

1. Title Page

Title	TARGET-EU: Effectiveness of BNT162b2 mRNA COVID-19 vaccine in healthy individuals or with stable pre-existing medical conditions against SARS-CoV-2 infection
Research question & Objectives	<p>What is the effectiveness of two-doses of BNT162b2 mRNA-based vaccine in preventing SARS-CoV-2 infection compared with no vaccination, among healthy adults (≥ 16 years old), including those with stable pre-existing medical conditions?</p> <p>Primary objective: To estimate the incidence of SARS-CoV-2 infection among individuals receiving the BNT162b2 vaccine, compared to unvaccinated individuals.</p>
Protocol version	1.0
Last update date	27 February 2026
Contributors	<p>Primary investigators contact information: Dr. Constanza Andaur Navarro (UMCU) c.l.andaurnavarro@umcutrecht.nl Dr. Oisín Ryan (UMCU) Dr. Sebastián Mildiner Moraga (UMCU)</p> <p>Contributor names: Prof. Ian Douglas (LSHTM) Dr. Emmy Manders (UMCU) Dr. Patrick Sourverein (UU-CPRD) Gabriel Sanfélix Gimeno (FISABIO-VID) Fran Llopis Cardona (FISABIO-VID) Aníbal García Sempere (FISABIO-VID)</p>
Study registration	<p>Site: https://catalogues.ema.europa.eu/node/4440/administrative-details</p> <p>Identifier: EUPAS1000000539</p>
Sponsor	<p>Organization: EU PE&PV research network Contact: eupepv@uu.nl</p>
Conflict of interest	CAN, OR, SMM and EM are currently salaried employees at University Medical Centre Utrecht, which receives institutional research funding from pharmaceutical companies and regulatory agencies and is administered by University Medical Centre Utrecht.

Table of Contents

1. Title Page	1
2. Abstract	3
3. Amendments and updates	4
4. Milestones	4
5. Rationale and background	5
6. Research questions and objectives	7
6.1 Primary Estimand 1	7
6.2 Supplementary Estimand 2.....	11
7. Research methods	14
7.1. Study design.....	14
7.2. Study design diagram.....	14
7.3. Setting	15
7.3.1 Definition of time 0 (and other primary time anchors) for entry to the study population – Context and rationale.....	15
7.3.2 Study inclusion criteria – Context and rationale.....	17
7.3.3 Study exclusion criteria – Context and rationale.....	19
7.4. Variables.....	21
7.4.1 Exposure(s) of interest – Context and rationale	21
7.4.2 Outcome(s) of interest – Context and rationale	25
7.4.3 Follow up – Context and rationale.....	27
7.4.4 Covariates (confounding variables and effect modifiers, e.g. risk factors, comorbidities, comedications) – Context and rational	29
7.5. Core Emulation Table – Design Summary	34
7.6. Data analysis.....	41
7.6.1 Analysis plan – Context and rationale.....	41
7.6.2 Primary Estimand Main Analysis.....	42
7.6.3 Supplementary Estimand (2) Main Analysis: Principal Stratum Effect.....	46
7.6.4 Sensitivity Analyses.....	48
7.6.5. Other Supplemental Analyses.....	51
7.6.6 Core Emulation Table – Estimation Summary	52
7.7. Data sources	57
7.7.1 Data sources – Context and rationale	57
7.8. Data management.....	59
7.9. Quality control	62
7.10. Study size and feasibility.....	68
8. Limitation of the methods	70
9. Protection of human subjects	72
10. Reporting of adverse events	72
11. References	73
12. Appendices	74

2. Abstract

Background: The rapid development of COVID-19 vaccines represents a critical milestone in mitigating the impact of the pandemic. However, ongoing evaluation of their real-world effectiveness is essential given the varied vaccination coverage.

Objectives: To estimate the effectiveness of BNT162b, an mRNA-based COVID-19 vaccine developed by Pfizer and BioNTech, in preventing SARS-CoV-2 infection among healthy individuals aged 16 years or older, including those with stable pre-existing medical conditions, accounting for deviations from standard dosing schedules or vaccination series.

Methods: We will conduct a matched cohort study using routinely collected healthcare data from the CPRD and VID databases from 01 December 2020 to 30 April 2022. Eligible participants will be ≥ 16 years old with at least 6 months of continuous database registration prior to vaccination or matched index date. Individuals with prior receipt of non-BNT162b2 vaccines, prior prophylactic COVID-19 treatments, or evidence of immunocompromised status will be excluded.

Each vaccinated individual will be matched 1:1 with an unvaccinated comparator on age, sex, geographic location, calendar time, and key risk factors. Time zero will be defined as the date of first BNT162b2 dose (or corresponding matched date for unvaccinated). Follow-up will start at time zero and continued for up to 90 days, until the earliest of SARS-CoV-2 infection, death, de-registration, or end of data availability.

Vaccine exposure and outcomes will be ascertained through linked pharmacy records, immunization registers, general practice records, and medical encounters. The primary outcome is laboratory-confirmed SARS-CoV-2 infection, identified using diagnostic and test codes.

The study will follow the target trial emulation framework in combination with the ICH E9(R1) estimand framework. A critical aspect is to identify and address intercurrent events (IEs). The IEs of interest include not receiving or being ineligible for a second vaccine dose; receipt of a BNT162b2 booster within three months of the second dose; receipt of a booster with a non-BNT162b2 vaccine; receipt of a non-COVID-19 vaccine after treatment completion; pre-exposure prophylactic COVID-19 treatment; and death. Different strategies will be applied to handle IEs. In the primary estimand, for missing a second dose, early receipt of a BNT162b2 booster and for receipt of a non-COVID-19 vaccine, a treatment policy strategy will be applied. For receipt of a non-BNT162b2 booster or other preventive COVID-19 treatment, a hypothetical strategy will be applied. Death will be handled using a composite strategy. In a supplementary estimand, a principal stratum strategy will be used to handle missing a second dose, and a while-alive strategy will be used to handle death, instead.

3. Amendments and updates

Version date	Version number	Section of protocol	Amendment or update	Reason
27 February 2026	1.0			

4. Milestones

Table 1. Milestones

Milestone	Date
Study protocol for RWD study	08 August 2025
Preliminary results RWD study	April 2026
Final Study report	10 June 2026

5. Rationale and background

What is known about the condition?

COVID-19 illness, caused by the severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), emerged in late 2019 in Wuhan, China, and quickly grew into a global pandemic, arriving in Europe in February 2020. SARS-CoV-2 virus primarily affects the respiratory system, with common symptoms including fever or shivering, cough, shortness of breath, sore throat, congestion, fatigue, muscle aches, headache, and loss or change in the sense of taste or smell. Symptoms usually appear 2 to 14 days after exposure to the virus, and their severity can range from mild to life-threatening. While many people recover within a few weeks, some develop severe illness requiring hospitalization, especially those with risk factors such as older age, obesity, or chronic diseases. COVID-19 spreads through respiratory droplets and close contact with infected individuals, making it highly contagious. In the early pandemic control efforts, preventing measures such as the use of masks, testing and vaccination became crucial in managing the early waves of the disease.

The first vaccines against COVID-19 were authorized for emergency use in late 2020. The Pfizer-BioNTech (BNT162b2) was among the earliest mRNA vaccines, utilizing novel technology to instruct cells to produce the SARS-CoV-2 spike protein, triggering an immune response. The vaccine rollout began with healthcare workers, elderly populations, and high-risk groups. To note, much of the early data was derived from populations and pandemic phases that may differ significantly from current conditions. Changes in circulating variants (e.g., Delta, Omicron), waning immunity over time, and differential uptake of booster doses have all introduced new complexities. These factors highlight the continued need for updated real-world assessments of vaccine effectiveness under evolving epidemiological and regulatory circumstances.

What is known about the exposure of interest?

The exposure of interest is BNT162b2, sold under the brand name Comirnaty.¹ It encodes the full-length SARS-CoV-2 spike protein stabilized in its prefusion conformation, delivered via lipid nanoparticles. Extensive clinical trials conducted the early pandemic period comparing individuals who are vaccinated to unvaccinated individuals demonstrated the efficacy of mRNA-based COVID-19 vaccine in preventing SARS-CoV-2 infection, hospital admission, and death from SARS-CoV-2.² These findings were further supported by observational studies confirming the vaccines' effectiveness in real world settings.³⁻⁵ Initial guidance recommended a two-dose primary series, with two 0.3 mL intramuscular injections administered 21 days apart.

Gaps in knowledge

Little or no evidence exists regarding the effectiveness of incomplete vaccination series, series mixed with other potentially available vaccines (e.g., heterologous prime-boost schedules) or altered dosing schedules (i.e., extended or shortened intervals) that deviate from the original clinical trial protocols. These variations often arise in real-world settings due to vaccine shortages, personal choice, or changing public health guidance. Such deviations may introduce uncertainty about the level and duration of protection conferred.

What is the expected contribution of this study?

This study aims to estimate the effectiveness of vaccination with BNT162b2 on preventing SARS-CoV-2 infection using target trial emulation (TTE) and the ICH E9(R1) estimand frameworks. TTE helps mitigate more structural biases that commonly arise in observational vaccine studies, particularly immortal time bias and selection (collider stratification) bias. By explicitly framing the

observational analysis as a trial analogue, we aim to reduce the risk of bias and strengthen causal inference and generate more reliable evidence on vaccine effectiveness under varied real-world conditions.

Specific research question:

Among healthy individuals aged 16 years or older, including those with stable pre-existing medical conditions, does receiving BNT162b2 vaccination – compared with no COVID-19 vaccination – reduce the risk of SARS-CoV-2 infection?

6. Research questions and objectives

The overall aim is to assess whether BNT162b2 mRNA COVID-19 vaccine reduce the risk of SARS-CoV-2 infection compared to no COVID-19 vaccination in healthy individuals or with pre-existing stable medical condition

6.1 Primary Estimand 1

Primary research question targeted by Estimand 1: What is the effectiveness (1-RR) of BNT162b2 mRNA COVID-19 vaccine (30 µg per dose) compared to no vaccine to prevent SARS-CoV-2 infection or death in healthy subjects aged 16 years or older regardless of incomplete dosing, a third booster, any non-COVID-19 vaccine, and as if any (non-target) COVID-19 vaccination or receipt of any other preventative COVID-19 treatment would not occur?

Table 2A. Primary estimand (estimand 1)

Attribute	Target Trial	Target Trial Emulation	Comment
Population	Individuals aged 16 years or older, healthy or with stable preexisting medical conditions.	Individuals aged 16 years or older, healthy or with stable preexisting medical conditions.	Identical in Target Trial and Emulation.
Treatment Conditions	Intervention: Two doses BNT162b2 mRNA COVID-19 vaccine separated by 21 days Control: placebo, with two matching doses.	Intervention: Two doses BNT162b2 mRNA COVID-19 vaccine separated by 21 days (with an allowable interval between 19 to 42 days) Control: No vaccine	Placebo interventions are not available in real world data sources. Instead, we use unvaccinated as the control condition.
Endpoint	Laboratory-confirmed SARS-CoV-2 infection or death from any cause.	Laboratory-confirmed SARS-CoV-2 infection or death from any cause.	Identical in Target Trial and Emulation. Emulation has access to population-scale health databases with information from PCR or antigen test, and mortality data.
Summary Measure	Vaccine Efficacy=1-RR where RR denotes rate ratio of disease or death in vaccinees compared to placebo recipients.	Vaccine Effectiveness =1-RR where RR denotes rate ratio of disease or death in vaccinees compared to control individuals.	The emulation will be considered an effectiveness study given the differences compared to the RCT setting to establish efficacy in terms of treatment conditions such as the use of an unvaccinated control in place of a placebo control, less control over and accurate measurements of

			participant behaviour as it relates to, e.g., timing of the second dose receipt in the intervention condition, testing for covid-19 when experiencing symptoms, etc.
Intercurrent events (IE) and strategies to handle them	<p>(1) missing or ineligible for second dose of target vaccine; (2) a third (booster) dose of target vaccine up to 3 months after second dose; (3) third (booster) dose of non-target COVID-19 vaccine; (4) receipt of any non-COVID-19 vaccine dose following treatment completion; (5) receipt of any other preventative COVID-19 treatment; (6) death</p> <p>Treatment policy (1), (2), (4), Hypothetical (3), (5) Composite (6)</p>	<p>(1) missing second dose of target vaccine (Exposed) ; (2) a third or booster dose of target vaccine up to 3 months after second dose (Exposed)(3) receipt of any non-target COVID-19 vaccine at any time (Exposed) or first dose of any COVID-19 vaccine at any time (Control); (4) receipt of any non-COVID-19 vaccine dose; (5) receipt of any other preventative COVID-19 treatment; (6) death</p> <p>Treatment policy (1), (2), (4), Hypothetical (3), (5) Composite (6)</p>	<p>The IEs have been altered to reflect the change from placebo to no vaccination in the control condition. For instance, in the target trial, individuals may miss or become ineligible to receive the second dose of placebo (IE1), following the same conditions as eligibility to receive the second dose of BNT162b2. This is not sensible to emulate in the no vaccine condition, and so the IE1 applies only to exposed. Furthermore, unlike in the trial setting it may be more difficult to assess when participants become ineligible to receive the second dose rather than simply failing to receive the second dose; however, since this IE is handled with a treatment policy strategy this distinction is also not strictly necessary. Finally, we change the period of intended second dose receipt from exactly 21 days in the trial to any time in the three months following receipt of the first dose. This reflects real world conditions, in which appointments may not be made for exactly 21 days after initial vaccination, or appointments may be missed and rescheduled for various reasons. The timeframe for IE4 has also been altered to reflect that this IE may occur at any time since enrolment and it would be handled with the same strategy (treatment policy), regardless of when it occurred. In practice this is</p>

			<p>not a deviation from the target trial but rather an omission at the time of developing the protocol of the target trial.</p> <p>Codes will be used to identify occurrences of IE5.</p> <p>For IE3 we will identify receipt of competitor vaccines (such as mRNA-1273[Spikevax] or AZD1222[Vaxzevria]) using local codes to retrieve manufacturer. This information will suffice to identify occurrence of IE3 in cases where the vaccine is obtained from sources available to the databases used (such as hospitals/GPs/public health services). We may fail to ascertain the receipt of other vaccines when this is done by participants travelling to a neighbouring country or by acquiring the vaccine privately, and this is not registered with the home data source.</p>
--	--	--	---

Rationale for selected strategies to handle intercurrent events are chosen

Treatment Policy

Missing second dose of target COVID-19 vaccine [IE1]

We handle this intercurrent event (IE) using the treatment policy strategy to allow for variations in treatment adherence commonly encountered in real-world vaccination practices. Although the intended treatment is the receipt of two doses separated by 21 days, in real life scenarios, patients may miss or become ineligible to receive the second dose for a number of reasons, such as missing the scheduled appointment for the second dose, receiving a dose of a COVID-19 vaccine earlier or later than planned beginning pre-exposure prophylaxis (e.g., monoclonal antibodies, antivirals) before the scheduled appointment for the second dose (likely postponing receipt of the second dose), or experiencing a SARS-CoV-2 infection after the first dose. In this case, some individuals may postpone the second dose or miss it entirely.⁶ With this strategy we can estimate the effectiveness of the vaccine regardless of not receiving the second dose, or receiving the second dose at another time than planned (any time between receipt of the first dose and end of the three-month follow-up), which is of regulatory interest.

Third (booster) dose of target COVID-19 vaccine [IE2] (same in Estimand 2)

This IE represents commonly encountered variations in real-world vaccination practices. Including them in the analysis reflects the overall efficacy of the vaccine as it would be used in routine care, enhancing the generalizability of the findings.

Hypothetical strategy

Receipt of any non-target COVID-19 vaccine (Exposed) or first dose of any COVID-19 vaccine at any time (Control) [IE3] & receipt of any other preventative COVID-19 treatment [IE5] (same in Estimand 2)

This IE could confound the direct evaluation of the target COVID-19 vaccine's efficacy. By estimating treatment effects as if these events had not occurred, this approach isolates the intrinsic protective effect of the BNT162b2 vaccine in a hypothetical scenario without interference from external interventions. This also applies for IE receipt of any preventative COVID-19 treatment (pre-exposure prophylaxis).

Composite

Death [IE6]

Death from any cause was included as part of a composite endpoint with SARS-CoV-2 infection to capture severe clinical outcomes comprehensively, especially in the context of COVID-19 where death due to COVID-19 might not be captured properly.

