID-19 report 34 [50] (left panel of Figure 12). The remaining



**Figure 13:** Rates of different infection outcomes for males and females living in households and care homes. For care home residents, we only show the rates for people aged over 50, as the younger ones are assumed to follow the general population rates.

rates just follow from the same methodology,

$$R_{I \rightarrow X}^{GP}(p) = \frac{N_X(p) - N_X^{ch}(p)}{N_{tot}(p) - N_{tot}^{ch}(p)}, \quad (5.8)$$

$$R_{I \rightarrow X}^{CH}(p) = \frac{N_X^{ch}(p)}{N_{tot}^{ch}(p)}, \quad (5.9)$$

$$(5.10)$$

where X refers to one of deaths or hospital admission in the hospital ward or ICU/ITU. The rate of resident deaths simply follows from subtracting the hospital death rates from the overall IFRs. Finally, the probability of having severe symptoms but recovering at home is, consistently,

$$R_{I \rightarrow S}(p) = 1 - \sum_{i \neq S} R_{I \rightarrow X_i}(p). \quad (5.11)$$

The results of computing the individual infection outcome rates by age, sex, and residence type are shown in Figure 13. The most important visible difference is the disparity on the fatality rates between care home residents and the general population. This could be the reflection of various reasons, including, for example, a generally poorer health condition of the care home population, or differences in admission policies to hospitals. Consistently with the ONS data [53], most of the care home deaths occur at the residence itself, while the probability of being admitted to the hospital decreases with age. Likewise, both for the general population and the care home population, people aged 55-70 years old are the group most likely to be admitted in the ICU/ITU. Females are less likely in general to develop a severe infection of COVID-19, with fatality rates roughly equivalent to those of a male 5 years younger.

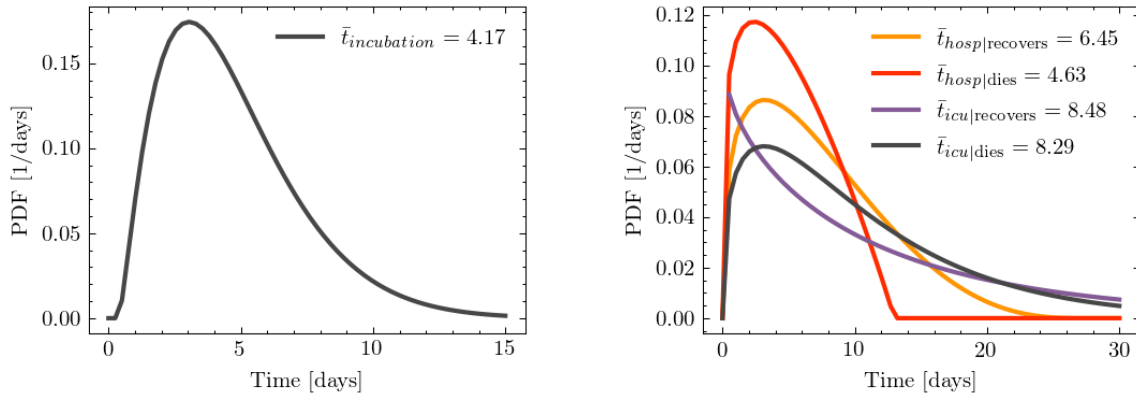
Once an infection outcome has been determined, the infected individual follows a symptoms trajectory composed of different stages. The time spent at each stage is sampled from different distributions derived from different data sources. In Table 3, we list the different stages per trajectory by infection outcome, and the details on the various timings are listed in Appendix C. In the left panel of Figure 14, we show the probability density functions for the incubation time, and the time to die or recover in hospital.

### 5.3 Seeding Infections

In the absence of sufficiently detailed knowledge of how epidemics arrive in a country, we seed infections using secondary information such as the number and regional distribution of observed cases. In the example of the simulating the spread of COVID-19 in England, we use the number of COVID-19-related deaths recorded in hospitals to estimate initial infection numbers and their regional distribution.

Trajectory	Stages					
asymptomatic	I[ $\beta_I$ ]	A[ $C_{14}$ ]				R
mild	I[ $\beta_I$ ]	M[ $C_{20}$ ]				R
severe	I[ $\beta_I$ ]	M[ $C_{20}$ ]	S[ $C_{20}$ ]			R
death at home	I[ $\beta_I$ ]	M[ $LN_M$ ]	S[ $C_3$ ]			D
ward	I[ $\beta_I$ ]	M[ $LN_M$ ]	H[ $\beta_H$ ]	M[ $C_8$ ]		R
death in ward	I[ $\beta_I$ ]	M[ $LN_M$ ]	H[ $\beta_D$ ]			D
ICU/ITU	I[ $\beta_I$ ]	M[ $LN_M$ ]	H[ $LN_{ICU}$ ]	ICU[ $e_{ICU}$ ]	H[ $e_H$ ]	M[ $C_3$ ]
death in ICU/ITU	I[ $\beta_I$ ]	M[ $LN_M$ ]	H[ $LN_{ICU}$ ]	ICU[ $e_D$ ]		D

**Table 3:** List of different trajectories through disease progression, with stages and, in brackets, the distribution from which corresponding timings are drawn. For their definition see Table 6. The available stages are **I**nfected, **A**symptomatic, **M**ild and **S**evere symptoms, admitted to a regular **H**ospital or an **ICU/ITU** ward, and, finally, as outcomes, **R**ecovered or **D**ead.



(a) Time taken for an infected individual to develop symptom. (b) Time spent in hospital by patients given their infection.

**Figure 14:** Probability density functions for symptom and progression timing.

Accounting for the time delay between infection and possible death, and for the probability of admitted patients to die, we have

$$N_{\text{tot}}(t, x) = \frac{1}{\bar{R}_{H \rightarrow D}(x)} N_{H \rightarrow D}(t + \Delta t_D, x), \quad (5.12)$$

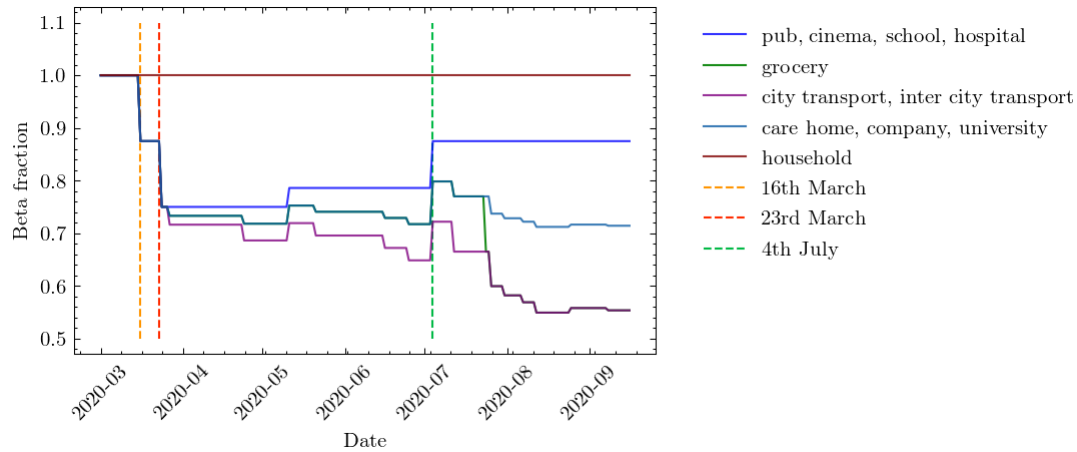
where  $N_{\text{tot}}(t, x)$  is the estimated number of cases in a region,  $x$ , on day,  $t$ ,  $N_{H \rightarrow D}(t, x)$  is the number of observed deaths in the region at date  $t$ , and  $\bar{R}_{H \rightarrow D}(x)$  is the rate for people dying in hospital in the region  $x$ , where the average is over the characteristics  $p$  is given by:

$$\bar{R}_{H \rightarrow D}(x) = \frac{1}{N_{H \rightarrow D}(t + \Delta t_D, x)} \sum_{i \in N_{H \rightarrow D}(t + \Delta t_D, x)} R_{H \rightarrow D}(p_i). \quad (5.13)$$

The relatively large statistical fluctuations in the initial phase of an epidemic, and possibly differing time profiles across regions, translate into the need for a region-specific seeding. This difference is highlighted by contrasting the seeding for London, where we introduce initial infections over two days only (28<sup>th</sup>-29<sup>th</sup> February 2020) with the North East of England and Yorkshire, where we seeded infections for a week, 28<sup>th</sup> February - 5<sup>th</sup> March 2020. We introduce the estimated number of daily cases in each of the regions until the following criterium is met,

$$N_{\text{tot}}(t < T(x), x) > 0.1 N_{\text{tot}}(t_{\text{max}}, x) \quad (5.14)$$

where  $T(x)$  is the number of days over which we seed new infections in region  $x$ , and  $N_{\text{tot}}(t_{\text{max}}, x)$  is the maximum number of cases that region  $x$  would reach in any given day, estimated from the maximum



**Figure 15:** Example scenario of different intensity parameters,  $\beta^{(L,g)}$ , over time normalised to unity (see Eq. (5.1)). The parameters change due to the effects of compliance with social distancing and mask wearing advice and regulations.

number of daily deaths in hospital. It is important to define the seeding for the infection based on the maximum number of cases each region will have, since the different regions are experiencing different stages of the epidemic at any given time.

## 6 Mitigation Policies and Strategies

Policies and interventions, often enacted by governing bodies, are introduced in an attempt to mitigate and control the spread of infectious diseases. In general, such policies are highly dependent on the type of infection and social norms in the affected population, and may include guidelines on how to change individual patterns of behaviour or the closure of certain venues where transmission is estimated to be highly likely. The modular nature of JUNE allows policies to be dynamically activated and deactivated at different points in time to allow for changes in policy decisions. Due to JUNE's granularity, these policies can be implemented at a highly localised level: by type and place of social interactions, by geographical region, by industry sector or venue type. JUNE can also model the population's compliance with the measures, again with high granularity. In this section we present a variety of policies which can be implemented in JUNE and exemplify their application through those measures that have been enacted by the UK Government to mitigate the spread of SARS-CoV-2.

### 6.1 Behavioral Changes

There are a variety of changes in behavioral patterns that are designed to reduce the probability of viral transmission, ranging from simple social distancing, increased hygiene and mask wearing, to quarantining of infected individuals or those who have been in sufficiently close contact with them, and the shielding of vulnerable parts of the population. We model the impact of the former set of measures, social distancing, increased hygiene and mask wearing, through multiplicative reductions in the location-specific contact-intensity parameters,  $\beta^{(L,g)}$ , see Fig. 15 for an example. The impact of compliance with social distancing and other, similar measures can be recorded both nationally and sometimes even in specific locations. This allows us to calculate the reduction in the corresponding intensity parameters as follows:

$$\beta^{(L,g)} = M^{(L,g)} \beta^{(L,g)} \quad (6.1)$$

$$= \left[ 1 - C^{(N)} \cdot C^{(L)} \cdot (1 - E) \right] \beta^{(L,g)}, \quad (6.2)$$

where  $M^{(L,g)}$  is the location and group specific modification factor,  $C^{(N)}$  is the national compliance (i.e. percentage of the population following guidelines),  $C^{(L)}$  is the compliance in a given location or

social setting  $L$ , and  $E$  denotes the efficiency of the measure. Quarantining is simulated by keeping the individuals in question in their homes for a certain amount of time, and allowing them to interact with members of their household in an otherwise unchanged household setting only. In JUNE, we have the ability to apply different policies to those with mild and severe symptoms, and to quarantine household members of symptomatic individuals. Similarly, JUNE also allows the definition of vulnerable individuals – typically by characteristics such as age – and of a prescription of how shielding policies are enacted relative to this group.

We will now turn to discuss our choices for specific measures. There have been a variety of studies on the effectiveness of social distancing with respect to COVID–19 and other infectious diseases. A comprehensive systematic review and meta-analysis [54] suggested that the relative risk of infection decreases by approximately a factor of 2 per meter distance. In practice, however, the efficiency of social distancing is highly dependent on external factors, in terms of both physical and social environment. We therefore use this literature as a benchmark, assuming on average 1 meter social distancing,  $E = 0.5$ , and fit the effects of social distancing to data where possible (see Section 6.3).

We simulate mask wearing according to Eq. (6.2), i.e. by multiplicatively reducing the  $\beta$  parameters in different locations. There is a significant body of literature on the effectiveness of mask wearing, including differences based on the material of the mask and the locations in which they are worn [54, 55, 56], as well as changes in efficiency due to re-using or washing them [57, 58]. In general, we focus on the wearing of masks by non-healthcare workers in settings outside the home and estimate mask effectiveness,  $E$ , to be 50% [59], irrespective of the specific of the location. However, after adjustments for compliance the actual, intensity parameter reduction may be much lower than this, which leads us to believe that this represents a conservative estimate.

In JUNE, quarantining of infected people with mild or severe symptoms is relatively straightforward: afflicted individuals do not leave their household for a pre-defined period of time – usually 7 to 14 days – but do not change interaction patterns with the residents in their household. In some versions of quarantine policies, household members must also stay at home and isolate themselves. This is modelled in JUNE along the same lines, only the possibly time-dependent compliance of the population with quarantine measures. Clearly, infected individuals with severe symptoms will always stay at home, until they are either recovered, moved to hospital, or died. It should be noted that quarantine sits on top of other, less individual-driven policy interventions included in JUNE which may restrict movement, such as the closure of companies and leisure venues.

Given the additional danger infectious diseases may pose to the more vulnerable and elderly populations, various policies, usually referred to as “shielding” can be introduced with an aim to protect these individuals. In JUNE, shielding is realised similar to quarantine: vulnerable individuals – usually defined by their age or other characteristics – stay at home and do not interact with. Apart from the definition of relevant characteristics, this only leaves a compliance probability to be introduced, which reduces the participation in any other social settings other than households.

## 6.2 Closure of Venues

Mitigation strategies that aim at reducing infection transmission through changes in individual behaviour may have to be further supplemented through partial or complete closure of certain parts of public life such as companies, transport, schools and universities.

Starting with the closure of companies, JUNE can realise this important measuring sector-specific way. JUNE allows the definition of “key” and “furloughed” workers, again in a sector-specific way. While the former represent those parts of the work force that continue with their essential work as usual, the latter never goes to work and stays at home during regular working hours. For the rest of the work force, JUNE allows the definition of flexible work patterns by assigning daily probabilities for workers go to their companies.

School and university closure is handled similarly to the closure of companies in JUNE. However, in the case of schools we are able to close individual year groups as well as entire schools, and we can identify the children of key workers and have them continue going to school. Since the return to in-person schooling may also be voluntary at certain points, or children may only go to school on certain days of the week, we also can apply a compliance factor at the year group level which is used to probabilistically determine which children attend school on any particular day.

