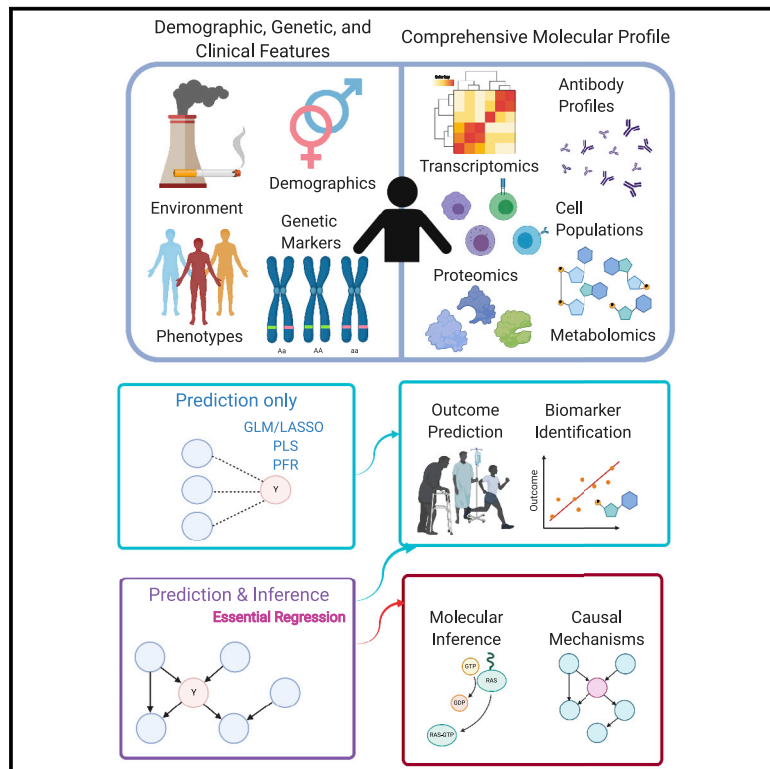


Patterns

Essential Regression: A generalizable framework for inferring causal latent factors from multi-omic datasets

Graphical abstract



Highlights

- ER is a novel interpretable machine-learning method for high-dimensional multi-omic data
- ER outperforms a wide range of state-of-the-art methods in terms of prediction
- Beyond prediction, ER identifies causal latent factors of groups/outcomes of interest
- ER generated novel immunological inferences, consistent with evidence in model organisms

Authors

Xin Bing, Tyler Lovelace, Florentina Bunea, ..., Harinder Singh, Panayiotis V. Benos, Jishnu Das

Correspondence

harinder@pitt.edu (H.S.),
benos@pitt.edu (P.V.B.),
jishnu@pitt.edu (J.D.)

In brief

Current analytical approaches for multi-omic datasets are limited by high dimensionality, differences in data distributions, and causal inference beyond prediction. Here, we present Essential Regression (ER), a novel latent-factor-regression-based interpretable machine-learning approach that integrates high-dimensional multi-omic datasets without distributional assumptions regarding the data and identifies significant latent factors and their causal relationships with system-wide outcomes/properties of interest. ER outperforms a range of state-of-the-art methods in terms of prediction and generates novel immunological inferences, consistent with evidence in model organisms.



Article

Essential Regression: A generalizable framework for inferring causal latent factors from multi-omic datasets

Xin Bing,^{1,7} Tyler Lovelace,^{2,3,7} Florentina Bunea,¹ Marten Wegkamp,^{1,4} Sudhir Pai Kasturi,⁵ Harinder Singh,^{6,*} Panayiotis V. Benos,^{2,*} and Jishnu Das^{6,8,*}

¹Department of Statistics and Data Science, Cornell University, Ithaca, NY, USA

²Department of Computational & Systems Biology, University of Pittsburgh, Pittsburgh, PA, USA

³Joint CMU-Pitt PhD Program in Computational Biology, Carnegie Mellon – University of Pittsburgh, Pittsburgh, PA, USA

⁴Department of Mathematics, Cornell University, Ithaca, NY, USA

⁵Division of Microbiology and Immunology, Yerkes National Primate Research Center, Emory University, Atlanta, GA, USA

⁶Center for Systems Immunology, Departments of Immunology and Computational & Systems Biology, University of Pittsburgh, Pittsburgh, PA, USA

⁷These authors contributed equally

⁸Lead contact

*Correspondence: harinder@pitt.edu (H.S.), benos@pitt.edu (P.V.B.), jishnu@pitt.edu (J.D.)

<https://doi.org/10.1016/j.patter.2022.100473>

THE BIGGER PICTURE Multi-omic technologies for deep cellular and molecular profiling from model organisms or humans have rapidly expanded. However, existing analytical approaches are constrained by the high dimensionality of these datasets, differences in data distributions, and the inability to generate causal inference beyond predictive biomarkers. To address these issues, we developed a novel interpretable machine-learning framework, Essential Regression (ER). ER integrates high-dimensional multi-omic datasets without distributional assumptions regarding the data and identifies significant latent factors and their causal relationships with system-wide outcomes/properties of interest. ER uses higher-order relationships encapsulated in the latent factors, rather than the individual observables, to home in on novel mechanistic insights. Our approach outperforms a range of state-of-the-art methods in terms of prediction and generates novel immunological inferences, consistent with evidence in model organisms.