6.2 Supplementary Estimand 2

Research question targeted by Estimand 2: What is the vaccine effectiveness (1-RR) of two doses of BNT162b2 mRNA COVID-19 vaccine (30 µg per dose) compared to no vaccine to prevent SARS-CoV-2 infection in healthy subjects aged 16 years or older, amongst those individuals who would receive both doses of BNT162b2 when assigned to the treated condition, regardless of receipt of a third booster, or any non-COVID-19 vaccine, while alive, and as if any (non-target) COVID-19 vaccination or receipt of any other preventative COVID-19 treatment would not occur?

Table 2B. Supplementary estimand (estimand 2)

Attribute	Target Trial	Target Trial Emulation	Comment
Population	Individuals aged 16 years or older, healthy or with stable preexisting medical conditions who would tolerate and receive both scheduled doses under the intervention treatment and the control conditions.	Individuals aged 16 years or older, health or with stable preexisting medical conditions who would tolerate and receive both scheduled doses under the intervention treatment.	In comparison to the target trial, a no-vaccine control is used instead of a placebo control. As a consequence, the emulation targets the principal stratum of those who would tolerate both scheduled doses under the intervention treatment, rather than under intervention and placebo. Emulation has access to population scale health databases in order to ascertain eligibility.
Treatment Conditions	Intervention: Two doses BNT162b2 mRNA COVID-19 vaccine separated by 21 days. Control: placebo, with two matched doses	Intervention: Two doses BNT162b2 mRNA COVID-19 vaccine separated by 21 days. Control: No vaccine	Placebo interventions are not available in real world data sources. Instead, we use unvaccinated as the control condition.
Endpoint	Laboratory-confirmed SARS-CoV-2 infection.	Laboratory-confirmed SARS-CoV-2 infection	Identical in Target Trial and Emulation. Emulation has access to population-scale health databases with information from PCR or antigen test.
Summary Measure	Vaccine Efficacy = 1-RR where RR denotes risk ratio of disease in vaccinees compared to placebo recipients.	Vaccine Effectiveness = 1-RR where RR denotes risk ratio of disease in vaccinees compared to control individuals.	The emulation will be considered an effectiveness study given the differences compared to the RCT settings in treatment conditions such

			as the use of an unvaccinated control in place of a placebo control, less control over and accurate measurements of participant behaviour as it relates to, e.g., timing of the second dose receipt in the intervention condition, testing for covid-19 when experiencing symptoms, etc.
Intercurrent events and strategies to handle them	<p>(1) missing or ineligible for second dose of target vaccine; (2) a third (booster) dose of target vaccine up to 3 months after second dose; (3) third (booster) dose of non-target COVID-19 vaccine; (4) receipt of any non-COVID-19 vaccine dose following treatment completion; (5) receipt of any other preventative COVID-19 treatment following treatment completion; (6) death</p> <p>Treatment Policy (2), (4)</p> <p>Hypothetical (3), (5)</p> <p>Principal Stratum (1)</p> <p>While-alive (6)</p>	<p>(1) missing second dose of target vaccine (Exposed); (2) a third or booster dose of target vaccine up to 3 months after second dose (Exposed); (3) receipt of any non-target COVID-19 vaccine at any time (Exposed) or first dose of any COVID-19 vaccine at any time (Control); (4) receipt of any non-COVID-19 vaccine dose; (5) receipt of any other preventative COVID-19 treatment (defined as in IE4); (6) death</p> <p>Treatment Policy (2), (4)</p> <p>Hypothetical (3), (5)</p> <p>Principal Stratum (1)</p> <p>While-alive (6)</p>	<p>Note that the same changes have been made to the IE definitions as in Table 2A above. Note here that while IE1 is treated with the principal stratum strategy, the actual principal stratum targeted in the emulation differs from that of the target trial; we now target the stratum of individuals who would receive the second dose in the treatment condition, in contrast to the target trial, which targets the stratum of individuals who would receive the second dose in both treatment and control conditions</p>

Rationale for selected strategies to handle intercurrent events are chosen

When compared with estimand 1 above, different strategies to handle IE1 and IE6 are proposed: Principal stratum and While-alive. The other IEs are handled with the same strategies as estimand 1, and therefore not described below.

Principal Stratum

Missing Second Dose of target COVID-19 vaccine [IE1]

We aim to estimate the effectiveness of the treatment under full adherence. For this reason, we handle the IE about missing the second dose using the principal stratum strategy. As described in Tables 2A and 2B, the definition of this intercurrent event differs from the corresponding IE in the target trial (in which the control condition is placebo rather than no vaccination). In the emulation, this IE applies only to the exposed group and represents failure to receive the second dose of the vaccine any time following receipt of the first dose and before the end of the 3-month follow-up period. Our aim is to estimate the effect of treatment in the stratum of the population who would receive both vaccine doses within a three-month period if they were to undergo vaccination.

While Alive

Death [IE6]

Death was addressed using a while-alive strategy, as post-mortem outcomes are undefined and the estimand pertains to outcomes while participants are living.

7. Research methods

7.1. Study design

Research design (e.g. cohort, case-control, etc.):

This study will employ a matched cohort design to assess the real-world effectiveness of the BNT162b2 mRNA COVID-19 vaccine against SARS-CoV-2 infection involving two databases, that is CPRD and VID. We will estimate the incidence of SARS-CoV-2 infection after receipt of the Pfizer-BioNTech COVID-19 vaccine compared to not receiving any COVID-19 vaccine over 90 days.

Rationale for study design choice:

The matched cohort approach is a robust observational design that enables the estimation of causal effects in the absence of randomization by pairing vaccinated individuals with unvaccinated comparators who share similar baseline characteristics.

Matched on key covariates—such as age, sex, geographic location, prior SARS-CoV-2 infection, pre-existing comorbidities (using the *Centers for Disease Control and Prevention [CDC] risk criteria*), and calendar time—helps to reduce confounding. Furthermore, the use of multiple large, population-based healthcare databases enhances external validity by capturing diverse demographic and clinical profiles, health system structures, and epidemic dynamics.

7.2. Study design diagram

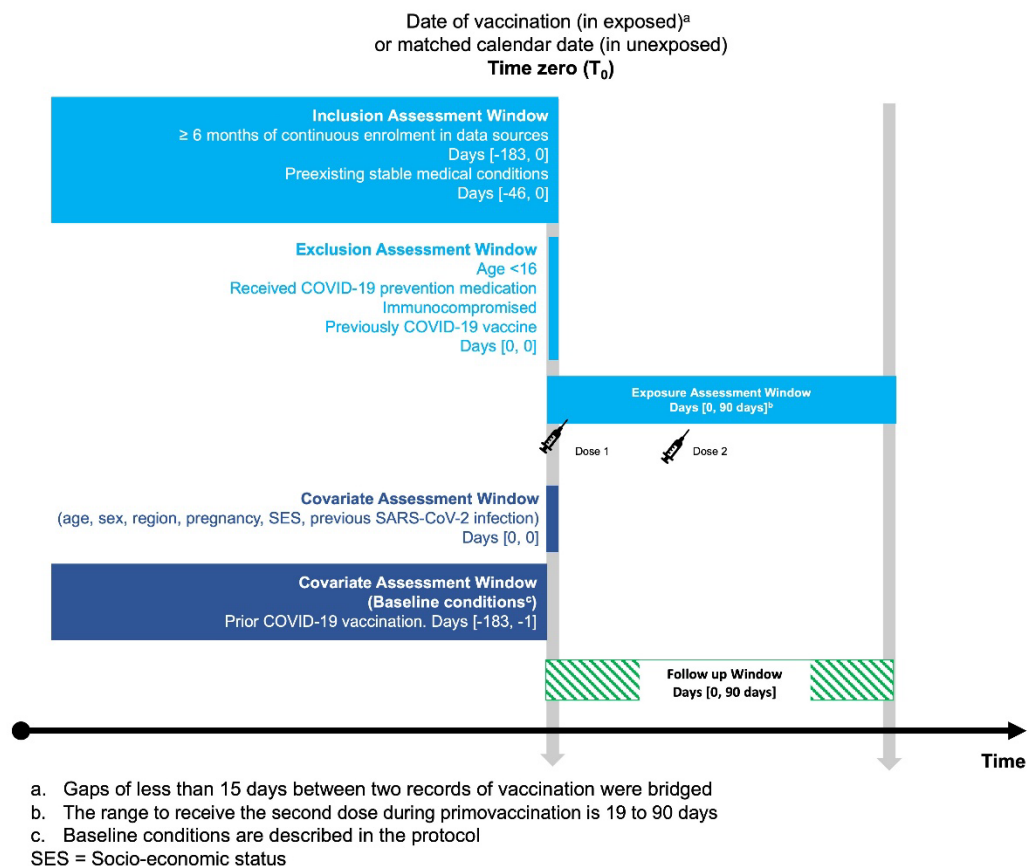


Figure 1. Study design diagram

7.3. Setting

The exposure assessment window is defined from December 1, 2020, to January 31, 2022. Additionally, to ensure sufficient baseline and follow-up data, individuals must have at least 6 months of available data prior to time zero (to capture covariates and prior health status) thus covering the start of the COVID-19 vaccine rollout through to the end of mass testing for both symptomatic and asymptomatic individuals in the UK. Consequently, the study period will be from 01 June 2020 to 30 April 2022 (i.e. source data range), with the latter date allowing for individuals enrolled on 31 January 2022 to enable 90 days of follow-up.⁷

7.3.1 Definition of time 0 (and other primary time anchors) for entry to the study population – Context and rationale

Time Zero (index date, T_0) is defined as the calendar day on which the eligibility criteria are fulfilled, exposure status is assigned, and follow-up begins. Time-zero will be:

- **Exposed group:** Calendar day on which a first dose of BNT162b2 mRNA COVID-19 vaccine is administered.
- **Unexposed group:** Calendar day on which the matched exposed comparator receives their first dose of BNT162b2 mRNA COVID-19 vaccine.

Table 3. Operational Definition of Time 0 (index date) and other primary time anchors

Study population name(s)	Time Anchor Description (e.g. time 0)	Number of entries	Type of entry	Washout window	Care Setting ¹	Code Type ²	Diagnosis position	Incident with respect to...	Measurement characteristics/validation	Source of algorithm
Exposed	Calendar day of vaccine administration	Single entry	Incident	[-183,-1]	OT	ATC/PROCEDUREID, local codes	n/a	Incident to vaccine administration	n/a	n/a
Unexposed	Matched calendar day	Single entry	Prevalent	[-183,-1]	OT	ATC/PROCEDUREID, local codes	n/a	Not required	n/a	n/a

¹ IP = inpatient, OP = outpatient, ED = emergency department, OT = other, n/a = not applicable

²See appendix for listing of clinical codes for each study parameter

Note: The washout window refers to the 183 days prior to the index date ([-183, -1]), during which no record of vaccination with the same ATC/PROCEDUREID/local code for manufacturer is allowed prior to index date for the exposed group. This criterion ensures that the vaccination capture at index date represents a new (incidence) vaccination event. For the unexposed (prevalent) group, the same window is listed for consistency, but no exclusion is applied because individuals in this arm have no vaccination at index date, and incident vaccination status is not required.

7.3.2 Study inclusion criteria – Context and rationale

Individuals must meet all the following inclusion criteria at index date:

1. Participants aged 16 years or older at index date will be identified from linked population databases, CPRD and VID.
2. Participants with continuous (or active) enrolment in the database for at least 6 months prior to index date (time zero) to allow for assessment of baseline characteristics, pre-existing conditions, and confounding factors. We assume *observable* participants are the ones who have not opted out of data use for research purposes. Our approach allows for inclusion without active consent but excludes recent entrants to the database (e.g., due to recent migration).
3. Participants with stable pre-existing conditions are eligible, provided there is no record of hospitalization within six weeks prior to the index date (time zero). Diagnostic codes and medication prescription or dispensation data will be used to define these conditions. This approach allows inclusion of individuals with managed chronic illnesses while excluding those with recent acute health events.

Table 4. Operational Definitions of Inclusion Criteria

Criterion	Details	Assessment window	Care Settings ¹	Code Type ²	Diagnosis position ³	Applied to study populations:	Measurement characteristics/ validation	Source for algorithm
Age	Age in years defined by (time 0 – year of birth)/365	[0,0]	OT	n/a	n/a	Exposed Unexposed	n/a	n/a
Registration	Presence healthcare encounters, or registration dates	[-183,0]	IP, OP, OT	ICD9-CM, ICD10-CM, SNOMED/READ /MEDCODEID	any	Exposed Unexposed	n/a	n/a
Pre-existing stable medical conditions ⁴	(1) ≥1 diagnosis record in the 6-month baseline period, and (2) no hospitalization with primary diagnosis in the 6 weeks prior to index. ⁵	[-46,0]	IP, OP, ED	ICD9-CM, ICD10-CM SNOMED/READ /MEDCODEID	Primary	Exposed Unexposed	n/a	n/a

¹ IP = inpatient, OP = outpatient, ED = emergency department, OT = other, n/a = not applicable

² See appendix for listing of clinical codes for each study parameter, ICD10 = International Classification of Diseases 10th edition

³ Specify whether a diagnosis code is required to be in the primary position (main reason for encounter)

⁴ The list of pre-existing stable medical conditions is presented in Appendix 4, Table 2.

⁵ Primary diagnosis of hospitalization refers to any acute events prompting admission. List is provided in Appendix 5.

7.3.3 Study exclusion criteria – Context and rationale

Individuals will be excluded if they meet any of these criteria at index date:

1. Participants who received any COVID-19 vaccine other than the Pfizer-BioNTech COVID-19 vaccine, as indicated by immunization records, to ensure a vaccine-naïve cohort for the purposes of this analysis.
2. Participants who received any COVID-19 prevention medication or pre-exposure prophylaxis (PrEP) (Pemivibart, Tixagevimab/Cilgavimab [Evusheld]) prior to or on the index date, as identified through prescription or dispensation records. Long-acting monoclonal antibody therapies used for PrEP can confer partial protection against SARS-CoV-2 infection independent of vaccination.
3. Participants with immunocompromised status based on diagnostic codes, medication use (e.g., immunosuppressants including ATC codes for L04A, L01X, L01B, H02AB, P01BA), or relevant clinical indicators in the 6 months prior to the index date. Immunocompromised individuals may have altered immune responses to vaccination and differential risk of infection and severe outcomes. Including such individuals could reduce the comparability of vaccinated and unvaccinated groups and affect the generalizability of results to the broader population. Excluding these individuals ensures a more homogeneous cohort for estimating vaccine effectiveness under typical immune function conditions.

Table 5. Operational Definitions of Exclusion Criteria

Criterion	Details	Assessment window	Care Settings ¹	Code Type ²	Diagnosis position ³	Applied to study populations:	Measurement characteristics/ validation	Source for algorithm
Any COVID-19 vaccine	Presence diagnostic/procedures codes, or registration	[-183,-1]	OT	ATC/PROCEDURE	n/a	Exposed Unexposed	n/a	n/a
PrEP	Presence of related codes	[-183,-1]	IP, OP, OT	ATC/PROCEDURE	n/a	Exposed Unexposed	n/a	n/a
Immunocompromised	(1) ≥1 diagnosis record, OR (2) had ≥1 prescription	[-183,0]	IP, OP, ED, OT	ICD9-CM, ICD10-CM, SNOMED/READ /MEDCODEID	any	Exposed Unexposed	n/a	n/a

¹ IP = inpatient, OP = outpatient, ED = emergency department, OT = other, n/a = not applicable

² See appendix for listing of clinical codes for each study parameter

³ Specify whether a diagnosis code is required to be in the primary position (main reason for encounter)

7.4. Variables

7.4.1 Exposure(s) of interest – Context and rationale

Study exposure: Identification of BNT162b2 mRNA COVID-19 vaccination will be based on recorded prescriptions, dispensing, or administration of the vaccine (as available). Vaccine receipt and date of vaccination will be obtained from all possible sources that capture COVID-19 vaccination, such as pharmacy dispensing records, general practice records, immunisation registers, vaccination records, medical records, or other secondary data sources (as available). Vaccination with the other existing COVID-19 vaccines will be identified for exclusion purposes, in a similar way as for the BNT162b2 mRNA COVID-19 vaccines.

The identification of the second dose of the primovaccination schema for BNT162b2 mRNA COVID-19 will be based on two approaches: (1) when available in data sources, it will be feasible to discern whether the vaccination record is dose one or dose two (2) when data sources do not distinguish between the first and second dose, we will consider only records of the same vaccine that occur at least 14 days after the first recorded vaccine (identified on the index date) to represent a second (separate) vaccination dose. This second approach is necessary to deal with the possibility of double-recordings of the same vaccination dose.