In addition to the partial or complete closure of companies in some industry sectors and of schools or universities, government policies may also close or limit the number or people attending leisure venues, such as restaurant and pubs, cinemas, or similar. In JUNE we are able to fully or partially close different types of leisure venues either nationally or at a more local level, down to super areas. Partial closure is enacted through a change in the probabilities that people attend different venues, which is both sex and age dis-aggregated. Modifications to other leisure activities, such as household visits, are also simple to realise in JUNE, by directly modifying the daily probabilities for such activities to take place.

### 6.3 Policies in the UK

We will now turn to present the mitigation strategies we have implemented in our simulation of the spread of COVID-19 in England and to discuss in more detail some of the parameter choices related to them.

To contextualise the discussion, it is worthwhile to list the operational policy interventions enacted by the UK government from the beginning of March 2020 in an effort to mitigate the spread of the virus, see Table 4.

Date (dd/mm/yy)	Policy
04/03/20	Encourage increased hand-washing
12/03/20	Case isolation at home
16/03/20	Voluntary household quarantine
16/03/20	Stop all non-essential travel
16/03/20	Stop all non-essential contact
16/03/20	Voluntary working from home
16/03/20	Voluntary avoidance of leisure venues
16/03/20	Encourage social distancing of entire population
16/03/20	Shielding of over 70s
20/03/20	Closure of schools and universities
21/03/20	Closure of leisure venues
21/03/20	Stopping of mass gatherings
23/03/20	'Stay at home' messaging
11/05/20	Multiple trips outside are allowed in England only
13/05/20	Encouraged to go back to work if they can while distancing
01/06/20	Meeting in groups of up to 6 outside allowed
01/06/20	Shielding of over 70s relaxed
01/06/20	School reopening for Early Year and Year 6 students
13/06/20	'Support bubbles' allowed
15/06/20	School reopening for Year 10 and 12 students for face-to-face support
04/07/20	Leisure venues allowed to reopen
04/07/20	Household-to-household visits permitted along with overnight stays
24/07/20	Mask wearing compulsory in grocery stores
01/08/20	Shielding is paused
01/08/20	'Eat Out to Help Out' scheme introduced
31/08/20	'Eat Out to Help Out' scheme ends
01/09/20	Schools and Universities allowed to reopen
01/09/20	'Rule of 6' introduced
14/10/20	Tiered local lockdown system introduced

**Table 4:** List of policies introduced in England by the UK Government at different points in time.

We already discussed the modifications of the  $\beta$ -factors related to the individual measures, which we fit to data, or set to values motivated by literature. Here we want to turn to a brief discussion of our modelling of social compliance with some of the measures. In order to estimate the effects of social distancing on the epidemiological development of COVID-19, we implement multiple staggered



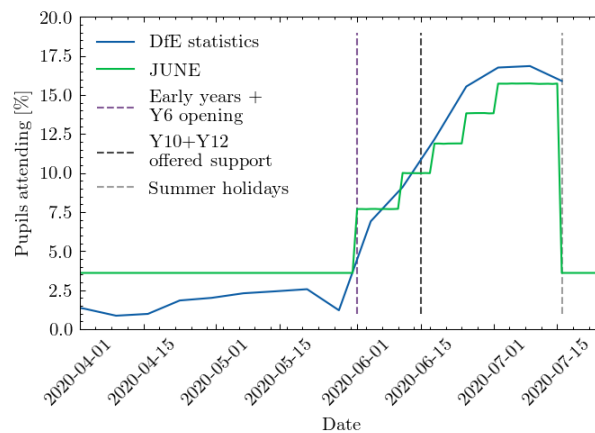
social distancing steps during the first wave of the pandemic between 16<sup>th</sup> March - 4<sup>th</sup> July 2020 and then again going into September 2020 as schools and universities begin to fully reopen. We fit the national compliance,  $C^{(N)}$ , with social distancing between 24<sup>th</sup> March - 11<sup>th</sup> May 2020 in the range 20-100% when fitting the rest the parameters (see Section 7). This is taken to be the harshest social distancing step against and others are determined relative to this fit. The location specific compliance,  $C^{(L)}$ , is set to be 100% in all locations during fitting to avoid parameter degeneracy and then altered manually thereafter. No social distancing is assumed between household members. We derived the compliance with mask wearing from a YouGov survey [60], and we further stratify the results by social environment or locations. Specifically, we assume complete (100%) compliance with mask wearing during commuting, 50% in care homes and no compliance in pubs, schools or in the household. Compliance with mask wearing in grocery stores is assumed to be at 50% before 24<sup>th</sup> July, after which we assume complete compliance given the change in government regulations. Since we already assume low intensity parameters in hospitals due to the significant amount of personal protective equipment (PPE) being worn in these scenarios, we do not apply any additional mask wearing in these settings.

On 16<sup>th</sup> March 2020, the UK Government encouraged people with COVID-19 symptoms to quarantine in their household for 7 days and all those in their household to quarantine for 14 days from symptom onset. We assume that compliance with this measure varies with time as people become more aware of the dangers of COVID-19. Between the 16<sup>th</sup> - 23<sup>rd</sup> March (i.e. the week leading up to the nationwide 'lockdown') we fit compliance with the quarantine policy of those symptomatic to be between 5-45%, and the probability that the rest of the household of a symptomatic individual complies is set to the same fitted value. After 'lockdown' comes into effect, the government tightened these rules to only leave the house for essential trips and one form of exercise per day. To account for this, we increase the symptomatic and household compliance with quarantine to be double their fitted value. In addition, the UK government strongly suggested that people over the age of 70 were encouraged to shield, from 16<sup>th</sup> March 2020. As in the case of quarantine, we assume people become more compliant with this policy over time and that the initial compliance with the shielding policy for this age bracket to increasing from 20% in the first week to 70% afterwards. Indeed, one of the reasons the compliance was set to only 70% even after lockdown is due to the fact that people in this age bracket already have a reduced mobility and interaction potential. A 70% compliance therefore still allows them a small chance to interact with others, e.g. in grocery stores, and any higher compliance figures would mean a complete and unrealistic decoupling of this critical population from any social interactions. The shielding policy initially runs until 1<sup>st</sup> August 2020 and after which the UK government paused the policy.

To model the partial or complete closure of industry sectors, it is important to understand the descriptions of key workers provided by the UK Government [61], and match these up with the relevant 5-digit SIC codes [26]. This ultimately allows us to deduce the proportion key workers in each sector and assign the corresponding key worker attribute probabilistically according to these proportions. In our simulation we encode findings from the ONS [61], reporting that 33% of the total workforce were key workers in 2019 with 14% able to work from home. We therefore set the proportion of key workers, i.e. those who go to work each day, at 19% of the workforce. We use the same logic to also decide which workers are furloughed in JUNE by identifying the 5-digit SIC codes of the relevant affected industries and proportionally assigning the relevant percentage of a given sector to be furloughed. We derive the relevant SIC codes from the Institute for Fiscal Studies in the UK [62], and we dynamically correct for any over or underestimation of furloughed workers by defining the proportion of the workforce who should be furloughed at any particular time, derived from Government reports [63]. A similar dynamic correction is also applied to the key work force. To model the more random work patterns of the remaining part of the work force, we derive a probability that a random worker goes into the company for work from a YouGov survey [64]. We note that in many surveys, including this and others undertaken (e.g. by the ONS [65]), the methodology does not explicitly state if key or furloughed workers were included. We believe, however, that our use of these surveys presents at least a conservative estimate of work attendance.

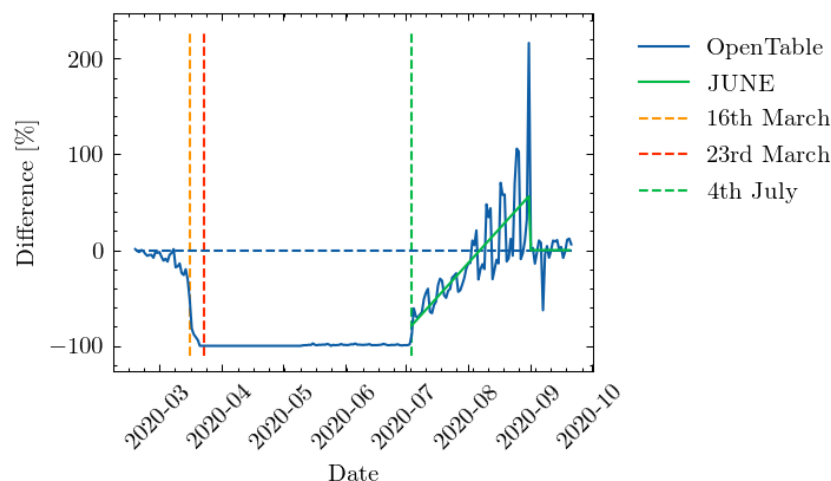
On 20<sup>th</sup> March 2020, all schools and universities in England were asked to close, with the exception that children of key workers could still attend school. To account for the partial school reopening of Early Years (nursery and reception age children) and Year 6 students on 1<sup>st</sup> June 2020, we open up these year groups in JUNE with an attendance compliance based on data derived from the Department for Education (DfE) [66]. While the government also asked schools to offer face-to-face support for Year 10 and 12 students from 15<sup>th</sup> June, we do not include this as the sessions were generally limited and

had an attendance rate around the 10% level [66]. Fig. 16 shows the good agreement in the number of children attending school as derived from JUNE compared with DfE data. The slight deviation from data after 15th June 2020, can be explained by not fully capturing the partial return of Year 10 and 12 students. The good agreement between JUNE and the DfE data before 1<sup>st</sup> June 2020, is of particular note



**Figure 16:** School attendance in JUNE compared to data collected by the UK’s Department for Education [66]

since this option was available only for children where all parents in the household were classified as key workers. This serves as an implicit partial validation of our method of selecting which individuals are key workers, as well as the household and company sector distribution algorithms. From 1<sup>st</sup> September 2020, we reopen schools fully in JUNE, while accounting for a closure for the national school holidays. While the timings of this week-long holiday varies across the country, we assume all schools share the same holiday period 26<sup>th</sup>-30<sup>th</sup> October 2020. Similarly, universities are opened from 1<sup>st</sup> September 2020, but with harsh social distancing measures in place. Given the modelling of where university students live, their inter-mixing is naturally captured in the household component of JUNE (see Section 3.2).



**Figure 17:** Year-on-year restaurant attendance from OpenTable [67] including a fit to the simulated re-opening change in probabilities used to derive the probability that people attend restaurants in JUNE.

On 16<sup>th</sup> March 2020, the UK Government encouraged people to avoid going to leisure venues such as bars and restaurants, although this rule was not imposed through the closure of such venues. However, on the 21<sup>st</sup> March 2020, this closure took place. We model these policies first by reducing the probability that people leave the house from 16<sup>th</sup> March followed by the closure of all relevant leisure venues included

Location	Value	Location	Value	Location	Value
Household	$\beta^{(H)} =$	School	$\beta^{(S)} =$	Work	$\beta^{(W)} =$

**Table 5:** Values for the various  $\beta^{(L,g)}$  in our simulation.

**Figure 18:** Example results for the first wave of COVID-19 in England, from March 1<sup>st</sup>, 2020, to July 31<sup>st</sup>, 2020: daily infections (upper right), daily hospital admissions (upper left), daily deaths in hospitals (lower left), and daily total deaths (lower right). Data from ONS and PHE in black contrasted with a typical JUNE run, in blue. The vertical lines indicate the soft and hard lockdown on March 16<sup>th</sup> and 23<sup>rd</sup>, respectively, and the reopening of leisure venues on July 4<sup>th</sup>.

in the simulation – cinemas, pubs, and restaurants – from 21<sup>st</sup> March 2020. Visits to care homes are also halted from this time. Since many of these venues were permitted to reopen from the July 4<sup>th</sup>, 2020, we assumed all venues reopen at this point. Additionally, data collected by OpenTable suggests that restaurant attendance after that date saw a significant increase likely encouraged by the UK Government’s ‘Eat Out To Help Out’ scheme which we capture in JUNE (see Fig. 17) [67, 68]. For the simulation of other leisure activities, and in particular household-to-household visits, we assume a drop in compliance and a consequently increasing number of such visits. In line with data collected by the ONS [69], we model this be gradually increasing the probability of visiting another household from mid-May until July 4<sup>th</sup>, 2020, when overnight visits were permitted.

## 7 Discussion of model outputs

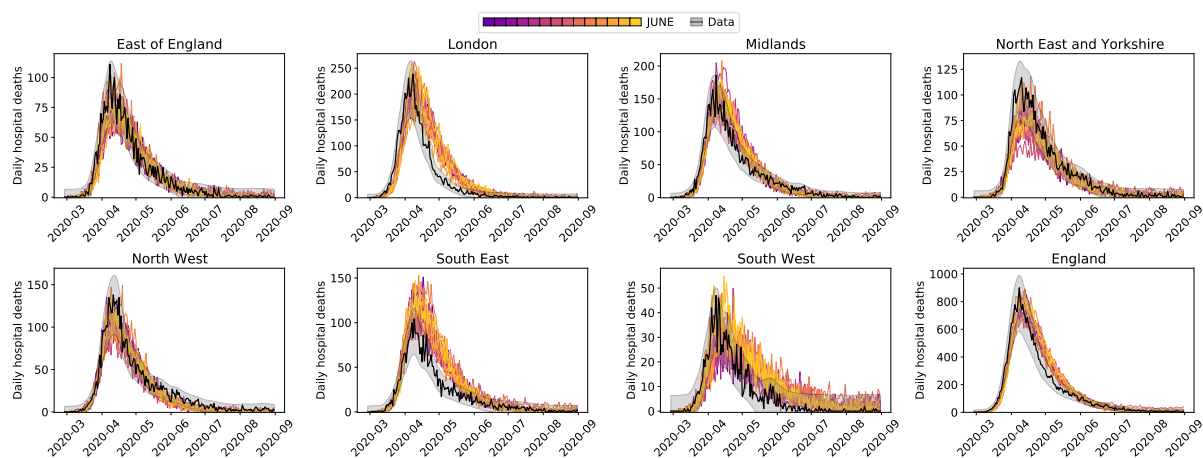
In this section we finally highlight the ability of JUNE to capture intricate social dynamics through a number of model outputs. It is worth stating that we did not yet undertake a systematic fit of model parameters to available data but rather show typical examples from first parameter scans of the code. We refer the reader to Table 5 for values for the various  $\beta^{(L,g)}$  and  $\alpha$ , cf. Eq. (5.1) as well as to ?? for the asymptomatic ratio  $R_A$  and other relevant inputs. In Fig. 18 we exhibit results for the number of daily infections, daily hospital admissions, daily deaths in hospitals and in total, and compare JUNE with corresponding data from NHSE’s SITREP [49] for admissions, the Covid Patient Notification System (CPNS) [70] for hospital deaths, and ONS data [42] for total deaths. The hospital admissions data has been corrected to take into account the delay between the actual date of admission and the date of a positive test result, which is when a patient’s admission is recorded in the SITREP data [49]. This is done by bootstrapping: we sample from daily test lags from the CHESS dataset [48] and subtract from the test result date to estimate the actual daily admissions.

