

<https://doi.org/10.1038/s43856-025-00956-x>

# A survey of pathogenic involvement in non-communicable human diseases



Michael Lape<sup>1,2,3</sup>, Daniel Schnell<sup>2</sup>, Sreeja Parameswaran<sup>3,4</sup>, Kevin Ernst<sup>3</sup>, Shannon O'Connor<sup>5</sup>, Nathan Salomonis<sup>1,2,6</sup>, Lisa J. Martin<sup>4,6</sup>, Brett M. Harnett<sup>1</sup>, Leah C. Kottyan<sup>3,4,6,7</sup> ✉ & Matthew T. Weirauch<sup>2,3,4,6,7,8</sup> ✉

## Abstract

**Background** Many relationships between pathogens and human disease are well-established. However, only a small fraction involve diseases considered non-communicable (NCDs). In this study, we sought to leverage the vast amount of newly available electronic health record data to identify potentially novel pathogen-NCD associations and find additional evidence supporting known associations.

**Methods** We leverage data from The UK Biobank and TriNetX to perform a systematic survey across 20 pathogens and 426 diseases, primarily NCDs. To this end, we assess the association between disease status and infection history proxies using a logistic regression-based statistical approach.

**Results** Our approach identifies 206 pathogen-disease pairs that replicate in both cohorts. We replicate many established relationships, including *Helicobacter pylori*, with several gastroenterological diseases and connections between Epstein-Barr virus and both multiple sclerosis and lupus. Overall, our approach identifies evidence of association for 15 pathogens and 96 distinct diseases, including a currently controversial link between human cytomegalovirus (CMV) and ulcerative colitis (UC). We validate the CMV-UC connection through two orthogonal analyses, revealing increased CMV gene expression in UC patients and enrichment for UC genetic risk signal near human genes that have altered expression upon CMV infection.

**Conclusions** Collectively, these results form a foundation for future investigations into mechanistic roles played by pathogens in the processes underlying NCDs. All results are easily accessible on our website, <https://tf.cchmc.org/pathogen-disease>.

## Plain Language Summary

Pathogens can hide in our bodies for years after initial infection. For example, the chickenpox virus can cause shingles decades after infection. Likewise, certain pathogens may contribute to the development or severity of subsequent conditions that are not contagious, referred to as non-communicable diseases (NCDs). Our study analyzed a vast number of electronic medical records to discover potential pathogen-NCD connections. We identified 206 such links. Future research based on these findings could revolutionize healthcare by enabling the development of vaccines targeted against these pathogens. For example, approaches similar to the introduction of the human papillomavirus vaccine which has led to declining cervical cancer rates. The development of vaccines for the pathogens identified in this study could potentially enable dramatic reductions in the NCDs linked to these pathogens.

Humans are exposed to infectious agents throughout life. Many of the communicable diseases associated with specific infectious agents are well-characterized<sup>1</sup>. In acute infection, virus-driven mechanisms are clearly and predominantly seen as the agent of disease. For example, respiratory syncytial virus (RSV) causes an upper-respiratory disease that can be life-threatening in young infants and causes cold-like symptoms in adolescents and adults<sup>2</sup>. Likewise, varicella zoster virus (VZV) causes varicella (chickenpox) upon primary infection. This infection typically occurs during

childhood and is well tolerated. However, if primary infection occurs in infants, adults, or the immunocompromised, the viral infection is less well contained, and the virally mediated pathology can be life-threatening<sup>3</sup>. Even after primary infection, VZV lies dormant for decades, reactivating in a portion of adults with the lytic virus directly leading to zoster (shingles)<sup>3</sup>.

The role of infectious agents in non-communicable diseases (NCDs) is much less well-explored, although several associations are well-known. For example, *Helicobacter pylori* infection is the strongest risk factor for gastric

<sup>1</sup>Department of Biomedical Informatics, University of Cincinnati College of Medicine, Cincinnati, OH, USA. <sup>2</sup>Division of Biomedical Informatics, Cincinnati Children's Hospital Medical Center, Cincinnati, OH, USA. <sup>3</sup>Center for Autoimmune Genomics and Etiology, Cincinnati Children's Hospital Medical Center, Cincinnati, OH, USA. <sup>4</sup>Division of Human Genetics, Cincinnati Children's Hospital Medical Center, Cincinnati, OH, USA. <sup>5</sup>Division of Rheumatology, Cincinnati Children's Hospital Medical Center, Cincinnati, OH, USA. <sup>6</sup>Department of Pediatrics, University of Cincinnati College of Medicine, Cincinnati, OH, USA. <sup>7</sup>Division of Allergy & Immunology, Cincinnati Children's Hospital Medical Center, Cincinnati, OH, USA. <sup>8</sup>Division of Developmental Biology, Cincinnati Children's Hospital Medical Center, Cincinnati, OH, USA. ✉e-mail: [Leah.Kottyan@cchmc.org](mailto:Leah.Kottyan@cchmc.org); [Matthew.Weirauch@cchmc.org](mailto:Matthew.Weirauch@cchmc.org)

cancers<sup>4</sup>. Likewise, multiple studies have demonstrated an important role for both hepatitis B (HBV) and C viruses (HCV) in the development of cirrhosis and other chronic liver diseases<sup>5,6</sup>. Further, certain strains of human papillomavirus (HPV) are known to cause a large proportion of cervical cancers and a smaller proportion of several other cancers<sup>7</sup>. Each of the aforementioned pathogens, as well as several others, have been classified as biological carcinogens by the International Agency for Research on Cancer (IARC)<sup>8</sup>. Indeed, it is estimated that up to 15% of all new cancer diagnoses are attributable to these and other infectious agents<sup>9</sup>. More recently, strong epidemiologic and mechanistic links have been identified between Epstein-Barr virus (EBV) infection and multiple sclerosis (MS)<sup>10–12</sup>, following decades of suggestive epidemiological and molecular evidence<sup>13,14</sup>. Many unknown pathogen-NCD connections likely remain yet to be discovered.

NCDs are often chronic diseases and have complex etiologies. Unlike acute viral infections, the virus itself might not contribute to the full etiology of virally associated NCDs. While many individuals might be exposed and mount an immune response to a virus, virally associated NCDs often occur only in a small subset of individuals in the context of specific genetic factors and other environmental exposures<sup>15</sup>.

Historically, connections between pathogens and diseases have been made one pair at a time. However, the recent establishment of large-scale national biobanks and the general shift to electronic health records (EHRs) enables the concurrent analysis of many pathogen-disease pairs. Here, we leverage biobank data from The UK Biobank (UKB) and EHR data from TriNetX, LLC (TNX), resources containing both serologic and diagnostic records, enabling the systematic detection of associations between multiple pathogens and multiple diseases simultaneously. Using a discovery-replication approach, our analysis reveals 206 replicated pathogen-disease associations, including previously established pathogen-disease links and many potentially novel or previously suggestive relationships. In particular, we identify strong evidence for the currently controversial connection between cytomegalovirus (CMV) infection and ulcerative colitis (UC). Orthogonal analyses of this relationship reveal that: (1) patients with UC have elevated CMV mRNA levels in intestinal tissue samples compared to healthy controls, and (2) UC genetic risk loci are enriched near human genes that change expression upon CMV infection in two independent datasets. Collectively, these results implicate multiple pathogens in dozens of non-communicable human diseases, providing a unique and powerful resource for future studies of the mechanistic roles played by these pathogens in disease development.

## Methods

### UK Biobank cohort

The UK Biobank (UKB) is a prospective cohort study containing medical, sociodemographic, and genetic data for nearly 500,000 adults from across the United Kingdom<sup>16</sup>. In this study, diagnoses for all participants were extracted from two UKB fields: (1) “First Occurrences” [UKB Category 1712], which is a synthetic field generated by UKB analysts that collates diagnoses from primary care, hospital inpatient, death registry, and self-reported records; and (2) “Type of Cancer – ICD10” [UKB Field 40006], which contains cancer diagnoses extracted from national cancer registries for linked participants. The first occurrences data are limited to 3-character International Classification of Diseases 10<sup>th</sup> revision (ICD10) codes, e.g., M32.9 is recorded as M32, so the cancer registry diagnoses were truncated to match. Analyses were limited to diseases with at least 17 cases and 187 total samples within the cohort (see power calculation below).

Data for 45 antibody titer levels representing immune responses to 20 unique pathogens [UKB Category 1307] were downloaded and Log10 transformed before analysis. Antibody titer measurements were performed using a multiplex serology approach based on the enzyme-linked immunosorbent assay (ELISA) concept as described by Mentzer et al.<sup>17</sup>. Similar to Mentzer et al., a series of 10 additional health-related and sociodemographic variables that could potentially be associated with both disease status and antibody titer level were collected to test for confounding during analyses.

The continuous covariates, age and body mass index (BMI), were scaled by a factor of 10, while all other covariates were either already categorical or were discretized (Supplementary Data 1). The scikit-learn IterativeImputer method, based on the Multivariate Imputation by Chained Equations (MICE) algorithm, was used to impute missing covariate values. Thirty-two participants were missing BMI values and eight were missing Townsend deprivation index values.

The UK Biobank data obtained, under application 47377, does not require separate clearance, instead operating under the approval granted to The UK Biobank by The North West Multi-centre Research Ethics Committee as a Research Tissue Bank. All participants in The UK Biobank provide informed consent for the collection of data directly by The UK Biobank, as well as consent for that data to be shared with researchers working on approved research projects.

### TriNetX cohort

TriNetX, LLC (TNX) is a private organization that has built a global medical research network that enables healthcare organizations to make their electronic health record (EHR) data more easily accessible to researchers in a de-identified manner, enabling Real World Data analyses. To prevent possible outlier values due to results being reported in different units or encoding errors, continuous laboratory test results were excluded, thereby limiting our analysis to only binary tests where the results were either positive or negative. A query built by combining a complete list of Logical Observation Identifiers Names and Codes (LOINC) codes that corresponded to binary clinical laboratory tests for our pathogens of interest was run on 02-14-2023 across 73 healthcare organizations on the TNX Research Network. The search resulted in a list of just over 11 million unique patients who had in their EHR a result for at least one of the laboratory tests in our query. Diagnoses for all participants were collected, and cohorts were automatically generated for each pathogen-disease pair. The use of GNU Parallel<sup>18</sup> made the processing of the immense amount of TNX data much more tractable. For a specific pathogen-disease pair, those with a test result for a particular LOINC code but without the diagnosis of interest were considered controls, and only those with a test result (positive or negative) before the earliest diagnosis for the disease of interest in their medical record were considered cases. We removed those participants with the diagnosis appearing in the EHR before the laboratory test, as the temporal relationship of infection prior to disease diagnosis could not be firmly established for them. We attempted to extract all potential confounding variables considered in the UKB analysis from the TNX data; however, data for only three of the ten were available. For situations where a covariate was included in the UKB model but was not available in the TNX data, the covariates were dropped from the TNX model before refitting. After finding few ICD10 B24 cases in TNX (a diagnosis used in UKB to indicate human immunodeficiency virus (HIV) infection, “Unspecified human immunodeficiency virus [HIV] disease”), it was determined that in the United States (the primary source of TNX data) the ICD10 code B20 appears to be primarily used to indicate HIV infection. Thus, the TNX results for B20 were merged with the UKB B24 results. However, this is the only ICD10 code for which this was done.

Use of the data obtained from the TriNetX research network was accessed through the University of Cincinnati using a previously approved IRB umbrella publishing protocol (2019-1238). The TriNetX data used in this study previously existed and were deidentified. Furthermore, our study involved no interaction or intervention with human subjects, indicating that it is a secondary analysis and not human subjects research. Finally, TriNetX requires each participating healthcare organization to obtain informed consent from its patients, including consent for the secondary analysis of the data collected.