Comparator: Eligible individuals in the comparator group (i.e., unvaccinated) will be identified based on the absence of any recorded receipt of a COVID-19 vaccine up to and including the index date (also described in section 7.3.1, table 3). This will be determined by verifying that there are no records of prescription, dispensing, or administration of any COVID-19 vaccine, including BNT162b2 and other authorized vaccines, across all available data sources. These sources include pharmacy dispensing data, general practice records, immunisation registries, vaccination records, medical records, and other relevant secondary databases. To ensure accurate classification, individuals with any evidence of COVID-19 vaccination at or prior to the index date will be excluded from the unvaccinated comparator group.

Algorithm to define duration of exposure effect:

Data on vaccines will be cleaned, as required. For example, if two doses of the same vaccine are recorded on the same date (or within a period of 14 days) this will be considered the same dose and the earliest date considered as time zero (i.e., Bridged).

The immunological effect – i.e., the development of protective immunity or measurable antibody response – is expected to begin or become noticeable approximately 14 days after the first vaccine dose.⁸

Matching procedure: The aim of the matching procedure is to construct a cohort with a 1:1 ratio of vaccinated individuals to control individuals (i.e., matched unit). Individuals will be evaluated for eligibility on each calendar day of the enrolment period (until 31 January 2022). In order to most directly or literally emulate a parallel-arms trial we specify a matching procedure such that each individual contributes data only once to the analysis in either the vaccinated or the unvaccinated group. This would emulate a true trial setting where it is an implicit inclusion criterion that a participant must not already have been enrolled in the study.

- On the first day of the study period, we identify vaccinated individuals as those eligible individuals who receive their first vaccination with BNT162b2 on that date, defined as time-zero for the vaccinated.

- We match each vaccinated individual to a control individual, that is, an individual who is eligible but unvaccinated on that calendar date, and who has the same values on the matching covariates. When more than one candidate match is available, one is sampled at random.
- Once a match has been selected, they are removed from the general pool of available control individuals for that calendar date so that they cannot be matched to another vaccinated individual; in other words, controls are sampled without replacement. Furthermore, once a control has been selected, they are removed from the pool of eligible exposed and control matched units for all future time points in the study period.
- We repeat the procedure for the second day of the study period based on the subset of vaccinated and control units who remain eligible.

The following variables will be considered as potential confounders and used as matching variables:

- Year of birth (Age)
- Sex
- Place of residence (as available – clinical practice [CPRD], neighbourhood [VID])
- Socio-economic status (as available – deprivation index [CPRD], income [VID])
- Pregnancy status (yes/no), and amongst pregnant individuals, trimester (first, second, third)
- Number of pre-existing conditions (categorised as shown below) considered by the US Center for Disease Control and Prevention (CDC) as high-risk conditions for severe COVID-19 disease⁹
 - CDC at risk group 0 = none of the conditions
 - CDC at risk group 1 = one of the conditions
 - CDC at risk group 2 = more than one condition
- Healthcare utilization (or proxy as available)
- Immunocompromised status

The matching procedure described above ensures that each individual is included only once in the analysis, which in turn allows us to use relatively simpler methods of statistical analysis which rely on the assumption of independence of observations.

However, in so-called “sequential trial emulation” settings it is typically more common to allow control units to be re-sampled, serve as controls at multiple points in time, and, unless the inclusion criteria of the target trial explicitly state so, to serve as exposed units at later points in time. There are two potential disadvantages of the above “parallel-arm” design when compared to the typical sequential trial emulation approach. First, it may be relatively statistically inefficient, as observations of vaccinated and control units are dropped later in the study period. Second, the sequential trial approach essentially ensures that all individual who are eligible and vaccinated while eligible are included in the analysis, thus ensuring that we can estimate the treatment effect marginalised over the distribution of baseline covariates of all vaccinated individuals who meet the eligibility criteria at the moment of vaccination. In the parallel-arm approach, by necessity, not all such individuals are included, as a) those who serve as matched controls earlier are removed from the analysis and b) vaccinated individuals may be dropped due to the absence of a suitable matched control with a higher probability.

On the other hand, the sequential trial approach has several disadvantages. It does not emulate the parallel-arm design of the target trial, in which individuals who previously served as a control would not at later dates be considered for enrolment in the exposed group. As such, the alignment of the design with the target estimand is unclear. Also, re-use of individual data artificially inflates the sample size in the sense that the number of exposed and control “units” in the analysis will outnumber unique individuals who serve as either exposed or controls in the data. Furthermore, there is substantial overlapping between the “units” of analysis from the same individual. Treating the non-independence of observations necessitates the use of additional assumptions, along with more complex statistical models that account for or model the correlation induced between observations from the same subject and between treatment groups (as the same subject can contribute to both). These complexities may hinder the interpretation of the results and introduce uncertainties that could negatively affect the regulatory assessment.

We will implement the typical “sequential trial” approach described above as a supplementary analysis and investigate the difference in distribution of baseline covariates of the matched vaccinated in both cases.

Table 6. Operational Definitions of Exposure

Exposure group name(s)	Details	Washout window	Assessment Window	Care Setting ¹	Code Type ²	Diagnosis position ³	Applied to study populations:	Incident with respect to...	Measurement characteristics/validation	Source of algorithm
Exposure	BNT162b2 mRNA COVID-19	[-183, 0]	[0, censor]	OT	ATC, PROCODEID, local codes	n/a	Exposed	Incident to vaccine administration	n/a	n/a
Comparator	No vaccination	[-183, 0]	[0, censor]	n/a	n/a	n/a	Unexposed	n/a	n/a	n/a

¹ IP = inpatient, OP = outpatient, ED = emergency department, OT = other, n/a = not applicable

² See appendix for listing of clinical codes for each study parameter

³ Specify whether a diagnosis code is required to be in the primary position (main reason for encounter)

Note: The washout window (-183 to 0 days) refers to the 183-day period prior to the index date.

- For the **exposed group**, this ensures inclusion of incident vaccination events (no previous COVID-19 vaccine during the window).
- For the **comparator group**, the same window is applied to confirm absence of vaccination before the matched index date.

7.4.2 Outcome(s) of interest – Context and rationale

SARS-CoV-2 infection, the causative agent of COVID-19, remains a central public health concern due to its ongoing transmission, potential for severe disease, and the emergence of new variants. In the context of evaluating vaccination strategies, incident SARS-CoV-2 infection serves as a primary outcome to assess the intervention's effectiveness in reducing viral acquisition. Capturing infection through laboratory-confirmed test results (positive PCR/antigen test) provides a direct measure of prophylactic benefit.

Table 7. Operational Definitions of Outcome

Outcome name	Details	Primary outcome?	Type of outcome	Washout window	Care Settings ¹	Code Type ²	Diagnosis Position ³	Applied to study populations:	Measurement characteristics/validation	Source of algorithm
SARS-CoV-2 infection	1 code within 90 days of follow-up	Yes	Binary	n/a	IP, OP, ED, OT	ICD10, SNOMED/R EAD/MEDC ODEID, Local codes	n/a	Exposed Unexposed	n/a	n/a

¹ IP = inpatient, OP = outpatient, ED = emergency department, OT = other, n/a = not applicable

² See appendix for listing of clinical codes for each study parameter

³ Specify whether a diagnosis code is required to be in the primary position (main reason for encounter)

7.4.3 Follow up – Context and rationale

Follow-up of the exposed and unexposed groups will begin at time zero (first dose vaccination for the treated and matched date for the controls) and will continue up to 90 days. This duration was selected based on evidence that COVID-19 vaccine-induced protection, particularly against symptomatic SARS-CoV-2 infection, wanes substantially after the first three months post-vaccination.⁸ End of follow up will differ depending on the estimand, as different strategies are used to handle intercurrent events.

Estimand 1

Follow-up ends at the earliest of:

- SARS-CoV-2 infection (outcome)
- Death (outcome)

- Loss to follow-up (i.e., de-registration from the data source)
- Administrative end of follow-up at day 90
- End of data availability in the data source (see note)
- Receipt of a third (booster) dose of a non-target COVID-19 vaccine in the exposed (IE3)
- Receipt of a first dose of any COVID-19 vaccine for the control (IE3)
- Receipt of any other preventative COVID-19 treatment (pre-exposure prophylaxis PrEP) (IE5)

Estimand 2

Follow-up ends at the earliest of:

- SARS-CoV-2 infection (outcome)
- Death (IE6)
- Loss to follow-up (i.e., de-registration)
- Administrative end of follow-up at day 90
- End of data availability in the data source (see note)
- Receipt of a third (booster) dose of a non-target COVID-19 vaccine in the exposed (IE3)
- Receipt of a first dose of any COVID-19 vaccine for the control (IE3)
- Receipt of any other preventative COVID-19 treatment (pre-exposure prophylaxis PrEP) (IE5)

Note that “End of data availability in the data source” is included here for reasons of completeness. As the study period ends in 2022, and data is extracted in 2025/2026 we would expect that data availability will end after the study period is complete.

Table 8. Operational Definitions of Follow Up

	Day 0	
	Select all that apply	Specify
Follow up start	Day 0	
Follow up end¹		
Date of outcome	Yes	First Occurrence of SARS-CoV-2 infection
Date of death	Yes	Included as part of the composite outcome (estimand 1) or censor at date of death (estimand 2)
End of observation in data	Yes	Censor at date of last healthcare contact or known de-registration, emigration or end of data availability (non-administrative censoring)
Day X following index date <i>(specify day)</i>	Yes	Day 90 (administrative censoring)
End of study period <i>(specify date)</i>	Yes	30 April 2022
End of exposure <i>(specify operational details, e.g. stockpiling algorithm, grace period)</i>	Yes	Used to define inclusion for principal stratum strategy in estimand 2. Specifically, receipt of the second dose (end of exposure) defines principal stratum membership for the exposed
Date of add to/switch from exposure <i>(specify algorithm)</i>	Yes	Only used to define follow-up end for hypothetical strategy in both estimand 1 and 2.

¹ Follow up ends at the first occurrence of the first event corresponding to criteria that end follow up.

7.4.4 Covariates (confounding variables and effect modifiers, e.g. risk factors, comorbidities, comedICATIONS) – Context and rationale

Including the following covariates is critical to control for confounding variables at baseline that can influence both vaccination likelihood and SARS-CoV-2 infection:

- Age, sex/gender, and ethnicity can affect immune responses and disease susceptibility.
- Socioeconomic status, smoking status, and BMI impact health behaviours and infection risk.
- Comorbidities (including immunocompromised status), prior SARS-CoV-2 infection, and vaccination history alter baseline immunity and disease severity risk, which are important to consider when measuring vaccine protection. Comorbidities will be summarised via (1) Number of pre-existing conditions at risk of developing severe COVID-19 considered by the US Centers for Disease Control and Prevention (CDC), (2) defining whether individuals have an immunocompromised status, and (3) through the Charlson Comorbidity Index (CCI) with 19 conditions based on the Glasheen adaptation for use with standard diagnostic codes such as ICD-10.^{10,11} using records from the last 10 years.
 - For the CDC, three categories will be used:
 - CDC at-risk group 0 = none of the conditions
 - CDC at-risk group 1 = 1 of the conditions
 - CDC at-risk group 2 = >1 conditions.
 - At-risk medical conditions will be identified by diagnostic codes and related medicinal products. Medical conditions are summarised below:
 - Cancer (with chemo/immuno/radiotherapy, cancer treatment, immunosuppressant; targeted cancer treatment (such as protein kinase inhibitors or PARP* inhibitors); blood or bone marrow cancer (such as leukaemia, lymphoma, myeloma))
 - Type 1 & 2 Diabetes
 - Obesity (BMI ≥ 30 kg/m²)
 - Cardiovascular disease / serious heart conditions including heart failure, coronary artery disease, cardiomyopathies
 - Chronic lung disease including COPD, asthma, bronchiectasis, interstitial lung disease, cystic fibrosis, tuberculosis.
 - Chronic kidney disease
 - HIV
 - Immunosuppression
 - Sickle cell disease
 - Hypertension
 - For the immunocompromised status, we will define it as having at least one of the following: immunodeficiencies, immunosuppressant medication use, human immunodeficiency virus infection (with and without AIDS) and other immunosuppressing conditions
- Healthcare access and utilization, along with medication use, impact testing and clinical management, affecting outcome ascertainment.
- Geographic location captures heterogeneity in exposure risk and healthcare resources.

- Pregnancy status represents a unique physiological condition associated with differences in immune response and COVID-19 risk.

For further details on operationalisation, see Table 9 below.

Although not a matching variable, the following covariate will be defined and used for as a covariate to predict receipt of the second dose in the exposed group (part of the principal stratum analysis for Supplementary Estimand 1):

Charlson Comorbidity Index (CCI): each condition will be assigned an integer weight from 1 to 6, with 6 representing the most severe comorbidity. The sum of the weighted comorbidity scores results in a summary score, for example a patient with chronic respiratory disease (score 1) and hemiplegia (score 2) would have a total score of 3. The following conditions and assigned weights will be use:

- Weight = 1:
 - Congestive heart failure
 - Myocardial infarction
 - Rheumatic disease
 - Chronic respiratory disease
 - Peripheral vascular disease
 - Cerebrovascular disease
 - Mild liver disease
 - Peptic ulcer disease
 - Diabetes without chronic complications
 - Dementia
 - Mild to moderate renal disease
- Weight = 2:
 - Hemiplegia or paraplegia
 - Diabetes with chronic complications
 - Any malignancy (excluding metastatic disease)
- Weight = 3:
 - Moderate or severe liver disease
 - Severe kidney disease
 - HIV infection without AIDS
- Weight = 6:
 - Metastatic solid tumor
 - AIDS

Table 9. Operational Definitions of Covariates

Characteristic	Details	Type of variable	Assessment window	Care Settings ¹	Code Type ²	Diagnosis Position ³	Applied to study populations:	Measurement characteristics/validation	Source for algorithm
Age*	age in years defined by (time 0 – year of birth)/365	Continuous, Categorized	[0, 0]	OT	n/a	n/a			n/a
Death	Date of death in administrative record in addition to diagnostic code (exclusion criteria)	Binary	[0, 0]	OT	ICD10, SNOMED/READ/MEDCODEID, Local codes	n/a			n/a
Sex/Gender*	As recorded in data instance	Categorical	[0, 0]	OT	n/a	n/a			n/a
Ethnicity	As recorded in data instance	Categorical	[0, 0]	OT	Local codes	n/a			n/a
Socioeconomic status*	As recorded in data instance. It might include deprivation index [CPRD] or income [VID]	Categorical ⁴	[0, 0]	OT	Local codes	n/a			n/a
Smoking status	Never, Former, Current smoker	Categorical	[-180, 0]	IP, OP, OT	ICD9-CM, ICD10-CM, SNOMED/READ/MEDCODEID, Local codes	any			n/a

Characteristic	Details	Type of variable	Assessment window	Care Settings ¹	Code Type ²	Diagnosis Position ³	Applied to study populations:	Measurement characteristics/validation	Source for algorithm
BMI	Most recent BMI prior to index date	Continuous, Categorized	[-180, 0]	IP, OP, OT	ICD9-CM, ICD10-CM, SNOMED/READ/MEDCODEID, Local codes or values	any			n/a
Comorbidities using CDC-at risk groups, Immunocompromised status, and Charlson Index	Presence of chronic conditions defined using diagnostic codes in the 12 months prior to index date	Binary	[-365, 0]	IP, OP	ICD9-CM, ICD10-CM, SNOMED/READ/MEDCODEID	any			n/a
Prior SARS-CoV-2 infection	Any recorded positive PCR or antigen test	Binary	[-365, 0]	IP, OP, OT	ICD9-CM, ICD10-CM, SNOMED/READ/MEDCODEID, Local codes	n/a			n/a
Prior vaccination history (e.g., influenza)	Any record of influenza vaccine in the 12 months prior to index date	Binary	[-365, 0]	IP, OP, OT	ATC/PROCEDURE, local codes	n/a			n/a
Other COVID-19 vaccines	Any record of non-target COVID-19	Binary	[-365, 90]	IP, OP, OT	Local codes	n/a			n/a

Characteristic	Details	Type of variable	Assessment window	Care Settings ¹	Code Type ²	Diagnosis Position ³	Applied to study populations:	Measurement characteristics/validation	Source for algorithm
	vaccine in the 12 months prior to index date and during follow-up								
Healthcare access/utilization*	Number of GP visits, hospital admissions or other medical encounters in the 12 months prior to index date	Continuous	[-365, 0]	IP, OP, ED, OT	Local codes	n/a			n/a
Medicines	Use of relevant medications defined via prescription or dispensation data (see Appendix 3, Table 3)	Binary	[-365, 0]	IP, OP, ED, OT	ATC/PROCOD EID	n/a			n/a
Geographic location*	Region, administrative area, or postal code	Categorical	[0, 0]	OT	Local codes	n/a			n/a
Pregnancy status*		Binary	[0,0]	IP, OP, OT	ICD10,	n/a			n/a

Characteristic	Details	Type of variable	Assessment window	Care Settings ¹	Code Type ²	Diagnosis Position ³	Applied to study populations:	Measurement characteristics/validation	Source for algorithm
					SNOMED/READ/MEDCODE ICD, Local codes				

¹ IP = inpatient, OP = outpatient, ED = emergency department, OT = other, n/a = not applicable

² See appendix for listing of clinical codes for each study parameter

³ Specify whether a diagnosis code is required to be in the primary position (main reason for encounter)

³ Origin variables can be continuous or categorical, may require harmonization across sources.