Proof-of-Concept: Data science output has been formulated, implemented, and tested for one domain/problem

SUMMARY

High-dimensional cellular and molecular profiling of biological samples highlights the need for analytical approaches that can integrate multi-omic datasets to generate prioritized causal inferences. Current methods are limited by high dimensionality of the combined datasets, the differences in their data distributions, and their integration to infer causal relationships. Here, we present Essential Regression (ER), a novel latent-factor-regression-based interpretable machine-learning approach that addresses these problems by identifying latent factors and their likely cause-effect relationships with system-wide outcomes/properties of interest. ER can integrate many multi-omic datasets without structural or distributional assumptions regarding the data. It outperforms a range of state-of-the-art methods in terms of prediction. ER can be coupled with probabilistic graphical modeling, thereby strengthening the causal inferences. The utility of ER is demonstrated using multi-omic system immunology datasets to generate and validate novel cellular and molecular inferences in a wide range of contexts including immunosenescence and immune dysregulation.



INTRODUCTION

Over the last decade, genomic, proteomic, metabolomic, and other technologies for generating deep molecular profiles of tissues and cells from model organisms or humans have rapidly expanded.^{1–4} However, the explosion in data, especially from a range of such “omic” technologies, has not been coupled to a proportional increase in our understanding of the underlying causal mechanisms. Existing analytical approaches have primarily focused on individual omic datasets, with relatively few attempts at integration of multi-omic datasets. In either case, we^{5–9} and others^{10–12} have primarily emphasized the delineation of predictive biomarkers with limited exploration of putative causal factors based on prior biological knowledge (Figure 1A). A key focus of these efforts has been to overcome the “curse of dimensionality” (a very large number of variables being measured in relation to a comparatively low number of samples) and the multiplicity of predictive signatures due to multi-collinear data, i.e., large correlated sets of variables. While there are several methods for reliably uncovering predictive markers from high-dimensional data, none of these analyze cause-effect relationships in relation to the outcomes/outputs of interest. This in turn has hampered efforts to undertake perturbative/translational experiments and/or clinical investigations that can test a functionally prioritized set of hypotheses generated by the large datasets.

In addition to the high dimensionality of datasets at any given scale of organization (e.g., cellular, molecular), biological systems, particularly in humans, manifest extreme complexity in terms of numbers of molecular components and their interaction rules as well as their hierarchical scales of organization, which include macromolecular complexes/condensates, organelles, cells, tissues, and organs. Each scale of organization in such a complex system has components and interaction rules that are unique to its level of organization. Thus, predicting changes in properties or behaviors of the system based on measuring components that are operating at different scales of organization represents a formidable challenge. Methods that make assumptions regarding data-generating mechanisms typically perform poorly at multi-scale integration as there are key differences in data distributions at each scale of organization.

We propose a novel framework, Essential Regression (ER), to address these key challenges and limitations of existing approaches by focusing on latent factors rather than observables in high-dimensional datasets that are significantly associated with a system-wide property or outcome that is of interest (Figure 1A). Critically, ER makes no assumptions regarding the underlying data distributions, enabling principled integration of multi-omic datasets. ER is also fundamentally different from three kinds of modern approaches. The first kind of approaches are designed specifically for multi-modal single-cell data,¹³ i.e., they require single-cell data as inputs. These are constrained by structural and/or distributional requirements. ER can work on any multi-omic datasets as there are no structural assumptions regarding the data; it can even combine bulk and single-cell multi-omic datasets. The second set of approaches require prior knowledge.¹⁴ However, ER works without the use of any priors, making it suitable across contexts even when prior knowledge is weak or unavailable. The third set of approaches^{15,16}

provide accurate prediction (i.e., predictive markers/correlates) from high-dimensional multi-collinear multi-omic datasets but not meaningful inference with provable statistical guarantees. ER uses regression on the latent factors rather than the observables, a novel statistical framework that comes with rigorous guarantees regarding both prediction and inference.

Overall, our analytical framework derives causal latent factors from thousands of variables from multi-omics datasets across various scales of biological organization (Figure 1A). After identifying significant latent factors, ER can be coupled with causal graphical-model analyses to examine the connectivity of these factors to the system-wide property or outcome of interest. In so doing, ER generates a high-confidence and prioritized set of latent factors comprised of known observables that are most proximal in the causal graph network to the system property/outcome of interest. We note that while causal discovery approaches have become popular over the last two decades, they have been confined to low-dimensional datasets due to the associated computational complexity. ER overcomes this fundamental conceptual limitation by first identifying latent factors from the observables (which achieves an inherent dimensionality reduction) and then identifying which latent factors are causally linked to the outcome/system-wide property of interest.