### Phecode analysis

Phecodes were generated using UKB and TNX ICD10 data with the PheWAS R library (PheWAS v0.99.6.1; R v4.0.2) and a minimum code count of one following previously published methods<sup>19–21</sup>. The resulting Phecodes were filtered to remove those with a disease group of infectious diseases,

injuries & poisonings, congenital anomalies, or symptoms. We applied the same statistical analysis to the Phecodes that we used for ICD10 codes, including the same study design involving both the UKB and TNX cohorts, requiring the same statistical thresholds to be met (UKB per-Phecode false discovery rate (FDR) < 0.3 and TNX per-Phecode FDR < 0.01). Due to technical issues, we were unable to include the requirement for a laboratory result to be present in the EHR before diagnosis in the TNX cohort. A set of infectious disease Phecodes was used as positive controls (Tier 1 only) or negative controls (Expected Negatives) – see Supplementary Data 2.

### Establishment of the minimum number of cases and total samples for analysis

To ensure that the comparisons made had sufficient statistical power to identify associations if present and minimize the multiple testing correction burden, we sought to determine what a minimum number of cases and controls would be. While our analytical strategy used a logistic regression model, to determine a preliminary power estimate, we opted to consider the power to detect differences in antibody levels between those with and without disease for known antibody-disease pairs. This approach enabled us to calculate the effect sizes (mean difference between groups divided by pooled standard deviation) and evaluate the power required to detect such differences. While the logistic regression models (especially with covariates) likely have different power estimates, our goal was to identify a pragmatic threshold for the minimum number of cases with a reasonable likelihood of success. We looked at effect sizes across 14 known positive antibody-disease pairs in the UKB data, which collapsed to eight pathogen-disease pairs. These pairs constitute our “Tier 1” positive controls and represent those infectious diseases that are directly caused by a pathogen, such as “herpes zoster (HZ)” (ICD10: B02) diagnosis and varicella-zoster virus, the virus that causes HZ. Effect sizes at the antibody-disease level ranged from 0.1 to 3.9. After collapsing the antibody results to the largest effect per pathogen-disease pair, effect sizes ranged from 0.2 to 3.9, with a median effect size of 0.72. Using the G\*Power 3.1 software package<sup>22</sup>, we calculated the number of cases required for 80% power with an alpha of 0.05, assuming ten controls per case, yielding a minimum case number of 17.

We recognize that more controls may be available, but higher numbers of controls impact power only slightly (data not shown). Further, using the median effect size of known positives might underestimate the power for some pathogen-disease associations; this conservative choice was made to balance the inclusion of diseases while ensuring sufficient statistical power.

### Model development and application

We modeled the association between a given disease and pathogen using a logistic regression model with disease status as the binary outcome and the pathogen proxy value (continuous for UKB antibody titers, categorical for TNX positive/negative binary lab tests) as the predictor (Supplementary Fig. S1a). Since each cohort uses a different type of predictor, we note that the odds ratio (OR) has a slightly different meaning between the two. In the UKB cohort, the OR represents the increase in odds of developing a disease per 10-fold increase in antibody titer level. Whereas in the TNX cohort, the OR represents the increase in odds of developing a particular disease in those subjects after infection by a given pathogen. All sex-specific diseases had controls limited to those participants at risk, e.g., cervical cancer (ICD10: C53) used only females without a cervical cancer diagnosis as controls. For all pregnancy and childbirth-related diseases, only patients with a record of a healthy birth (ICD10 codes O80–O84) were used as controls (i.e., patients with records of both a healthy birth and the diagnosis of interest were removed from the analysis to prevent them being used as both cases and controls).

After adjusting our logistic regression model for all covariates found to be significantly associated with both disease status and pathogen proxy in separate two-sided univariate statistical tests (Supplementary Table 1), we ran a backward elimination procedure (stepAIC, MASS R library), removing any non-significant covariates from the final model, to mitigate possible overfitting. We attempted to replicate models that showed

statistical significance in the UKB cohort using the replication cohort by simply refitting the same model with the TNX data.

In the rare situation where a categorical pathogen test (TNX) had five or fewer patients in one of the cells of the pathogen-disease contingency table, we altered our model slightly. For example, for the diagnosis of “herpesviral [herpes simplex] infections” (ICD10: B00) and the herpes simplex virus 1 (HSV1) test with LOINC code 93439-8, there is only one disease control with a positive test result for HSV1, the cause of this disease. This represents a strong association that would be lost with a standard logistic regression model. Thus, in these situations, we instead employed Firth’s bias-reduced logistic regression method<sup>23</sup> (R library *logistf*), which is appropriate in such situations.

A permutation procedure was performed to verify the robustness of the UKB model results. Briefly, 10,000 permutations were run for each antibody-disease pair (45 total pairs per disease), where disease status was randomly shuffled amongst the participants while keeping the number of cases and controls and the model itself, i.e., confounders adjusted for, constant. Next, all permutation results for a particular disease were pooled into a larger per-disease null distribution now containing 450,000 permutation results. Empirical p-values for each antibody-disease pair were calculated by comparing the nominal p-value from our analytical model to the respective per-disease null distribution.

We applied a per-disease Benjamini–Hochberg (BH) false discovery rate<sup>24</sup> at the pathogen-disease level to control for multiple testing. Since we were using a discovery-replication model, we used a lenient FDR threshold of 0.3 for our discovery cohort to reduce the likelihood of false negatives that might occur due to the smaller size of the discovery cohort. Despite the smaller number of subjects, the UKB cohort has the advantage of systematic measurements of antibody titers. We therefore applied a much more stringent FDR threshold of 0.01 to our larger replication cohort to minimize false positives. Our Tier 1 analysis drove the choice of a precise 0.3 cutoff in our discovery cohort. All eight Tier 1 associations were significant in the UKB data at this cutoff, which is expected given their well-accepted causal relationships. We emphasize that many Tier 2 associations identified by a semi-automated literature search might be incorrect or might represent weaker associations that are harder to detect. We therefore did not require all Tier 2 associations to be significant in the UKB data and instead expected an intermediate replication rate for the Tier 2’s between the Tier 1 and Expected Negatives.

### Model assessment

To assess the model’s performance, we calculated associations for a set of positive and negative control pairs. We used the Tier 1 controls as described above for the positive controls. Six infectious diseases are included in the Tier 1 controls, two of which can be caused by two different pathogens included in this study (“herpesviral [herpes simplex] infections” by herpes simplex virus 1 or 2 and “unspecified viral hepatitis” by hepatitis B (HBV) or hepatitis C (HCV)). The negative control set, termed “Expected Negatives”, represents the complement of the Tier 1 controls, e.g., a herpes zoster diagnosis paired with HBV instead of the causal agent, varicella zoster virus (VZV). As an additional assessment, a second set of positive controls using only non-communicable diseases (NCDs) (“Tier 2”) was collected using a semi-automated literature mining approach. In brief, we employed the log product frequency (LPF), a previously published method for quantifying the co-occurrence of search terms in PubMed<sup>25,26</sup>. Specifically, we employ a negated form of LPF to rank pathogen-disease pairs by the number of PubMed co-citations (disease and pathogen), normalized by the number of citations of each separately (Supplementary Fig. S1b). The negation rotates the LPF values around the origin, allowing us to deal with LPF values greater than zero, whereas regular LPF values are all less than or equal to zero. The closer the LPF value is to zero, the more the disease and pathogen were co-cited as opposed to cited individually.

Searches of PubMed were conducted using the Entrez functions from the Python library Biopython using default settings on 8/11/2020. The top 175 results ranked by LPF (Supplementary Data 3) were manually reviewed

by M.L. The evidence column in Supplementary Data 3 contains one of four levels used to gauge the literature support found for a particular pathogen-disease pair during the manual review. Briefly, a rating of “High” indicates that at least one published meta-analysis supporting the association was found, or the connection is textbook knowledge (e.g., Hepatitis B causes hepatitis). “Medium” signifies that some cross-sectional or longitudinal studies were found supporting the pair’s association. An evidence level of “Low” indicates that only a handful of case reports pointing to a possible association were found. “None” means no papers were found indicating an association.

After limiting pairs to only those with a “high” evidence grade, we were left with 83 pathogen-disease pairs, making up our “Tier 2” positive controls. Pairs with less evidence and thus not rated as “high” were classified as either “medium”, “low”, or “none”, for several reasons, all of which were excluded from the Tier 2 list. For example, D50, “iron-deficiency anemia” paired with Merkel cell polyomavirus (MCV), ranked 22<sup>nd</sup> by LPF, was marked as “none” after manual review. Notably, the acronym MCV can also stand for “mean corpuscular volume”, a blood measurement that can indicate iron-deficiency anemia, which is likely the cause for the high ranking. Further, K58, “irritable bowel syndrome (IBS)” paired with *H. pylori*, ranked 138<sup>th</sup> by LPF, was classified as having low evidence by M.L. after the manual literature review. Although many studies have examined this association, the two most recent, i.e., published before our literature search, meta-analyses<sup>27,28</sup> have found no significant association between the two. Further, since the initial literature review, two additional meta-analyses found non-significant associations between IBS and *H. pylori* as well<sup>29,30</sup>.

This semi-automated approach was required because a manual literature review for (20 pathogens x 426 diseases = 8520) pathogen-disease pairs was intractable. The raw automated queries and results, including the number of citations and PubMed ID (PMID) lists for pathogen-disease co-citations, pathogen only, and disease only, can be found in Supplementary Data 4.

### Orthogonal validation of CMV-UC association

We attempted to validate two of our replicated findings using two distinct ‘omics-based methods.

First, to capture differences in viral gene expression levels between cases and controls for the diseases of interest, we used publicly available RNA-seq data from case/control cohorts (Supplementary Data 5) and the bioinformatics tool VIRTUS v1.2.1<sup>31</sup> using default parameters except reducing the default “hit\_cutoff” from 400 to 0 since we were only investigating a small number of pathogens. Briefly, VIRTUS is a pipeline that takes an input RNA-seq FASTQ file and aligns the reads to a reference human genome (hg38; GCF\_000001405.39) using STAR<sup>32</sup>. It then attempts to align any unaligned reads to a second pre-compiled index file containing a user-specified set of viral genomes, here Epstein-Barr virus (EBV) (NC\_007605.1) and human cytomegalovirus (CMV) (NC\_006273.2). The resulting number of mapped reads for each virus was first normalized by the pathogen’s genome length, then normalized by the number of mapped human reads in the RNA-seq sample, providing the flexibility to compare results across both different samples and experiments as well as different pathogens. Finally, the non-parametric Mann–Whitney U test was used to statistically compare the normalized read counts for a particular pathogen between cases and controls.

Second, to examine if genome-wide association study (GWAS) loci for the diseases of interest were enriched near human genes that have altered expression levels upon infection by the pathogens of interest more so than unchanged genes, we used our RELI tool<sup>11</sup>. RELI estimates the statistical significance of the overlap between a set of input genomic loci and a peak file, typically GWAS loci and ChIP-seq peaks. It does this through a permutation-based procedure by first counting the number of overlaps between the input loci and peaks, then permuting the input loci around the genome and again counting the number of overlaps with the peaks, building a null distribution. Finally, to calculate significance, RELI compares the total number of overlaps seen in the input loci set to the constructed null distribution.