*These variables are used during the matching procedure

7.5. Core Emulation Table – Design Summary

Table 10 A. Comparison of Target Trial and Proposed Target Trial Emulation Design Elements for the Primary Estimand

	Target Trial	Target Trial Emulation	Comment
Inclusion criteria	<p>Participants 16 years of age or older.</p> <p>Participants must demonstrate willingness and ability to comply with all scheduled visits, the vaccination plan, laboratory tests, lifestyle considerations, and other study-related procedures.</p> <p>Participants with preexisting stable medical conditions may be included, provided their condition does not</p>	<p>Participants 16 years of age or older.</p> <p>Participants must be enrolled in the population databases for 6 months to assess pre-existing conditions or any other medical conditions relevant for confounding adjustment.</p> <p>Participants with preexisting stable medical conditions may be included, provided their condition does not require hospitalization within six weeks prior to enrolment.</p>	<p>Emulation is more inclusive as willingness to participate and providing consent are not required. However, requirement in the emulation is to be present in the database for 6 months meaning we miss individuals who have recently entered the data source (e.g., due to migration). We assume such individuals are not systematically different with respect to their response to the vaccine.</p>

	<p>require hospitalization within six weeks prior to enrolment.</p> <p>Participants must be capable of providing personally signed, informed consent.</p>	<p>Participants' data is pre-approved for usage in the population databases. Included participants are those who do not avail of any opt-out clause to have their data removed from all research.</p>	<p>Operationalized using diagnostic codes prescription/dispensation, which are subject to inaccuracies</p> <p>Although the target population is the same in the trial and emulation, in practice, individuals who choose to enrol in a vaccine efficacy trial may display different person level characteristics and behaviour to those who remain unvaccinated in the general population</p>
Exclusion criteria	<p>Participants who have received any medication intended to prevent COVID-19.</p> <p>Immunocompromised individuals with known or suspected immunodeficiency, as determined by medical history, laboratory tests, or physical examination.</p> <p>Participants who have previously received any COVID-19 vaccine.</p>	<p>Participants who have received any medication intended to prevent COVID-19.</p> <p>Immunocompromised individuals with known or suspected immunodeficiency, as determined by medical history, laboratory tests, or physical examination.</p> <p>Participants who have previously received any COVID-19 vaccine.</p> <p>Participants have not already been enrolled in the study as either a vaccinated or control unit (primary analysis)</p>	<p>Identical in Target Trial and Emulation. Operationalized using diagnostic codes, prescription/dispensation, and vaccine administration which are subject to inaccuracies. To reflect the specific parallel-arms emulation procedure (i.e. the primary matching procedure), we here explicitly state that in the emulation we exclude individuals who have already been enrolled in the study. In the real trial, this exclusion criteria is implicit. In the supplementary analysis, this exclusion criteria is not applied.</p>
Setting	Multicentre	Routine care data sources (e.g., primary care and/or administrative databases)	Real-world data captures care as delivered.

<p>Study treatment conditions</p>	<p>Intervention: Two doses BNT162b2 mRNA COVID-19 vaccine separated by 21 days</p> <p>Control: placebo, with two matched doses.</p>	<p>Intervention: Two doses BNT162b2 mRNA COVID-19 vaccine separated by 21 days (with an allowable interval between 19 to 42 days)</p> <p>Control: No vaccine</p>	<p>Placebo interventions are not available in real world data sources. Instead, we use unvaccinated as the control condition. This assumes a lack of coded COVID-19 vaccination accurately reflects non-vaccination.</p>
<p>Method of Assignment to Trial Intervention</p>	<p>Participants will be randomly assigned to either strategy at baseline</p>	<p>On each calendar day, we will identify eligible vaccinated and unvaccinated units. All eligible vaccinated individuals will be “assigned” to the vaccinated condition. Eligible unvaccinated individuals will be assigned to the control condition if matched to an eligible vaccinated individual. Matching will be based on covariates with a 1:1 ratio.</p>	<p>Randomisation cannot be directly emulated. Emulation of randomization relies on assumption of conditional exchangeability between groups given observed covariates. 1:1 matching is used to emulate randomization, matching on year of birth, sex at birth, place of residence, socio-economic status (when available), pregnancy status, and risk factors for COVID-19 (e.g., immunocompromised status, CDC high-risk conditions). Observed covariates identified using codelists, register information in the data source. Information on some covariates may be missing.</p>
<p>Time (when follow up begins and ends):</p>	<p>For each participant, follow-up starts at the time of randomisation (which coincides with the first dose of the intervention) and ends at diagnosis of SARS-CoV-2 infection, death, study withdrawal, loss to follow-up, any intercurrent event after which observed follow-up is not of interest (see above section), or administrative</p>	<p>For treated participants meeting eligibility criteria, follow-up starts at the time they receive the first dose of the target vaccine. For the control participants, follow-up starts on index date determined by the start of follow-up of their matched treated unit. Follow-up ends at the first record of COVID-19 infection, death, loss to follow-up, any intercurrent event after which observed follow-up is not of</p>	<p>The end of follow-up is identical in Target Trial and Emulation. The start of follow-up is also the same for the treated group. For individuals in the control group, the start of follow-up cannot be determined by a clinical event (e.g. start of intervention). To address this, it will be set to the date of the first vaccine dose of the matched treated individual.</p>

	end of study period (3-months), whichever occurs earlier.	interest (see above section) or administrative end of study period (3 months), whichever occurs earlier.	
Outcome (including operational definition)	Laboratory-confirmed SARS-CoV-2 infection (Individual who test positive for SARS-CoV-2 using an approved laboratory test, such as reverse transcription polymerase chain reaction (RT-PCR) test, nucleic acid amplification test (NAAT) and rapid antigen test, only confirmed by laboratory), or death from any cause	Laboratory-confirmed SARS-CoV-2 infection (Individual who test positive for SARS-CoV-2 using an approved laboratory test, such as reverse transcription polymerase chain reaction (RT-PCR) test, nucleic acid amplification test (NAAT) and rapid antigen test, only confirmed by laboratory and results are documented in the individual's medical or public health records during the study period, or death from any cause as recorded in the data source.	Identical in Target Trial and Emulation for laboratory-confirmed SARS-CoV-2 infection. Code lists and algorithms will be used to identify death (i.e., a combination of ICD10 codes together with records of death dates for individual persons) However, ascertainment of SARS-CoV-2 infection is likely to be incomplete, as both testing practices and recording of results in routine care are variable. This limits the accuracy of infection identification. Although death is generally more reliably captured than infection, completeness may still be constrained by the extent and quality of linkage to vital statistics or mortality registries. It is assumed that absence of records indicating COVID-19 infection reflect no infection occurred. And equally with death. In both cases this could lead to misclassification of the outcome.
Intercurrent events and strategies to handle them	Intercurrent events are handled using three strategies: - Treatment policy for (1) missing or ineligible second	Intercurrent events are handled using three strategies: - Treatment policy for (1) missing or ineligible second	Identical in Target Trial and Emulation. To implement the treatment policy , explicit identification of IE1 is not required. . IE2 and IE4 can also be

	<p>dose, (2) receipt of third dose of target vaccine within 3 months after the second dose, and (4) non-COVID-19 vaccine post-treatment</p> <ul style="list-style-type: none"> - Hypothetical for (3) third dose of a non-target COVID-19 vaccine and (5) post-treatment preventive therapies - Composite for (6) death. 	<p>dose, (2) receipt of third dose of target vaccine within 3 months after the second dose, and (4) non-COVID-19 vaccine post-treatment</p> <ul style="list-style-type: none"> - Hypothetical for (3) third dose of a non-target COVID-19 vaccine and (5) post-treatment preventive therapies - Composite for (6) death. 	<p>identified using vaccine administration dates and batch numbers to identify early third doses, but again, as they are considered part of the treatment policy, this is not strictly necessary to implement the hypothetical strategy, IE3 can be identified filtering by vaccine manufacturer and timing post-second dose. Similarly, IE5 can be identified through prescription records for preventive medications initiated post-vaccination.</p> <p>To implement the composite strategy, IE6 can be retrieved using diagnostic codes for death as well as administrative information.</p>
Loss to follow up	A participant will be considered lost to follow-up if they repeatedly fail to return for scheduled visits and they cannot be contacted by the study staff.	A participant will be considered lost to follow-up when the individual is deemed to have administratively exited the data source (e.g. as directly recorded in the patient record)	Emulation differs from the target trial as follow-up is passive and relies on the administrative enrolment or presence of healthcare interactions in the database.

Table 10 B. Comparison of Target Trial and Proposed Target Trial Emulation Design Elements for the Supplementary Estimand

	Target Trial	Target Trial Emulation	Comment
Inclusion criteria	Same as estimand 1	Same as estimand 1	Same as Estimand 1
Exclusion criteria	Same as Estimand 1	Same as Estimand 1	Same as Estimand 1

Setting	Same as Estimand 1	Same as Estimand 1	Same as Estimand 1
Study treatment conditions	Same as Estimand 1	Same as Estimand 1	Same as Estimand 1
Method of Assignment to Trial Intervention	Same as Estimand 1	Same as Estimand 1	Same as Estimand 1
Time (<i>when follow up begins and ends</i>):	For each participant, follow-up starts at the time of randomisation (which coincides with the first dose of the intervention) and ends at diagnosis of SARS-CoV-2 infection, death, study withdrawal, loss to follow-up, any intercurrent event which censors follow-up (see above section), or end of study period (3-months), whichever occurs earlier.	For treated participants meeting eligibility criteria, follow-up starts at the time they receive the first dose of the target vaccine. For the control participants, follow-up starts on index date determined by the start of follow-up of their matched treated unit. Follow-up ends at the first record of COVID-19 infection, death, loss to follow-up, any intercurrent event which censors follow up (see above section) or end of study period (3 months), whichever occurs earlier.	The end of follow-up is identical in Target Trial and Emulation. The start of follow-up is also the same for the treated group. For individuals in the control group, the start of follow-up cannot be determined by a clinical event (e.g. start of intervention). To address this, it will be set to the date of the first vaccine dose of the matched treated individuals.
Outcome (including operational definition)	Same as Estimand 1	Same as Estimand 1	Same as Estimand 1
Intercurrent events and strategies to handle them	Intercurrent events are handled using three strategies: <ul style="list-style-type: none"> - Principal Stratum for (1) missing or ineligible second dose 	Intercurrent events are handled using three strategies: <ul style="list-style-type: none"> - Principal Stratum for (1) missing second dose 	Identical in Target Trial and Emulation. <p>To implement the principal stratification, IE1 can be identified via absence of a second dose within the expected interval which can be flagged using date fields and vaccine</p>

	<ul style="list-style-type: none"> - Treatment policy for (2) receipt of third dose of target vaccine within 3 months, and (4) non-COVID-19 vaccine post-treatment - Hypothetical for (3) third dose of a non-target COVID-19 vaccine and (5) post-treatment preventive therapies - While alive for (6) death. 	<ul style="list-style-type: none"> - Treatment policy for (2) receipt of third dose of target vaccine within 3 months, and (4) non-COVID-19 vaccine post-treatment - Hypothetical for (3) third dose of a non-target COVID-19 vaccine and (5) post-treatment preventive therapies - While alive for (6) death. 	<p>codes. To implement the hypothetical strategy, IE3 can be identified filtering by vaccine manufacturer and timing post-second dose. Similarly, IE5 can be identified through prescription records for preventive medications initiated post-vaccination.</p> <p>To implement the while-alive, IE5 can be retrieved using diagnostic codes for death as well as administrative information.</p>
Loss to follow up	Same as Estimand 1	Same as Estimand 1	Same as Estimand 1

7.6. Data analysis

7.6.1 Analysis plan – Context and rationale

For Estimand 1, the main estimand supporting decision making, we will estimate the incidence rate ratio (RR) under a hypothetical scenario in which participants in the control arm would not experience the intercurrent event of receiving the BNT162b2 vaccine. For the primary analysis, the rate ratio will be estimated using a Poisson regression model including treatment group and log of time at risk as offset. Sensitivity analyses will be performed to investigate the robustness of the results to assumptions made in the primary estimand regarding potential violations of the missing-at-random assumption for non-administrative censoring. The IPCW analysis will allow censoring to depend on observed baseline and time-varying covariates, while best- and worst-case scenario analyses assess the impact of extreme departures from this assumption.

For Estimand 2, a weighted Poisson regression model to estimate the RR in the principal stratum of individuals who would tolerate and receive the two doses of the vaccine, if they were to undergo vaccination.

The BNT162b2 effectiveness in preventing the SARS-Cov-2 infection will be estimated as $1 - RR$.

Supplementary analyses including diagnostic and descriptive assessments to support the main analysis will be presented to contextualise data and results from the primary and sensitivity analyses. These will include number of infections per treatment group, crude incidence rates, covariate distributions of matched cohorts, number of individuals censored and censoring patterns (intercurrent events, loss to follow-up) as well as model diagnostics.

An additional supplementary analysis will be produced whereby individuals in the control arm can contribute data multiple times to the control arm and also to the BNT162b2 arm after reception of the vaccine.

The results will be presented separately for each data source. In addition, we will explore the feasibility of combining the estimated Incidence Rate Ratios (IRR) from CPRD Aurum and VID using a random-effects meta-analysis of logarithmic transformation of the rate ratios. The combined IRRs and 95% CI will be calculated, and heterogeneity will be assessed using I².

7.6.2 Primary Estimand Main Analysis

i. Objective

No hypothesis testing will be performed in this study. The focus is on estimation of the Incidence Rate Ratio to quantify vaccine effectiveness and the corresponding 95% to reflect the uncertainty due to sampling variation.

ii. Exposure contrast

Exposed vs unexposed to BNT162b2 mRNA COVID-19

iii. Outcome

Occurrence of laboratory-confirmed SARS-CoV-2 infection or death

iv. Software

R (including libraries: data.table, duckdb, dplyr, fst, geeM, msm, logger, lubridate, rlang, survival, survminer)

v. Handling of intercurrent events (explaining how follow-up is handled post intercurrent event)

- **Missing or Ineligible for Second Dose**
Strategy used: *Treatment Policy*
Handling: Data after the occurrence of the IE will be included in the analysis as per the treatment policy strategy
- **Third (Booster) Dose of Target Vaccine**
Strategy used: *Treatment Policy*
Handling: Data after the occurrence of the IE will be included in the analysis as per the treatment policy strategy
- **Third Dose of Non-Target Vaccine**
Strategy used: *Hypothetical*
Handling: Non-administratively censored at the time of receiving the non-target vaccine dose.
- **Receipt of any non COVID-19 vaccine**
Strategy used: *Treatment policy*
Handling: Data after the occurrence of the IE will be included in the analysis as per the treatment policy strategy.
- **Pre-exposure Prophylaxis After First Dose** (e.g., *monoclonal antibodies, antivirals*)
Strategy used: *Hypothetical*
Handling: Non-administratively censored at the time of receiving the PrEP treatment.
- **Death**
Strategy used: *Composite endpoint with SARS-CoV-2 Infection*
Handling: Considered as outcome.

Diagnostics for IE handling:

As a descriptive supplementary analysis, we will tabulate for both matched-exposed and matched-control how often each IE occurs, and which IE leads to censoring of the follow-up. Note that for IEs under the treatment policy strategy (IE1, IE2, IE4), the IE may occur but not censor follow-up.

vi. Outcome Model

Relative risk parameterized with the Incidence Rate Ratio (IRR).

The IRR will be estimated using a Poisson regression model, in which the occurrence of the outcome is predicted by group membership, with an offset term equal to the log person time at risk. This model assumes

Matching will be used to adjust for baseline confounding (see below).

For the primary analysis, model-based standard errors will be used to estimate the 95%CI.

Standard Vaccine effectiveness as $VE = (1 - IRR) \times 100$.

Assumptions

Event occurrence follows a Poisson distribution with constant risk over time. VE is assumed constant over time within the period under consideration. Each individual contributes independent person-time to the analysis. In Poisson regression, censoring is implicitly assumed independent of the outcome (not informative) conditional on treatment assignment and no infection up to the time of censoring.

vii. Confounding adjustment

Cohort Matching (for comparator selection)

To construct comparable exposure groups, a 1:1 matching approach without replacement will be used to select unvaccinated individuals as comparators. This step ensures comparability at baseline on the matching variables listed in section 7.4.1 in the sense that, when every eligible vaccinated individual can be matched to an eligible control with the same covariate profile, the distribution of the matching variables matches in the control group matches that of the vaccinated.