By analyzing both simulated and real-world immunological multi-omic datasets, we demonstrate that ER and the associated causal graphical-model analyses significantly outperform a wide range of state-of-the-art approaches in predicting outcomes and provide multi-scale inferences not afforded by the existing methods. The novel causal predictions are corroborated by biological findings in relevant experimental systems.

RESULTS

ER: A novel data-distribution-free statistical regression framework for inferring causal latent factors

We present ER, a novel data-distribution-free latent factor regression approach that integrates high-dimensional multi-omic datasets and identifies latent factors that are significantly associated with a system property/outcome (Figures 1B and 1C; experimental procedures; Note S1). ER is a paradigm-altering concept in regression analysis for high-dimensional datasets i.e., datasets where the number of features exceeds the number of samples. Existing regression methods use techniques including regularization (e.g., L1 regularization: least absolute selection and shrinkage operator [LASSO],¹⁵ L1 + L2 regularization: Elastic Net), bootstrap aggregation (e.g., random forest¹⁶), or the incorporation of pre-specified group structures (e.g., group LASSO) to avoid overfitting. However, biomarkers/features identified using these approaches are merely predictive/correlative and may have no connection with the underlying mechanisms driving the system property/outcome of interest. ER, on the other hand, uses a two-step approach that allows for the identification of latent factors that can be used to infer causal structures underlying the system property/outcome. ER first finds latent factors in a data-dependent fashion, without the need for a pre-specified group structure. Of all latent factors, ER identifies a specific subset of latent factors that can be used to infer causal associations with the property/outcome of interest. Critically, ER makes no assumptions regarding the