We obtained GWAS data for ulcerative colitis (UC) and systemic lupus erythematosus (SLE) from the NHGRI-EBI GWAS catalog (v1.0.2-associations\_e96\_r2019-05-03)<sup>33</sup>. A genome-wide significance cutoff of  $5 \times 10^{-8}$  was used, and we considered only data from European populations due to the prevalence of GWAS data for this ancestry group. Independent loci were identified for each phenotype using linkage disequilibrium (LD)-based pruning with PLINK<sup>34</sup> (window size 300,000 kb, SNP shift size 100,000 kb, and  $r^2 < 0.2$ ). These independent loci were then expanded to incorporate variants in strong LD ( $r^2 > 0.8$ ) again using PLINK. Next, we downloaded lists of genes with expression changes in response to CMV or separately EBV infection from the VExD database (<https://vexd.cchmc.org>)<sup>35</sup> (Supplementary Data 6). Differentially expressed genes (DEGs) are identified in VExD as those genes with an adjusted  $p$ -value  $< 0.05$  and absolute fold change  $> 2$ , when comparing infected and uninfected cells of the same type within a single study. As a null model, for each study, we also identified sets of genes that, while expressed, do not significantly change upon infection (adjusted  $p \geq 0.05$ , and fold change  $< 1.2$ ) and randomly selected the same number of genes to match the corresponding differentially expressed gene set. We then ran RELI using the GWAS risk loci for each disease as input against the genomic regions defined by a 200 kb window centered on the transcription start site of each input gene.

### Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

## Results

### Pathogen-disease data collection in two large independent cohorts

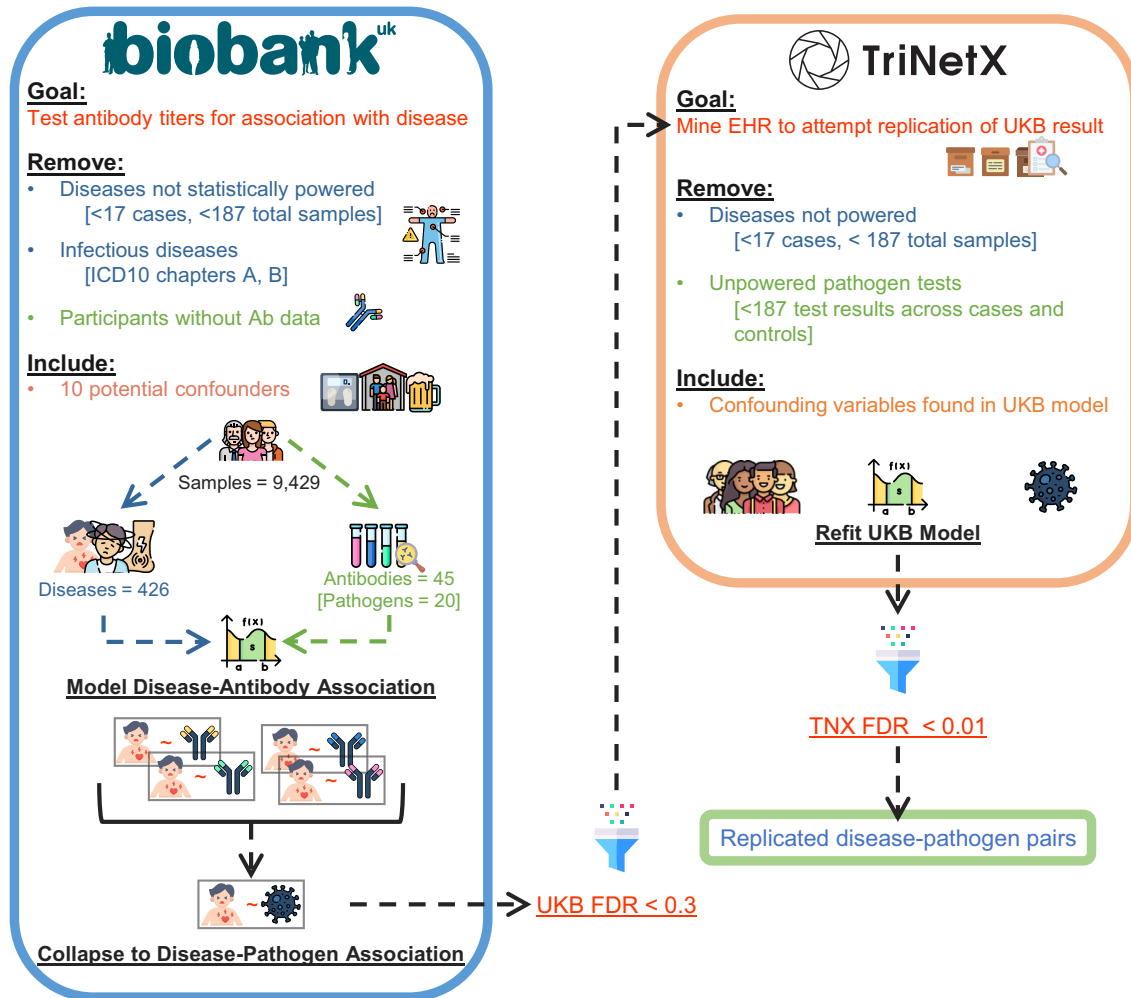
We employed a discovery-replication experimental design using two newly available resources, The UK Biobank (UKB) and TriNetX (TNX). For our discovery cohort, we used the portion of UKB participants that had 45 antibody titers systematically measured, representing immune responses to 20 unique pathogens (9429 UKB participants). Within the set of UKB subjects with serologic data, setting aside International Classification of Diseases 10<sup>th</sup> revision (ICD10) codes A00–B99 covering “Certain infectious and parasitic diseases”, we found 399 non-communicable diseases (NCDs) and 27 communicable diseases (ICD10 codes between C00 and O99) for which we were statistically powered to test (Supplementary Data 7). The primary focus of this study is on the 399 NCDs with a secondary assessment of communicable diseases that are associated with non-causative pathogens. This discovery cohort (UKB) has the advantage of a carefully controlled experimental design with many covariates for the detection of confounding effects, but it has the disadvantage of a relatively small number of subjects.

We used data extracted from The TriNetX Research Network as our independent replication cohort. These data contain both the clinical diagnoses and the serologic test results required for our model for over 11 million individuals. A total of 209 separate binary (positive/negative) clinical laboratory tests targeting the 20 UKB pathogens were identified in the medical records provided by TNX. Utilizing the same statistical power requirements, TNX was powered to test nearly 75% of the 8616 tested UKB pathogen-disease pairs. This replication cohort (TNX) has the disadvantage of a less well-controlled study design involving data from different clinical tests and sites, with the advantage of a very large number of subjects. The larger number of subjects enabled the use of an additional restriction requiring infection status to appear in a participant’s medical record prior to disease diagnosis (see Methods).

### Development of a statistical model with high sensitivity and specificity

We developed a workflow for systematically detecting and replicating pathogen-disease pairs within our discovery and replication cohorts (Fig. 1). To this end, we used a logistic regression model to test for association between one of the 426 disease statuses of interest and a proxy for a history of infection by a given pathogen: either one of the 45 continuous antibody titer





**Fig. 1 | Overview of study design.** The UK Biobank (UKB) was used as the discovery cohort (left). 426 diseases with sufficient sample counts in UKB (along with positive and negative controls) were tested for association with the 45 antibody titers representing immune responses to 20 unique pathogens across 9429 UKB participants. Models were adjusted for any of ten additional health-related and socio-demographic variables that were determined to be confounding for a particular

antibody-disease pair. The antibody-disease pair results were collapsed to the pathogen-disease level by selecting the most significant antibody-disease result to represent that pathogen's association with the disease. Significant pathogen-disease pairs identified in UKB were tested in the independent TriNetX (TNX) cohort (right). Pathogen-disease pairs that were significant in both the discovery and the replication cohort were considered replicated pairs.

values (UKB) or binary clinical laboratory test results (TNX). To mitigate potential confounding, we adjusted each model for any of ten additional sociodemographic and health-related variables (Supplementary Data 1) that were found to be significantly associated with both disease status and the pathogen proxy in the discovery cohort (Supplementary Table S1). To prevent overfitting, a backward elimination procedure was used to prune non-significant covariates before fitting the final model (see Methods). Six infectious diseases (ICD10 codes between A00 and B99) were also included as part of a control set. Within the resulting 19,289 separate antibody-disease tests, we collapsed instances of multiple antibodies to the same pathogen by selecting the antibody with the most significant association. In total, this procedure resulted in 8616 pathogen-disease tests.

The most commonly adjusted-for covariates were age, sex, and body mass index (BMI), included in 35.5%, 31.4%, and 24.4% of the UKB antibody-disease models, respectively (Supplementary Fig. S2). As expected, age was most significant in models for diseases that are more common in geriatric patients, such as “disorders of lipoprotein metabolism and other lipidaemias”, “other arthrosis”, and “other cataract”. Likewise, the strongest BMI associations were seen in models for metabolic syndrome diseases such as “obesity”, “essential primary hypertension”, and “diabetes mellitus”. Roughly a quarter of models did not require adjustment for any covariates.

We confirmed the robustness of our analytical model in our UKB discovery cohort using a permutation-based approach. To this end, we calculated an empirical  $p$ -value for each antibody-disease pair by permuting disease status across individuals (see Methods). We observed exceptionally strong correlation (Pearson's  $r > 0.99$ ) between the nominal  $p$ -value obtained from our analytical model and the empirically derived permutation-based  $p$ -value (Supplementary Fig. S3).

To reduce the multiple testing burden, we only tested pathogen-disease pairs in the replication cohort that were statistically significant in the discovery cohort. During replication, we refit the UKB model using the replication cohort data, adjusting for the same covariates when data were available. To reduce false negatives during discovery, a lenient per-disease Benjamini–Hochberg (BH) false discovery rate (FDR) threshold of 0.3 was applied to the discovery cohort results (see Methods). Then, to reduce false positives during replication, a much more stringent per-disease FDR threshold of 0.01 was used in the replication stage. Only pairs with significant associations and effects (odds ratios) in the same direction across both cohorts were considered replicated pathogen-disease relationships.

We first assessed the sensitivity of our model on a set of positive controls. “Tier 1” positive controls were identified as infectious disease diagnoses paired with their causal pathogen, e.g., hepatitis C (HCV) paired with a diagnosis of “unspecified viral hepatitis” (ICD10: B19). As expected,