The matching covariates are selected based on prior knowledge of factors which determine both the likelihood of receiving the vaccine and the probability of SARS-CoV-2 infection: age, sex, place of residence, socio-economic status, pregnancy status, number of pre-existing conditions considered by the US Centre for Disease Control and Prevention (CDC) as high-risk conditions, healthcare utilization. ⁹⁸

Assumptions Underlying Matching

- **No unmeasured confounding** (all relevant baseline confounders are matched on).
- **Positivity** (each individual has a non-zero probability of receiving either treatment, given their covariates).
- **Stable Unit Treatment Value (SUTVA)**. Each individual's potential outcome under the observed treatment equals their actual outcome, and the outcome of one individual is not influenced by the exposure of the other.

Diagnostics for Matching. Supplementary analyses will be presented concerning:

- **Covariate balance:** Check that both matching covariates and additionally available baseline characteristics (see section 7.4.4) are balanced across treatment groups after matching.

- Evaluate standardized mean differences (SMDs): SMDs < 0.1 will be considered acceptable.
- Characterisation of the exposed, matched-exposed and matched-control populations: We will describe the overall number of eligible exposed units (individuals) and their baseline characteristics (including the components of algorithms used for matching such as the CDC score). We will then describe these same baseline characteristics for the matched-exposed and matched-control units. This will allow us to check (a) how representative the matched-exposed population is to the entire exposed population, and (b) how well balanced the exposed and control units are on components of the matching variables.

viii. Missing Data Handling

Missing Exposure Data

Missing information on prescription or dispensation of the target vaccine will be interpreted as the exposure not occurring, and not incomplete data capture.

Missing Outcome Data

The Poisson model implicitly assumes non-informative censoring, meaning that censored participants contribute time at risk up to the time of censoring and their censoring is unrelated to the outcome, **conditional** on model covariates and survival up to the time of censoring (i.e., outcome data is missing at random under these assumptions). It is also assumed that participants who do not experience the outcome before censoring are correctly classified as having not had the event—that is, the available data provide complete outcome coverage with respect to the defined endpoint. This assumption concerns the correct classification of the outcome variable, and violations (e.g., missed or delayed event recording) could lead to outcome misclassification.

Missing Baseline/Matching Covariates

Missing data in baseline (matching) covariates is handled using a combination of complete case analysis and multiple imputation. Note that, for variables that are assessed using diagnostic codes or records of prescription/dispensation of a medicine, we assume complete coverage, and absence of a record will be taken to reflect a true absence of the corresponding event/medicine prescription. As such, missingness in this context largely concerns lifestyle factors such as socioeconomic status and BMI. We will first quantify the extent and pattern of missingness to evaluate which missing data strategy is most appropriate

- If low numbers of individuals (< 5%) have missing values on these matching variables we will conduct a complete case analysis.
- Otherwise, missing values for lifestyle factors will be imputed using multiple imputation with chained equations (MICE) under the Missing at Random (MAR) assumption. Information from individuals who have not yet experienced the IE will be used instead.
 - If a covariate has more than 40% missing data, we will consider alternative approaches (e.g., exclusion of the variable, sensitivity analyses) and justify the decision. Thresholds of 40% have been cited because effect estimates begin to be less reliable as the level of missingness increases beyond this threshold.¹²

Imputation Model

The MICE procedure will include all matching covariates, as well as predictors of missingness (if there are any additional factors not covered by the covariates in the treatment and outcome models). The treatment and outcome of interest will also be included.

Key covariates included in the imputation model will be:

- Demographics (age, sex, ethnicity [if available])
- Clinical history and comorbidities (see section 7.4.4 and Table 9)
- Lifestyle factors (e.g., smoking, BMI)
- Neighbourhood
- Marital status, educational level, employment status, or household size (as available)
- Education level

Full Conditional Distributions

MICE will use variable-specific conditional models:

- Logistic regression for binary variables (e.g., smoking yes/no).
- Multinomial logistic regression for categorical variables with >2 categories.
- Predictive mean matching for continuous variables

Number of Imputations and Diagnostics

We will generate at least 10 imputed datasets (to ensure stable estimates given the level of missingness) and pool results across imputations using Rubin's rules. Diagnostics will include:

- Checking whether imputed values are plausible and consistent with observed distributions.
- Evaluating convergence of the chained equations.
- Assessing stability and consistency of results across imputed datasets.

Effect Estimation Under Multiple Imputation

The imputation model will be applied prior to matching and effect estimation. Cohorts will be matched and outcome models will then be fitted in each imputed dataset, producing treatment effect estimates and corresponding variances. These estimates will be combined across the imputed datasets using Rubin's rules, which account for both within-imputation variance (the average estimation error within each imputed dataset) and between-imputation variance (the variability in estimates across imputations). The total variance therefore reflects uncertainty from both the imputation process and the effect estimation, producing valid confidence intervals.

Based on the feasibility assessment, we expect BMI to have from 5 to 30 percent missing values. Socioeconomic Status is expected to be partially missing, but we cannot specify a priori an expected missingness percentage.

ix. Subgroup Analyses

N/A

7.6.3 Supplementary Estimand (2) Main Analysis: Principal Stratum Effect

i. Objective

No hypothesis testing will be performed in this study. The focus is on estimation of the Incidence Rate Ratio to quantify vaccine effectiveness and the corresponding 95% to reflect the uncertainty due to sampling variation.

ii. Exposure contrast

Exposed vs unexposed to BNT162b2 mRNA COVID-19

iii. Outcome

Occurrence of laboratory-confirmed SARS-CoV-2 infection

iv. Software

R (including libraries: data.table, duckdb, dplyr, fst, geeM, msm, logger, lubridate, rlang, survival, survminer)

v. Handling of intercurrent events (explaining how follow-up is handled post intercurrent event)

Same as table 7.6.2, except:

Missing second dose of target vaccine

- **Strategy Used:** Principal Stratum
- **Handling:**

We consider the trial population to be divided into two strata based on whether they would (stratum (a)) or would not (stratum (b)) receive the second vaccine dose (within 90 days after T0) under the treated condition.

The principal stratum of interest is stratum (a) as these are the individuals in both the vaccinated and control groups who would receive the second dose of the vaccine if vaccinated. Within the vaccinated group, we directly observe membership of stratum (a): All vaccinated individuals who receive both doses of the vaccine are members of stratum (a). Note here that all exposed individuals who receive the second dose are considered here, regardless of whether this occurs after the occurrence of the outcome or any other IE.

Members of the control group may belong to stratum (a) or (b). However, in this group, strata membership is latent. In order to identify the causal effect from observed information, we must assume principal ignorability: that potential outcomes are independent of principal strata membership given observed covariates. This implies that we can sufficiently model strata membership using observed information. In order to do this, we will use principal score weighting. This technique allows us to estimate the principal stratum effect under two additional but closely related assumptions: That membership to stratum (a) is independent of treatment assignment conditional on observed covariates, and that the principal score model estimated in the vaccinated treatment arm is transportable to the control treatment arm (i.e., can predict the principal stratum membership in an unbiased way). Under these assumptions, we first fit a model predicting the probability of receiving the second dose during follow up among the matched treated members. The predictors in this model are the baseline matching covariates and

additional characteristics that might predict the probability of receiving the second dose. We consider here that adverse reactions to the vaccine, covid-19 infection itself, and medical history / frailty may be strong predictors of the likelihood to obtain the second vaccine. Thus, in addition to the matching variables (which already capture many of these predictors), we additionally include the Charlson Comorbidity Index (CCI) as a more fine-grained index of medical history / frailty, also including several factors which may make one more susceptible to experiencing an adverse reaction to the first dose.

In a second step, we use this estimated model to predict probabilities of stratum (a) membership (i.e. the principal score) amongst the matched control members. These Control group members receive a weight of $\text{PrincipalScore}/(1 - \text{PrincipalScore})$, reflecting Inverse Probability of Stratum Membership Weights (IPSMW).

To estimate the principal stratum effect, the analysis will be restricted to matched treated units who receive the second dose at any point in the 90 days following the first dose (described above) but include all matched control units (weighted by the IPSMW).

Death

- **Strategy Used:** While-Alive
- **Handling:** Participants who die are administratively censored; the estimand reflects the treatment effect while being alive.

vi. Outcome Model

As in A, except a weighted Poisson regression model will be used on the sub-selection of matched control and treated units described above, with IPSM weights. Robust standard errors will be used to compute confidence intervals.

vii. Confounding adjustment

As in 7.6.2. Note that as we take a selection of the matched units, pre-weighting the distribution of baseline confounders may be unequal between groups. However, as the baseline matching variables are included in the IPSMW model, these will be balanced in the weighted groups. To check this we will perform balance checks as described in 7.6.2 using the matched-weighted controls and selected exposed units.

viii. Missing Data Handling

Same as 7.6.2.

ix. Subgroup Analyses

N/A

7.6.4 Sensitivity Analyses

(1a,2a) Censoring Not Missing At Random: Inverse Probability of Censoring Weights

To estimate Estimand (1) and (2) we perform non-administrative censoring of follow up when an individual de-registers from the data source, or when they experience IE3 or IE5 (see 7.4.3 for details). This means we have missing time-at-risk and missing outcomes following time of non-administrative censoring. In the main analysis we assume that this data is Missing at Random (MAR) conditional on the treatment assignment.

Analysis Methods:

In this sensitivity analysis we relax this assumption, assuming instead that data which is non-administratively censored for the reasons detailed above is MAR conditional on baseline and time-varying confounders and treatment assignment. In order to do this, we must model the probability of being censored vs uncensored at each day of follow up, predicted by baseline matching covariates (e.g., age, sex, prior SARS-CoV-2 infection) and intermediate time-varying covariates that may predict non-administrative censoring (i.e., hospitalization, diagnosis with acute/incident immunosuppressing conditions).⁶

To do this, we predict, for each person-day of follow-up, the probability of being censored on that day. IPC weights can be calculated as: $1/(\text{the estimated probability of remaining uncensored given the covariates})$. The weight is calculated separately for each day and then multiplied together across all days of follow-up to give each participant's cumulative weight on each day.

In order to utilize IPCWs we must also alter other elements of our analytic strategy for estimating the IRR. In the main analysis, our analytic dataset consists of one row per individual, and the estimation proceeds by fitting a standard poisson regression model with an offset equal to the total person time of follow-up for each individual. However, this approach does not allow us to weight individual person-days with IPCWs. Instead, in the sensitivity analysis, we restructure the data such that each person-day of follow-up is represented by a single row, indexed by a person ID. For Estimand 1, each row obtains a weight given by the (time-varying) IPCW weights. For Estimand 2, each row obtains a weight equal to the product of the IPSM weight (obtained at baseline and time-invariant) and the (time-varying) IPCW weight. We estimate the IRR using generalized estimating equations (GEE) with a poisson link function, using the person-level identifier as a clustering variable to estimate the standard error accounting for the fact that we have repeated observations from the same individuals.

Assumptions:

- Censoring is conditionally independent of the outcome given baseline and time-varying observed covariates.
- Correct model specification
- The distributions of the censoring probabilities among censored and uncensored individuals must overlap
- All reasons for non-administrative censoring can be treated with a single model (same model for the censoring holds for all censoring reasons)

(1b, 2b) Non-administrative censoring: best/worst case scenarios

Analysis Methods:

- Best/worst case scenario

- For non-administrative censored individuals, repeat the analysis under four scenarios, assuming a) all censored individuals had the outcome of interest at the censoring date and b) no censored individuals had the outcome of interest at the censoring date, respectively for each treatment group. This equates to the following scenarios

Scenario	Exposed Group	Unexposed Group	Interpretation
1	Best case (lowest event rate)	Worst case (highest event rate)	Maximally favors vaccinated
2	Worst case (highest event rate)	Best case (lowest event rate)	Maximally favors control
3	Best case	Best case	Optimistic for both groups
4	Worst case	Worst case	Pessimistic for both groups

Assumptions:

- All or none of the non-administratively censored individuals had the outcome of interest

Table 11. Sensitivity analyses – rationale, strengths and limitations

	What is being varied? How?	Why? (What do you expect to learn?)	Strengths of the sensitivity analysis compared to the primary	Limitations of the sensitivity analysis compared to the primary
(1a, 2a) Informative censoring (IPCW)	The assumption that censoring is random conditional on treatment assignment is varied in the sense that it is assumed censoring is random conditional on a larger set of baseline and time-varying covariates. IPCW explicitly models the censoring mechanism based on observed covariates.	To assess the robustness of the treatment effect estimate to violations of the non-informative censoring assumption. If results are stable across scenarios, confidence increases that findings are not driven by bias (due to censoring).	IPCW adjusts for measured predictors of censoring. Offers a principled method for recovering unbiased estimates under informative censoring given the censoring mechanism can be modelled with observed covariates (MAR conditionally on observed information)	IPCW is sensitive to model misspecification. Cannot account for unmeasured factors affecting censoring, such as MNAR. Weighting can increase variance, especially if weights are unstable due to positivity violations.
(1b,2b) Informative censoring (best/worst case scenario)	The assumption that censoring is non-informative is varied.	To assess the robustness of the treatment effect estimate to violations of the non-informative censoring assumption. If results are stable across scenarios, confidence increases that findings are not driven by bias (due to censoring).	Offers an empirical method for exploring alternative assumptions about censoring. Sets bound on the extent of maximum possible bias due to informative censoring	Best/worst case scenarios are extreme assumptions. It does not allow for any variability in the process; the most likely values for the treatment effect under non-informative censoring may be very far away from the bounds.

7.6.5. Other Supplemental Analyses

The supplemental analyses aim to probe the degree to which re-using individual patient data (i.e., to serve as controls multiple times, and/or serve as an exposed unit after contributing to the control arm) yields different effect estimates than the primary approach of using patient data only once. To that end we replicate the analyses of the primary and supplementary estimand with this alternative “sequential trial emulation” design, described in 7.4.1.

Table 12. Supplemental analysis 1: Primary Estimand with Sequential Trial Design

Hypothesis:	No hypothesis testing will be performed in this study. The focus is on estimation of the Incidence Rate Ratio to quantify vaccine effectiveness and the corresponding 95% to reflect the uncertainty due to sampling variation.
Exposure contrast:	Exposed vs unexposed to BNT162b2 mRNA COVID-19
Outcome:	Occurrence of laboratory-confirmed SARS-CoV-2 infection or death
Analytic software:	R (including libraries: data.table, duckdb, dplyr, fst, geeM, msm, logger, lubridate, rlang, survival, survminer)
Handling of Intercurrent Events:	As in Primary Estimand section 7.6.2
Outcome Model(s): (provide details or code)	As in Primary Estimand section 7.6.2, except the IRR will be estimated using a generalized estimating equation (GEE) Poisson regression model, with the person-level identifier used as a clustering variable, to correct the standard error estimation for the repeated participation of individuals both within the control arm and between the control and exposed arm.
Confounding adjustment method	<i>Name method and provide relevant details, e.g. bivariate, multivariable, propensity score matched (specify matched algorithm ratio and caliper), propensity score weighting (specify weight formula, trimming, truncation), propensity score stratification (specify strata definition), other.</i>
	<i>Cohort Matching (for comparator selection)</i> To construct comparable exposure groups, a 1:1 matching approach with replacement will be used to select unvaccinated individuals as comparators. Otherwise, as in Primary Estimand section 7.6.2
Missing data methods	<i>Name method and provide relevant details, e.g. missing indicators, complete case, last value carried forward, multiple imputation (specify model/variables), other.</i>
	As in Primary Estimand section 7.6.2
Subgroup Analyses	<i>List all subgroups</i>
	N/A

Table 13. Supplemental analysis 2: Supplementary Estimand with Sequential Trial Design