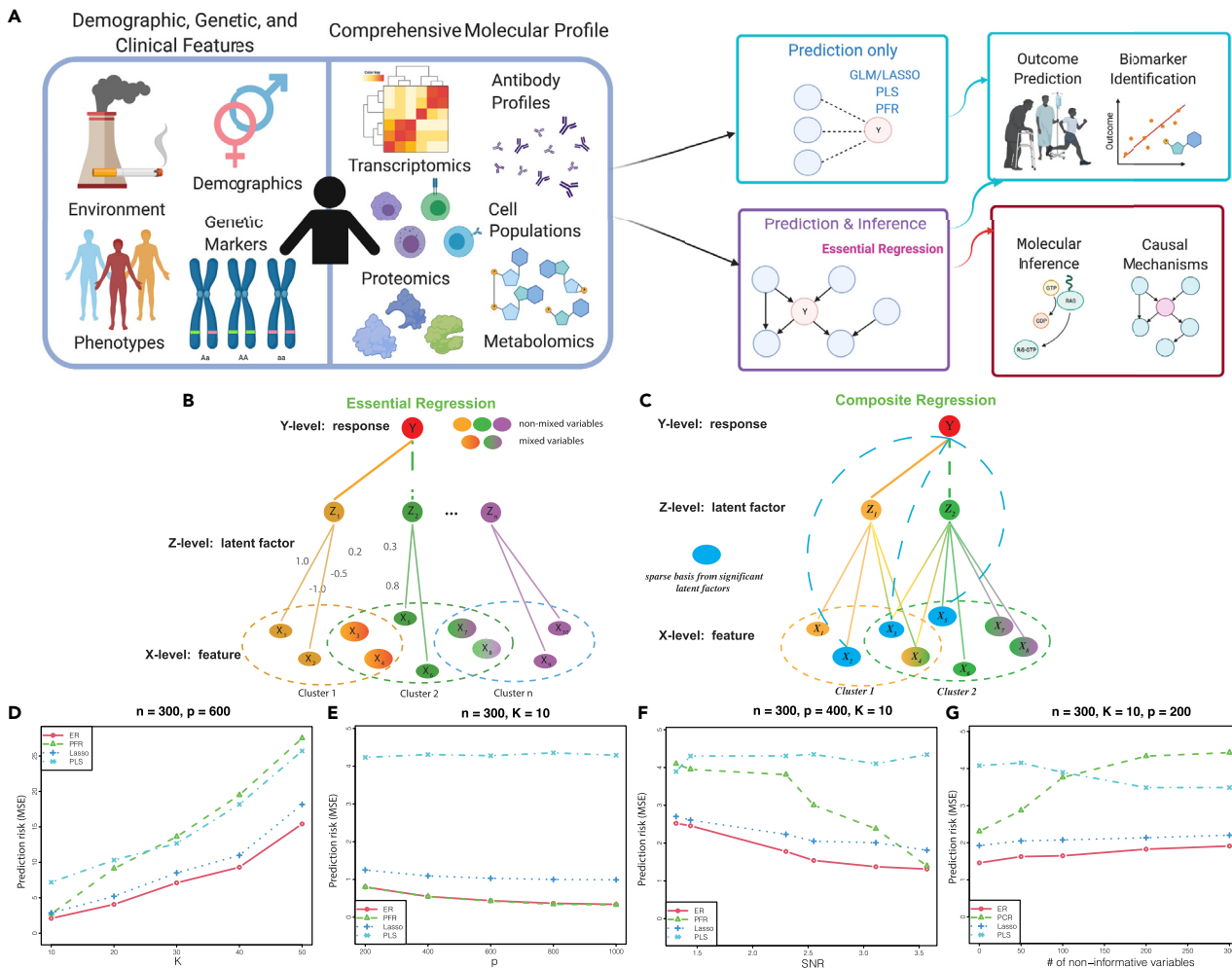


Figure 1. Essential Regression: A novel interpretable machine-learning approach to uncover causal latent factors from high-dimensional multi-omic datasets

(A) Schematic illustrating the different kinds of multi-omic datasets typically used in systems analyses and the key advantages of the methods introduced in this study over existing approaches.

(B) Schematic summarizing the steps in ER.

(C) Schematic summarizing the steps in Composite Regression.

(D–G) Comparison of the predictive performance of PLS, PFR, LASSO, and ER on simulated datasets across a range of parameter settings.

underlying data generating mechanisms and can be broadly used across multi-omic datasets (experimental procedures; Note S1).

Formally, ER is a latent-factor regression model in which the unobservable factor Z influences linearly both the response Y and the data X . Its novelty lies in the formulation that enables the latent factor Z to be meaningfully interpreted.

$$X = AZ + W$$

$$Y = \beta^T Z + \epsilon$$

Here, X is the matrix of observables and belongs to a high (p) dimensional space (dimensionality of X is $p \times n$, where n is the number of observations/samples). X is decomposed into the allocation matrix A of dimension $p \times K$, and Z is the latent-factor

matrix of dimension $K \times n$, i.e., it reduces X from a p to a K dimensional space (Note S1). The matrix Z is used to regress to Y , i.e., the regression coefficients correspond to Z s. W and ϵ are independent error terms (Note S1). This formulation helps cluster the observables (X s) into overlapping clusters/latent factors (Z s) in a data-dependent fashion and then identify which of the latent factors are significantly associated with and can be used to infer the outcome.

ER comes with two key provable statistical guarantees. The first step is to decompose the matrix of observables X into the latent factor matrix Z . To do this, the membership matrix A needs to be identifiable, up to a $K \times K$ signed permutation matrix. The first guarantee ensures this: we prove that the allocation matrix (A) is indeed identifiable up to a $K \times K$ signed permutation matrix under the assumption that there are at least 2 observables anchoring each latent factor (Note S1).¹⁷ This is a reasonable

assumption as the model only requires each latent factor to be defined by two observables that are not associated with other latent factors; all other observables may or may not be associated with multiple latent factors. This allows for the identification of a group structure from the observables entirely in a data-dependent fashion without the need to incorporate any prior knowledge. The second guarantee relates to the identifiability of the regression coefficients. We also prove that the coefficient matrix is identifiable up to a signed permutation matrix (Note S1),¹⁸ ensuring that the model can rigorously infer significant latent factors driving the outcome.

It is important to note that our current model assumes linearity at 2 different levels: (1) between the observables (X s) and the latent factors (Z s) and (2) between the outcome/response variable of interest (Y) and the latent factors (Z s). However, this does not necessarily translate into a linearity assumption between the X s and Y (it only translates into a linearity assumption when X and Y are Gaussian). Thus, the current model can incorporate non-linear relationships between X and Y , when X and Y are not Gaussian. Further, the linearity assumption between Y and Z is reasonable even for moderately large K (number of latent factors) as $K \ll p$ (p is the dimensionality of the original dataset). Further, both our algorithm and associated theoretical guarantees are still valid for a moderate and large K , even when $K \sim n$ (n = number of samples). Thus, ER is a first-in-class interpretable machine-learning framework that can uncover significant latent factors associated (linearly or, in many instances, non-linearly) with any system-wide property/outcome of interest.

ER outperforms state-of-the-art approaches on simulated high-dimensional datasets

We investigated the performance of ER on simulated data (Note S1), comparing it with a suite of state-of-the-art approaches including LASSO,¹⁵ partial least squares (PLS) regression,¹⁹ and principal components/factors regression (PFR).²⁰ We evaluated the performance of ER, PFR, PLS, and LASSO changes across a range of parameters for original dimensionality (p), reduced dimensionality of the dataset (i.e., K), and the signal-to-noise ratio (SNR). We varied these parameters one at a time and computed the corresponding prediction risk (mean squared error) on data not used to build the model (Note S1). We did a grid search on the relevant parameters and found that the prediction error for all four methods deteriorates as K increases or the SNR decreases (Figures 1D–1F). This indicates that prediction becomes more difficult for large K and a small SNR. On the other hand, ER, PFR, and Lasso perform better as p increases.