The agreement with data is satisfying as a first check of the validity of the simulation and motivates a more differential look into regional differences. In Fig. 20 we show results of JUNE for four example regions, namely, London, the Midlands, the North East and Yorkshire, and the South West of England. While we start to see minor discrepancies, in particular in the case of the South West, where visibly JUNE overshoots data, the overall agreement is quite satisfying given the very different age and employment profiles in the different regions. In particular the agreement in the slight delay in the onset of the infection suggests that JUNE not only captures global effects but also provides a good and detailed understanding of the spatio-temporal structure of the infection spread.

Another way to gain insights into the character of the infection and its progression is to look at its impact in different bins of age. We therefore look at hospital admissions in England, in different age

**Figure 19:** Hospital deaths for four regions, North-East and Yorkshire (upper left), the Midlands (upper right), South-West (lower left), and London (lower right). Data from ONS and PHE in black contrasted with a typical JUNE run, in blue. The vertical lines indicate the soft and hard lockdown on March 16<sup>th</sup> and 23<sup>rd</sup>, respectively, and the reopening of leisure venues on July 4<sup>th</sup>.

**Figure 20:** Hospital admissions for four age bins, -18 (upper left), 18-65 (upper right), 65-75 (lower left), 75+ (lower right). Data from ONS and PHE in black contrasted with a typical JUNE run, in blue. The vertical lines indicate the soft and hard lockdown on March 16<sup>th</sup> and 23<sup>rd</sup>, respectively, and the reopening of leisure venues on July 4<sup>th</sup>.



**Figure 21:** Daily hospital deaths for each region in England, and England itself, for a number of realisations of JUNE as described in this section. Observed data in black with 3 standard deviation error bands. Data from CPNS [70].

bins in ?? . Again, while not perfect, the agreement of the simulation with data is very satisfying and indicates that JUNE broadly captures how the infection spreads through the population. The agreement

It is worth noting that in JUNE the details of all interactions resulting in infections are stored, enabling future detailed analyses on the sociological nature of SARS-CoV-2 infections.

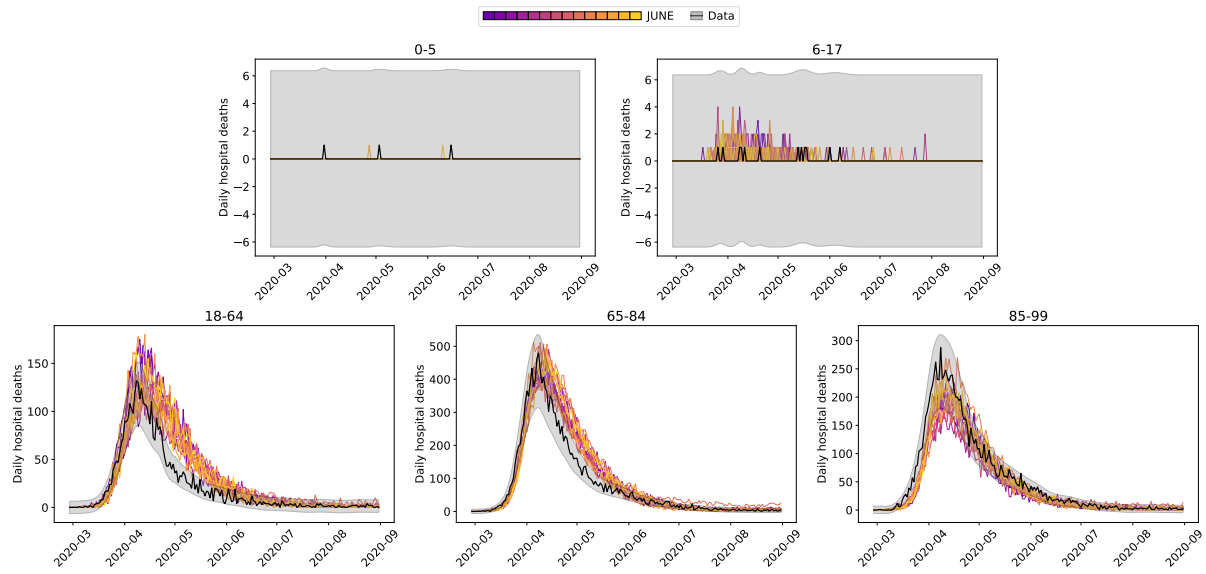
In this section we finally highlight the ability of JUNE to capture intricate social dynamics through a number of model outputs. It is worth noting that the realisations of JUNE presented in the following are ran with parameter sets from the “non-implausible” region of the global parameter space. See Table Table 7 for the ranges of the global parameter space. We refer the reader to Section 8 and Appendix D for more information, and to [71] for a complete exploration of the global parameter space, along with a full uncertainty analysis.

In Figure Fig. 21 we exhibit results for the number of daily deaths in hospital for regions of England and England itself. In addition, in Figure Fig. 22 we show the same realisations for daily deaths in England stratified by age. The agreement with data is satisfying and while there are minor discrepancies for certain outputs, we would like to stress that all of these outputs are simultaneously fit by JUNE.

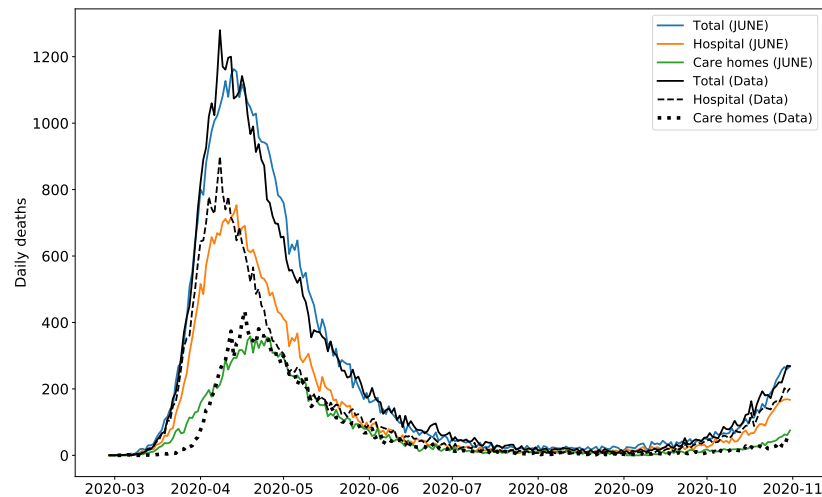
Along with deaths in hospitals, there have been a non-negligible number of fatalities in care homes in England during this pandemic. JUNE successfully models both deaths in hospitals, and deaths within care homes as illustrated in Figure Fig. 23 where there is good agreement with data even into the second wave of the pandemic.

We would like to emphasise that the outputs shown here are illustrative of the capabilities of JUNE to capture the social dynamics of a heterogeneous population giving rise to large differences in disease spread to different age strata and regions.

Deeper analysis of the sociological nature of disease spread can be undertaken with JUNE as details of all interactions resulting in infections are stored, as well as for hospitalisations and deaths. In Figure xx), JUNE displays good agreement with data for percentage cumulative infections by household size and also by ethnicity. We note that these are emergent outputs of JUNE, and are a virtue of the granular virtual population. These are but a few examples of the types of analyses that can be undertaken with JUNE due to its detailed recording of events during the simulation.



**Figure 22:** Daily hospital deaths in England stratified by age, for the same realisations as in Fig. 21. Observed data in black with 3 standard deviation error bands. Data from CPNS [70].



**Figure 23:** Deaths in England illustrated as different lines for total deaths, hospital deaths, and deaths within care homes. Data from ONS [42].

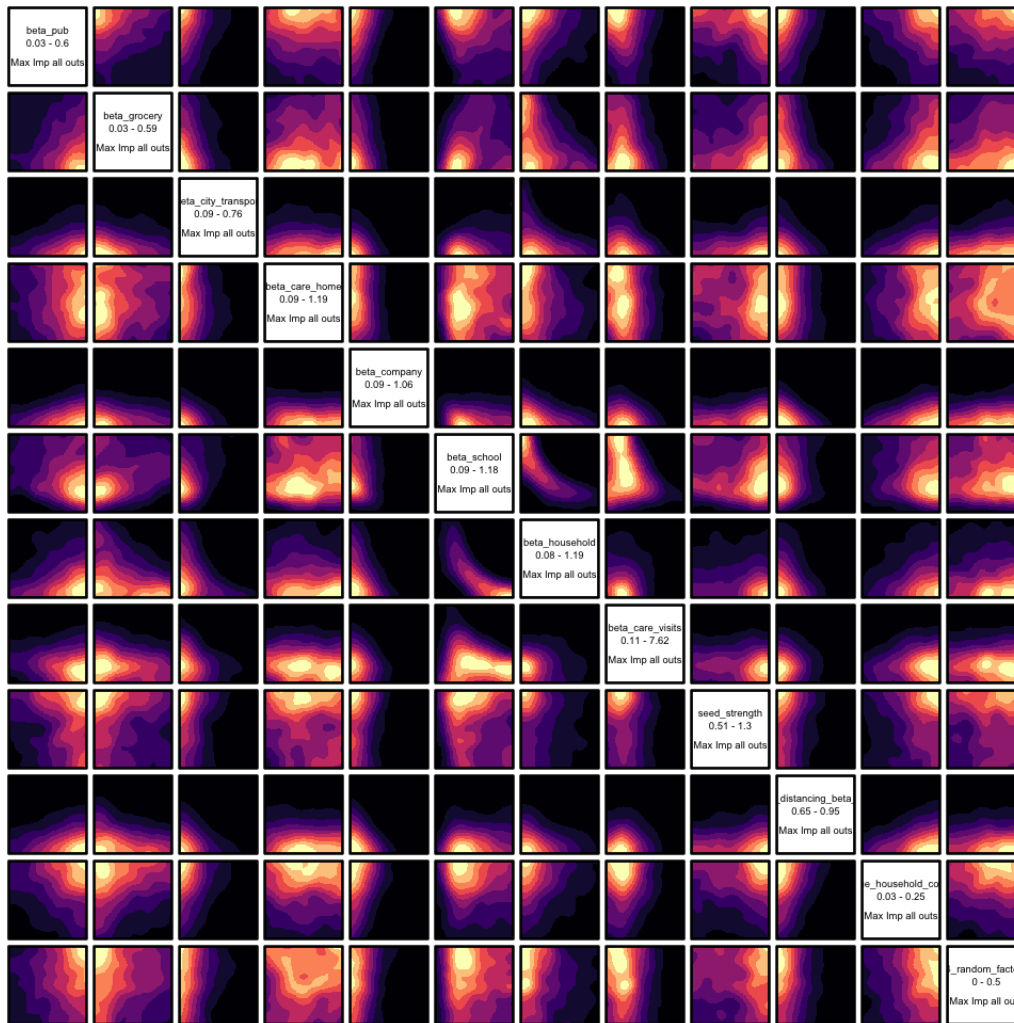
## 8 Fitting via Bayesian Emulation

We now discuss efficient calibration strategies which form a critical part of our ability to extract core insights from JUNE. Fitting a complex model such as JUNE to observed data presents a challenging task. This is mainly due to a) the detailed nature of June and the inevitable computational expense of performing model evaluations, b) the large number of input parameters that we may wish to explore, c) the stochastic nature of the output of JUNE and d) the various uncertainties present in the comparison between model and data. The combination of computational expense and high dimensional input parameter space precludes the use of many parameter exploration methods that rely upon large numbers of model evaluations (including many standard optimisers, sampling approaches such as MCMC etc.). The stochastic output, which implies we will be exploring a much more complex surface, requires methods developed to deal with stochastic functions. Even more challenging is that the substantial uncertainties present imply that we may not even want to optimise for a single “best fit to data” as it may have limited statistical relevance, but instead search for the set of all input parameter values that give acceptable matches between model output and observed data, thereby fully capturing the induced parametric uncertainty.

We hence employ the Bayes linear emulation and history matching methodology [72, 73, 74], a widely applied uncertainty quantification approach designed to facilitate the exploration of large parameter spaces for expensive to evaluate models of deterministic or stochastic form. This approach centres around the concept of an *emulator*: a statistical construct that mimics the slow to evaluate scientific model in question, providing predictions of the model outputs with associated uncertainty, at as yet unevaluated input parameter settings. In contrast to the model, the emulator is extremely fast to evaluate: for example, in the case of JUNE, the emulator exhibited a speed increase of nine orders of magnitude. The emulator provides insight into the model’s structure and, thanks to its speed, can be used to perform the global parameter search far more efficiently than approaches that attempt to use the comparatively slow scientific model itself. Here we give a brief overview of emulation and history matching, but for more details see Appendix D. See also [74, 75, 76, 77] for further examples of its application within epidemiology, [78] for a comparison to Approximate Bayesian Computation in an epidemiological setting, and [79] for a tutorial introduction in the context of systems biology. For an extensive treatment see [73] along with the discussion in [80]. See also [81] for a general introduction to emulation.

Initially, we identify a large set of input parameters to search over, primarily composed of interaction intensity parameters at the group level, along with associated broad ranges, as given in Table 7. We then identify a set of particular model outputs to match to corresponding observed data. Here we focus on hospital deaths (CPNS [70]) and total deaths (ONS) at well-spaced time points throughout the period of the first wave of the epidemic. We then construct Bayes linear emulators for each of the model outputs at each of the chosen time points. The emulators are trained using a set of JUNE runs, initially designed using a 18-dimensional Latin hypercube, and seek to mimic the behaviour of each of the JUNE outputs as a function over the 18-dimensional parameter space. The emulators provide, at each unevaluated input location, an expectation for the possible JUNE model output value and a position dependent variance representing the emulator’s uncertainty about this estimate. Close to known runs the emulator’s uncertainty will be low, however it will increase appropriately as we move to less well explored regions of the parameter space [79]. Note that we deliberately choose to emulate the direct physical outputs of the model as this has multiple benefits for emulator construction, in contrast to emulating a combined metric such as the likelihood (for discussion of this point see [73, 79, 80]).

Due to the emulators’ speed, they are ideal for global parameter exploration. This is performed by constructing an implausibility measure, that gives the distance between the emulator’s expected JUNE model output and the observed data we are trying to match, standardised by all the major uncertainties present: observational errors, emulator uncertainty and structural model discrepancy, the latter being a direct acknowledgement that the model is an imperfect representation of reality (see Appendix D for details). The implausibility measures are used to rule out large regions of the input parameter space that will not provide acceptable matches, and the analysis then proceeds in iterations: a second batch of JUNE runs is performed over the remaining region of parameter space, new emulators constructed (which are only defined over this region), new implausibility measures formed and more parameter space removed. This process is referred to as iterative History Matching [73, 74]. See for example [76] where it was successfully applied to a stochastic disease model with 96 input parameters.



**Figure 24:** 2-dimensional projections of the 18-dimensional input space, coloured by the optical depth of the non-implausible region, which gives the depth of the region of interest conditioned on the two given inputs [79]. These plots are formed from 500,000 emulator evaluations over the input space. The emulators were trained on 3 iterations of 125 JUNE model evaluations.