Hypothesis:	No hypothesis testing will be performed in this study. The focus is on estimation of the Incidence Rate Ratio to quantify vaccine effectiveness and the corresponding 95% to reflect the uncertainty due to sampling variation.
Exposure contrast:	Exposed vs unexposed to BNT162b2 mRNA COVID-19
Outcome:	Occurrence of laboratory-confirmed SARS-CoV-2 infection or death
Analytic software:	R (including libraries: data.table, duckdb, dplyr, fst, geeM, msm, logger, lubridate, rlang, survival, survminer)
Handling of Intercurrent Events:	As in Primary Estimand section 7.6.3
Outcome Model(s): (provide details or code)	As in section 7.6.2, except the IRR will be estimated using a generalized estimating equation (GEE) Poisson regression model, with the person-level identifier used as a clustering variable, to account for the fact that individuals may contribute to both vaccinated and unvaccinated groups at different times.
Confounding adjustment method	<i>Name method and provide relevant details, e.g. bivariate, multivariable, propensity score matched (specify matched algorithm ratio and caliper), propensity score weighting (specify weight formula, trimming, truncation), propensity score stratification (specify strata definition), other.</i>
	<i>Cohort Matching (for comparator selection)</i> To construct comparable exposure groups, a 1:1 matching approach with replacement will be used to select unvaccinated individuals as comparators. Otherwise, as in Primary Estimand section 7.6.2
Missing data methods	<i>Name method and provide relevant details, e.g. missing indicators, complete case, last value carried forward, multiple imputation (specify model/variables), other.</i>
	As in Primary Estimand section 7.6.2
Subgroup Analyses	<i>List all subgroups</i>
	N/A

7.6.6 Core Emulation Table – Estimation Summary

Table 14(A). Estimation Summary for Primary Estimand 1

	Target Trial	Target Trial Emulation	Comment

Analysis Method	RR is parameterized with the Incidence Rate Ratio (IRR). The IRR is estimated using a Poisson regression model, in which the occurrence of the outcome is modelled conditional on treated vs placebo group membership, with an offset term equal to the log person time at risk. Vaccine Efficacy = $(1-IRR) \times 100$.	RR is parameterized with the Incidence Rate Ratio (IRR). The IRR is estimated using a Poisson regression model, in which the occurrence of the outcome is predicted by treated vs control group membership, with an offset term equal to the log person time at risk. Vaccine Effectiveness = $(1-IRR) \times 100$.	Identical in the Target Trial and Emulation, with the exception that the exposed and control group members are constructed based on a matching scheme for covariate adjustment. As a consequence the distribution of covariates at baseline in the matched controls is equal to that of the (matched) exposed.
Missing Data Assumptions and Methods to Handle	For intercurrent events handled using the hypothetical strategy, that is with non-administrative censoring, individuals have missing time at risk. This assumes that missing time at risk (i.e., time post-censoring) is missing at random conditional on treatment assignment and on not being infected up to the censoring time.	For intercurrent events handled using the hypothetical strategy, that is with non-administrative censoring, individuals have missing time at risk. This assumes that missing time at risk (i.e., time post-censoring) is missing at random conditional on treatment assignment and on not being infected up to the censoring time. Missing covariate values are handled using either complete case, multiple imputation or variable omission strategies depending on the degree of missingness	Identical in Target Trial and Emulation for censoring. In the trial, information on any important covariates is measured directly. In the emulation we both need to identify information on more covariates (for the confounding adjustment model) and to identify this information using codes. Lifestyle characteristics such as weight/BMI and smoking behaviour may not be measured when these values are in the normal range, or the individual does not present with relevant medical concerns.
Statistical Model Assumptions	<ul style="list-style-type: none"> - Poisson model assumes constant hazard rate over time. - Non-informative censoring. - Correct model specification (i.e., 	<ul style="list-style-type: none"> - Poisson model assumes constant hazard rate over time. - Non-informative censoring. 	Identical in Target Trial and Emulation.

	correct functional form and inclusion of relevant covariates).	- Correct model specification.	
Sensitivity Analyses	A sensitivity analysis under the censoring not at random assumption will be conducted to assess the sensitivity of our risk estimates to the censoring at random (conditional on treatment assignment) assumption.	<p>A sensitivity analysis will be conducted to assess the sensitivity of our risk estimates to censoring at random (conditional on treatment assignment) assumption using Inverse Probability of Censoring Weights (IPCW) under a censoring at random assumption conditioning on treatment assignment and additional covariates.</p> <p>In addition, four sensitivity analyses under the censoring not at random assumption will be run. These analyses represent a best-case and a worst-case scenario where all non-administratively censored subjects within each arm are respectively assumed to not experience SARS-CoV-19 infection by the end of the study (best case scenario) or to experience the SARS-CoV-19 infection at the time of censoring (worst case scenario).</p>	Identical in Target Trial and Emulation for censoring. Sensitivity analyses under the best/worse case scenarios are necessary in the emulation setting, as we would assume such information is available to the study team directly (SARS-CoV-19 infections) in the trial.

Table 14(B). Estimation Summary for Primary Estimand 2

	Target Trial	Target Trial Emulation	Comment
Analysis Method	RR is parameterized with the Incidence Rate Ratio (IRR). The IRR is estimated using a Poisson regression model, in which the occurrence of the outcome is predicted by treated vs placebo group membership, with an	RR is parameterized with the Incidence Rate Ratio (IRR). The IRR is estimated using a weighted Poisson regression model, in which the occurrence of the outcome is predicted by treated vs control group membership, with an	Poisson regression with IPSM weighting amongst the matched cohort is used in the emulation to adjust for baseline confounding and model membership of stratum (a). To account for the uncertainty introduced

	<p>offset term equal to the log person time at risk. Observations are weighted using Inverse Probability of Stratum Membership (IPSM) weights. Analysis is restricted to individuals who receive both doses of their assigned treatment. Vaccine Efficacy = $(1-IRR) \times 100$.</p>	<p>offset term equal to the log person time at risk.</p> <p>Observations are weighted using Inverse Probability of Stratum Membership (IPSM) weights. Analysis is restricted to individuals who receive both doses of the treatment in the vaccinated group. Vaccine Effectiveness = $(1-IRR) \times 100$.</p>	<p>by the estimation of inverse probability weights, outcome models will use robust (sandwich) variance estimators.</p>
<p>Missing Data Assumptions and Methods to Handle</p>	<p>For intercurrent events handled using the hypothetical strategy, that is with non-administrative censoring, individuals have missing time at risk. This assumes that missing time at risk (i.e., time post-censoring) is missing at random conditional on treatment assignment and on not being infected up to the censoring time.</p>	<p>For intercurrent events handled using the hypothetical strategy, that is with non-administrative censoring, individuals have missing time at risk. This assumes that missing time at risk (i.e., time post-censoring) is missing at random conditional on treatment assignment and on not being infected up to the censoring time.</p> <p>Missing covariate values are handled using either complete case, dummy coding, or variable omission strategies depending on the level of missingness</p>	<p>Identical in Target Trial and Emulation for censoring.</p> <p>In the trial, information on any important covariates is measured directly. In the emulation we both need to identify information on more covariates (for the confounding adjustment model) and to identify this information using codes. Lifestyle characteristics such as weight/BMI and smoking behaviour may not be measured when these values are in the normal range, or the individual does not present with relevant medical concerns.</p>
<p>Statistical Model Assumptions</p>	<ul style="list-style-type: none"> - Poisson model assumes constant hazard rate over time. - Non-informative censoring. - Correct model specification (i.e., correct functional form and inclusion of relevant covariates). 	<ul style="list-style-type: none"> - Poisson model assumes constant hazard rate over time. - Non-informative censoring. - Correct model specification. - Exchangeability assumed conditional 	<p>Identical in Target Trial and Emulation, though exchangeability is achieved through randomization in the trial and through adjustment (matching and weighting) in the emulation.</p>

	- Exchangeability between treatment groups due to randomization.	on covariates used in the matching and principal score model.	
Sensitivity Analyses	<p>A sensitivity analysis under the censoring not at random assumption will be conducted to assess the sensitivity of our risk estimates to the censoring at random (conditional on treatment assignment) assumption.</p> <p>A sensitivity analysis relaxing the assumption of monotonicity will be conducted, modelling the probability of stratum (a) membership amongst vaccinated as well as placebo rather than assuming stratum (a) membership to be observed in the vaccinated.</p>	<p>A sensitivity analysis will be conducted to assess the sensitivity of our risk estimates to censoring at random (conditional on treatment assignment) assumption using Inverse Probability of Censoring Weights (IPCW) under a censoring at random assumption conditioning on treatment assignment and additional covariates.</p> <p>In addition, four sensitivity analyses under the censoring not at random assumption will be run. These analyses represent a best-case and a worst-case scenario where all non-administratively censored subjects within each arm are respectively assumed to not experience SARS-CoV-19 infection by the end of the study (best case scenario) or to experience the SARS-CoV-19 infection at the time of censoring (worst case scenario).</p>	<p>Identical in Target Trial and Emulation for censoring. Sensitivity analyses under the best/worse case scenarios are necessary in the emulation setting, as we would assume such information is available to the study team directly (SARS-CoV-19 infections) in the trial. No sensitivity analysis for the monotonicity assumption is used in the emulation: Due to the change in the nature of the estimand enforced by the absence of a placebo group, no monotonicity assumption is needed to identify the principal stratum effect.</p>

7.7. Data sources

7.7.1 Data sources – Context and rationale

Reason for selection / Rationale for selection and feasibility:

The Clinical Practice Research Datalink (CPRD) and The Valencia health system integrated database (VID) were selected as data sources for this study because they offer large, high-quality, population-based electronic health records with national coverage in the UK and Spain, respectively. Both sources provide the required data elements to operationalize the study design, including demographics, diagnoses, prescriptions, laboratory test results, hospitalizations, and mortality data. In addition, both databases have previously been used in studies of vaccines effectiveness and safety.

Strengths of data source(s):

CPRD covers over 19 million patients in total, with 16.5 million current acceptable patients in Aurum and nearly 3 million in GOLD. The data are updated frequently (monthly for GOLD, quarterly for Aurum), and mortality, prescribing, and diagnosis information are available with high completeness. CPRD includes validated linkages to secondary care, cancer registries, and national death registries. Both databases support validated outcome definitions and have been widely used in pharmacoepidemiological studies.

VID is a set of multiple, public, population-wide electronic databases for the Valencia Region, with ~5 million inhabitants. The VID provides exhaustive longitudinal information including sociodemographic and administrative data (sex, age, nationality, etc.), clinical (diagnoses, procedures, diagnostic tests, imaging, etc.), pharmaceutical (prescription, dispensation) and healthcare utilisation data from hospital care, emergency departments, specialised care, primary care and other public health services. It also includes a set of associated population databases and registries of significant care areas such as cancer, rare diseases, vaccines, congenital anomalies, microbiology and others, and public health databases from the population screening programmes.

All electronic health systems in the VID use the ICD-9-CM and the ICD-10-CM. All the information in the VID databases can be linked at the individual level through a single personal identification code. The databases were initiated at different moments in time, but all in all the VID provides comprehensive individual-level data fed by all the databases from 2009 to date.

Limitations of data source(s) (with potential impact on study results):

For CPRD, only the year of birth is available for adults, which may slightly impact precision in age-based eligibility criteria. Ethnicity is missing for approximately 20% of patients overall (higher in earlier years). Linked hospital and mortality data are available only for practices in England and may be delayed.

Data source quality

CPRD is managed by the Medicines and Healthcare products Regulatory Agency (MHRA) in the UK. Data from participating general practices are extracted monthly (GOLD) or quarterly (Aurum) and undergo multi-level validation and quality checks. The CPRD provides data through its secure Trusted Research Environment or via multi-study licenses. Full metadata and SOPs are available at <https://www.cprd.com/data-access> and <https://www.cprd.com/data-quality>

VID is managed by the Valencia regional government. Research teams at FISABIO, such as the HSRP Unit, can be considered as providers, as they are granted access to VID data on a project-basis, after 1. ethics committee approval of a research protocol and 2. data commission approval of the data extraction. More information can be found at

<https://www.san.gva.es/es/web/investigacio/solicitud-datos-sia-gaia> and <https://fisabio.san.gva.es/es/registro-de-actividades-de-tratamiento-de-datos/>

CPRD and VID was evaluated for its feasibility in supporting this study. CPRD database was found to be well-suited for this purpose, offering high-quality data on vaccination (including dose, type, and batch number), diagnoses, and mortality. CPRD includes a large and diverse population, with sufficient follow-up time and coverage to meet the estimated sample size of 44,000 participants. Limitations include the absence of direct data on informed consent though these are not expected to significantly impact study validity. VID provides comprehensive coverage of the Valencia region and includes detailed vaccination records (dose, manufacturer, batch number), diagnostic codes, and hospitalization data. It supports robust longitudinal analyses and has the population size needed to meet the target sample size. Limitations include delayed access due to regulatory approvals and missing inpatient medication data.

Further details on the feasibility assessment for data sources in this study can be found in appendix 2.

Table 15. Metadata about data sources and software

	Data 1	Data 2
Data Source(s):	CPRD Aurum (via UU)	VID (via FISABIO)
Study Period:	Expected: 2020-2022	Expected: 2020-2022
Eligible Cohort Entry Period:	Expected: 2020–2022	Expected: 2020–2022
Data Version (or date of last update):	CPRD GOLD: 06/2025 CPRD Aurum: 06/2025	31/12/24
Data sampling/extraction criteria:	tbd	tbd
Type(s) of data:	Primary care EHR, linked secondary care data (HES)prescriptions, lab values, demographics, outcomes	Primary and specialist care EHR, emergency and hospitalization care, prescriptions and dispensing, lab values, demographics, outcomes, microbiological surveillance network (REDMIVA)
Data linkage:	Yes, deterministic linkage with HES.	Yes, deterministic linkage.
Conversion to CDM*:	Yes, ConcePTION CDM v2.2	Yes, ConcePTION CDM v2.2
Software for data management:	EMIS Web for Aurum, Vision for GOLD; analyses via local TRE	R version 3.6.0.

*CDM = Common Data Model

*UU = Utrecht University

*FISABIO = Foundation for the Promotion of Health and Biomedical Research of Valencia Region

7.8. Data management

The study will be conducted in a distributed manner using the UMCU, ARS Toscana and VAC4EU tools, procedures, and pipeline. Figure 1 specifies the data sets (D) and transformation processes (T), programming follows this pipeline, with involvement of different types of experts.

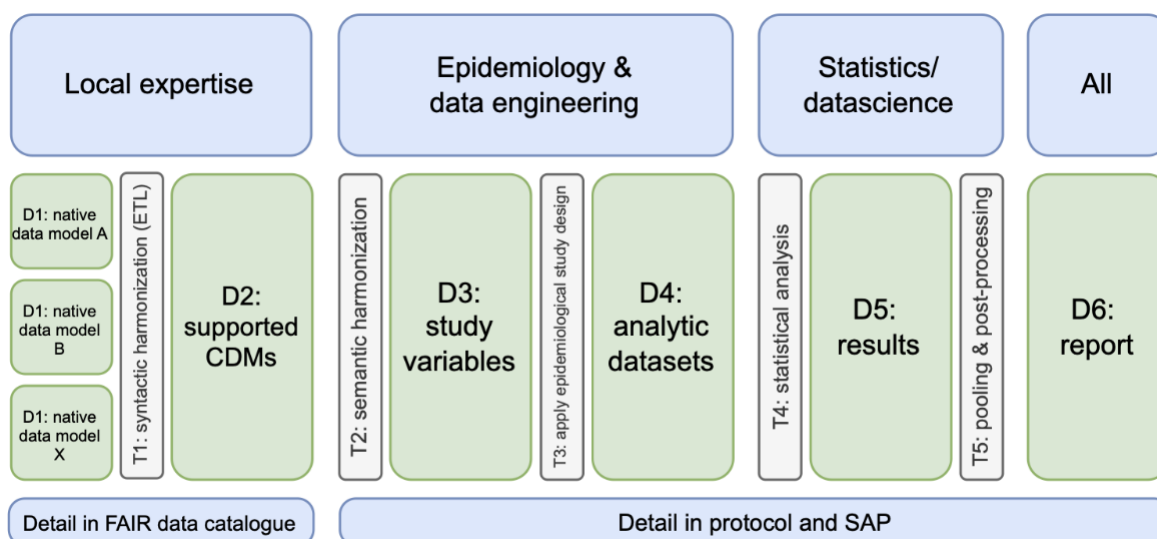


Figure 2. Data Management from the data transformation perspective

D1: Original data can be in any native format

The RWD-RWE pipeline used by VAC4EU starts with data banks that are controlled by the Data Expert and Access Partner (DEAP) and can be in any format. Data always stays local and never leaves the secure environments of the DEAPs. The ETL (extract, transform, load, see below for more details under 'T1') design is shared in a searchable FAIR VAC4EU catalogue. The VAC4EU FAIR Molgenis data catalogue is a meta-data management tool designed to contain searchable meta-data describing organisations that can provide access to specific data sources.

T1: Syntactic harmonisation (ETL)

T1: Syntactic harmonisation is conducted through an extraction, transformation, and loading (ETL) process of native data into the ConcePTION common data model (CDM) (see section 'D2: Common data model'). To harmonise the structure of the data sets stored and maintained by each data partner, a shared syntactic foundation is used. The ETL process has various structured steps as described by Thurin et al.¹³

- DEAPs are asked to share the data dictionaries of their data banks (selected tables and variable names/structure)
- Metadata (descriptive data about the data sources and databanks) & data dictionaries, are uploaded in FAIR data catalogue (Molgenis).