Among the four methods, ER systematically had the smallest prediction error in all settings, and PLS had the worst performance in most settings. Furthermore, PFR failed to accurately identify K and tended to a very low and sub-optimal K in most scenarios (Figures 1D–1F). This also indicates that, for principal-component regression approaches, detecting K requires a larger SNR, i.e., the other approaches are able to accurately detect K at lower SNRs. In a moderate SNR regime, PFR has comparable performance to ER (Figure 1D). However, as K increases, the advantage of ER becomes considerable, which supports the fact that PFR only has guarantees for fixed K (Figure 1E). Further, the performance of PFR is more sensitive to the SNR compared with the other three methods (Figure 1F). Finally,

when increasing the number of uninformative variables, ER has the best performance (Figure 1G). Overall, ER worked very well for very high p , was able to accurately identify K , and did not have a significant reduction in performance at lower SNR regimes or with a higher number of uninformative variables (Figures 1D–1G), outperforming state-of-the-art approaches on one or more of these fronts. ER functions counterintuitively when challenged by the curse of dimensionality (i.e., having higher dimensionality is worse as it induces higher variance and can lead to overfitting). The higher dimensionality of the datasets generates more features that provide additional information, which are used by ER to predict the latent factors (Z) more accurately, thereby overcoming the curse of dimensionality.

Extension of ER as composite regression enables uncovering of observables, within significant latent factors, that underlie outcomes

While the significant latent factors uncovered by ER provide insights into the interplay of the different observables driving outcome, in some contexts their complexity can prove challenging. In these instances, smaller sets of observables underlying outcome are desirable. Currently, regularization is widely used to identify a sparse set of observables (biomarkers). However, regularization-based approaches such as LASSO or Elastic Net uncover predictive biomarkers that may simply be correlative. Given that ER identifies latent factors significantly driving outcome, we sought to develop an approach to identify a sparse set of observables from the significant latent factors identified by ER (Figure 1C). Using L1-regularization on the significant latent factors identified by ER allows us to identify a sparse set of observables, within these factors, tied to outcome. We term this ER-derivative-approach Composite Regression (CR) (Figure 1C). As the sparse set of observables delineated by CR are selected from those that lie within the significant latent factors, unlike LASSO-based biomarkers, these are no longer simply predictive but capture causal relationships that can be used to infer the underlying mechanistic basis of outcome. Together, ER and CR provide a highly prioritized set of significant latent factors and associated observables, which can be used both for inference of underlying cellular/molecular mechanisms as well as corresponding biomarkers.

Inferring causal factors underlying immunosenescence in a vaccine response

A recent study comprehensively profiled cellular and molecular responses induced by the shingles Zostavax vaccine in a cohort comprising both younger adults and elderly individuals.²¹ The high-dimensional multi-omic analysis included immune-cell frequencies and phenotypes, as well as transcriptomic, metabolomic, cytokine, and antibody analyses. The vaccine induced robust antigen-specific antibody titers as well as CD4⁺, but not CD8⁺, T cell responses.²¹ Using a multi-scale, multifactorial response network, the authors identified associations between transcriptomic, metabolomic, cellular phenotypic, and cytokine datasets that pointed to immune and metabolic correlates of vaccine immunity.²¹ Interestingly, differences in the quality of the vaccine-induced responses by age were also noted.²¹ We hypothesized that a method based on latent factors rather than measurables

would improve the delineation of components that underlie the quality and magnitude of the vaccine-induced responses. If so, then such a method would be able to leverage the differences in vaccine-induced responses and accurately predict age as the system-wide property of interest. The latent factors identified in this manner could then provide insights into the cellular and molecular basis of age-induced immunosenescence manifested by diminished responses to the Zostavax vaccine.

To explore the above formulation of immunosenescence as a predictor of age, we first applied a suite of state-of-the-art approaches, LASSO, PLS, and PFR, on the entire spectrum of multi-omic vaccine-induced responses (including transcriptomic, metabolomic, cytokine, antibody, and cellular phenotypic data) to predict age (Figure 2A). As most individuals in the cohort were in 2 distinct age groups, adults under 40 and elderly people over 60, we first sought to explore the performance of LASSO, PLS, and PFR in predicting the two age groups as binary categorical variables, i.e., younger adults and elderly people. The predictive performance of all methods was evaluated in a stringent leave-one-out cross-validation (LOOCV) framework (experimental procedures). We have previously demonstrated that on such multi-omic datasets, cross validation is a gold standard to evaluate model performance with data held out.^{5,6,8} In an LOOCV framework, we found that PFR had no predictive power (area under the curve [AUC] <0.5), while LASSO and PLS had weak predictive power, in predicting age as a categorical variable (Figure 2B, AUCs = 0.63 and 0.60, respectively). The receiver operating characteristic (ROC) curve for LASSO had an interesting shape. It attained a true positive rate of ~0.4 at a false positive rate of ~0.15, but beyond that it was essentially no better than random (Figure 2B). This observation is consistent with the observation that differences in an age-associated multiscale multifactorial response network (MMRN) were driven by only a subset of elderly vaccinees.²¹ Thus, a purely predictive modeling approach like LASSO can leverage these relatively straightforward differences to accurately predict age for a subset of the vaccinees but fails to predict age for others. We then compared these methods with the performance of ER and CR. In a matched, LOOCV framework, ER and CR were very accurate at predicting age (Figure 2B, AUCs = 0.79 and 0.77, respectively, $p < 0.01$).