**Table 1 | Results for all Tier 1 pathogen-disease pairs and other discussed pairs**

Disease name	ICD10	Path	Group	UKB FDR	TNX FDR	UKB OR	TNX OR	Rep status
Anogenital Herpesviral Infect	A60	HSV-2	Tier 1	1.4E−03	2.6E−298	4.10	11.76	REP
Herpesviral Infect	B00	HSV-1	Tier 1	9.5E−04	8.9E−266	2.08	2.93	REP
Herpesviral Infect	B00	HSV-2	Tier 1	2.2E−01	1.6E−298	1.36	5.18	REP
Herpes Zoster	B02	VZV	Tier 1	6.4E−02	4.1E−259	1.43	14.43	REP
Unspec Viral Hepatitis	B19	HBV	Tier 1	2.4E−04	1.9E−299	2.28	7.20	REP
Unspec Viral Hepatitis	B19	HCV	Tier 1	6.7E−06	1.9E−299	2.96	11.57	REP
Unspec HIV	B24	HIV	Tier 1	1.3E−06	8.0E−300	38.21	635.37	REP
Infectious Mono	B27	EBV	Tier 1	5.0E−02	8.2E−212	2.14	8.39	REP
Iron-Deficiency Anemia	D50	<i>H. pylori</i>	Tier 2	9.9E−03	1.1E−08	1.17	1.24	REP
Multiple Sclerosis	G35	EBV	Tier 2	2.5E−01	2.5E−03	2.55	4.45	REP
Duodenal Ulcer	K26	<i>H. pylori</i>	Tier 2	1.3E−03	1.7E−04	1.73	1.43	REP
Peptic Ulcer Site Unspec	K27	<i>H. pylori</i>	Tier 2	9.0E−02	5.6E−138	1.60	3.68	REP
Gastritis and Duodenitis	K29	<i>H. pylori</i>	Tier 2	1.5E−01	7.9E−134	1.13	1.97	REP
Fibrosis and Cirrhosis of Liver	K74	HCV	Tier 2	7.3E−02	6.7E−299	2.69	4.31	REP
Other Dis of Liver	K76	HBV	Tier 2	1.6E−02	7.8E−299	1.44	2.97	REP
Other Dis of Liver	K76	HCV	Tier 2	1.6E−02	7.8E−299	1.47	2.03	REP
Asthma	J45	<i>H. pylori</i>	Unk	1.8E−02	2.8E−04	0.89	0.80	REP
GERD	K21	<i>H. pylori</i>	Unk	2.1E−02	2.9E−07	0.88	0.66	REP
Other Dis of Esophagus	K22	<i>H. pylori</i>	Unk	5.1E−03	1.4E−09	0.68	0.72	REP
Diaphragmatic Hernia	K44	<i>H. pylori</i>	Unk	2.7E−01	1.4E−06	0.90	0.76	REP
Inflammatory Bowel Dis	K58	<i>H. pylori</i>	Unk	1.9E−01	6.10E−29	0.91	0.48	REP
Crohn's Dis	K50	CMV	Unk	8.7E−01	–	0.86	–	DNAR
Crohn's Dis	K50	EBV	Unk	8.7E−01	–	1.65	–	DNAR
Ulcerative Colitis	K51	CMV	Unk	2.0E−01	8.6E−06	1.36	2.78	REP
Ulcerative Colitis	K51	EBV	Unk	3.7E−01	–	0.80	–	DNAR
Unspec Noninfective Gastroenteritis/Colitis	K52	CMV	Unk	8.4E−01	–	0.97	–	DNAR
Unspec Noninfective Gastroenteritis/Colitis	K52	EBV	Unk	8.4E−01	–	1.04	–	DNAR
Systemic Lupus Erythematosus	M32	EBV	Tier 2	1.8E−01	1.4E−21	3.98	4.96	REP
Systemic Lupus Erythematosus	M32	CMV	Unk	4.0E−01	–	1.89	–	DNAR

The disease name with its 3-character International Classification of Diseases 10th revision (ICD10) code is listed next to the paired pathogen, followed by the electronic health record statistical analysis results, including the UK Biobank (UKB) per-disease false discovery rate (FDR) and odds ratio (OR), as well as the TriNetX (TNX) per-disease FDR and OR. All eight of the Tier 1 pathogen-disease pairs and any other pairs discussed are listed.

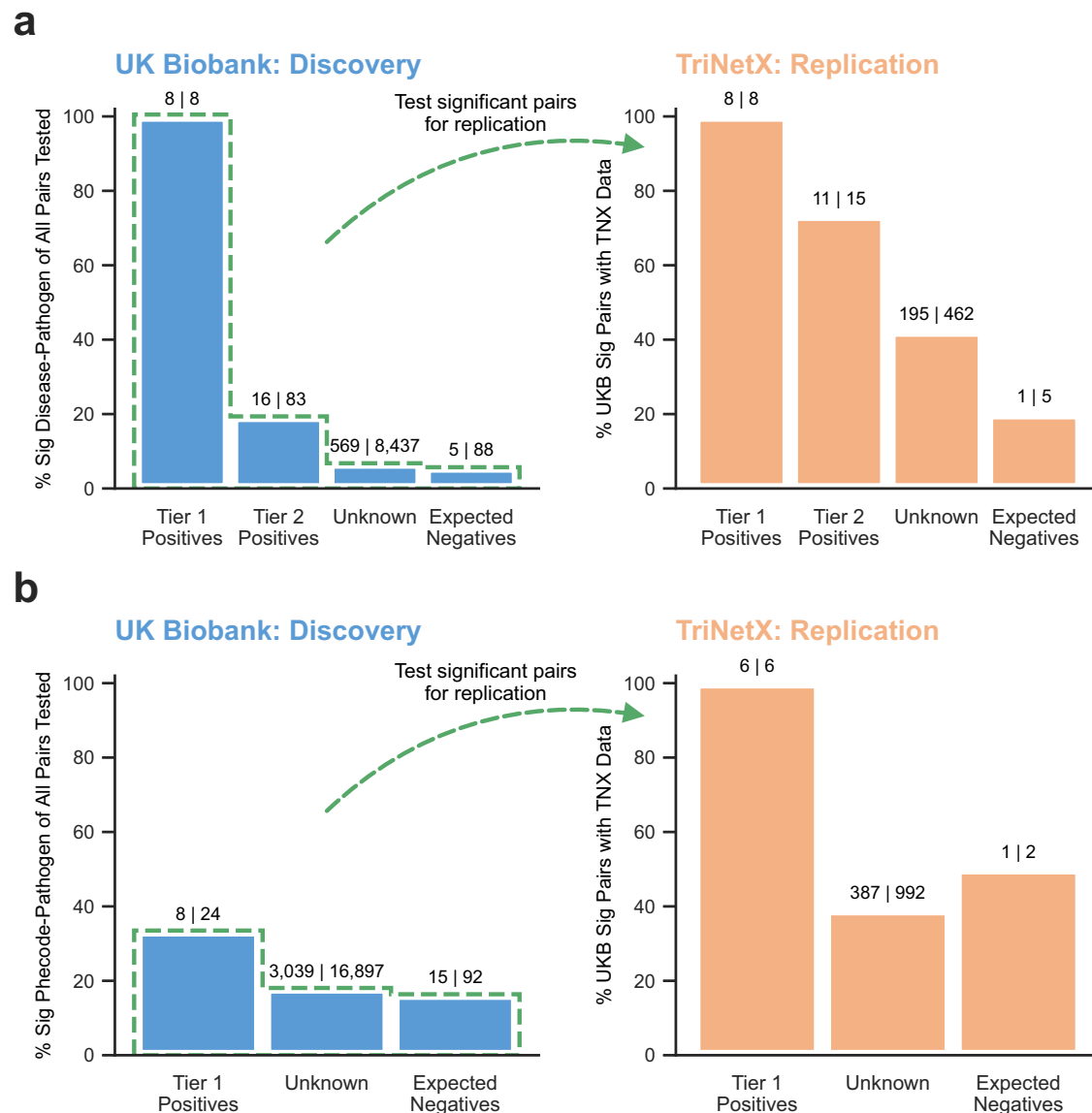
Column abbreviations: *Path* Pathogen, *Group* Control group (Tier 1 control, Tier 2 control, or previously Unknown association (Unk)). Disease name abbreviations: *Infect* Infection, *Unspec* Unspecified, *Mono* Mononucleosis, *Dis* Disease, *GERD* Gastroesophageal reflux disease. Replication Status (Rep Status) abbreviations: *REP* Replicated (UKB per-disease FDR < 0.3 AND TNX per-disease FDR < 0.01), *DNAR* Did not attempt replication (UKB per-disease FDR > 0.3).

our model found significant associations between the disease “unspecified viral hepatitis” and both hepatitis B (HBV) and HCV pathogens. It also identified an association with the BK virus (BKV) at the more lenient discovery cohort threshold. However, upon testing for replication, only the HBV and HCV results remained significant (Supplementary Fig. S4). To assess the specificity of the model, we used a set of “Expected Negatives”, which we identified as the complement of the Tier 1 positive controls, i.e., an infectious disease diagnoses with a pathogen that does not cause the disease, such as “unspecified viral hepatitis” with Epstein-Barr virus (EBV). We excluded human immunodeficiency virus (HIV) from the Expected Negative set due to its indirect involvement in many immune-mediated diseases.

Overall, our model identified significant association for all eight (100%) of the Tier 1 positive control pathogen-disease pairs in the UKB cohort, and all eight replicated in TNX (Table 1). Conversely, the model identified only five (5.68%) of the Expected Negatives as significant in UKB, only one of which replicated in TNX, a diagnosis of “infectious mononucleosis” (ICD10: B27) and human herpesvirus 6 (HHV-6). Upon further investigation, this replicated Expected Negative pair represents a previously established

relationship: HHV-6 infection can account for up to five percent of infectious mononucleosis (IM)-like syndrome diagnoses in adult patients<sup>36</sup>. Considering only the fully replicated pairs as predicted positives, and those pairs that were either identified as not significant in UKB or failed replication as predicted negatives, our model has a sensitivity of 1.0, a specificity of 0.80, and a precision of 0.89.

We next assessed a set of 83 pathogen-NCD pairs with suggestive literature evidence, the “Tier 2” positive controls, collected via a semi-automated literature search approach (see Methods). Sixteen (19.28%) of these Tier 2 pairs were significantly associated in UKB, and 11 of the 15 with available data (68.75%) replicated in TNX. Included in these are well-known associations such as *H. pylori* with several gastroenterological diseases such as “duodenal ulcer”, “peptic ulcer, site unspecified”, and “gastritis and duodenitis”<sup>37–39</sup>. We also replicated connections between particular pathogens and hepatic diseases, such as “fibrosis and cirrhosis of liver” with HCV and “other diseases of liver” with both HBV and HCV<sup>40</sup>. Finally, our model validated the now well-established associations of EBV with multiple sclerosis (MS)<sup>10,12,41</sup> and with systemic lupus erythematosus (SLE)<sup>42,43</sup>. Taken together, these results establish that our approach can capture both



**Fig. 2 | Summary of pathogen-disease pairs identified across both cohorts.** Bar charts show the percent of pathogen-disease pairs in each subgroup (“Tier 1 Positives”, “Tier 2 Positives” (ICD10 analysis only), “Expected Negatives”, and “Unknown”) that were significant in the discovery cohort (blue bars) and replication cohort (orange bars). The numbers used to calculate the percentages are indicated above each bar. All discovery cohort (blue bars) significant pairs were assessed for

replication in the independent replication cohort (orange bars). Note that the total number of pairs tested for replication may not match the number of significant UK Biobank pairs due to insufficient data available in the replication cohort. Such cases are not considered replication successes or failures. **a** Results using International Classification of Diseases 10<sup>th</sup> revision (ICD10) codes. **b** Results using Phecodes.