D2: Common data model

For this project, the CDM (D2) is the ConcePTION common data model. The CDM version that is used is v2.2, which is available as an open-source CDM. In this CDM, data are represented in a

common structure, but the values of the data remain in their original language (e.g. codes will have either ICD9/10/ICPC/SNOMED or MEDCODEID values).

T2: Semantic harmonisation

During the T2 step, many data transformations occur related to the completion of missing features in the data. Based on the relevant diagnostic medical codes and keywords, as well as other relevant concepts (e.g., medications), one or more phenotype algorithms are constructed (typically one sensitive, or broad, algorithm and one specific, or narrow, algorithm) to operationalise the identification and measurement of each event. In this step we conduct time anchoring (observation periods, look back periods), clean the data such as the dose of vaccines, sort on record level, aggregate across multiple records, and combine concepts for implantation of algorithms, and rule-based creation of study variables.

In this phase of the creation of study variables, semantic mapping is conducted. This semantic mapping across different vocabularies is conducted as part of the R-study script using different functionalities. To reconcile differences between different terminologies and native data availability, machine-readable code lists are used that comprise the terminologies that are used in the network (e.g. ICD-9, ICD10, SNOMED, ICPC and DEAP specific adaptations). This is combined with the BRIDGE metadata file that defines risk windows, look-back periods, and algorithms for each study variable. ¹⁴

D3: Study variables

D3 datasets are interim data sets with information on study variables for each study participant, the unit may be a person, a medicine, or episode of time. The design of these datasets is described in codebooks. Examples of D3 datasets are the outputs of the ConcePTION pregnancy algorithm, and outputs of functions that define smoking. Multiple functions/packages have been developed for previous studies for different study variables.

T3: Application of epidemiological design

In the T3 step epidemiological designs are applied such as sampling, matching (on specific variables and/or propensity scores), and selection based on inclusion and exclusion criteria using the study variables in the D3 datasets. The designs will be implemented for the various study objectives using R-scripts, and these may use the existing functions (R-cran) or functions that have been developed in the VAC4EU community (e.g. matching).

D4: Analytical data set

D4 is an analytical dataset, and multiple D4 data sets may be produced based on the objectives of the study. The format is described initially in a code book for communication between programmers and statisticians.

T4: Statistical analysis

This step in the data transformation pipeline will produce statistical estimates such as descriptives (counts, percentages), distributions (mean, percentiles), rates (prevalence, incidence), regression coefficients, or other relevant estimates. This will be conducted using R.

D5: Results

D5 is the set of estimates, tables or aggregate data that is transferred from the DEAPs to the Digital Research Environment (DRE). The aggregated results produced by these scripts at the DEAP's site will be uploaded to the UMCU DRE for post-processing, pooling and visualisation (Figure 1). The DRE is a cloud-based, globally available research environment where data are stored and organised securely and where researchers can collaborate. The DRE is made available through UMCU. The DRE applies double authentication where researchers can collaborate using data that are stored and organised securely [ref]. UMCU is responsible for data processing and data security.

All researchers who need access to the DRE will be granted access to study-specific secure workspaces by UMCU. Access to the workspaces will be possible only after double authentication using an identification code and password together with the user's mobile phone for authentication.

Uploading files will be possible for all researchers with access to the workspace within the DRE. Downloading of files will be possible only after requesting and receiving permission from a workspace member with an "owner" role, who will be a UMCU team member.

T5: Post-processing/pooling

In this step, the result from different DEAPs is pooled and converted into tables and figures for reporting.

7.9. Quality control

All key study documents such as the hypothetical trial protocol, target trial emulation protocol and study reports will undergo senior scientific and editorial review.

Data quality

For all data sources and for each data instance we will conduct *INSIGHT* level 1-2 quality checks, detailed statistical analysis plans for the indicators are available on the public repositories:

- <https://github.com/UMC-Utrecht-RWE/INSIGHT-Level1> Hoxhaj, V. (2023). UMC-Utrecht-RWE/INSIGHT-Level1: <https://doi.org/10.5281/zenodo.10035167>
- <https://github.com/UMC-Utrecht-RWE/INSIGHT-Level2> Hoxhaj, V., & van den Bor, R. (2023). UMC-Utrecht-RWE/INSIGHT-Level2: <https://doi.org/10.5281/zenodo.10035169>

Briefly, level 1 verifies Data Completeness and level 2 Data Consistency.

Level 1 – Data Completeness

The purpose of the level 1 check is to verify the completeness of the ETL process and the data in the variables. Examples of tests are:

- Presence of variables in each of the CDM tables in D2
- Checks for misspellings and letter case in variable names in each of the CDM tables
- Verification of vocabularies
- Check date formats
- Check conventions of values
- Missing data analysis
- Frequency tables for categorical variables

Level 2 – Data Consistency

Real data is not random but follows certain logical constraints that reflect rules governing real-world situations. Examples of indicators generated by level 2 checks are:

- Event dates before date of birth
- Event dates after date of death
- Event dates out of observation periods
- Subjects having an observation but not present in the PERSONS table
- Observations associated with a visit id and occurred before/after the visit start/end date
- Subjects younger than 12 years old reported as parents
- Age at the observation period older than 115 y old Data

Code Quality

These coding practices define how the TARGET programming team collaborates to write clean, reliable, and reproducible code for the VAC4EU Real-World Evidence (RWE) Analytical Pipeline. They aim to ensure clarity, consistency, and maintainability across all case studies within the project.

Coding conventions

To ensure clarity, consistency, and maintainability across the project, the following conventions will be applied to all codebases within the project:

- Consistent style: Code follows a consistent and readable style (see the [tidyverse style guide](#) for R).
- Meaningful names: Use clear, descriptive names for variables, functions, and files to convey their purpose.
- Modular code: Break down code into small, reusable functions where possible.
- No hardcoded paths: Use configuration files or relative paths to ensure portability.

Following these conventions makes the code easier to understand, test, and reuse across case studies and teams.

Documenting Code

Code documentation is used to promote good coding practices and ensure our work is understandable, maintainable, and reproducible. To achieve this, we will:

- Use descriptive comments that explain the purpose and rationale behind code sections, focusing on why something is done, not just what.
- Clearly document function inputs, outputs, and side effects, using standardized formats (e.g., `roxygen2` in R) where appropriate and supported.
- Write meaningful variable and function names to make the code as self-explanatory as possible.

Version Control

We use Git and GitHub to manage version control. These tools support good coding practices by enabling collaboration, tracking changes, accessing a project's history, and ensuring code quality through review and documentation.

A dedicated GitHub organisation has been created for the project (<https://github.com/target-roc19>). Each case study is managed in its own repository within this organisation. Repositories are structured consistently across case studies, to reinforce modularity. Access to repositories is controlled through teams.

During development, all repositories remain private to ensure confidentiality. Once the project is finalised, relevant repositories will be made public and assigned a digital object identifier (DOI) via Zenodo to support transparency, reproducibility, and reuse by the wider research community.

To maintain code quality and clarity, we follow the git and GitHub guidelines below.

- Always use pull requests (PRs): never push directly to the main branch.
- Open an issue before creating a new branch. Ideally, one PR resolves one issue to keep changes focused and reviewable.
- Every PR must be reviewed by at least one other person before merging.
- The PR author merges the PR after it has been reviewed and approved.
- Write clear, descriptive commit messages.
- Write informative PR descriptions, including:
 - A concise title
 - Links to related issues
 - A summary of the changes

Continuous Integration

Continuous Integration (CI) is set up to automatically check code quality and run tests whenever changes are pushed to the repository or submitted through a pull request (PR). The CI workflow ensures that the package adheres to predefined style guidelines and that all automated tests pass before changes are merged.

Coding Template

Every case study follows the general coding template used across all code in the TARGET project. The folder structure is organised as follows:

```
case-study-template
|___data
| |___D2_cdm
| |___D3_study_variables
| |___D4_analytic_datasets
| |___D5_results
| |___D6_report
|___docs
|___logs
|___run
|___tests
|___transformations
| |___T2_semantic_harmonization
| |___T3_study_design
| |___T4_statistical_analysis
| |___T5_processing_results
|___CHANGELOG.md
|___LICENSE
|___README.md
```

Project Data Structure and Storage

The data folder follows the Real-World Evidence pipeline structure. Data conforming to the common data model is stored in the D2_cdm folder.

Results from transformations T2, T3, T4, and T5 are saved in the respective folders:

- D3_study_variables
- D4_analytic_datasets
- D5_results
- D6_report

Each dataset is associated with a codebook, explained in more detail below.

All data remain securely stored on the Data Expert and Access Partners (DEAPs) servers and are never transferred externally. For testing purposes, dummy datasets are created. These fall into two categories:

- Unit test data: Small, predefined input and output pairs used to test individual transformation steps. These are stored in the tests folder, not in data, and can support automated testing.
- Pipeline test data: Larger, more complex dummy datasets used to test whether the full pipeline runs as expected. These may be included in the repository only if they remain below GitHub's 100 MiB file size limit and will otherwise be shared via SharePoint.

Logging System

When the pipeline is executed, log files are saved in the logs folder. These logs are especially helpful when running the code in the DEAPs environment, as they help trace and diagnose potential errors. We recommend using the logger R package to handle logging throughout the pipeline. A sample logging setup can be found in the logger.R script located at the root of the project directory.

Executing the Analytical Pipeline

The run folder contains scripts used to execute each transformation step in the pipeline.

- A central script, run_pipeline.R, orchestrates the full pipeline from start to finish.
- Subscripts (e.g., run_T2.R or similar) are available to run individual transformation steps separately.

Typically, the run_pipeline.R script is the main entry point used by a DEAP to execute the full pipeline. Before running it in the DEAP environment, the pipeline may need to be adapted to local settings. This can be done using a configuration file that defines variables required to tailor the pipeline to a specific DEAP. Please note that configuration files should not include sensitive information.

Such a file might include variables like:

- The name of the DEAP
- The path to the local data instance
- The path to any required external resources

Testing and Quality Assurance

The tests folder contains scripts to test the analytical pipeline. Tests will be used to ensure code behaves as expected and remains stable over time. By systematically checking inputs, outputs, and edge cases, tests help catch errors early and make future changes safer. We use the `testthat` R package to structure and run unit tests.

Continuous integration (CI) is used to automate testing. With CI, tests are automatically run each time code is pushed to the repository (e.g., via GitHub Actions). This helps identify issues immediately, ensures that new changes do not break existing functionality, and supports better collaboration by enforcing consistent code quality across contributors.

Modular Data Transformation Workflow

The transformations folder follows the Real-World Evidence pipeline structure. It contains the source code for all transformation steps, which is typically written in R. Each subfolder corresponds to a specific step in the pipeline (e.g., `T2_semantic_harmonization`, `T3_study_design`, `T4_statistical_analysis`, `T5_processing_results`) and includes the relevant scripts and helper functions for that step.

During the T2 step, a database is usually created (e.g., using DuckDB). This database can be queried using SQL, and it is recommended that all SQL queries be saved as clearly named, standalone SQL script files to ensure readability and reusability.

The purpose of the transformations folder is to structure and modularise the processing logic, making it easier to maintain, test, and reuse across different case studies. By organising code by transformation step, teams can work in parallel, increasing efficiency.

Changelog

A changelog will be kept for all notable changes in the project. Changelogs help track the evolution of the project over time, making it easier for collaborators to understand what has changed between versions. We follow the structure and best practices outlined in [Keep a Changelog](#).

Codebooks

Before developing code, codebooks are created to describe each dataset (D) within the pipeline. A codebook is a comprehensive document that outlines the structure, contents, and metadata of a dataset. It serves as a detailed reference guide for anyone working with the data and plays a crucial role in guiding the development of the analytical pipeline by clearly defining both the inputs and expected outputs.

All codebooks are summarized in a central index file, which provides a high-level overview of the pipeline's structure. For each codebook, the index file includes:

- A brief description of its purpose,
- A list of the scripts used to generate the corresponding dataset,
- A description of the input datasets and input parameters required.

The datasets D2, D3, D4, and D5 are typically subdivided into multiple smaller transformation steps, each detailed within their respective codebooks. These smaller transformation steps ensure that each part of the pipeline is clearly scoped and well-documented.

In addition to supporting development, codebooks help ensure quality control by making transformation logic transparent and verifiable, and they enhance reproducibility by documenting exactly how data is structured and used throughout the analytical pipeline.

Deployment

The analytical pipeline is delivered to DEAPs as a GitHub release, tagged with a version number. Versioning follows the format: YYYYMMDD.XX, where the date indicates the release date and XX denotes the sub-version or revision number.

Any deployment issues can be reported via the GitHub repository using the issues feature, where the programming team responsible for the R code will collaborate with the local DEAP to resolve them as needed.

Reproducibility

It is recommended to locally use the `renv` R package to maintain the R version and version of packages for reproducibility purposes.

At this time, however, using `renv` reliably across different systems and environments remains challenging. For this reason, we currently recommend its use only in local development setups.

We are actively monitoring developments in the R ecosystem related to cross-platform reproducibility. As soon as a more stable and portable solution becomes available, we will revisit this guidance and promote broader adoption.

Open-Source Licensing

The code will be made available under an open source license.

README Guidelines

Each case study repository includes a README that covers the following points:

- **Project Overview:** brief summary of the study goals and key research questions.
- **Background:** context and rationale for the study.
- **Repository Structure:** Outline of main folders and their contents.
- **Data Overview:** Description of data sources, formats, and data privacy considerations.
- **How to Run:** Instructions for running the pipeline and key scripts, plus where outputs are saved.
- **Testing:** How to run tests to verify code functionality.
- **Contributing:** Guidelines for code contributions and issue tracking.
- **License:** Information about the code license.
- **Contact:** Who to reach out to for help or questions.

7.10. Study size and feasibility

We computed target sample sizes analytically for a two-arm Poisson GLM with robust standard errors with sandwich estimator. Precision was defined as the relative precision of the vaccine efficacy (VE): the ratio of the upper 95% confidence interval (CI) bound for VE to the assumed VE. For a target relative precision p (e.g. 30%, $p = 0.30$), we solve for the required standard error of $\log(\text{IRR})$ and then the required expected number of events. IPTW weights (lognormal-like) inflate the variance by $R_w = 1 + (\text{weight_sd}/\text{weight_mean})^2$, so required events are multiplied by R_w . Person-time and sample size follow from expected events and the control incidence rate; final N is inflated to account for loss to follow-up.

Given:

r = IRR (vaccine / control)
 p = precision_target (decimal, e.g. 0.30 for 30%)
 z = 1.96
 $R_w = 1 + (\text{weight_sd} / \text{weight_mean})^2$
 λ_c = control incidence (events per person-year)
 f = mean follow-up (years)
loss = loss proportion (e.g. 0.20)

1) Compute required SE for $\log(\text{IRR})$:

$\text{rhs} = (1 + p) - p / r$
 $\text{SE}_{\text{req}} = - (1 / z) * \log(\text{rhs})$

2) Required expected events in control (unweighted):

$\text{EO}_{\text{unw}} = (1 + 1/r) / (\text{SE}_{\text{req}}^2)$

3) Account for weighting (IPTW variance inflation):

$\text{EO} = R_w * \text{EO}_{\text{unw}}$
 $\text{E1} = r * \text{EO}$
 $\text{Total_events} = \text{EO} + \text{E1}$

4) Convert to person-time (equal allocation assumption):

$T_{\text{total}} = 2 * \text{EO} / \lambda_c$ # total person-years

5) Convert to number of participants:

$N_{\text{no_loss}} = T_{\text{total}} / f$
 $N_{\text{adj}} = \text{ceiling}(N_{\text{no_loss}} / (1 - \text{loss}))$

For the precision-based sample size calculation the following assumptions were made: VE = (60%, 70%, 80%, 90%), an incidence rate in the controls $\lambda_c = 0.013$ (Polack et al., 2020 used a higher rate of 0.017), mean follow-up = 3 months, 0.55 and 2% loss to follow up, no clustering and propensity score weights (MEAN = 1, SD = 0.20, inflated variance by $R_w = 1.04$).

Weight distribution justification: based on previous Post Authorization Safety Studies using CPRD we know the distribution of stabilized propensity score weights can be approximated with a lognormal distribution with MEAN = 1.0 and SD = 0.18. For this sample size calculation, we decided to use a 10% safety margin and assume a SD = 0.2, leading to a variance inflation factor of $R_w = 1.04$.

Loss to follow up justification: Post Authorization Safety Studies we expect a low loss to follow up due to participants leaving the data source of about 0.5%. In addition, a 2% loss to follow up has been adopted as a safety margin.