For the JUNE model, we constructed emulators for hospital deaths and total deaths at 8 time points over the period March to June, for England and for each of the 7 regions, and for the age bins (defined by the SITREP dataset) 0–5, 6–17, 18–64, 65–84, 85<sup>+</sup>. The emulators were trained in three iterations formed from 125 JUNE evaluations each. The emulators were then evaluated at 500,000 locations across the 18-dimensional input space, taking ten minutes on a single processor. The results of the global parameter search are given in the optical depth plots [73] of Figure 24, which shows the location of the “non-implausible” region of interest in various 2-dimensional projections of the 18-dimensional parameter space. The JUNE runs discussed in the preceding section were sampled from this region. Note the various joint constraints on the input parameter space imposed by the matching process, for example the reciprocal relationship between  $\beta_{school}$  and  $\beta_{household}$ . For more details of this approach, of emulator diagnostics, and further output plots see Appendix D.

We can see that the Bayes linear emulation and history matching methodology facilitates the efficient exploration, development and calibration of the highly complex JUNE model using a modest number of runs, a process which would be extremely challenging to perform directly. While here we have performed a provisional exploration of the parameter space as part of the model development, for a full uncertainty analysis of the JUNE model, including the emulator driven generation of full probabilistic forecasts incorporating all major sources of uncertainty, see [71].

## 9 Summary

In this paper we introduced the new JUNE model to simulate the spread of epidemics through a population. JUNE is an individual-based model (IBM) enabling a highly granular geographical and sociological resolution. The frequent and persisting perception that IBMs such as JUNE are heavily parameterised and therefore lack predictive power is misleading. As noted in [82], many of the properties and building blocks of these types of model are not globally fitted to observed cases or fatalities, as is the case for deterministic and stochastic models built from differential equations. Instead, the JUNE framework separates the uncertainty arising from unknown disease dynamics from uncertainties in the population structure, where the latter is informed by demographic statistics and other available data.

The model is formulated and encoded in four distinct layers, **population**, **interaction**, **disease**, and **policy**. Its modular structure allows not only the flexible and seamless addition of many details and novel features, but it also lends itself to application to other populations with different sociological setups. As a first example we discuss its application to the case of the spread of COVID-19 in England, with convincing results underlining the quality of the model and its ability to understand the spread of an epidemic in great detail and with high geographical and sociological resolution.

Studies where JUNE is applied to different settings are forthcoming [83]. One of the strengths of the model is its ability to capture differences in geographical and sociological structure with unprecedented resolution, facilitated through the hierarchical structure in which the **population** is organised. JUNE also allows a flexible yet detailed modelling of daily activities of the virtual population, by combining the geographical position of buildings and other structures with the social interactions taking place. In contrast to other models this enables a very granular understanding of work patterns, leisure activities, etc.. In forthcoming publications we will exploit this high level of detail to try and answer pertinent questions relating to social imbalances in the impact of COVID-19.

## Acknowledgements

This work was undertaken as a contribution to the Rapid Assistance in Modelling the Pandemic (RAMP) initiative, coordinated by the Royal Society.

We are indebted to a number of people who shared their insights into various aspects of the project with us: We would like to thank Sinclair Sutherland for his patience and support in using the ONS database of the census data - without his help it would have been near impossible for us to produce our virtual population. James Nightingale and Richard Hayes provided valuable insights into the construction of efficient algorithms in the initial phase of the project. We are grateful to Bryan Lawrence, Grenville Lister, Sadie Bartholomew and Valeriu Predoi from the National Centre of Atmospheric Science and the University of Reading for assistance in improving the computational performance of the model. We gratefully acknowledge the generous provision of computing time on the Hartree and JASMIN facilities. We would like to thank the GridPP team at Durham and Manchester for their support and computing time spent on their systems. We would also thank Michael Goldstein and TJ McKinley for their statistical and epidemiological advice. Christina Pagel and Rebecca Shipley provided invaluable advice in producing this publication and looking for holes in our arguments.

Map data copyrighted OpenStreetMap contributors and available from <https://www.openstreetmap.org>

**Todo** [Richard]:  
**need to acknowledge compute time on cosma**

J.B., A.C., C.C.L., E.E., M.I.L., A.Q.B., A.S., and H.T. thank the STFC-funded Centre for Doctoral Training in Data-Intensive Science <sup>7</sup> for financial support. F.K. gratefully acknowledges funding as Royal Society Wolfson Research fellow. I.V. gratefully acknowledges Wellcome funding (218261/Z/19/Z).

This work used the DiRAC@Durham facility managed by the Institute for Computational Cosmology on behalf of the STFC DiRAC HPC Facility ([www.dirac.ac.uk](http://www.dirac.ac.uk)). The equipment was funded by BEIS capital funding via STFC capital grants ST/K00042X/1, ST/P002293/1, ST/R002371/1 and ST/S002502/1, Durham University and STFC operations grant ST/R000832/1. DiRAC is part of the National e-Infrastructure.

---

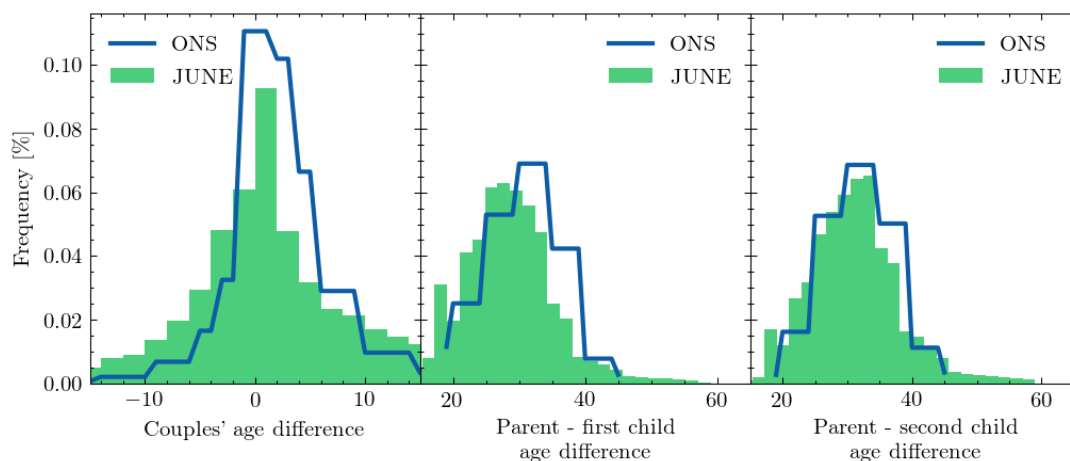
<sup>7</sup><https://ddis.physics.dur.ac.uk/>



## A Algorithms

### A.1 Constructing Credible Households

The ONS divides households into the following broad categories: single, couple, family, student, communal, and other [20]. We populate the households in this ordering, giving preference to those types for which we have the most precise and unambiguous data.



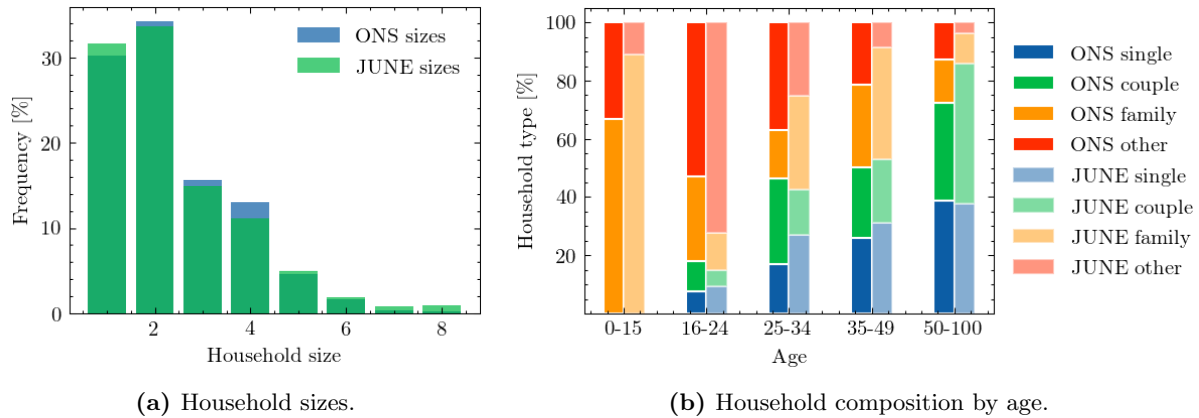
**Figure 25:** Distributions of age differences between partners (left), between parents and their first (middle) and second child (right): outputs of JUNE compared with the input data from the ONS database.

We define and construct households types as follows:

1. Single: These are households with a single person living in it. The census data differentiate single households occupied by an adult or an older adult ( $\geq 65$  years old), and we fill the households accordingly.
2. Couple: These are households occupied by a couple without children. Again, the census differentiates between household with adults or older adults living in them. We preferentially fill these households with two people of different sex, with an age difference sampled from the corresponding UK distribution of age differences at the time of marriage [84] (see also the left panel Fig. 25).
3. Family: These households are defined by the number of adults (singles or couples) and the number of children. A difficulty here is that the census data does not stratify beyond “two or more” children. To compensate for this, we introduce a distribution to select the number of children in these households. To fill a family household, we allocate a female adult first. If there are no female adults available (because they have already been allocated somewhere else), we chose a male adult. In case of families with two adults, we match the person with a partner, preferentially with different sex, and an age difference sampled from the same dataset we use for couples. The census data provides us with the number of dependent children for each OA (area), and we add a suitable number of children according to the age difference between the mother and the  $n$ -th child as given by ONS data collected on birth characteristics [85] (see also the middle and right panels of Fig. 25).
4. Students: From the census data, we know how many student households there are and how many students live in a given OA (area). We uniformly distribute students among their households, assuming a constant ratio of the number of students per household. Students are selected from the population aged between 18 and 25 years old.
5. Communal: We use census data on the number of people in an OA (area) living in a communal establishment, as well as the number of such establishments, such as care homes [21]. The communal establishments are filled last after the types described above their residents will be those who do

not live in any of the other household types. As in the case of student households, we assume a constant ratio of the number of communal residents per establishment.

6. Other: This category encapsulates the uncertain household compositions given by ONS. These may include groups of adults living together, multi-family or multi-generational families. In a similar manner to the communal households, these are filled last with those people that have not yet been allocated.

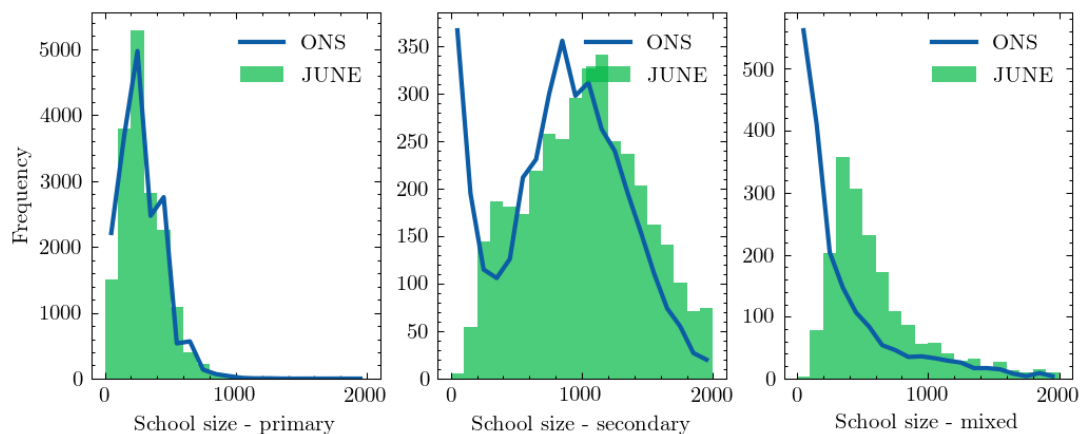


**Figure 26:** Comparisons between outputs of JUNE and data from the ONS database for all England.

As a further test of our household populating algorithm against available data, we compare the JUNE household size distribution and age dependence of people living in different household types with that given by ONS [20]. Fig. 26 demonstrates that the JUNE household composition algorithm clearly produces a household size distribution in good agreement with the census data. We also observe the impact of our assumptions on the composition of families and more complex household compositions (“other”). Given the unknown specifics of certain household composition types, we believe our overall household composition by age to be in reasonable agreement with the data.

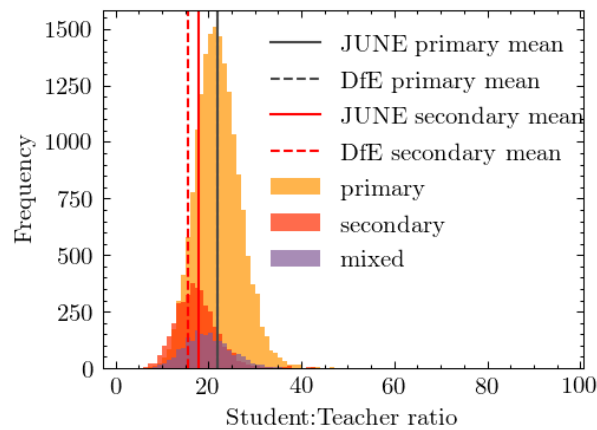
## A.2 Schools

The procedure for assigning children and teachers to schools throughout England is specified in Section 3.3.



**Figure 27:** Distribution of school sizes (left) and distances students have to travel to their schools (right), comparing the JUNE simulation with data

Following our algorithm, we arrive at a distribution of school sizes displayed in Fig. 27, which we see to be in reasonable agreement with the data. Similarly Fig. 28 shows the full distribution of class sizes



**Figure 28:** Distribution of student to teacher ratios for primary schools, secondary schools, and mixed schools.

in JUNE. In the case of COVID-19, most countries have prioritised the return of children to school from younger age brackets. Therefore, recovering good agreement with data particularly in these age brackets is crucial.

### A.3 Work places

We use ONS data on industries and companies in England categorised according to 21 sectors following the Standard Industrial Classification (SIC) code convention (see Table 1) [26] as our framework for differentiating between different types of work.

Companies are initialised according to ONS data on company sizes and sectors at the MSOA (super area) level [27]. We use data on the geographical distribution of company sizes to fix the number of companies at the MSOA (super area) level and use the data on the distribution of sectors to probabilistically assign an industry sector to these companies at the same geographical level. Since the ONS provides information on company sizes by binned size ranges, we take the median size of each bin and assign this to each company. The largest bin is 1000+ employees which we assume to be 1500. It should be noted that companies are not assigned a sector based on their size, but purely on their geography. This does not mean there is no correlation between company size and their sector in JUNE, but that this would arise implicitly based on the geographical distributions, rather than explicitly from data input.