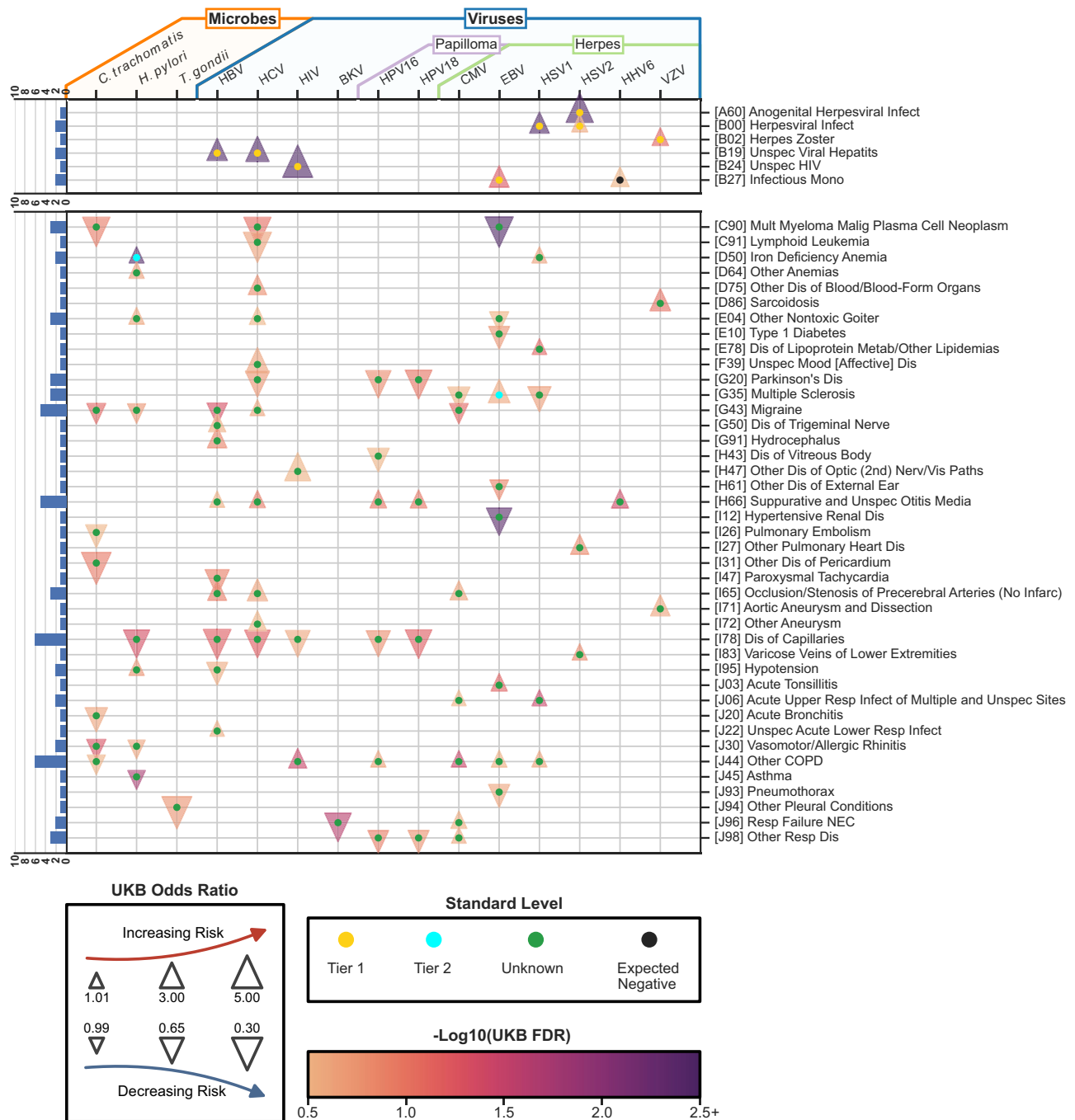
well-established and suggestive pathogen-disease relationships while maintaining a substantial degree of specificity.

### Identification of 206 replicated pathogen-disease relationships

Encouraged by the performance of our model on our positive and negative control sets, we next sought to identify, to the best of our knowledge, novel pathogen-disease relationships. In total, of the 8437 “unknown” (non-Tier 1, Tier 2, or Expected Negatives) pathogen-disease tests that met our requirements in UKB, 569 were significant. 462 of these pairs had sufficient data in the TNX cohort to test for replication, and 195 of these were replicated (2.3% of the total 8437 “unknown” pairs that were initially tested) (Fig. 2a). The 195 replicated pairs represent a diverse collection of diseases and pathogens, with 15 of the 20 tested pathogens connected to at least one of 96 distinct diseases, 89 of which are NCDs (Figs. 3, 4, and Supplementary Data 8). Altogether, the 11 Tier 2 and the 195 “unknown” pathogen-disease relationships equate to 206 replicated associations. While these relationships

are correlations, due to the size of the TNX cohort, we could restrict our study population to only those with a pathogen test result (positive or negative) before disease development while remaining statistically powered. Thus, these “temporal correlations” provide stronger evidence that the pathogen plays a role in disease development than a simple correlation. In contrast, the antibody titer data present in the UKB were obtained from the near-simultaneous measurement of all 45 antibody titers across all 9429 pilot project participants. Thus, they cannot provide the same temporal correlations.

Outside of our Tier 1 positive controls, the largest odds ratio (OR) obtained in the UKB is for systemic lupus erythematosus (SLE) with EBV (UKB OR = 3.98), which also has a large TNX odds ratio of 4.96. We also replicated the well-established connection between EBV and multiple sclerosis (UKB OR = 2.55; TNX OR = 4.45). Overall, HCV and HBV have the most replicated associations (25) of the pathogens tested (Supplementary Fig. S5). EBV, HBV, and HCV all have replicated associations in over



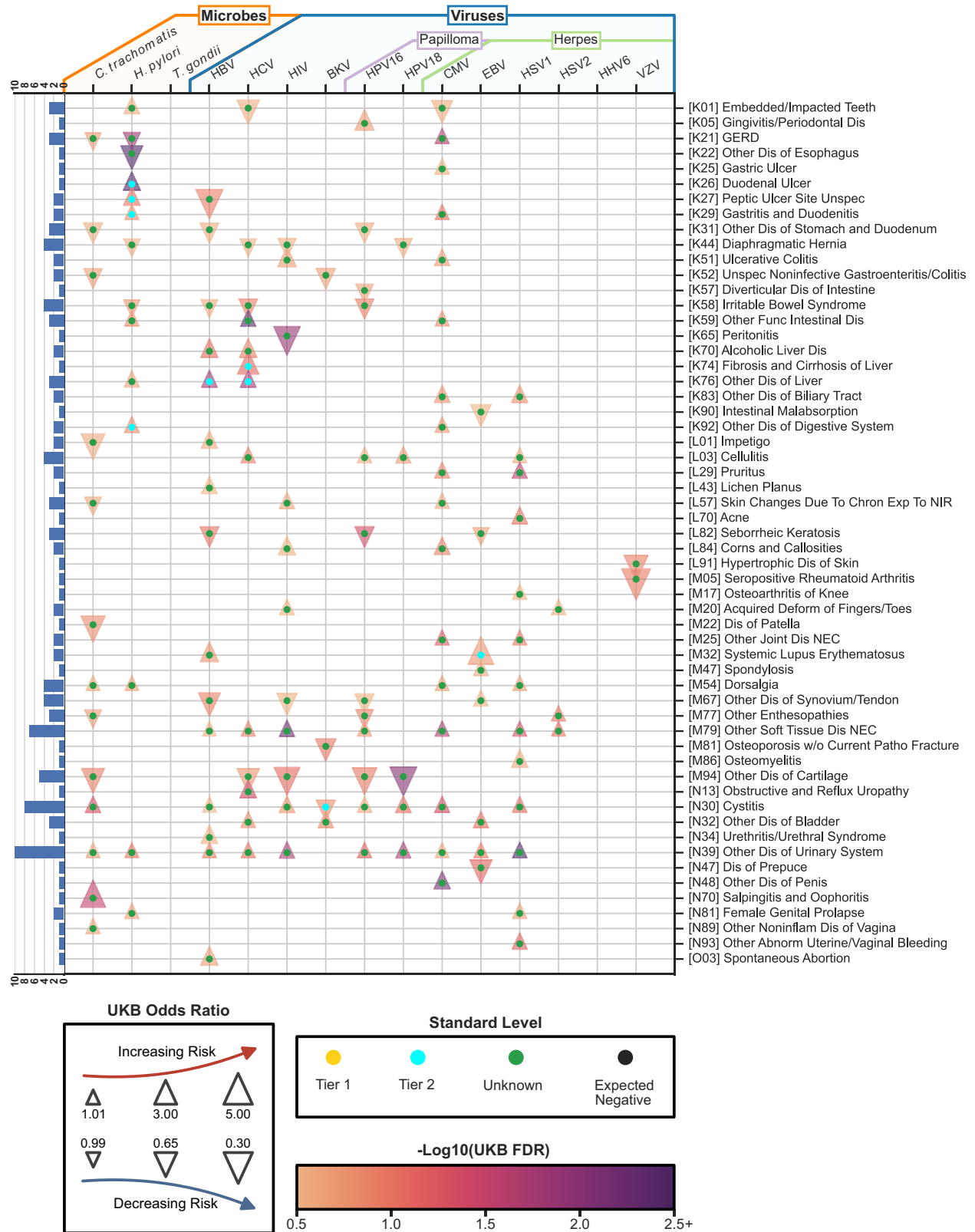
**Fig. 3 | Overview of replicated Tier 1 controls and ICD10 codes C00–J99 pathogen-disease pairs.** Heatmap containing all pathogens and diseases, either Tier 1 or with an ICD10 code between C00 and J99, with at least one replicated result. Pathogens are first grouped by type (microbe or virus) and then sub-grouped by family if more than one member is present. Diseases are ordered by International Classification of Diseases 10<sup>th</sup> revision (ICD10) code, with the Tier 1 positive controls at the top. A histogram opposite to each disease indicates the total number of replicated associations between that disease and all pathogens investigated. Each replicated result for a pathogen-disease pair is represented by a triangle either pointing up (indicating an odds ratio greater than one) or down (indicating an odds

ratio less than one). The size of each triangle represents the discovery cohort (UK Biobank, UKB) odds ratio (OR) and is capped at a maximum of 5 to account for the very large effect size between “unspecified hiv disease” and human immunodeficiency virus (HIV) (OR: 38.2). The color of each triangle represents the UKB negative log base-10 per-disease FDR, capped at a maximum of 2.5 to enable better visual distinction in the region covering all but the eight most significant pathogen-disease pairs. Each triangle is marked with a central dot, the color of which indicates whether the pair is a Tier 1 (gold) or Tier 2 (cyan) positive control, Expected Negative (black) control, or an unknown relationship (green).

ten different ICD10 chapters, reflecting the often systemic effects of infections by these pathogens. Human cytomegalovirus (CMV) has the most associations with an OR indicating risk (21 of its 24 replication associations have an OR greater than 1). In contrast, *Chlamydia trachomatis* has the largest number of protective relationships (15 of its 20 with an OR less than

1). “Other diseases of the urinary system” (ICD10: N39) had the most replicated associations across all pathogens (10), all of which are predicted to increase the risk of disease. The 3-character ICD10 code N39 includes both the communicable diagnosis “urinary tract infection” (N39.0) as well as several non-communicable forms of incontinence (N39.3, N39.4).





**Fig. 4 | Overview of all replicated ICD10 codes K00–O99 pathogen-disease pairs.** Heatmap similar to Fig. 3, containing all diseases with at least one replicated result for ICD codes in chapters XI–XV (K00–O99). One pathogen with no replicated

results, HHV-6, is included to better facilitate interpretation alongside Fig. 3. Please refer to Fig. 3's legend for a detailed description of the data.

To characterize the replicated pathogen-disease pairs more broadly, we examined our results at the ICD10 chapter level, normally representing distinct body systems. *H. pylori* has the most replicated associations in a particular chapter: chapter XI (K00–K95), which contains diseases of the digestive system (Supplementary Fig. S6). In addition to the three Tier 2 positive controls discussed above, *H. pylori* is predicted to increase the risk of four additional digestive system diseases and to decrease the risk of “irritable bowel syndrome”, “diaphragmatic hernia”, “other diseases of esophagus”, and “gastro-esophageal reflux disease” (GERD).

### Phecode analysis of pathogen-disease associations

Phecodes are a phenotype encoding scheme developed originally for PheWAS studies<sup>44,45</sup>. Phecodes are sets of ICD codes bundled together in an attempt to represent a single phenotype. They also include the use of exclusion criteria, which helps to reduce the presence of cases in a control cohort. To demonstrate the robustness of our pathogen-disease associations, we next repeated our ICD10-based analysis using Phecodes as the outcome variable instead of individual ICD10 codes. This analysis included 18,820 pathogen-Phecode unknown pairs, 10 Tier 1 pairs, and 104 Expected Negatives.