We then coupled ER to causal-inference analyses on the ER-identified significant latent factors using directed graphical models.²² Directed acyclic graphs (DAGs) are sometimes referred to as causal graphs because under certain assumptions the learned DAGs from observational data (Markov equivalence classes) asymptotically represent the true data-generating causal graph. Although these algorithms have shown considerable success in analyzing many biological processes and biomedical problems,^{23–27} including biomarker selection and classification,^{28–30} scalability limits the datasets to which they can be applied.^{31,32} Here, we use the causal-learning algorithm for mixed data, CausalMGM,^{23,33} only on the significant latent factors delineated by ER to overcome the scale limitation. By applying CausalMGM only on the significant latent factors, we greatly reduce the dimensionality of the input dataset while preserving the information of individual (correlated) variables in the latent factors. Thus, CausER (CausalMGM on the significant latent factors from ER) prioritizes further within the significant latent factors (experimental procedures; Note S1)

by virtue of their direct connections to the outcome in the graphical model. Furthermore, it predicts potential cause-effect relationships between the latent factors and the property/outcome of interest, which leads to hypotheses generation, while CausER was the best predictor of age as a categorical variable (AUC = 0.86, $p < 0.01$). Together, these results demonstrate that while LASSO, PLS, and PFR fail to accurately predict age from Zostavax-induced vaccine responses, ER, CR, and CausER can overcome this challenging problem by leveraging non-trivial differences in latent factors comprised of discrete sets of measurables.

Next, we evaluated whether these methods could predict actual age as a continuous variable beyond the categorical classifiers of younger adults and elderly individuals. As before, performance was measured in a rigorous cross-validation framework (experimental procedures). Using the vaccine-induced responses, PFR was not at all predictive of age (Figure 2C, Pearson $r = -0.71$; Figure S1, Spearman $r = -0.82$). LASSO and PLS had poor performance in predicting age as a continuous variable (Figure 2C, Pearson $r = 0.29$ and 0.13 , respectively; Figure S1, Spearman $r = 0.25$ and 0.09 , respectively). In fact, the predictive powers of PLS and PFR were not significantly different from a negative control model built on permuted data (Figure 2C). However, both ER and CR were significantly predictive of age as a continuous variable (Pearson $r = 0.48$ for both, Spearman $r = 0.44$ and 0.49 , respectively, $p < 0.01$; Figures 2C and S1), and as in the previous instance, CausER had the best performance in predicting age as a continuous variable (Pearson $r = 0.61$, Spearman $r = 0.59$, $p < 0.01$; Figures 2C and S1). Together, these results demonstrate that while state-of-the-art methods including LASSO, PLS, and PFR fail to predict age either as a categorical or a continuous variable, all three of the new approaches that are based on latent factors, ER, CR, and CausER, are able to do so reasonably accurately based on the multi-omic profiles of vaccine-induced responses.

We next explored the likely causal relationships among the latent factors that lead to age-induced immunosenescence and diminished responses to the Zostavax vaccine. CausalMGM was used to construct a causal graph with all latent factors identified in the latent-model-identification step of ER (Figure 2D). Notably, the majority of significant latent factors identified by ER were seen to be proximal to the outcome variable (age) in the causal graph. Importantly, all 4 latent factors in the Markov blanket generated by CausalMGM were also identified as significant by ER (Figure 2D). Overall, the significant latent factors revealed by ER had significantly lower network distances (i.e., they had stronger cause-effect relationships) from age compared with the non-significant latent factors (Figure 2E, $p < 0.05$). These results demonstrate that the cause-effect relationships identified by ER are validated by CausalMGM. Importantly, while CausER hits are identified via the sequential application of ER and CausalMGM, respectively, the order is critical, with ER being the key first step. Without the two-stage dimensionality reduction (first from observables to latent factors and then the identification of significant latent factors) afforded by ER, running CausalMGM or other allied causal graphical models on the initial set of observables would be computationally intractable.