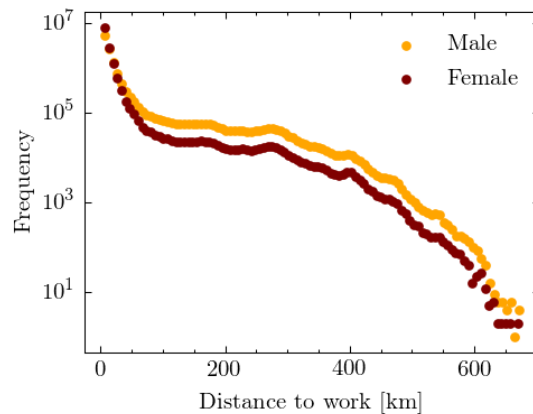
Individuals are assigned a sector attribute probabilistically, following the distributions of sectors disaggregated by sex at the MSOA (super area) level [28]. We determine the MSOA (super area) in which they work according to the ONS commuting origin-destination matrix (or ‘flow’ data) [29] which provides information on the number of people by sex travelling from one MSOA (super area) to another for work. Finally, a matching is carried out between people who work in a certain MSOA (super area), and the companies available to them based on their respective sector attributes. In future work we plan to use additional demographic attributes to assign individuals their sectors and companies.

### A.4 Commuting

The commuting structure in JUNE is built upon the national transport network constructed from nodes representing cities, and edges representing possible transit routes. Commuters are defined as either ‘internal’, i.e. they live and work in a city, or ‘external’, i.e. they live outside a city’s metropolitan boundary and commute into it (see Section 3.4 for more details on how people are assigned locations of work). The metropolitan boundary of each city is defined using data collected by the ONS [35] which specifies the MSOAs (super areas) belong to the cities.

The following procedure is used to determine the groups within which people have the chance to mixing during a commute.

1. For each city, we seed several additional nodes which act as ‘gateway stations’ outside the metropoli-



**Figure 29:** Distance travelled to work by sex according to JUNE. Here we see that men are more likely to travel further to work than women. This is in broad agreement with data presented in [30]

tan area boundary. These serve as funnels into the city and determine the mixing of external commuters. In the case of London we seed 8 stations which are placed evenly around the boundary of the metropolitan area. For all other stations we seed 4 evenly spaced stations North, South, East and West of the city boundaries. These figures are informed by the approximate and number of lines entering each city, and the proportional differences between the number of London public transport links and those of other cities [36].

2. We model the commuting of all people who travel by public transport into a city’s metropolitan area. We assign all external commuters to their nearest gateway station. During each commuting time-step in the simulation, people commuting through the same gateway station are split into ‘carriages’ containing people with whom they have the potential to interact. Similarly, internal commuters are also split into carriages and able to interact with each other.
3. During a commute time-step, each carriage is assigned to be travelling at ‘peak time’ with an 80% probability.
4. The default number of people in an average carriage is fixed to 50 people. For each city this number is adjusted in proportion to data from the UK Department for Transport (DfT) data on overcrowding in trains [33]. This data also dis-aggregates at the level of peak or off-peak travel which is used to further adjust the filling of carriages.
5. The commuting time-step is run twice a day in order to simulate commuting in each direction.

We calculate the distance travelled to work by sex, in Fig. 29, and we see that men are more likely to travel further to work in our model than women. Our findings are in reasonably good agreement with the survey [30] and serve as an independent validation of our model.

## B Contact matrices

We use the contact matrices from the BBC Pandemic survey [15] and supplement them with the `Po1yMod` matrices [10] for interactions of children with other children in the age bracket of 0-12 years. When comparing the matrices that capture interactions in all settings given in the BBC study, an anomaly appears in the matrix describing physical contacts - the original `Po1yMod` data approximately a factor of 3 higher than the BBC matrices in neighboring age bins. We account for that by a simple scaling of the physical contacts by 1/3 before using these data. To arrive at matrices including interactions at home, in school, or in other settings for the age brackets 5-12, missing in the BBC study, we scale the `Po1yMod` setting-inclusive results by a ratio of contacts in the the respective setting for the age bin of 13-14 year old kids, while we assume that the interactions of 0-4 year old children are concentrated at home.

To extract mixing matrices that are suitable for our context-specific simulation, we have to correct for the fact that the reported matrices average over the corresponding age bins in the UK population. For example, contacts teachers and school children are normalised to the full UK population in the respective age bin instead of the number of teachers in schools that actually participate in the interaction. This necessitates rescaling to the number of people in the social context to arrive at corrected social interaction matrices  $\bar{\mathcal{M}}_{ij}^{(H,W,S,O)}$ . This correction step will be detailed in the relevant subsections below.

## B.1 Social mixing at work

The matrices for the age-dependent interaction frequency at the work place show only a very mild correlation with age, typically favouring interactions of workers with a similar age by about a factor of 2. We will therefore not include age effects at the work place into the matrices used in JUNE. To minimize effects due to early retirement, students etc. we average over the ages of 25-60 and we compare this to the average over the working age, 18-64, but correct for an employment rate of 75%. In so doing we arrive at the number of daily contacts for adults at work:

$$n_{AA}^{(W)} = \begin{cases} 4 & (0.35 \text{ physical}) & \text{for ages 25-60} \\ 4.8 & (0.35 \text{ physical}) & \text{for ages 18-64, corrected for employment rate} \end{cases} \quad (\text{B.1})$$

In JUNE we will use  $n_{AA}^{(W)} = 4.8$ , with a ratio of about 7% physical contacts. While it is obvious that different industrial sectors will in reality have very different numbers of daily contacts, with corresponding impact on their vulnerability towards infections, we have not made any attempt to account for such a sector-dependent modulation, apart from effects that naturally arise from different sizes of work forces in different companies.

## B.2 Social mixing in schools

We decompose school populations into year groups labelled with indices  $i \in \{1, 2, \dots, N\}$  for a school with  $N$  year groups and denote teachers with  $T$ . Starting with the interaction of pupils in various year groups an apparent large asymmetry emerges between the summed number of interactions of pupils with adults in the school and of adults with pupils in the BBC data set. This, however, is easily explained by realising that the number of interaction in a given context is normalised to the fraction of the population in a given age bin, irrespective of whether they can participate in the interaction or not. This means that the number of interactions between teachers and pupils have to be renormalised to the ratio of teachers in the adult population – about 500,000 teachers out of 36,300,000 adults, with about 216,000 working in primary and 208,000 working in secondary schools.

Summing the number of interactions of children in the age range of 5-17 with adults in the range 25-65 in schools, and assuming the latter are all teachers yields an average of 0.75 pupil-teacher interactions ( $0.06 = 8\%$  of them physical) per day with very little dependence on the children's age. Conversely, adults have about 0.2 ( $0.02 = 10\%$  of them physical) interactions per day with children in schools, again, relatively independent of the age of the children. Normalising this to the number of teachers in the population, we arrive at about 15 teacher-pupil interactions per day, which fits very well to approximate teacher-pupil ratios of 1:20-1:25<sup>8</sup>. We therefore assume that the individual interaction frequency of one specific teacher-pupil pair is consistently described with 0.75/day. For interactions among adults in the school setting we including the interaction of parents with teachers and of parents among themselves, thereby blurring the picture. We therefore assume that teachers inherit the daily contact frequencies from the work place mixing above. Turning finally to the interactions amongst children, we see a very dominant correlation in age. In order to capture this, we assume that per year of age-difference the number of interactions among children in school,  $n_{KK}^{(S)}$ , will be reduced by a factor  $\xi$ . By fitting to the combination of BBC and PolyMod studies we find, to good approximation,  $\xi = 0.3$  and  $n_{KK}^{(S)} = 2.5$ , with on average 15% of the interactions being physical.

As a consequence we obtain the following social interaction frequency matrix for individual pairings

<sup>8</sup>In fact, for primary schools, the average class size is about 21 pupils, while for secondary schools it is about 16 pupils [25].

at schools:

$$\bar{\mathcal{M}}_{ij}^{(S)} \approx \begin{pmatrix} 4.8 & 0.75 & 0.75 & 0.75 & \dots \\ 15 & 2.50 & 0.75 & 0.25 & \dots \\ 15 & 0.75 & 2.50 & 0.75 & \dots \\ 15 & 0.25 & 0.75 & 2.50 & \dots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix}, \quad (\text{B.2})$$

where the first row and the first column specify the interactions between teachers and students in different year groups, and the second and following row and columns are populated by interactions of the pupils with other pupils across year groups ordered by age.

### B.3 Social mixing at home

In our model we decompose the household population into 4 subgroups, namely children ( $K$ , ages 0-19), young adults ( $Y$ , 18-24), adults ( $A$ , 25-64), and older adults ( $O$ , 65+). We therefore arrive at a  $4 \times 4$  matrix of corrected social interactions at home,  $\bar{\mathcal{M}}_{ij}^{(H)}$ , where the indices  $i, j \in \{K, Y, A, O\}$ . In the following we will detail how we arrive at the various matrix elements. When correcting for the impact of social environment, *i.e.* the household compositions, we will ignore household compositions which are listed as “other” in the ONS database, due to a lack of detailed information (see Section 3.2 for more details). When using these data, we will use numbers in units of millions,  $H_{OAYK}$  of households with a composition of  $O$  older adults,  $A$  adults,  $Y$  independent children or young adults living at home, and  $K$  children aged 0-19.

- $\bar{\mathcal{M}}_{OO}^{(H)}$ : we ignore the case of care homes or other facilities with more than two residents. Then the average interaction frequency from the BBC data is given by  $n_{OO}^{(H)} = 0.78$  (0.44 physical) and 0.62 at weekends<sup>9</sup>. With  $H_{2000} = 2.131$  and  $H_{1000} = 3.294$ <sup>10</sup>

$$\bar{\mathcal{M}}_{OO}^{(H)} = 0.78 \cdot \frac{2H_{2000} + H_{1000}}{2H_{2000}} \approx 1.4. \quad (\text{B.3})$$

- $\bar{\mathcal{M}}_{AA}^{(H)}$ : the interaction frequency between adults aged 20-65 at home from the BBC data is given by  $n_{AA}^{(H)} = 1.2$  (0.74 = 62% of them physical).

$$\bar{\mathcal{M}}_{AA}^{(H)} = 1.2 \cdot \frac{\sum_{x,y} (2H_{02xy} + H_{01xy})}{\sum_{x,y} 2H_{02xy}} \approx 1.34, \quad (\text{B.4})$$

where  $\sum_{x,y} H_{02xy} = 8.751$  and  $\sum_{x,y} H_{01xy} = 7.644$ .

- $\bar{\mathcal{M}}_{YY}^{(H)}$ : the interaction frequency between young adults age 18-26 at home from the BBC data is given by  $n_{YY}^{(H)} = 1.3$  (0.4 = 34% of them physical). There is no obvious household correction that we can apply, but the number of contacts is relatively close to the value of  $\bar{\mathcal{M}}_{AA}^{(H)} = 1.34$ , so we will assume that young adults interact with each other with a frequency similar to that of adults:

$$\bar{\mathcal{M}}_{YY}^{(H)} = \bar{\mathcal{M}}_{AA}^{(H)}. \quad (\text{B.5})$$

It is worth noting that the age range for young adults is relatively narrow, and that there will be edge effects that may effectively increase the interaction frequency.

- $\bar{\mathcal{M}}_{YA}^{(H)}$  and  $\bar{\mathcal{M}}_{AY}^{(H)}$ : we have  $n_{YA}^{(H)} \approx 0.7$  with a relatively steep decline with the age of the young adults, which we attribute to the fact that with increasing age young adults move out of their parents’ home. To obtain some better understanding of the situation, we look at the interaction of

<sup>9</sup>One may speculate in how far this drop is a reflection of uncertainties in the data or a true “physical” effect, for example due to visitors, travel, or similar.

<sup>10</sup>Here and in the following, the numbers of different household configurations are taken from [86].

adults in the age range 40-65 with young adults, aged 18-24. From this we arrive at an average of  $n_{AY}^{(H)} = 0.17$  ( $0.07 = 40\%$  of them physical).

To relate this to a corrected value we must make an assumption concerning the number of young adults in the three age bins that still live with their parents, which we take as 75%, 50%, and 40% for the three age bins. To correct the AY number we assume that the majority of households with young adults living as non-dependent children with their parents is composed of households with one young adult adult. Therefore:

$$\begin{aligned}\bar{\mathcal{M}}_{YA}^{(H)} &= \frac{1}{3} \left[ \frac{0.87}{0.75} + \frac{0.65}{0.5} + \frac{0.55}{0.4} \right] \approx 1.3 \\ \bar{\mathcal{M}}_{AY}^{(H)} &= 0.17 \cdot \frac{\sum_{xy} (2H_{02xy} + H_{01xy})}{\sum_y (2H_{021y} + H_{011y})} \approx 1.47,\end{aligned}\tag{B.6}$$

where  $\sum_y H_{021y} = 1.514$  and  $\sum_y H_{011y} = 0.946$ .

- $\bar{\mathcal{M}}_{KK}^{(H)}$ : the average number of daily contacts at home between children age 0-17 is  $n_{KK}^{(H)} = 0.47$  (79% of them physical). Assuming all children live as dependents with their parents, and demanding that households with “2 or more children” (ONS classification) have, on average, 2.3 children to account for the UK reproduction rate, we arrive at:

$$\bar{\mathcal{M}}_{AA}^{(H)} = 0.87 \cdot \frac{\sum_x (H_{02x1} + H_{01x1}) + 2.3(H_{02x2} + H_{0.1x2})}{\sum_x 2.3(H_{02x2} + H_{0.1x2})} \approx 1.2.\tag{B.7}$$

- $\bar{\mathcal{M}}_{KA}^{(H)}$  and  $\bar{\mathcal{M}}_{AK}^{(H)}$ : to account for contacts of children with adults we will use sliding age windows in dependence on the age of the child, using that parents are usually between 20-40 years older than their children. We then arrive at  $n_{KA}^{(H)} = 1.27$  (70% of them physical) and  $n_{AK}^{(H)} = 0.67$ , the former with an only mild dependence on the age of the child, while the latter shows clear edge effects for the first and last bins of the adult age distribution. These numbers translate into:

$$\bar{\mathcal{M}}_{KA}^{(H)} = 1.27\tag{B.8}$$

$$\bar{\mathcal{M}}_{AK}^{(H)} = 0.67 \cdot \frac{\sum_{x,y} (2H_{02xy} + H_{01xy})}{\sum_{[x,y]} 2(H_{02x1} + H_{02x2}) + (H_{01x1} + H_{01x2})} \approx 1.69.^{11}\tag{B.9}$$

We will also assume that the interaction frequency and intensity of children and young adults living in the same household is determined by

$$\bar{\mathcal{M}}_{KY}^{(H)} = \bar{\mathcal{M}}_{KA}^{(H)} \quad \text{and} \quad \bar{\mathcal{M}}_{YK}^{(H)} = \bar{\mathcal{M}}_{AK}^{(H)}\tag{B.10}$$

- $\bar{\mathcal{M}}_{O,KYA}^{(H)}$  and  $\bar{\mathcal{M}}_{KYA,O}^{(H)}$ : we assume that interactions of children, young adults, and adults with older adults at home have three different realizations:
  1. as regular contacts in a multi-generational household, where we assume that older adults behave like adults in terms of interaction frequency and intensity;
  2. as regular contacts between children and their grand-parents who act as child-minders while the parents are at work;
  3. through regular or sporadic visits, where we again assume that interactions of older adults follow the pattern of adults.