We first examined the positive control results, finding that 8 of the 10 Tier 1 results are significant in the discovery cohort (UKB data), all of which are also significant in the replication cohort (TNX data). As expected, this fraction of pairs that replicated pairs was much higher than the fraction observed for the Expected Negatives, 6 of the 104 pairs (Fig. 2b). These results demonstrate that our Phecode-based model also has the discriminatory capacity to separate positive controls from negative controls.

After assessing the model performance, we next analyzed the unknown pathogen-Phecode pairs. Of the 18,820 unknown pairs, 3289 (589 unique Phecodes) were significant in the discovery cohort. We were powered to test 2784 of these pairs for replication in the TriNetX data, where we found that 1341 pairs (449 unique Phecodes) fully replicated (Supplementary Data 2). Nearly all of the pairs with a Phecode corresponding to an ICD10 in Table 1 replicate, including multiple sclerosis with EBV, systemic lupus erythematosus with EBV, and ulcerative colitis with CMV, indicating agreement between the ICD10 results and the Phecode results.

### Orthogonal validations of the cytomegalovirus – ulcerative colitis relationship

We next sought orthogonal evidence supporting the 206 replicated associations identified by our approach. Virus-disease relationships are often reflected by higher expression levels of virus-encoded genes in patients compared to controls<sup>46–50</sup>. Likewise, many viruses manipulate host gene expression patterns<sup>51,52</sup>, and the molecular processes of most complex diseases are now appreciated to be impacted by alterations to human gene expression levels<sup>53,54</sup>. We thus hypothesized that causative pathogen-disease relationships would be reflected in publicly available gene expression data. To test this hypothesis, we performed two complementary analyses. First, we examined viral gene expression levels in patients compared to controls. Second, we asked if genome-wide association study (GWAS) risk loci were enriched near human genes with altered expression levels following viral infection.

As a positive control, we first examined the well-established link between EBV and SLE. To this end, we identified six publicly available SLE case/control RNA-seq data sets performed in blood and B cell subsets (Supplementary Data 5). Collectively, these data contain 378 SLE cases and 74 control subjects. We used the VIRTUS software package<sup>31</sup> to identify and quantify viral read counts in these data. As expected, this analysis revealed significantly higher EBV transcript levels in SLE cases compared to controls ( $p$ -value =  $4.9 \times 10^{-3}$ ) (Fig. 5a, Supplementary Data 9).

We next considered ulcerative colitis (UC), a disease with a suggestive but still unclear role for viral infection<sup>55</sup>. Our main analyses implicated both HIV and CMV in UC disease processes, with the replication cohort showing infection occurred before disease diagnosis. We thus examined RNA-seq data for a set of 669 UC cases and 59 controls obtained from seven studies

using intestinal biopsies (Supplementary Data 5). Similar to EBV and SLE, we observe significantly higher levels of CMV transcripts relative to controls in these samples ( $p$ -value =  $2.2 \times 10^{-2}$ ) (Fig. 5b, Supplementary Data 9), providing additional evidence of a role for CMV in UC disease processes.

As a second orthogonal analysis, we asked if GWAS disease risk loci are enriched near human genes with virus-induced expression level changes. To this end, we used our RELI algorithm<sup>11</sup> to relate GWAS risk loci and public RNA-seq experiments examining virally infected and uninfected cells. In brief, this procedure uses a permutation-based method to estimate the significance of the overlap between the genomic coordinates of GWAS-identified risk loci and 200 kilobase windows around the transcription start site for genes with virus-altered gene expression levels (see Methods).

As a positive control, we first examined EBV and SLE. As expected, our analyses revealed significant overlap between SLE risk loci and genes that are differentially expressed upon EBV infection, in two independent studies performed in B lymphocytes and peripheral blood mononuclear cells (Fig. 5c, purple bars; Supplementary Data 10). In contrast, when considering genes that did not change significantly upon infection (Fig. 5c, gray bars), SLE risk loci are not enriched. These results are consistent with our previous observation that the genomic binding events of the EBNA2 regulatory protein, encoded by EBV, coincide with approximately half of all SLE risk loci<sup>11</sup>. Encouraged by these results, we next compared CMV-altered genes to UC genetic risk loci. Similar to the EBV-SLE results, we again observe highly significant overlap between CMV-altered genes and UC risk loci and insignificant overlap for expressed but unchanged genes (Fig. 5d, purple and gray bars, respectively; Supplementary Data 10) in two different cell types (monocytes and dendritic cells). Similar results were obtained for the highly related Crohn’s disease and inflammatory bowel disease (Supplementary Data 10).

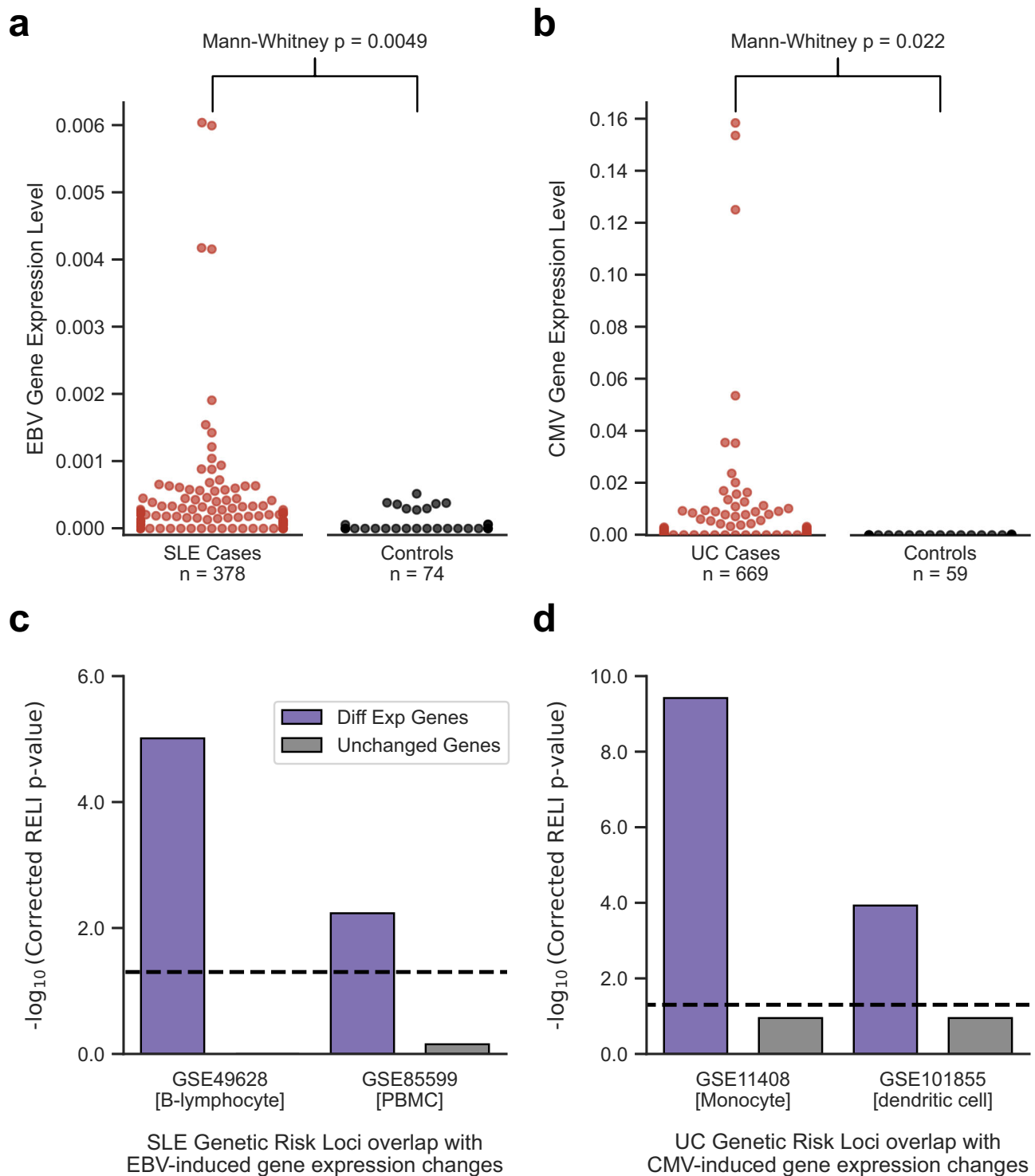
Collectively, these analyses provide compelling orthogonal evidence that the CMV-UC connection identified in both of our independent cohorts might represent a causative relationship.

## Discussion

In this study, we sought a broader understanding of the role played by pathogens in what are traditionally considered non-communicable diseases (NCDs). To this end, we developed a logistic regression model and applied the model to two large, independent biobank resources. Our model showed strong discriminatory performance on a set of positive and negative controls and replicated many additional well-documented pathogen-NCD associations. Overall, our results were robust regardless of the choice of International Classification of Diseases 10<sup>th</sup> revision (ICD10) codes or Phecodes. We report evidence for over 200 new, to the best of our knowledge, or previously tenuous pathogen-disease connections, including a role for cytomegalovirus (CMV) in ulcerative colitis (UC), which was supported by two orthogonal genomics-based validations, and provide the corresponding data on a freely accessible and easily browsable web server, <https://tf.cchmc.org/pathogen-disease>.

The relationship between *H. pylori* and gastroesophageal reflux disease (GERD) continues to be debated, with one recent meta-analysis reporting that the eradication of *H. pylori* increases the risk of GERD, thus indicating a possible protective effect<sup>56</sup> and another recent systematic review rated the “evidence grade” for the association between *H. pylori* and GERD as low<sup>56,57</sup>. Our replicated result indicating that *H. pylori* has a protective effect adds additional evidence of this debated association. Outside of ICD10 chapter XI, *H. pylori* also has a risk relationship with “iron-deficiency anemia”, an association that has been published previously<sup>58</sup>, and a protective relationship with asthma, which has also been previously reported<sup>59,60</sup>.

Our results implicated both human immunodeficiency virus (HIV) and CMV in ulcerative colitis (UC). HIV has recently been linked to UC<sup>61</sup>. For CMV, a possible role is much less clear, perhaps due to the historical difficulty of differentiating between CMV colitis and inflammatory bowel diseases such as UC<sup>55</sup>. Although the hypothesis that CMV may be causal of UC remains contentious<sup>62</sup>, the ability of CMV to cause UC flare-ups is still heavily debated<sup>63</sup>. Our results should aid in these ongoing debates, including



**Fig. 5 | Orthogonal validation of the EBV/SLE and CMV/UC associations.** Orthogonal validation of the Epstein-Barr virus (EBV)/systemic lupus erythematosus (SLE) (positive control) and cytomegalovirus (CMV)/ulcerative colitis (UC) (new prediction) associations using virus gene expression levels (top) and enrichment of disease risk loci near virus-induced differentially expressed human genes (bottom). **a, b** Swarm plots showing viral read counts normalized by the size of the viral genome and the total number of human mapped reads in that sample, providing the final ‘Normalized Hit Rate’ calculated by the VIRTUS software package, referred to on the plots as ‘Gene Expression Level’. Normalized hit rates were

compared between cases (red dots) and controls (black dots) using a Mann–Whitney test with the  $p$ -values annotated on the plots: SLE versus controls (panel **a**) and UC versus controls (panel **b**). Bar plots indicating the enrichment as a  $-\log_{10}(\text{Corrected RELI } p\text{-value})$  of SLE (left, **c**) or UC (right, **d**) genome-wide association study (GWAS) loci proximal to genes with altered expression (purple bars) or unaltered expression (gray bars) after infection by the viruses EBV and CMV, respectively. Gene Expression Omnibus (GEO) ID and cell type are provided below each plot. The black dashed line indicates statistical significance ( $p = 0.05$ ).

the replicated associations based on serologic and diagnostic data as well as our two orthogonal analyses, which all suggest a role for CMV in UC processes.