Based on these assumptions and values for target relative precision targets of 30%, 10% and 5% we obtained the following expected sample sizes:

Table 16. Precision-based sample size estimates for vaccine efficacy (VE)

Assumed VE	Required SE (logIRR)	Total events	Events (control)	Events (vaccine)	Sample size		VE 95% CI lower	VE 95% CI upper	VE precision target ^a (%)
					0.5% loss	2% loss			
60	0.04	3221	2301	921	1422854	1444633	56.75	63.00	5
60	0.08	742	530	212	327427	332438	52.94	66.00	10
60	0.13	302	216	87	133110	135147	48.39	69.00	15
70	0.06	1463	1125	338	695766	706416	66.04	73.50	5
70	0.14	319	246	74	151664	153986	60.87	77.00	10
70	0.22	122	94	28	57698	58581	53.85	80.50	15
80	0.11	578	482	97	297739	302297	75.00	84.00	5
80	0.26	111	92	19	56815	57685	66.67	88.00	10
80	0.47	35	29	6	17658	17929	50.00	92.00	15
90	0.30	136	123	13	76047	77211	81.82	94.50	5
90	1.17	10	9	1	5127	5205	0.00	99.00	10
90	-	-	-	-	-	-	-	-	15 ^b

Notes: ^aVE precision target (%) = 100 × (VE upper limit / assumed VE - 1). For example, a 5% precision target corresponds to a precision ratio of 1.05.

^b The sample size could not be calculated for VE = 90% with a 15% precision target because the specified precision is not achievable with these parameter values (i.e., the required standard error is undefined).

These values assume no clustering and equal follow-up; clustering or informative/heavy-tailed weights would require simulation and typically increase sample-size requirements. However, the assumptions made for the VE (60%) and the control infection rate lambda_c (1.3%) are very conservative. The primary trial of BNT162b2 vaccine reported VE = 95.5% (90.3%–97.6%) in participants without prior infection and a control infection rate lambda_c of 1.7% (Polack et al., 2020).

8. Limitation of the methods

A limitation of the approach of comparing the original RCT with the TTE RWD study is the potential that the TTE cannot fully emulate all aspects of the original RCT, e.g. due to incomplete information on variables to assess in-/exclusion criteria, limited validity of outcomes, incomplete or flexible emulation of treatment strategies. These limitations are inherent to the TTE approach, and an assessment will be made regarding the impact of these limitations on potential differences between RCT estimates and TTE RWD estimates. This will help inform where potential data failings are and improve fit-for-purpose assessment on future emulations.

Exposure Misclassification: In this study, exposure is defined based on administration records, which do not capture whether patients actually filled or took their medications. As a result, some individuals categorized as “exposed” may not have received the intended pharmacological treatment. This limitation is acknowledged as a source of non-differential exposure misclassification, which would bias effect estimates toward the null. However, control units may also be misclassified in the sense that they may have received vaccinations abroad or in regions where they are not resident, which may not be captured in the data source.

Violation of the assumption of constant incidence rate: The Poisson model assumes that the relative treatment effect is constant over time, which may not be appropriate. There is no sensitivity analysis specified to relax this assumption.

Positivity Violations: Some subgroups of patients may be very unlikely to receive one of the treatments, leading to limited overlap in covariate distributions between groups. The propensity score model is examined for evidence of positivity violations by assessing the overlap of propensity scores between treatment groups. Patients in non-overlapping regions are excluded, and stabilized inverse probability of treatment weights (IPTW) are truncated at the 1st and 99th percentiles to limit the influence of extreme weights. These steps reduce bias and variance due to poor overlap while preserving generalizability within the area of clinical equipoise.

Parallel Arms and Re-sampling individuals: The primary analysis may be statistically inefficient, as individuals become ineligible to serve as exposed units at later time points, and ineligible to serve as controls on other calendar dates or to similar individuals. This may lead to a change in distribution of covariates among the selected vaccinated relative to the population of all eligible individuals who receive vaccination. It may also lead to positivity violations, in particular for exposed individuals who receive their vaccination later in calendar time during the study period, as few controls may be available to match. This potential limitation is assessed using the supplementary analysis in which individuals may serve as controls multiple times, or as vaccinated later in the study period.

Need for modelling of baseline conditional exchangeability: The lack of random assignment in the emulation means that the confounding structure must be modelled in order to ensure (conditional) exchangeability. In the current study we do this using a matching approach. This comes with the disadvantage of introducing additional assumptions, such as conditional exchangeability given observed covariates, and positivity (in the sense that all exposed should be matched to a control with the exact same profile of covariates). Another perhaps under-appreciated disadvantage is that the choice of any confounding adjustment method changes the definition of the estimand. The estimand in a randomized trial may be interpreted as a local-average treatment effect (LATE), that is, the average treatment effect amongst the subset of the population who meets the eligibility criteria and enrolls in a randomized trial. In that setting, there is no need to distinguish between the average treatment effect amongst the treated (ATT), the

untreated (ATU), or the target population as a whole (ATE), as random assignment ensures that all three are equal. Implicitly, the estimand in the randomized trial may be stated as the ATE, with the assumption that the LATE is equivalent to the ATE. In an observational setting, we must choose between the ATE, ATU and ATT (amongst others) as target estimands, and we should be aware that different methods of adjusting for confounding may yield estimates of subtly different estimands. For instance, using the matching scheme described in fact yields an estimate of the ATT. On the one hand this may not necessarily reflect the treatment effect in the population as a whole. However, on the other hand, in the current context the ATT it is a more pragmatic choice than the alternatives, and more likely that this treatment effect can be estimated with greater precision than the ATE or ATU. First, at any given moment in time, all individuals in a country were eligible to receive the target vaccine, and so from a purely computational standpoint estimating the ATE would be challenging. Second, those individuals who remain untreated with the target vaccine are less likely to engage with healthcare services, and have high-quality data available in the data source, and so any analytic strategy which aims to estimate the treatment effect amongst a population that includes the unvaccinated may be prone to bias. Finally, if in the end almost all individuals in the population are treated, or if the treatment sufficiently covers the population such that the distribution of baseline covariates among the treated is approximately equal to the distribution of baseline covariates in the eligible population as a whole, then we can consider the matching procedure described above to estimate the ATE.

Precision of timing for second dose in real world data sources: In the target trial we use cut-off of 21 days after the index date for individuals to receive the second dose of the target vaccine, as this reflects the intended gap between first and second dose of the treatment. However, in real world data sources, it is unlikely that all individuals who receive and intend to receive their second dose schedule an appointment exactly 21 days after receiving the first dose. As such, this strict cut-off applied to real world data may fail to estimate the effect of full adherence to the treatment strategy in real world settings, where the second dose is likely received somewhere between 19 and 42 days after the first dose. In the emulation we consider all those who receive a vaccination in the 90 days following the first dose to have received the second dose as scheduled. In particular in the case of the second estimand, this may produce differences in the estimand targeted in the trial and the emulation, if, for instance, those who receive the second dose later in this admissible period are systematically different (or receive a systematically different level of protection against the virus) than those who received the second dose close to the intended 21 day period.

Informative censoring: In the context of this study, if patients who discontinue their treatment are systematically different in terms of their risk to contract SARS-CoV-19 (e.g., sicker patients transferring out of practice or dying without timely data recording), this could bias the estimated treatment effect. To account for the potential of informative censoring, the study implements inverse probability of censoring weighting (IPCW). IPCW estimates the probability of remaining uncensored over time conditional on baseline covariates. These probabilities are then used to reweight individuals in the outcome model, thereby creating a pseudo-population where censoring is independent of the outcome, conditional on the included covariates.

Data Source Heterogeneity: CPRD and VID differ in population coverage, healthcare systems, coding practices, and linkage availability, which may introduce heterogeneity in effect estimates. Analyses will be performed separately within each data source using harmonized definitions under the Conception CDM framework. If pooled estimates are produced, heterogeneity will be assessed using meta-analytic methods (e.g., I^2 statistics).

9. Protection of human subjects

This is a non-interventional study using secondary data collection and does not pose any risks for individuals. Each data source research partner will apply for an independent ethics committee review according to local regulations. Data protection and privacy regulations will be observed in collecting, forwarding, processing, and storing data from study participants.

Patient information

This study involves data that exists in an anonymized structured format and contains no patient personal information. All parties will comply with all applicable laws, including laws regarding the implementation of organisational and technical measures to ensure the protection of patient personal data. Such measures will include omitting patient names or other directly identifiable data in any reports, publications, or other disclosures, except where required by applicable laws. Patient personal data will be stored at DEAPs in encrypted electronic form and will be password protected to ensure that only authorised study staff have access. DEAPs will implement appropriate technical and organisational measures to ensure that personal data can be recovered in the event of a disaster. In the event of a potential personal data breach, DEAPs shall be responsible for determining whether a personal data breach has in fact occurred and, if so, providing breach notifications as required by law.

Patient consent

As this study does not involve data subject to privacy laws according to applicable legal requirements, obtaining informed consent from individuals is not required.

10. Reporting of adverse events

For studies in which the research team uses only data from automated healthcare databases, according to the International Society for Pharmacoepidemiology Guidelines for Good Pharmacovigilance Practices: *“Aggregate analysis of database studies can identify an unexpected increase in risk associated with a particular exposure. Such studies may be reportable as study reports but typically do not require reporting of individual cases. Moreover, access to automated databases does not confer a special obligation to assess and/or report any individual events contained in the databases. Formal studies conducted using these databases should adhere to these guidelines.”* For non-interventional study designs that are based on secondary use of data, such as studies based on medical chart reviews or electronic health records, systematic reviews, or meta-analyses, reporting of adverse events/adverse drug reactions is not required. Reports of adverse events/adverse drug reactions should only be summarized in the study report, where applicable. According to the EMA Guideline on good pharmacovigilance practices, Module VI – Management and Reporting of Adverse Reactions to Medicinal Products, *“All adverse events/reactions collected as part of [non-interventional post-authorization studies with a design based on secondary use of data], the submission of suspected adverse reactions in the form of [individual case safety reports] is not required. All adverse events/reactions collected for the study should be recorded and summarized in the interim safety analysis and in the final study report.”* Module VIII – Post-Authorization Safety Studies echoes this approach. Legislation in the EU further states that for certain study designs such as retrospective cohort studies, particularly those involving electronic health records, it may not be feasible to make a causality assessment at the individual case level.

11. References

1. European Medicines Agency (EMA). Comirnaty EPAR. December 21, 2020. Accessed July 31, 2025. <https://www.ema.europa.eu/en/medicines/human/EPAR/comirnaty>
2. Polack FP, Thomas SJ, Kitchin N, et al. Safety and Efficacy of the BNT162b2 mRNA Covid-19 Vaccine. *New England Journal of Medicine*. 2020;383(27):2603-2615. doi:10.1056/nejmoa2034577
3. Dagan N, Barda N, Kepten E, et al. BNT162b2 mRNA Covid-19 Vaccine in a Nationwide Mass Vaccination Setting. *New England Journal of Medicine*. 2021;384(15):1412-1423. doi:10.1056/nejmoa2101765
4. Piechotta V, Siemens W, Thielemann I, et al. Safety and effectiveness of vaccines against COVID-19 in children aged 5–11 years: a systematic review and meta-analysis. *Lancet Child Adolesc Health*. 2023;7(6):379-391. doi:10.1016/S2352-4642(23)00078-0
5. Kow CSC, Ramachandram DS, Hasan SS. The effectiveness of BNT162b2 mRNA vaccine against COVID-19 caused by Delta variant of SARS-CoV-2: a systematic review and meta-analysis. *Inflammopharmacology*. 2022;30(1):149-157. doi:10.1007/s10787-021-00915-7
6. Sartorão-Filho CI, Zoqui MC, Duarte DO, et al. Prediction and reasons for COVID-19 second dose vaccine hesitation: a cross-sectional study in a municipality of Brazil. *Sao Paulo Medical Journal*. 2023;141(3). doi:10.1590/1516-3180.2022.0095.R1.06072022
7. Schneeweiss S, Rassen JA, Brown JS, et al. Graphical depiction of longitudinal study designs in health care databases. *Ann Intern Med*. 2019;170(6):398-406. doi:10.7326/M18-3079
8. Menegale F, Manica M, Zardini A, et al. Evaluation of Waning of SARS-CoV-2 Vaccine-Induced Immunity: A Systematic Review and Meta-analysis. *JAMA Netw Open*. 2023;6(5):E2310650. doi:10.1001/jamanetworkopen.2023.10650
9. Centers for Disease Control and Prevention. Underlying Medical Conditions Associated With Higher Risk for Severe COVID-19: Information for Healthcare Professionals. May 15, 2025. Accessed October 21, 2025. <https://www.cdc.gov/covid/hcp/clinical-care/underlying-conditions.html>
10. Glasheen WP, Cordier T, Gumpina R, Haugh G, Davis J, Renda A. *Charlson Comorbidity Index: ICD-9 Update and ICD-10 Translation*. Vol 12. www.ICD10Data.com,
11. Charlson ME, Pompei P, Ales KL, Mackenzie CR. *A NEW METHOD OF CLASSIFYING PROGNOSTIC COMORBIDITY IN LONGITUDINAL STUDIES: DEVELOPMENT AND VALIDATION*. Vol 40.; 1987.
12. Jakobsen JC, Gluud C, Wetterslev J, Winkel P. When and how should multiple imputation be used for handling missing data in randomised clinical trials - A practical guide with flowcharts. *BMC Med Res Methodol*. 2017;17(1). doi:10.1186/s12874-017-0442-1
13. Thurin NH, Pajouheshnia R, Roberto G, et al. From Inception to ConcePTION: Genesis of a Network to Support Better Monitoring and Communication of Medication Safety During Pregnancy and Breastfeeding. *Clin Pharmacol Ther*. 2022;111(1):321-331. doi:10.1002/cpt.2476
14. Cid Royo A, Elbers JHJ R, Weibel D, et al. Real-World Evidence BRIDGE: A Tool to Connect Protocol With Code Programming. *Pharmacoepidemiol Drug Saf*. 2024;33(12). doi:10.1002/pds.70062

12. Appendices

Appendix 1. List of pre-defined pre-existing stable medical conditions.

Appendix 2. Feasibility assessment

Appendix 3. Clinical codes

Appendix 4. Table Shells

Appendix 5. List of acute events prompting hospitalization.

Appendix 5. List of acute events prompting hospitalization.

1. Cardiovascular

- Acute myocardial infarction (heart attack)
- Unstable angina
- Acute decompensated heart failure
- Stroke / cerebrovascular accident (ischemic or hemorrhagic)
- Arrhythmias requiring inpatient management (e.g., atrial fibrillation with rapid ventricular response)
- Hypertensive crisis

2. Respiratory

- Pneumonia (bacterial, viral, including COVID-19)
- Acute exacerbation of COPD
- Severe asthma attack
- Pulmonary embolism
- Acute respiratory distress syndrome (ARDS)

3. Infectious / Sepsis

- Sepsis / septic shock
- Severe urinary tract infection (e.g., pyelonephritis)
- Meningitis / encephalitis
- Acute gastroenteritis with dehydration
- Influenza or other acute viral infections requiring hospitalization

4. Gastrointestinal / Hepatic

- Acute appendicitis
- Gastrointestinal bleeding (upper or lower)
- Acute pancreatitis
- Bowel obstruction
- Acute liver failure or acute-on-chronic liver failure
- Severe cholecystitis

5. Renal / Metabolic

- Acute kidney injury
- Hyperosmolar hyperglycemic state
- Diabetic ketoacidosis
- Severe electrolyte disturbances (e.g., hyponatremia, hyperkalemia)

6. Neurological

- Seizure requiring inpatient care
- Acute traumatic brain injury
- Acute neurologic deficits (e.g., Guillain-Barré, transverse myelitis)
- Intracranial hemorrhage

7. Trauma / Surgery

- Fractures requiring inpatient management (hip, femur, pelvis)
- Major burns
- Severe lacerations requiring surgical intervention
- Post-operative complications leading to re-hospitalization
- Acute surgical emergencies (e.g., perforated viscus)

8. Hematologic / Oncologic

- Severe anemia requiring transfusion
 - Acute leukemic crises
 - Thrombocytopenia with bleeding
 - Febrile neutropenia in oncology patients
9. Obstetric / Gynecologic
- Complicated labor or delivery requiring inpatient care
 - Ectopic pregnancy
 - Severe postpartum hemorrhage
10. Other Severe Acute Events
- Acute psychiatric emergencies requiring hospitalization (suicidal intent, psychosis)
 - Acute allergic reactions / anaphylaxis
 - Severe dehydration or malnutrition requiring IV therapy
 - Acute liver or kidney transplant complications