As a result we obtain the following social mixing matrix

$$\bar{\mathcal{M}}_{ij}^{(H)} = \begin{pmatrix} 1.2 & 1.69 & 1.69 & 1.69 \\ 1.27 & 1.34 & 1.47 & 1.50 \\ 1.27 & 1.30 & 1.34 & 1.34 \\ 1.27 & 1.50 & 1.34 & 2.00 \end{pmatrix} = \begin{array}{c|cccc} & \text{K} & \text{Y} & \text{A} & \text{O} \\ \hline \text{K} & 1.2 & 1.69 & 1.69 & 1.69 \\ \text{Y} & 1.27 & 1.34 & 1.47 & 1.50 \\ \text{A} & 1.27 & 1.30 & 1.34 & 1.34 \\ \text{O} & 1.27 & 1.50 & 1.34 & 2.00 \end{array}, \quad (\text{B.11})$$

where, for convenience, we have made the entries explicit.

## B.4 Social mixing in other venues

Social venues (“pubs”, “cinemas”, and “groceries”) in JUNE are assumed to have only one subgroup, “attendees”, meaning that the social mixing matrix for these interactions is a single-element;  $\bar{\mathcal{M}}^{(P)} = 3$ ,  $\bar{\mathcal{M}}^{(C)} = 3$ ,  $\bar{\mathcal{M}}^{(G)} = 1.5$  for “pubs”, “cinemas”, and “groceries” respectively. These values were chosen heuristically according to the estimated number of contacts in each location relative to the the number of contacts set elsewhere (as discussed above). Given that we do not consider different subgroups in these locations, making the matrix single-valued, these numbers only serve the purpose of intuitively introducing a hierarchy of contact intensities ( $\beta$  parameters) into the model structure. Since the intensity parameters are fitted to data (see Section 8), the form of Eq. (5.2) ensures that the choice these social mixing matrices values will not significantly affect the probabilities of transmission.

Hospitals have three subgroups: medical staff, ward patients, and ICU/ITU patients. The social mixing matrix for hospitals (where the superscript  $M$  refers to “medical facility”) is

$$\bar{\mathcal{M}}_{ij}^{(M)} = \begin{pmatrix} 5 & 10 & 10 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix} \quad \text{and} \quad \phi_{ij}^{(M)} = \begin{pmatrix} 0.05 & 1 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix}, \quad (\text{B.12})$$

where  $(i, j) \in \{S, W, I\}$  denoting the three subgroups, medical staff, ward patients, and ICU/ITU patients, respectively. The number of contacts between a medic and patients 10 represents the average number of patients per medic. We assume that a patient is visited by a medic once per characteristic time, set to 8 hours for hospitals. The number of contacts between patients is irrelevant, as patients are by definition already infected, but is set to zero.

Social mixing in care homes considers three subgroups: workers, residents, and visitors, with matrix

$$\bar{\mathcal{M}}_{ij}^{(CH)} = \begin{pmatrix} 15 & 15 & 1 \\ 1.5 & 4 & 20 \\ 1.5 & 6 & 0 \end{pmatrix} \quad \text{and} \quad \phi_{ij}^{(CH)} = \begin{pmatrix} 0.05 & 1 & 0 \\ 1 & 0.4 & 1 \\ 0 & 0.5 & 0 \end{pmatrix}, \quad (\text{B.13})$$

with a characteristic time of 24 hours. The seemingly-high number of contacts between workers and visitors, and residents and visitors is to compensate for the characteristic time of 24 hours; if visitors were to be present in a care home for a full characteristic time, they would experience this many contacts, but visits are a day time activity which take only a few hours, resulting in fewer contacts.

Finally, universities are modelled as having 6 groups to represent professors and 5 distinct groups of students (for the moment based only on age 19-23), with diagonal elements  $\bar{\mathcal{M}}_{i=j}^{(U)} = 2$  and off-diagonal elements  $\bar{\mathcal{M}}_{i \neq j}^{(U)} = 0.75$ , and all  $\phi_{ij}^{(U)} = 0.25$ .

## B.5 Deriving contact matrices from JUNE

We derive the contact matrices in Fig. 9 by simulating a week of pre-lockdown activity. For each person, in each subgroup  $i$ , in each venue, we choose the required  $N_{ij}$  people (with replacement) for all (non-empty) subgroups  $j$  in that venue (where  $N_{ij}$  from the relevant social mixing matrix). We populate “raw” contact matrices using these selected people. As these contacts are then uni-directional, we make the same corrections as in [15] to account for reciprocal contacts. We hope to produce contact matrices derived from constructing self-consistent (reciprocal) networks of contacts within groups in future work.



Name	Function	Source
$C_T$	constant with time $T$	
$\beta_I$	$\beta_{2.29,19.05,0.39,39.8}(t)$	[39]
$LN_M$	$LN_{0.83,5.7}(t)$	*
$\beta_H$	$\beta_{1.35,3.68,0.05,27.1}(t)$	[87]
$\beta_D$	$\beta_{1.21,1.97,0.08,12.9}(t)$	[87]
$LN_{ICU}$	$LN_{1.41,0.9}(t)$	[88]
$e_{ICU}$	$e_{1.06,0.89,12}(t)$	[88]
$e_D$	$e_{1.23,1,9.69}(t)$	[88]

**Table 6:** Characteristic functions and their parameters. \*We constrained the time from symptom onset to hospitalisation through private communication with hospital physicians at early stages of the first COVID-19 wave of infections. We later checked this assumed profile against published data and found our values to be broadly consistent.

## C Details on modelling health trajectories

For the times spent in different stages of disease progression we use a variety of functions, namely intervals of constant length, scaled and shifted  $\beta$  functions, scaled log-normal distributions, and exponential Weibull distributions, given by

$$\begin{aligned}
 C_{t_{\text{end}}}(t) &= \Theta(t_{\text{end}} - t) \\
 \beta_{a,b,l,S}(t) &= \beta_{a,b} \left( \frac{t-l}{S} \right) \\
 LN_{s,S}(t) &= LN_s \left( \frac{t}{S} \right) \\
 e_{a,c,S}(t) &= e_{a,c} \left( \frac{t}{S} \right).
 \end{aligned} \tag{C.1}$$

The trajectories and their building blocks to construct the corresponding time intervals infected individuals spend in various stages of the disease are listed in Table 3.

## D Calibration via Bayes Linear Emulation and History Matching

We now provide more details of the Bayes linear emulation and history matching process outlined in section 8. To set up the history matching problem, we identify a large set of 18 input parameters to the JUNE model to explore. This set is composed mainly of contact intensity parameters, but also contains such parameters governing social distancing effects, compliance and physical contact, with each parameter specified along with associated broad ranges, in Table 7. We denote this set of parameters by the 18-dimensional vector  $x$  and denote the initial search region defined by their combined ranges as  $\mathcal{X}_0$ . We identify a set of outputs to match to observed data, specifically the deaths and total deaths for England and for each of the 7 regions, and for the age bins 0–5, 6–17, 18–64, 65–84, 85+, at the time points of the 20th, 28th March, 5th, 13th, 21st and 29th April, 12th, 26th May, and 9th, 23rd June. We represent the list of all these outputs as the vector  $f$ .

We note that the JUNE model can now be viewed as a function that maps the inputs  $x$  to the vector of all outputs of interest  $f(x)$ . As we cannot evaluate the model  $f(x)$  exhaustively over the full parameter space  $\mathcal{X}_0$  due to computational expense, we mimic it using a fast to evaluate (but uncertain) Bayesian emulator. For an individual output  $f_i(x)$ , representing for example, the total deaths in England on the 28th March, we construct an emulator of the form

$$f_i(x) = \sum_j b_{ij} g_{ij}(x_{A_i}) + u_i(x_{A_i}) + v_i(x) \tag{D.1}$$

where the first term on the RHS is a regression term designed to capture global behaviour, composed of known deterministic functions  $g_{ij}$  (with a common choice being low order polynomials) of the active

Input Parameter ( $x_i$ )	Type	Range
$\beta_{pub}$	contact intensity	[0.02,0.6]
$\beta_{grocery}$	.	[0.02,0.6]
$\beta_{cinema}$	.	[0.02,0.6]
$\beta_{university}$	.	[0.02,0.6]
$\beta_{city\ transport}$	.	[0.08,0.77]
$\beta_{inter\ city\ transport}$	.	[0.08,1.2]
$\beta_{hospital}$	.	[0.08,1.2]
$\beta_{care\ home}$	.	[0.08,1.2]
$\beta_{company}$	.	[0.08,1.2]
$\beta_{school}$	.	[0.08,1.2]
$\beta_{household}$	.	[0.08,1.2]
$\beta_{care\ visits}$	.	[0.1,8]
$\beta_{household\ visits}$	.	[0.08,1.2]
$\alpha_{physical}$	physical contact factor	[1.8,3]
$\alpha_{seedstrength}$	seeding	[0.5,1.3]
$M_{quarantine\ household\ compliance}$	compliance	[0.034,0.26]
$M_{social\ distancing\ \beta\ factor}$	social distancing	[0.65,0.95]
$M_{sd4randomfactorall}$	social distancing	[0.004,0.5]

**Table 7:** The input parameters that form the 18-dimensional vector  $x$  explored in the global parameter search, their type and their ranges that define the search region  $\mathcal{X}_0$ .

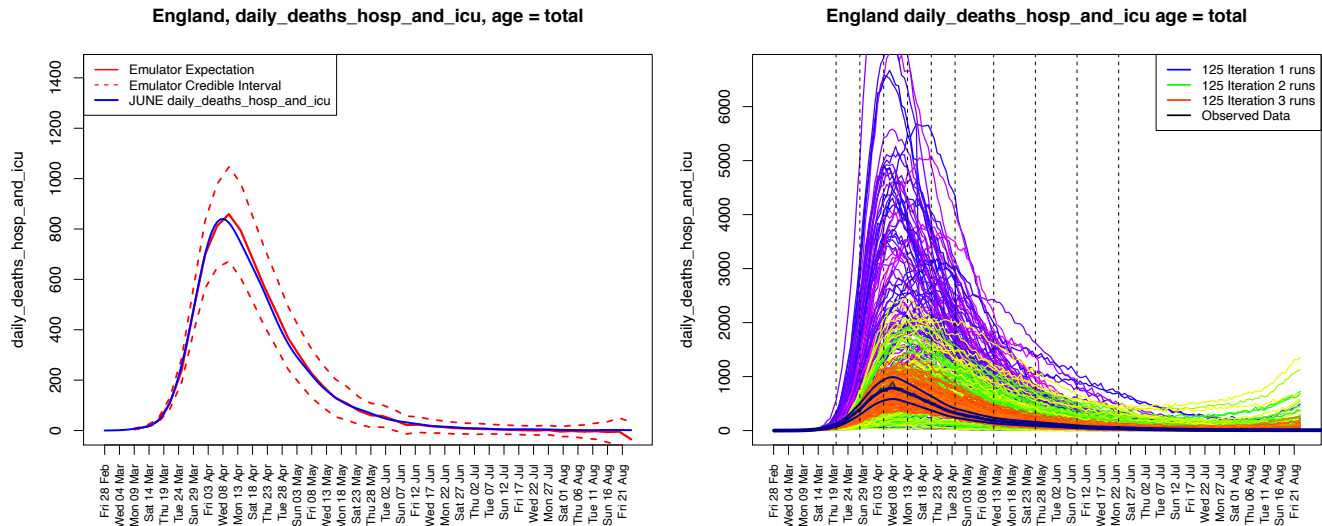
variables  $x_{A_i}$ , which are a subset of the inputs that are found to be most influential for output  $f_i(x)$ , and of the  $b_{ij}$  which are unknown regression coefficients. The middle term  $u_i(x_{A_i})$  is a Gaussian process with various forms of correlation structure available, capable of mimicking large classes of functions, which has the flexibility to capture more local behaviour of  $f_i(x)$ , and  $v_i(x)$  is an uncorrelated nugget that represents the effect of the remaining inactive input variables, and/or any stochasticity exhibited by the model.

We perform an initial space filling set of  $n = 125$  runs  $D = (f(x^{(1)}), f(x^{(2)}), \dots, f(x^{(n)}))$  with the  $x^{(i)} \in \mathcal{X}_0$  chosen using a maximin Latin hypercube design. The emulators are updated by the runs  $D$  using the Bayes linear update equations [73], and hence can give a prediction with corresponding uncertainty, of the unobserved  $f(x)$  at a new, previously unevaluated input point  $x$ , in the form of the adjusted expectation  $E_D(f_i(x))$  and the adjusted variance  $\text{Var}_D(f_i(x))$  respectively. The emulators have to satisfy extensive diagnostics [73, 89], an illustrative example of which is given in Figure 30, left panel, which shows the emulator prediction  $E_D(f_i(x))$  for  $f_i(x)$  across several time points (the solid red line) and the prediction interval  $E_D(f_i(x)) \pm 3\sqrt{\text{Var}_D(f_i(x))}$  (the red dashed lines) along with the held out run output  $f(x)$  (the blue line) which the emulator has not previously seen, showing excellent agreement between emulator and model. The emulator evaluation takes a fraction of a second, and mimics the JUNE model well.

By confronting the emulators with the observed data vector  $z$  corresponding to the outputs in  $f$ , and incorporating major sources of uncertainty (e.g. observation error, structural model discrepancy, stochasticity), we can rule out large parts of the input parameter space  $\mathcal{X}_0$  as implausible. We do this using an Implausibility measure, for which the univariate version  $I_i(x)$ , is defined for each output as

$$I_i^2(x) = \frac{(E_{D_i}(f_i(x)) - z_i)^2}{\text{Var}_D(f_i(x)) + \sigma_{e_i}^2 + \sigma_{\epsilon_i}^2} \quad (\text{D.2})$$

where  $E_D(f_i(x))$  and  $\text{Var}_D(f_i(x))$  are the emulator expectation and variance as before,  $z_i$  is the observed data point corresponding to model output  $f_i$ ,  $\sigma_{e_i}^2$  is the variance of the observation error  $e_i$  (a random quantity representing the imperfections of the measurement process), and  $\sigma_{\epsilon_i}^2$  is the variance of the model discrepancy  $\epsilon_i$  (an often neglected random quantity representing the imperfections of the model [72, 73, 90]). If  $I_i(x)$  is large, it is because the emulator expectation for  $f_i(x)$  is very far from the observed data  $z_i$ , even given all the major sources of uncertainty, and therefore the input parameter  $x$  is highly unlikely to yield model output similar to observed data were we to evaluate JUNE there, and hence  $x$  could be discarded from further analysis. A typical cutoff maybe  $I_i(x) < c$  where  $c = 3$ , motivated by



**Figure 30:** Left Panel: an example diagnostic showing the emulator prediction  $E_D(f_i(x))$  for  $f_i(x)$  across several time points (the solid red line) and the prediction interval  $E_D(f_i(x)) \pm 3\sqrt{\text{Var}_D(f_i(x))}$  (the red dashed lines) along with the held out smoothed run output  $f(x)$  (the blue line). The emulator captures the behaviour of the JUNE model well. Right panel: daily hospital deaths for all of England, showing the progression of the runs from iterations 1, 2 and 3 used in the history matching process (in purple, green and red respectively). Observed data (smoothed and original) in black. Vertical dashed lines: emulated outputs.