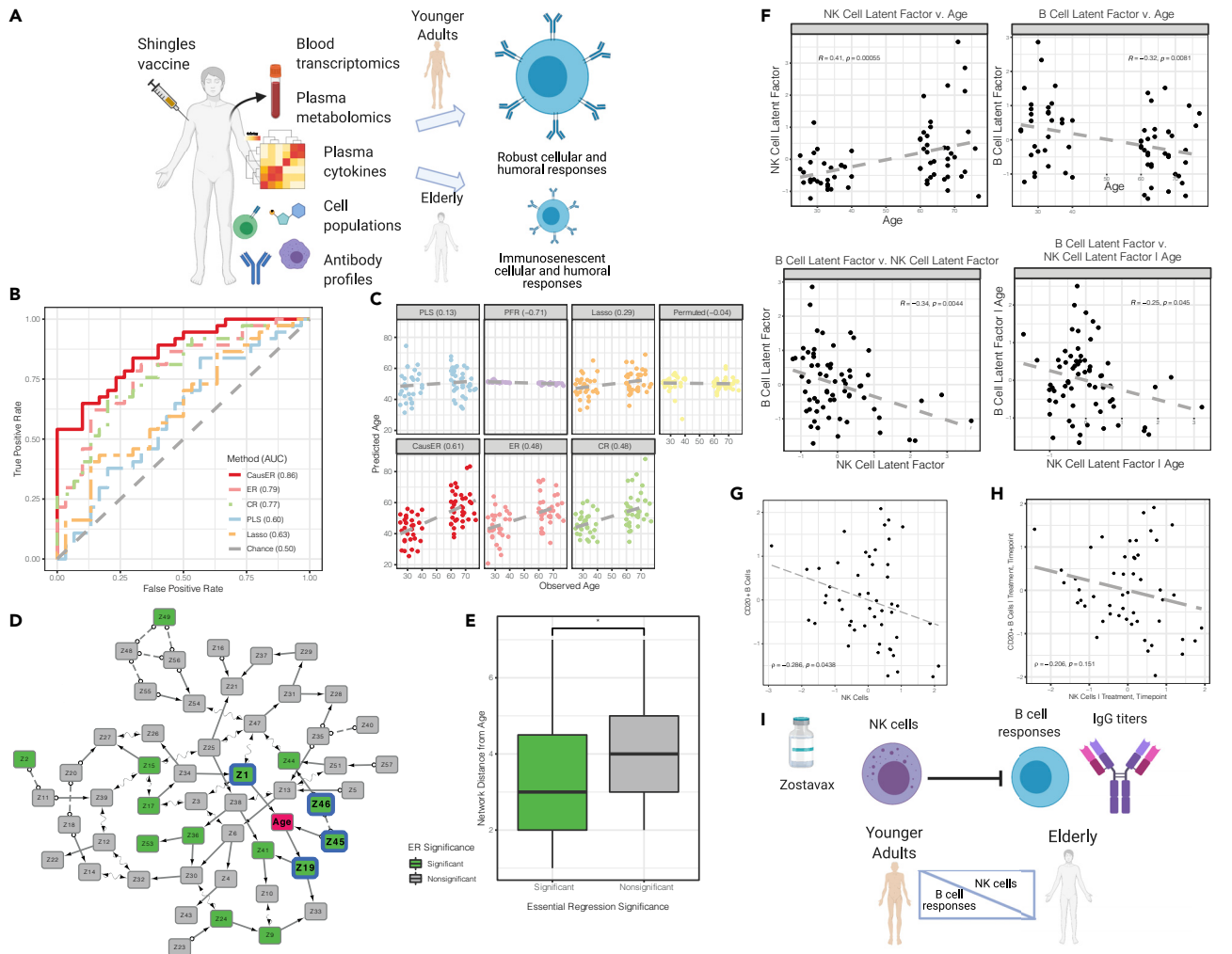


Figure 2. Identifying causal signatures of age-induced immunosenescent responses to the Zostavax vaccine

(A) Schematic summarizing the input data and the problem of interest.

(B) ROC curves for the different methods at discriminating between elderly people and younger adults in an LOOCV framework.

(C) Pearson correlations of the different methods at predicting age as a continuous variable, as measured in an LOOCV cross-validation framework.

(D) CausalMGM on all Zs identified by ER. The Markov blanket is highlighted with a blue border and bolder fonts. A directed edge $X \rightarrow Y$ indicates X is a cause of Y, while a bidirected edge $X \leftrightarrow Y$ indicates the presence of a latent confounder that is a common cause of X and Y. A partially oriented edge $X \circ \rightarrow Y$ indicates that Y is not a cause of X but that either X or a latent confounder causes Y. Unoriented edge indicates directionality could not be inferred for that edge.

(E) Network distances in the causal graph generated by CausalMGM of the significant and non-significant Zs identified by ER from the outcome variable of interest. p value calculated using a Mann-Whitney U test

(I) Mechanistic insights obtained from ER.

(F) Correlations involving the NK cell latent factor, B cell latent factor, and age. Top panels show correlations between the NK cell latent factor and age (top left), and the B cell latent factor and age (top right). Bottom panels show correlations between the NK cell latent factor and the B cell latent factor without correcting for age (bottom left) and after correcting for age (bottom right).

(G) Correlations between NK cells and B cells in the context of vaccination against SARS-CoV2 in a NHP model.

(H) Correlations between NK cells and B cells in the context of vaccination against SARS-CoV2 in a NHP model, after correcting for treatment (vaccination arm) and timepoint.