A recent study by Levine et al. took a similar approach to ours, with a specific focus on six neurodegenerative diseases<sup>64</sup>. Herein, we report results from an analysis across the disease spectrum. In addition to the comparatively limited scope of the Levine et al. study, looking at just six neurodegenerative diseases, there are several additional key differences between the 2023 Levine et al. study and ours. The Levine et al. study used hospital diagnosis codes as a pathogen proxy, some of which link to multiple pathogens, such as viral encephalitis. In comparison, we were able to pinpoint specific pathogens due to our use of serology data. Further, by restricting to inpatient hospital databases, the Levine et al. study focused on patients with infections sufficiently severe enough to require hospitalization. In contrast, our analyses included systematically measured titers for select UK Biobank (UKB) participants<sup>17</sup>, along with data from TriNetX (TNX), which pulls all available clinical laboratory test results for each patient, encompassing standard preventative screenings, outpatient diagnostic workups, as well as panels ordered during hospital stays. This is likely one reason why the odds ratios we report are much more modest than those reported by Levine et al. Thus, although both studies are of great utility, the results of the two studies are not directly comparable.

Roughly 40% (81/206) of the replicated associations in our study have odds ratios of less than one, indicating a potentially protective pathogen-disease relationship. For example, while high-risk strains of human papillomavirus (HPV) such as 16 and 18 are known to cause over 70% of cervical cancer, our results also suggest that HPV-16 and -18 can be “protective” for diseases such as “seborrheic dermatitis” and “other dermatitis”. Indeed, viral infections that increase the risk of one phenotype or disease can reduce the risk for others<sup>65–68</sup>. More generally, the role of viral infection in shaping the human immune system and subsequent immune responses has been studied extensively. It is well appreciated that viral infection can rewire the chromatin of immune cells and shape subsequent responses of a person toward additional inflammatory insults<sup>69,70</sup>. For example, a viral response that results in an interferon-based immune response could be protective in the context of diseases driven by T cell helper-2 type inflammation<sup>71</sup>. Vaccination against a particular pathogen will protect against infection and, thus, the disease risks associated with that pathogen. However, further studies will be required to investigate whether vaccination will confer the same protective effects against NCDs we report in this study.

The availability of the large datasets from the UK Biobank and TriNetX enabled this research. The associations identified in our study depend upon a sufficient number of subjects that have both accompanying diagnostic and serology data. As additional larger datasets are released, it will be critical to validate this study’s findings and use the additional statistical power to examine NCDs for which we were not powered. Identifying potentially causal etiological mechanisms driving these pathogen-disease associations will also be important, as recently exemplified by the Epstein-Barr virus (EBV) – multiple sclerosis (MS) field<sup>11,12,69,72,73</sup>. Furthermore, the goal of this study was to attempt to identify connections between pathogens and non-communicable diseases. Combinations of pathogens can also have impacts on human disease etiology. For example, Plasmodium falciparum and EBV have been shown to increase the risk of endemic Burkitt lymphoma synergistically<sup>74–76</sup>. Likewise, genetics likely plays a vital role in pathogen-host interactions. Future studies applying new methodologies to much larger cohorts than those presented here will be important for identifying novel combinations of pathogens and host genetic variant-pathogen interactions that impact human disease.

Although attempts were made to minimize limitations in this study, some remain. For example, by their very nature, electronic health records and biobank data are noisy. We attempted to address this noise by examining two independent datasets covering vastly separate geographic locations and requiring a pathogen-disease pair to be significant in both, to be considered replicated. Another limitation is that the datasets used in this

study are based on cohorts chiefly from only two countries (the United States and the United Kingdom). Accordingly, the 20 pathogens investigated are largely prevalent and of most importance to the people of those regions as compared to the rest of the global population. An additional limitation of the study includes the possibility of cross-reactivity occurring while testing for a particular pathogen. While each of the serological tests was approved by the Clinical Laboratory Improvement Amendments program and used to inform clinical care of patients, it is possible that some non-viral human protein epitopes cross-reacted with the viral antigens in the serological tests<sup>12,72,77</sup>. Finally, not all confounders in the UK Biobank models could be adjusted for in the replication cohort because the data were unavailable.

In summary, we present the largest systematic assessment to date of pathogens in the context of non-communicable human disease. Using complementary discovery and replication datasets, we identified 206 replicated pathogen-disease relationships, including additional orthogonal evidence strongly supporting a relationship between CMV infection and ulcerative colitis. We anticipate that this rich data resource will form the foundation for future characterization of the many currently unknown pathogen-disease relationships.

## Data availability

The data used for the main analysis from The UK Biobank (UKB) can be accessed via an application for Tier 2 UKB data. Further, the data from TriNetX (<https://live.trinetx.com>, accessed on 14 February 2023) can be requested directly from TriNetX. In both cases, costs may be incurred, and a material transfer agreement or data-sharing agreement is required. The data sources for each figure are as follows: Fig. 1 - None (pictorial overview of the methods); Fig. 2a - Supplementary Data 8; Fig. 2b - Supplementary Data 2; Fig. 3 - Supplementary Data 8; Fig. 4 - Supplementary Data 8; Fig. 5a, b - Supplementary Data 9; Fig. 5c, d - Supplementary Data 10; Supplementary Fig. S1 - None (two formulas); Supplementary Fig. S2 - Supplementary Data 8; Supplementary Fig. S3 - Supplementary Data 8; Supplementary Fig. S4 - Supplementary Data 8; Supplementary Fig. S5 - Supplementary Data 8; Supplementary Fig. S6 - Supplementary Data 8.

## Code availability

All code used in this project is available from the Weirauch Lab’s Research GitHub page, [https://github.com/WeirauchLab/pathogen\\_ncd](https://github.com/WeirauchLab/pathogen_ncd)<sup>78</sup>. The code is also stored in Zenodo and can be accessed via the <https://doi.org/10.5281/zenodo.14985960><sup>78</sup>.

Received: 20 October 2023; Accepted: 9 June 2025;

Published online: 20 June 2025

## References

1. Mandell, Douglas, and Bennett’s principles and practice of infectious diseases (Elsevier, 2020).
2. Borchers, A. T., Chang, C., Gershwin, M. E. & Gershwin, L. J. Respiratory syncytial virus—a comprehensive review. *Clin. Rev. Allergy Immunol.* **45**, 331–379 (2013).
3. Gershon, A. A. et al. Varicella zoster virus infection. *Nat. Rev. Dis. Prim.* **1**, 15016 (2015).
4. Hardbower, D. M., Peek, R. M. & Wilson, K. T. At the bench: *Helicobacter pylori*, dysregulated host responses, DNA damage, and gastric cancer. *J. Leukoc. Biol.* **96**, 201–212 (2014).
5. Lauer, G. M. & Walker, B. D. Hepatitis C virus infection. *N. Engl. J. Med.* **345**, 41–52 (2001).
6. Trépo, C., Chan, H. L. Y. & Lok, A. Hepatitis B virus infection. *Lancet* **384**, 2053–2063 (2014).
7. zur Hausen, H. Papillomaviruses in the causation of human cancers — a brief historical account. *Virology* **384**, 260–265 (2009).
8. Bouvard, V. et al. A review of human carcinogens—Part B: biological agents. *Lancet Oncol.* **10**, 321–322 (2009).



9. Plummer, M. et al. Global burden of cancers attributable to infections in 2012: a synthetic analysis. *Lancet Glob. Health* **4**, e609–e616 (2016).
10. Bjornevik, K. et al. Longitudinal analysis reveals high prevalence of Epstein–Barr virus associated with multiple sclerosis. *Science* **375**, 296–301 (2022).
11. Harley, J. B. et al. Transcription factors operate across disease loci, with EBNA2 implicated in autoimmunity. *Nat. Genet.* **50**, 699–707 (2018).
12. Lanz, T. V. et al. Clonally expanded B cells in multiple sclerosis bind EBV EBNA1 and GlatIRAM. *Nature* **603**, 321–327 (2022).
13. Thacker, E. L., Mirzaei, F. & Ascherio, A. Infectious mononucleosis and risk for multiple sclerosis: a meta-analysis. *Ann. Neurol.* **59**, 499–503 (2006).
14. Belbasis, L., Bellou, V., Evangelou, E., Ioannidis, J. P. A. & Tzoulaki, I. Environmental risk factors and multiple sclerosis: an umbrella review of systematic reviews and meta-analyses. *Lancet Neurol.* **14**, 263–273 (2015).
15. Virolainen, S. J., VonHandorf, A., Viel, K. C. M. F., Weirauch, M. T. & Kottyan, L. C. Gene–environment interactions and their impact on human health. *Genes Immun.* **24**, 1–11 (2023).
16. Bycroft, C. et al. The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203 (2018).
17. Mentzer, A. J. et al. Identification of host–pathogen–disease relationships using a scalable multiplex serology platform in UK Biobank. *Nat. Commun.* **13**, 1818 (2022).
18. Tange, O. GNU Parallel 20220122 (‘20 years’). *Zenodo* <https://doi.org/10.5281/zenodo.5893336> (2022).
19. Wei, W.-Q. et al. Evaluating phecodes, clinical classification software, and ICD-9-CM codes for phenome-wide association studies in the electronic health record. *PLoS One* **12**, e0175508 (2017).
20. Wu, P. et al. Mapping ICD-10 and ICD-10-CM codes to phecodes: workflow development and initial evaluation. *JMIR Med. Inform.* **7**, e14325 (2019).
21. Denaxas, S. et al. Mapping the Read2/CTV3 controlled clinical terminologies to Phecodes in UK Biobank primary care electronic health records: implementation and evaluation. *AMIA Annu. Symp. Proc. AMIA Symp.* **2021**, 362–371 (2021).
22. Faul, F., Erdfelder, E., Buchner, A. & Lang, A.-G. Statistical power analyses using G\*Power 3.1: tests for correlation and regression analyses. *Behav. Res. Methods* **41**, 1149–1160 (2009).
23. Firth, D. Bias reduction of maximum likelihood estimates. *Biometrika* **80**, 27–38 (1993).
24. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. R. Stat. Soc. Ser. B Methodol.* **57**, 289–300 (1995).
25. Hu, Y. et al. Analysis of genomic and proteomic data using advanced literature mining. *J. Proteome Res.* **2**, 405–412 (2003).
26. Febbo, P. G. et al. Literature Lab: a method of automated literature interrogation to infer biology from microarray analysis. *BMC Genomics* **8**, 461 (2007).
27. Ng, Q. X. et al. Is there an association between Helicobacter pylori infection and irritable bowel syndrome? A meta-analysis. *World J. Gastroenterol.* **25**, 5702–5710 (2019).
28. Kim, Y.-A., Cho, Y. J. & Kwak, S. G. The association between Helicobacter pylori infection and irritable bowel syndrome: a meta-analysis. *Int. J. Environ. Res. Public Health* **17**, 2524 (2020).
29. Li, C., Shuai, Y., Zhou, X. & Chen, H. Association between Helicobacter pylori infection and irritable bowel syndrome: a systematic review and meta-analysis. *Medicine* **99**, e22975 (2020).
30. Wang, Z., Liu, Y., Peng, Y. & Peng, L. Helicobacter pylori infection—a risk factor for irritable bowel syndrome? An updated systematic review and meta-analysis. *Med. Kaunas. Lith.* **58**, 1035 (2022).
31. Yasumizu, Y., Hara, A., Sakaguchi, S. & Ohkura, N. VIRTUS: a pipeline for comprehensive virus analysis from conventional RNA-seq data. *Bioinformatics* **37**, 1465–1467 (2021).
32. Dobin, A. et al. STAR: ultrafast universal RNA-seq aligner. *Bioinform. Oxf. Engl.* **29**, 15–21 (2013).
33. Buniello, A. et al. The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res.* **47**, D1005–D1012 (2019).
34. Purcell, S. et al. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
35. Dexheimer, P. J., Pujato, M., Roskin, K. M. & Weirauch, M. T. VExD: a curated resource for human gene expression alterations following viral infection. *G3 GenesGenomesGenetics* jkad176, <https://doi.org/10.1093/g3journal/jkad176> (2023).
36. Naito, T. et al. Causes of infectious mononucleosis-like syndrome in adult patients. *Intern. Med.* **45**, 833–834 (2006).
37. Ahmed, S. & Belayneh, Y. M. Helicobacter pylori and duodenal ulcer: systematic review of controversies in causation. *Clin. Exp. Gastroenterol.* **12**, 441–447 (2019).
38. Suerbaum, S. & Michetti, P. Helicobacter pylori Infection. *N. Engl. J. Med.* **347**, 1175–1186 (2002).
39. Kandulski, A., Selgrad, M. & Malfertheiner, P. Helicobacter pylori infection: a clinical overview. *Dig. Liver Dis.* **40**, 619–626 (2008).
40. Smith, A., Baumgartner, K. & Bositis, C. Cirrhosis: diagnosis and management. *Am. Fam. Phys.* **100**, 759–770 (2019).
41. Bar-Or, A. et al. Epstein–Barr virus in multiple sclerosis: theory and emerging immunotherapies. *Trends Mol. Med.* **26**, 296–310 (2020).
42. James, J. A. et al. An increased prevalence of Epstein–Barr virus infection in young patients suggests a possible etiology for systemic lupus erythematosus. *J. Clin. Invest.* **100**, 3019–3026 (1997).
43. Li, Z.-X., Zeng, S., Wu, H.-X. & Zhou, Y. The risk of systemic lupus erythematosus associated with Epstein–Barr virus infection: a systematic review and meta-analysis. *Clin. Exp. Med.* **19**, 23–36 (2019).
44. Denny, J. C. et al. PheWAS: demonstrating the feasibility of a phenome-wide scan to discover gene–disease associations. *Bioinformatics* **26**, 1205–1210 (2010).
45. Denny, J. C. et al. Systematic comparison of phenome-wide association study of electronic medical record data and genome-wide association study data. *Nat. Biotechnol.* **31**, 1102–1110 (2013).
46. Gross, A. J., Hochberg, D., Rand, W. M. & Thorley-Lawson, D. A. EBV and systemic lupus erythematosus: a new perspective. *J. Immunol. Baltim. Md.* **174**, 6599–6607 (2005).
47. Poole, B. D. et al. Aberrant Epstein–Barr viral infection in systemic lupus erythematosus. *Autoimmun. Rev.* **8**, 337–342 (2009).
48. Piroozmand, A., Haddad Kashani, H. & Zamani, B. Correlation between Epstein–Barr virus infection and disease activity of systemic lupus erythematosus: a cross-sectional study. *Asian Pac. J. Cancer Prev. APJCP* **18**, 523–527 (2017).
49. Moon, U. Y. et al. Patients with systemic lupus erythematosus have abnormally elevated Epstein–Barr virus load in blood. *Arthritis Res. Ther.* **6**, R295–R302 (2004).
50. Agostini, S. et al. HLA alleles modulate EBV viral load in multiple sclerosis. *J. Transl. Med.* **16**, 80 (2018).
51. Paschos, K. & Allday, M. J. Epigenetic reprogramming of host genes in viral and microbial pathogenesis. *Trends Microbiol.* **18**, 439–447 (2010).
52. Liu, X. et al. Human virus transcriptional regulators. *Cell* **182**, 24–37 (2020).
53. Maurano, M. T. et al. Systematic localization of common disease-associated variation in regulatory DNA. *Science* **337**, 1190–1195 (2012).