Pukelsheim’s powerful 3-sigma rule<sup>12</sup>. There are various ways to combine implausibility measures for each of the individual outputs, the simplest being to maximise:  $I_M(x) = \max_i I_i(x)$ , although other more nuanced and/or robust versions are available, that capture more of the multivariate behaviour [73].

We now employ iterative history matching [73], a parameter search method that seeks to identify all parts of parameter space that would give rise to acceptable matches between model output and observational data. This proceeds at the  $j$ th iteration (or wave), by constructing emulators using the current set of runs, removing the implausible parts of the input space to define the new non-implausible region  $\mathcal{X}_j = \{x \in \mathcal{X}_0 : I_M(x) < c\}$ , designing and performing a new space filling set of runs across the reduced input space  $\mathcal{X}_j$  and re-emulating, but now with a more accurate emulator defined only over the reduced region  $\mathcal{X}_j$ . For further discussion see [73, 79, 80], but it suffices to note that the iterative nature of history matching is key, as it allows later iteration emulators to become far more accurate as they are only employed over far smaller parts of the input space, and are hence informed by a much higher density of runs.

The observed data for total deaths was obtained from the ONS, while the hospital deaths data is taken from CPNS - the Covid Patient Notification System [70]. For each output corresponding to the element of  $f$ , the data was first smoothed slightly with a standard kernel smoother, to reduce the day-to-day stochasticity. The observation error and model discrepancy variances for each output were each decomposed into multiplicative and additive components to represent possible systematic biases, in addition to a scaled  $\sqrt{n}$  component for the observation error only, to model the noisy count process. For example, we have the decompositions  $\sigma_{\epsilon_i}^2 = \alpha_{mult,\epsilon_i}^2 z_i^2 + \gamma_{add,\epsilon_i}^2$ , with  $\alpha_{mult,\epsilon_i} = 0.06$  and  $\gamma_{add,\epsilon_i}^2 = 3/2$ , and  $\sigma_{e_i}^2 = \alpha_{mult,e_i}^2 z_i^2 + \gamma_{add,e_i}^2 + (\delta_{corr,e_i} \sqrt{z_i})^2$ , with  $\alpha_{mult,e_i} = 0.06$  and  $\gamma_{add,e_i}^2 = 3/2$  and  $\delta_{corr,e_i} = 0.25$  governed by the mitigation of the smoothing process.

As described in section 8, we performed 3 waves of the history match with 125 runs each wave, finding that the emulators were of sufficient accuracy after the third wave. Figure 30, right panel, shows the progression of the runs from iterations 1, 2 and 3 used in the history matching process (in purple, green and red respectively) for the daily hospital deaths in England output, with the data (original and

<sup>12</sup>Pukelsheim’s 3-sigma rule is the powerful, general, but somewhat underused result that states for any continuous unimodal distribution, 95% of the probability must lie within  $\mu \pm 3\sigma$ , regardless of its asymmetry or skew.

smoothed) in black. We can see that the third iteration runs are vastly improved and surround the observed data. These allow accurate emulators to be constructed that can identify the region of input space of interest, which were used to construct Figure 24, as discussed in section 8.

## References

- [1] R. E. Russell, R. A. Katz, K. Richgels, D. P. Walsh, and E. Grant. A framework for modeling emerging diseases to inform management. *Emerging Infectious Diseases*, 23(1):1–6, 2017.
- [2] Namdi Brandon, Kathie L Dionisio, Kristin Isaacs, Rogelio Tornero-Velez, Dustin Kapraun, R Woodrow Setzer, and Paul S Price. Simulating exposure-related behaviors using agent-based models embedded with needs-based artificial intelligence. *Journal of exposure science & environmental epidemiology*, pages 1–10, 2018.
- [3] Amy H Auchincloss, Samson Y Gebreab, Christina Mair, and Ana V Diez Roux. A review of spatial methods in epidemiology, 2000–2010. *Annual review of public health*, 33:107–122, 2012.
- [4] Abdulrahman M El-Sayed, Peter Scarborough, Lars Seemann, and Sandro Galea. Social network analysis and agent-based modeling in social epidemiology. *Epidemiologic Perspectives & Innovations*, 9(1):1, 2012.
- [5] Rebecca J Rockett, Alicia Arnott, Connie Lam, Rosemarie Sadsad, Verlaine Timms, Karen-Ann Gray, John-Sebastian Eden, Sheryl Chang, Mailie Gall, Jenny Draper, et al. Revealing covid-19 transmission in australia by sars-cov-2 genome sequencing and agent-based modeling. *Nature medicine*, 26(9):1398–1404, 2020.
- [6] Neil M Ferguson, Derek AT Cummings, Christophe Fraser, James C Cajka, Philip C Cooley, and Donald S Burke. Strategies for mitigating an influenza pandemic. *Nature*, 442(7101):448–452, 2006.
- [7] Dennis L Chao, M Elizabeth Halloran, Valerie J Obenchain, and Ira M Longini Jr. Flute, a publicly available stochastic influenza epidemic simulation model. *PLoS Comput Biol*, 6(1):e1000656, 2010.
- [8] Elizabeth Hunter, Brian Mac Namee, and John D Kelleher. A taxonomy for agent-based models in human infectious disease epidemiology. *Journal of Artificial Societies and Social Simulation*, 20(3), 2017.
- [9] Sameera Abar, Georgios K Theodoropoulos, Pierre Lemarinier, and Gregory MP O’Hare. Agent based modelling and simulation tools: A review of the state-of-art software. *Computer Science Review*, 24:13–33, 2017.
- [10] Joël Mossong, Niel Hens, Mark Jit, Philippe Beutels, Kari Auranen, Rafael Mikolajczyk, Marco Massari, Stefania Salmaso, Gianpaolo Scalia Tomba, Jacco Wallinga, et al. Social contacts and mixing patterns relevant to the spread of infectious diseases. *PLoS Med*, 5(3):e74, 2008.
- [11] Elizabeth J. Williamson, Alex J. Walker, Krishnan Bhaskaran, Seb Bacon, Chris Bates, Caroline E. Morton, Helen J. Curtis, Amir Mehrkar, David Evans, Peter Inglesby, Jonathan Cockburn, Helen I. McDonald, Brian MacKenna, Laurie Tomlinson, Ian J. Douglas, Christopher T. Rentsch, Rohini Mathur, Angel Y. S. Wong, Richard Grieve, David Harrison, Harriet Forbes, Anna Schultze, Richard Croker, John Parry, Frank Hester, Sam Harper, Rafael Perera, Stephen J. W. Evans, Liam Smeeth, and Ben Goldacre. Factors associated with covid-19-related death using opensafely. *Nature*, 584(7821):430–436, 2020.
- [12] Claire Zoellner, Rachel Jennings, Martin Wiedmann, and Renata Ivanek. Enable: An agent-based model to understand listeria dynamics in food processing facilities. *Scientific reports*, 9(1):1–14, 2019.
- [13] Albert M Lund, Ramkiran Gouripeddi, and Julio C Facelli. Stham: an agent based model for simulating human exposure across high resolution spatiotemporal domains. *Journal of Exposure Science & Environmental Epidemiology*, 30(3):459–468, 2020.
- [14] Songzhu Mei, Hongtao Guan, and Qinglin Wang. An overview on the convergence of high performance computing and big data processing. In *2018 IEEE 24th International Conference on Parallel and Distributed Systems (ICPADS)*, pages 1046–1051. IEEE, 2018.

- [15] Petra Klepac, Adam J Kucharski, Andrew JK Conlan, Stephen Kissler, Maria Tang, Hannah Fry, and Julia R Gog. Contacts in context: large-scale setting-specific social mixing matrices from the bbc pandemic project. *medRxiv*, 2020.
- [16] Office for National Statistics. QS103EW (age by single year). <https://www.nomisweb.co.uk/census/2011/qs103ew>.
- [17] Office for National Statistics. Sex by age. <https://www.nomisweb.co.uk/census/2011/lc1117ew>.
- [18] Office for National Statistics. DC2101EW (ethnic group by sex by age). <https://www.nomisweb.co.uk/census/2011/dc2101ew>.
- [19] Ministry of Housing, Communities & Local Government. English indices of deprivation 2019. <https://www.gov.uk/government/statistics/english-indices-of-deprivation-2019>.
- [20] Office for National Statistics. KS105EW (household composition). <https://www.nomisweb.co.uk/census/2011/ks105ew>.
- [21] Office for National Statistics. KS405UK (communal establishment residents). <https://www.nomisweb.co.uk/census/2011/ks405uk>.
- [22] Office for National Statistics. DC1104EW (residence type by sex by age). <https://www.nomisweb.co.uk/census/2011/dc1104ew>.
- [23] Education and Skills Funding Agency. UK register of learning providers. <https://www.ukrlp.co.uk/>.
- [24] Department for Transport. National travel survey 2014: Travel to school. [https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment\\_data/file/476635/travel-to-school.pdf](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/476635/travel-to-school.pdf).
- [25] Department for Education. Class size and education in england evidence report. [https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment\\_data/file/183364/DFE-RR169.pdf](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/183364/DFE-RR169.pdf).
- [26] Office for National Statistics. Uk sic 2007. <https://www.ons.gov.uk/methodology/classificationsandstandards/ukstandardindustrialclassificationofeconomicactivities/uksic2007,2007>.
- [27] Office for National Statistics. UK Business Counts - enterprises by industry and employment size band. <https://www.nomisweb.co.uk/datasets/idbrent>, 2011.
- [28] Office for National Statistics. KS605EW-KS607EW (industry by sex). <https://www.nomisweb.co.uk/census/2011/ks605ew>.
- [29] Office for National Statistics. WU01EW (location of usual residence and place of work by sex). <https://www.nomisweb.co.uk/census/2011/wu01ew>.
- [30] Petra Klepac, Stephen Kissler, and Julia Gog. Contagion! the bbc four pandemic – the model behind the documentary. *Epidemics*, 24:49 – 59, 2018.
- [31] J.I. Gershuny and O. Sullivan. United Kingdom time use survey, 2014-2015. Technical Report SN: 8128, UK Data Service, 2015.
- [32] Age UK. Briefing: Health and care of older people in england 2019. [https://www.ageuk.org.uk/globalassets/age-uk/documents/reports-and-publications/reports-and-briefings/health--wellbeing/age\\_uk\\_briefing\\_state\\_of\\_health\\_and\\_care\\_of\\_older\\_people\\_july2019.pdf](https://www.ageuk.org.uk/globalassets/age-uk/documents/reports-and-publications/reports-and-briefings/health--wellbeing/age_uk_briefing_state_of_health_and_care_of_older_people_july2019.pdf). Accessed: 14 December 2020.
- [33] UK Department for Transport. RAI0201 (city centre peak and all day arrivals and departures by rail on a typical autumn weekday, by city). <https://www.gov.uk/government/statistical-data-sets/rai02-capacity-and-overcrowding>, 2011.

- [34] Office for National Statistics. QS701EW (method of travel to work). <https://www.nomisweb.co.uk/census/2011/qs701ew>.
- [35] Office for National Statistics. Output Area (2011) to Major Towns and Cities (December 2015) Lookup in England and Wales. <https://geoportal.statistics.gov.uk/datasets/78ff27e752e44c3194617017f3f15929,2015>.
- [36] National Rail. Maps of the national rail network of great britain. [https://www.nationalrail.co.uk/stations\\_destinations/maps.aspx](https://www.nationalrail.co.uk/stations_destinations/maps.aspx), 2015.
- [37] Yuanyuan Dong, Xi Mo, Yabin Hu, Xin Qi, Fan Jiang, Zhongyi Jiang, and Shilu Tong. Epidemiology of covid-19 among children in china. *Pediatrics*, 145(6), 2020.
- [38] Ping-Ing Lee, Ya-Li Hu, Po-Yen Chen, Yhu-Chering Huang, and Po-Ren Hsueh. Are children less susceptible to covid-19? *Journal of Microbiology, Immunology, and Infection*, 2020.
- [39] Xi He, Eric Lau, Peng Wu, Xilong Deng, Jian Wang, Xinxin Hao, Yiu Lau, Jessica Y Wong, Yujuan Guan, Xinghua Tan, Xiaoneng Mo, Yanqing Chen, Baolin Liao, Weilie Chen, Fengyu Hu, Qing Zhang, Mingqiu Zhong, Yanrong Wu, Lingzhai Zhao, and Gabriel Leung. Temporal dynamics in viral shedding and transmissibility of covid-19. *Nature Medicine*, 26, 05 2020.
- [40] OpenABM. Openabm-covid19: Agent-based model for modelling the covid-19 epidemic. [https://github.com/roberthinch/OpenABM-Covid19/blob/master/documentation/parameters/infection\\_parameters.md](https://github.com/roberthinch/OpenABM-Covid19/blob/master/documentation/parameters/infection_parameters.md).
- [41] Helen Ward, Christina J Atchison, Matthew Whitaker, Kylie E. C. Ainslie, Joshua Elliott, Lucy C Okell, Rozlyn Redd, Deborah Ashby, Christl A. Donnelly, Wendy Barclay, Ara Darzi, Graham Cooke, Steven Riley, and Paul Elliott. Antibody prevalence for sars-cov-2 in england following first peak of the pandemic: React2 study in 100,000 adults. *medRxiv*, 2020.
- [42] Office for National Statistics. Deaths registered weekly in england and wales, provisional. <https://www.ons.gov.uk/peoplepopulationandcommunity/birthsdeathsandmarriages/deaths/datasets/weeklyprovisionalfiguresondeathsregisteredinenglandandwales>.
- [43] Department of Health & Social Care. Vivaldi 1: Covid-19 care homes study report. <https://www.gov.uk/government/publications/vivaldi-1-coronavirus-covid-19-care-homes-study-report/vivaldi-1-covid-19-care-homes-study-report#fn:1>.
- [44] Office for National Statistics. Population estimates for the uk, england and wales, scotland and northern ireland: mid-2019. <https://www.ons.gov.uk/peoplepopulationandcommunity/populationandmigration/populationestimates/bulletins/annualmidyearpopulationestimates/latest>.
- [45] Office for National Statistics. Lc1105ew residence type by sex by age. <https://www.nomisweb.co.uk/census/2011/lc1105ew>.
- [46] Office for National Statistics. Deaths involving covid-19 in the care sector, england and wales: deaths occurring up to 12 june 2020 and registered up to 20 june 2020 (provisional). <https://www.ons.gov.uk/peoplepopulationandcommunity/birthsdeathsandmarriages/deaths/articles/deathsinvolvingcovid19inthecaresectorenglandandwales/deathsoccurringupto12june2020andregisteredupto20june2020provisional>.
- [47] Scientific Advisory Group for Emergencies. Dynamic co-cin report to sage and nervtag - 30 june 2020. [https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment\\_data/file/903395/S0612\\_Dynamic\\_CO-CIN\\_report\\_to\\_SAGE\\_and\\_NERVTAG.pdf](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/903395/S0612_Dynamic_CO-CIN_report_to_SAGE_and_NERVTAG.pdf).

- [48] Public Health England (PHE). Covid-19 hospitalisation in england surveillance system (chess) – daily reporting. <https://www.england.nhs.uk/coronavirus/wp-content/uploads/sites/52/2020/03/phe-letter-to-trusts-re-daily-covid-19-hospital-surveillance-11-march-2020.pdf>. Accessed: 14 December 2020.
- [49] NHS. Covid-19 situation reports. <https://digital.nhs.uk/about-nhs-digital/corporate-information-and-documents/directions-and-data-provision-notice/data-provision-notice-dpns/covid-19-situation-reports>. Accessed: 14 December 2020.
- [50] S Jenks et al. NF Brazeau, R Verity. Covid-19 infection fatality ratio: Estimates from seroprevalence, 2020.
- [51] Marina Pollán, Beatriz Pérez-Gómez, Roberto Pastor-Barriuso, Jesús Oteo, Miguel A Hernán, Mayte Pérez-Olmeda, Jose L Sanmartín, Aurora Fernández-García, Israel Cruz, Nerea Fernández de Larrea, Marta Molina, Francisco Rodríguez-Cabrera, Mariano Martín, Paloma Merino-Amador, Jose León Paniagua, Juan F Muñoz-Montalvo, Faustino Blanco, Raquel Yotti, Faustino Blanco, Rodrigo Gutiérrez Fernández, Mariano Martín, Saturnino Mezcua Navarro, Marta Molina, Juan F. Muñoz-Montalvo, Matías Salinero Hernández, Jose L. Sanmartín, Manuel Cuenca-Estrella, Raquel Yotti, José León Paniagua, Nerea Fernández de Larrea, Pablo Fernández-Navarro, Roberto Pastor-Barriuso, Beatriz Pérez-Gómez, Marina Pollán, Ana Avellón, Giovanni Fedele, Aurora Fernández-García, Jesús Oteo Iglesias, María Teresa Pérez Olmeda, Israel Cruz, Maria Elena Fernandez Martinez, Francisco D. Rodríguez-Cabrera, Miguel A. Hernán, Susana Padrones Fernández, José Manuel Rumbao Aguirre, José M. Navarro Marí, Begoña Palop Borrás, Ana Belén Pérez Jiménez, Manuel Rodríguez-Iglesias, Ana María Calvo Gascón, María Luz Lou Alcaine, Ignacio Donate Suárez, Oscar Suárez Álvarez, Mercedes Rodríguez Pérez, Margarita Cases Sanchís, Carlos Javier Villafáfila Gomila, Lluís Carbo Saladrigas, Adoración Hurtado Fernández, Antonio Oliver, Elías Castro Feliciano, María Noemí González Quintana, José María Barrasa Fernández, María Araceli Hernández Betancor, Melisa Hernández Febles, Leopoldo Martín Martín, Luis-Mariano López López, Teresa Ugarte Miota, Inés De Benito Población, María Sagrario Celada Pérez, María Natalia Vallés Fernández, Tomás Maté Enríquez, Miguel Villa Arranz, Marta Domínguez-Gil González, Isabel Fernández-Natal, Gregoria Megías Lobón, Juan Luis Muñoz Bellido, Pilar Ciruela, Ariadna Mas i Casals, Maria Doladé Botías, M. Angeles Marcos Maeso, Dúnia Pérez del Campo, Antonio Félix de Castro, Ramón Limón Ramírez, Maria Francisca Elías Retamosa, Manuela Rubio González, María Sinda Blanco Lobeiras, Alberto Fuentes Losada, Antonio Aguilera, German Bou, Yolanda Caro, Noemí Marauri, Luis Miguel Soria Blanco, Isabel del Cura González, Montserrat Hernández Pascual, Roberto Alonso Fernández, Paloma Merino-Amador, Natalia Cabrera Castro, Aurora Tomás Lizcano, Cristóbal Ramírez Almagro, Manuel Segovia Hernández, Nieves Ascunce Elizaga, María Ederra Sanz, Carmen Ezpeleta Baquedano, Ana Bustinduy Bascaran, Susana Iglesias Tamayo, Luis Elorduy Otazua, Rebeca Benarroch Benarroch, Jesús Lopera Flores, and Antonia Vázquez de la Villa. Prevalence of SARS-CoV-2 in Spain (ENE-COVID): a nationwide, population-based seroepidemiological study. *The Lancet*, 396(10250):535–544, August 2020. Publisher: Elsevier.
- [52] Flavia Riccardo, Marco Ajelli, Xanthi D Andrianou, Antonino Bella, Martina Del Manso, Massimo Fabiani, Stefania Bellino, Stefano Boros, Alberto Mateo Urdiales, Valentina Marziano, Maria Cristina Rota, Antonietta Filia, Fortunato Paolo D’Ancona, Andrea Siddu, Ornella Punzo, Filippo Trentini, Giorgio Guzzetta, Piero Poletti, Paola Stefanelli, Maria Rita Castrucci, Alessandra Ciervo, Corrado Di Benedetto, Marco Tallon, Andrea Piccioli, Silvio Brusaferrero, Giovanni Rezza, Stefano Merler, and Patrizio Pezzotti. Epidemiological characteristics of covid-19 cases in italy and estimates of the reproductive numbers one month into the epidemic. *medRxiv*, 2020.
- [53] Office for National Statistics. Deaths involving covid-19 in the care sector, england and wales. <https://www.ons.gov.uk/peoplepopulationandcommunity/birthsdeathsandmarriages/deaths/articles/deathsinvolvingcovid19inthecaresectorenglandandwales/>



deathsoccurringupto12june2020andregisteredupto20june2020provisional#  
characteristics-of-care-home-residents-who-died-from-covid-19.

- [54] Derek K Chu, Elie A Akl, Stephanie Duda, Karla Solo, Sally Yaacoub, and Holger J Schünemann. Physical distancing, face masks, and eye protection to prevent person-to-person transmission of sars-cov-2 and covid-19: a systematic review and meta-analysis. *Lancet*, (395):1973–1987, 2020.
- [55] Emma P Fischer, Martin C Fischer, David Grass, Isaac Henrion, Warren S Warren, and Eric Westman. Low-cost measurement of face mask efficacy for filtering expelled droplets during speech. *Science Advances*, 6(36):eabd3083, 2020.
- [56] Jeremy Howard, Austin Huang, Zhiyuan Li, Zeynep Tufekci, Vladimir Zdimal, Helene-Mari van der Westhuizen, Arne von Delft, Amy Price, Lex Fridman, Lei-Han Tang, et al. Face masks against COVID-19: an evidence review. 2020.
- [57] C.Y. Suen, H.H. Leung, K.W. Lam, et al. Feasibility of reusing surgical mask under different disinfection treatments. *medRxiv*, 2020.
- [58] E.C. Toomey, Y. Conway, C. Burton, S. Smith, M. Smalle, X.S. Chan, A. Adishes, S. Tanveer, L. Ross, I. Thomson, D. Devane, and T. Greenhalgh. Extended use or reuse of single-use surgical masks and filtering face-piece respirators during the coronavirus disease 2019 (covid-19) pandemic: A rapid systematic review. *Infect Control Hosp Epidemiol*, 8:1–9, 2020.
- [59] Mingming Liang, Liang Gao, Ce Cheng, Qin Zhou, John Patrick Uy, Kurt Heiner, and Chenyu Sun. Efficacy of face mask in preventing respiratory virus transmission: A systematic review and meta-analysis. *Travel medicine and infectious disease*, 36:101751–101751, 2020.
- [60] YouGov. YouGov COVID-19 behaviour changes tracker: Wearing a face mask when in public places. <https://yougov.co.uk/topics/international/articles-reports/2020/03/17/personal-measures-taken-avoid-covid-19>, 2020.
- [61] Office for National Statistics. Coronavirus and key workers in the uk. <https://www.ons.gov.uk/employmentandlabourmarket/peopleinwork/earningsandworkinghours/articles/coronavirusandkeyworkersintheuk/2020-05-15>, 2020.
- [62] Institute for Fiscal Studies. Sector shutdowns during the coronavirus crisis: which workers are most exposed? <https://www.ifs.org.uk/publications/14791>, 2020.
- [63] HM Revenue & Customs. HMRC coronavirus (COVID-19) statistics. <https://www.gov.uk/government/collections/hmrc-coronavirus-covid-19-statistics#Coronavirus-Job-Retention-Scheme-Management-information>, 2020.
- [64] YouGov. YouGov COVID-19 behaviour changes tracker: Avoiding going to work. <https://yougov.co.uk/topics/international/articles-reports/2020/03/17/personal-measures-taken-avoid-covid-19>, 2020.
- [65] Office for National Statistics. Coronavirus and the latest indicators for the uk economy and society: 29 october 2020. <https://www.ons.gov.uk/peoplepopulationandcommunity/healthandsocialcare/conditionsanddiseases/bulletins/coronavirustheukeconomyandsocietyfasterindicators/29october2020>, 2020.
- [66] UK Department for Education. Attendance in education and early years settings during the coronavirus (COVID-19) outbreak. <https://www.gov.uk/government/collections/attendance-in-education-and-early-years-settings-during-the-coronavirus-covid-19-outbreak>, 2020.
- [67] OpenTable. The state of the restaurant industry. <https://www.opentable.com/state-of-industry>, 2020.
- [68] Government of the United Kingdom. Get a discount with the eat out to help out scheme. <https://www.gov.uk/guidance/get-a-discount-with-the-eat-out-to-help-out-scheme>, 2020.

- [69] Office for National Statistics. Coronavirus and the social impacts on great britain: 25 september 2020. <https://www.ons.gov.uk/peoplepopulationandcommunity/healthandsocialcare/healthandwellbeing/bulletins/coronavirusandthesocialimpactsongreatbritain/25september2020>, 2020.
- [70] Covid-19 patient notification system (cpns) user guide. <https://www.england.nhs.uk/statistics/wp-content/uploads/sites/2/2020/09/CPNS-User-Guide-20200831.pdf>. Accessed: 14 December 2020.
- [71] Joseph Bullock, Carolina Cuesta-Lazaro, Arnau Quera-bofarull, Miguellcaza-Lizaola, Aidan Sedgewick, Henry Truong, Tristan Caulfield, Kevin Fong, Ian Vernon, Julian Williams, Richard Bower, and Frank Krauss. June: a bayesian uncertainty analysis. (*in preparation*), 2021.
- [72] P S Craig, M Goldstein, A H Seheult, and J A Smith. Pressure matching for hydrocarbon reservoirs: a case study in the use of bayes linear strategies for large computer experiments (with discussion). In C Gatsonis, J S Hodges, R E Kass, R McCulloch, P Rossi, and N D Singpurwalla, editors, *Case Studies in Bayesian Statistics*, volume 3, pages 36–93. New York, 1997.
- [73] I Vernon, M Goldstein, and R G. Bower. Galaxy formation: a bayesian uncertainty analysis. *Bayesian Analysis*, 5(4):619–670, 2010.
- [74] I. Andrianakis, I. Vernon, N. McCreesh, T.J. McKinley, J.E. Oakley, R. Nsubuga, M. Goldstein, and R.G. White. Bayesian history matching of complex infectious disease models using emulation: A tutorial and a case study on HIV in uganda. *PLoS Comput Biol.*, 11(1):e1003968, 2015.
- [75] I. Andrianakis, I. Vernon, N. McCreesh, T. J. McKinley, J. E. Oakley, R. N. Nsubuga, M. Goldstein, and R. G. White. History matching of a complex epidemiological model of human immunodeficiency virus transmission by using variance emulation. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 66(4):717–740, 2017.
- [76] I. Andrianakis, N. McCreesh, I. Vernon, T.J. McKinley, J.E. Oakley, R. Nsubuga, M. Goldstein, and R.G. White. Efficient history matching of a high dimensional individual based hiv transmission model. *SIAM/ASA Journal of Uncertainty Quantification*, 5(1):694–719, 2017.
- [77] Nicky McCreesh, Ioannis Andrianakis, Rebecca N. Nsubuga, Mark Strong, Ian Vernon, Trevelyan J. McKinley, Jeremy E. Oakley, Michael Goldstein, Richard Hayes, and Richard G. White. Universal test, treat, and keep: improving art retention is key in cost-effective hiv control in uganda. *BMC Infectious Diseases*, 17(1):322, May 2017.
- [78] T.J. McKinley, I. Vernon, I. Andrianakis, N. McCreesh, J.E. Oakley, R.N. Nsubuga, M. Goldstein, and R.G. White. Approximate bayesian computation and simulation-based inference for complex stochastic epidemic models. *Statistical Science*, 33(1):4–18, June 2018.
- [79] Ian Vernon, Junli Liu, Michael Goldstein, James Rowe, Jen Topping, and Keith Lindsey. Bayesian uncertainty analysis for complex systems biology models: emulation, global parameter searches and evaluation of gene functions. *BMC Systems Biology*, 12(1):arXiv:1607.06358 [q-bio.MN], 2018.
- [80] I Vernon, M Goldstein, and R G. Bower. Rejoinder for Galaxy formation: a bayesian uncertainty analysis. *Bayesian Analysis*, 5(4):697–708, 2010.
- [81] A O’Hagan. Bayesian analysis of computer code outputs: A tutorial. *Reliability Engineering and System Safety*, 91:1290–1300, 2006.
- [82] Madhav V Marathe and Naren Ramakrishnan. Recent advances in computational epidemiology. *IEEE intelligent systems*, 28(4):96–101, 2013.
- [83] Joseph Bullock, Carolina Cuesta-Lazaro, Arnau Quera-bofarull, Anjali Katta, Katherine Hoffmann Pham, Benjamin Hoover, Hendrik Strobel, Rebeca Moreno Jimenez, and Miguel Luengo-Oroz. Operational response simulation tool for epidemics within settlements. (*in preparation*), 2020.

- [84] Office for National Statistics. Marriages in England and Wales.  
<https://www.ons.gov.uk/peoplepopulationandcommunity/birthsdeathsandmarriages/marriagecohabitationandcivilpartnerships/datasets/marriagesinenglandandwales2013,2017>.
- [85] Office for National Statistics. Birth characteristics in England and Wales: 2017.  
<https://www.ons.gov.uk/releases/birthcharacteristicsinenglandandwales2017,2017>.
- [86] Office for National Statistics. LC1109EW (household composition by age by sex).  
<https://www.nomisweb.co.uk/census/2011/lc1109ew>.
- [87] Scientific Advisory Group for Emergencies. Dynamic co-cin report to sage and nervtag, 13 may 2020. <https://www.gov.uk/government/publications/dynamic-co-cin-report-to-sage-and-nervtag-13-may-2020>.
- [88] ICNARC. Icnarc report on covid-19 in critical care 10 july 2020.  
<https://www.icnarc.org/Our-Audit/Audits/Cmp/Reports>.
- [89] T S Bastos and A O'Hagan. Diagnostics for gaussian process emulators. 51:425–438, 2008.
- [90] M C Kennedy and A O'Hagan. Bayesian calibration of computer models. 63(3):425–464, 2001.