54. Oh, E. S. & Petronis, A. Origins of human disease: the chrono-epigenetic perspective. *Nat. Rev. Genet.* **22**, 533–546 (2021).
55. Mourad, F. H., Hashash, J. G., Kariyawasam, V. C. & Leong, R. W. Ulcerative colitis and cytomegalovirus infection: from A to Z. *J. Crohns Colitis* **14**, 1162–1171 (2020).
56. Mou, W.-L., Feng, M.-Y. & Hu, L.-H. Eradication of *Helicobacter pylori* infections and GERD: a systematic review and meta-analysis. *Turk. J. Gastroenterol.* **31**, 853–859 (2020).
57. Cheng, Y. et al. Systematic assessment of environmental factors for gastroesophageal reflux disease: an umbrella review of systematic reviews and meta-analyses. *Dig. Liver Dis.* **53**, 566–573 (2021).
58. DuBois, S. & Kearney, D. J. Iron-deficiency anemia and *Helicobacter pylori* infection: a review of the evidence. *J. Am. Coll. Gastroenterol. ACG* **100**, 453 (2005).
59. Amedei, A., Codolo, G., Del Prete, G., de Bernard, M. & D'Elios, M. M. The effect of *Helicobacter pylori* on asthma and allergy. *J. Asthma Allergy* **3**, 139–147 (2010).
60. Zuo, Z. T. et al. The protective effects of *Helicobacter pylori* infection on allergic asthma. *Int. Arch. Allergy Immunol.* **182**, 53–64 (2021).
61. Elmahdi, R. et al. Development of inflammatory bowel disease in HIV patients: a Danish cohort study (1983–2018) with American validation (1999–2018). *Gastro. Hep. Adv.* **1**, 1114–1121 (2022).
62. Beswick, L., Ye, B. & van Langenberg, D. R. Toward an algorithm for the diagnosis and management of CMV in patients with colitis. *Inflamm. Bowel Dis.* **22**, 2966–2976 (2016).
63. Jentzer, A. et al. Cytomegalovirus and inflammatory bowel diseases (IBD) with a special focus on the link with ulcerative colitis (UC). *Microorganisms* **8**, 1078 (2020).
64. Levine, K. S. et al. Virus exposure and neurodegenerative disease risk across national biobanks. *Neuron* **111**, 1086–1093.e2 (2023).
65. Bach, J.-F. The effect of infections on susceptibility to autoimmune and allergic diseases. *N. Engl. J. Med.* **347**, 911–920 (2002).
66. Cooper, P. J. Interactions between helminth parasites and allergy. *Curr. Opin. Allergy Clin. Immunol.* **9**, 29–37 (2009).
67. Garn, H. & Renz, H. Epidemiological and immunological evidence for the hygiene hypothesis. *Immunobiology* **212**, 441–452 (2007).
68. Kilpeläinen, M., Terho, E. O., Helenius, H. & Koskenvuo, M. Farm environment in childhood prevents the development of allergies. *Clin. Exp. Allergy J. Br. Soc. Allergy Clin. Immunol.* **30**, 201–208 (2000).
69. Hong, T. et al. Epstein-Barr virus nuclear antigen 2 extensively rewires the human chromatin landscape at autoimmune risk loci. *Genome Res.* **31**, 2185–2198 (2021).
70. Wang, R. et al. SARS-CoV-2 restructures host chromatin architecture. *Nat. Microbiol.* **8**, 679–694 (2023).
71. Okada, H., Kuhn, C., Feillet, H. & Bach, J.-F. The ‘hygiene hypothesis’ for autoimmune and allergic diseases: an update. *Clin. Exp. Immunol.* **160**, 1–9 (2010).
72. Thomas, O. G. et al. Cross-reactive EBNA1 immunity targets alpha-crystallin B and is associated with multiple sclerosis. *Sci. Adv.* **9**, eadg3032 (2023).
73. Keane, J. T. et al. The interaction of Epstein-Barr virus encoded transcription factor EBNA2 with multiple sclerosis risk loci is dependent on the risk genotype. *EBioMedicine* **71**, 103572 (2021).
74. Chene, A. et al. Endemic Burkitt’s lymphoma as a polymicrobial disease: New insights on the interaction between *Plasmodium falciparum* and Epstein-Barr virus. *Semin. Cancer Biol.* **19**, 411–420 (2009).
75. Thorley-Lawson, D., Deitsch, K. W., Duca, K. A. & Torgbor, C. The Link between *Plasmodium falciparum* malaria and endemic Burkitt’s lymphoma – new insight into a 50-year-old enigma. *PLOS Pathog.* **12**, e1005331 (2016).
76. Quintana, M., del, P., Smith-Togobo, C., Moormann, A. & Hviid, L. Endemic Burkitt lymphoma – an aggressive childhood cancer linked to *Plasmodium falciparum* exposure, but not to exposure to other malaria parasites. *APMIS* **128**, 129–135 (2020).
77. Singh, D. et al. Antibodies to an Epstein Barr Virus protein that cross-react with dsDNA have pathogenic potential. *Mol. Immunol.* **132**, 41–52 (2021).
78. Lape, M. & Ernst, K. WeirauchLab/pathogen\_ncd: communications medicine publication. *Zenodo*, <https://doi.org/10.5281/ZENODO.14985960> (2025).

## Acknowledgements

We thank Ivan Marazzi, Chris Benner, Brad Rosenberg, Emily Miraldi, Juan Fuxman-Bass, and Artem Babaian for their insightful comments on the manuscript. We also thank Matthew Skowronek, Eric Young, Jeff Warnick, and the rest of the engineering team at TriNetX for additional assistance working with the TriNetX data. We thank Monika Grabowska, Wei-Qi Wei, and Spiros Denaxas for their invaluable help with the Phecode analysis. We thank The University of Cincinnati Center for Clinical and Translational Science and Training (CCTST) for assistance with accessing the TriNetX data. This research was conducted using the UK Biobank Resource under application number 47377. TriNetX, LLC is compliant with the Health Insurance Portability and Accountability Act (HIPAA), the US federal law which protects the privacy and security of healthcare data, and any additional data privacy regulations applicable to the contributing healthcare organization. CCTST is supported by the National Center for Advancing Translational Sciences of the National Institutes of Health (NIH), under Award Number 2UL1TR001425-05A1. The content is solely the responsibility of the authors and does not necessarily represent the official views of the NIH. This research was funded by National Institutes of Health (NIH) R01 DK107502, R01 AI148276, U01 HG011172, U19 AI070235, and P30 AR070549 to L.C.K.; R01 HG010730, R01 GM055479, and U01 AI130830, to M.T.W.; R01 AR073228, R01 NS099068, and R01 AI024717 to M.T.W. and L.C.K.; R01 CA226802 to N.S. Funding was also provided by Cincinnati Children’s Hospital Medical Center ARC Award 53632 to M.T.W. and L.C.K.

## Author contributions

Conceptualization: M.L., L.C.K., M.T.W.; Study design: M.L., L.C.K., M.T.W.; Statistical analysis design: M.L., D.S., L.M., M.T.W.; Computational analysis: M.L., S.P.; Writing: M.L., L.C.K., M.T.W.; Review and approve final manuscript: M.L., D.S., S.P., N.S., B.H., L.M., L.C.K., M.T.W.; Funding: N.S., L.C.K., M.T.W.

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s43856-025-00956-x>.

**Correspondence** and requests for materials should be addressed to Leah C. Kottyan or Matthew T. Weirauch.

**Peer review information** *Communications Medicine* thanks Olamide Arege and the other, anonymous, reviewer(s) for their contribution to the peer review of this work. [Peer review reports are available].

**Reprints and permissions information** is available at <http://www.nature.com/reprints>

**Publisher’s note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Open Access** This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

© The Author(s) 2025