The prioritized CausER hits (Figure 2D), i.e., significant latent factors identified by ER that are also in the Markov blanket of the outcome variable (age) in the causal graph generated by CausalMGM, comprised antigen-specific immunoglobulin G (IgG) titers (Z1), a metabolic module (Z19), and B (Z46) and natural killer (NK; Z45) cell frequencies. CausER provides both prioritized cause-effect relationships and directions of these relationships. While the latter relates to mathematical conditional-indepen-

dence relationships (experimental procedures), the former provides prioritized mechanistic insights. Notably, the discovery and labeling of causal latent factors are completely unbiased and not based on prior knowledge. These significant latent factors are those that were identified as significant by ER and in the Markov blanket of outcome; neither step used any prior knowledge.

However, to evaluate the quality of these discoveries, we examined the uncovered latent factors in light of previously

elucidated bases of immunosenescence. The lowering of titers with age is expected and has been previously reported,²¹ so this corresponds to a recapitulation of known relationships. However, CausER also revealed a novel cause-effect relationship between altered B and NK cell numbers and immunosenescence. To further dissect the nature of this relationship, we examined correlations between NK cells, B cells, and age. We found that NK cells significantly increased, while the numbers of B cells significantly decreased, with age (Figure 2F). More interestingly, there was a significant negative correlation between NK and B cells (Figure 2F), and the correlation remained significant even after correcting for age (Figure 2F).

Notably, these causal inferences are supported by perturbation experiments involving biologically relevant organisms. Our findings relate to a previously described mechanistic linkage between NK cells and a weaker germinal center (GC) response in a murine model³⁴ in the context of vaccination with a model antigen (NP-KLH). NK cells can inhibit CD4 T cell responses, including those of T follicular helper cells, in a perforin-dependent manner; this leads to a weaker GC response and diminished antibody titers and affinity maturation.^{34,35} Furthermore, in the context of vaccination against severe acute respiratory syndrome coronavirus 2 (SARS-CoV2) in a non-human primate (NHP) model, we leveraged cell-subset-frequency data from a recent study³⁶ to examine the relationship between NK and B cells. We found a significant negative relationship between NK and B cells spanning multiple time points and vaccination arms corresponding to different adjuvants (Figure 2G). These relationships remained unaltered even after correcting for time point and vaccination arm using a linear model (Figure 2H). Together, these results demonstrate that a novel relationship uncovered solely by ER and CausER, without the use of any prior knowledge, from a human-systems vaccinology study have strong support in vaccination studies both in mice and NHPs. Notably, these studies use different antigens and adjuvants, suggesting that the uncovered novel relationship between NK and B cells is highly robust, and can be broadly extrapolated across vaccination strategies. Our results suggest a novel basis of human immunosenescence in the context of vaccine responses (Figure 2I). This discovery is especially striking as ER converged on this mechanism without the use of any prior knowledge.

Analyzing latent factors potentially reflective of trained immunity in a vaccine response

To test whether ER is applicable to datasets generated using alternate technological platforms, we applied it to analyze the temporal dynamics of transcriptional responses (microarray data) induced by the malaria RTS,S vaccine.³⁷ RTS,S has a standard regimen of 3 doses separated by a month and is currently the most advanced malaria-vaccine candidate that has consistently demonstrated 40%–80% protective efficacy in malaria-naïve individuals in controlled human challenge studies.⁵ There has been intense interest over the last decade in uncovering molecular signatures induced by the RTS,S vaccine and corresponding correlates of protection.^{5,38,39} In a controlled human-infection setting, differential expression of immunoproteasome genes was identified as a pre-challenge correlate of protection.³⁷ After the third dose, as expected, there was a striking but transitory shift in inflammatory gene expression followed a convergence of the

majority of gene signatures back to pre-vaccination levels within 2 weeks after the third dose.³⁷ We reasoned that aspects of trained immunity induced by the vaccine may be reflected in the transcriptomic signatures that do not converge after 2 weeks. Thus, a sensitive method such as ER would be able to discriminate between expression profiles at the following time points, pre-vaccination (G1), the day after the third dose (G2), and 14 days after the third dose (G3) (Figure 3A), and reveal candidate genes and molecular pathways that could contribute to trained immunity. In this instance, the use of a microarray dataset also afforded the opportunity to explore how ER performs with noisier but nevertheless valuable datasets generated using older technologies.

As before, the ability of the different methods to discriminate between G1, G2, and G3 transcriptional profiles was measured in a rigorous cross-validation framework (experimental procedures). We found that there were significant differences in the ability of the different methods to discriminate between the three kinds of expression profiles, with ER and CausER (CausalMGM on the significant latent factors from ER) having the best performance, significantly better than the other methods ($p < 0.01$, Figures 3B and 3C). Next, we chose to focus on the ability of the different methods to specifically distinguish the G3 profile from the other two (Figure 3D) or just the G1 profile (Figure 3E). This constituted the most “difficult” discrimination as there are broad differences in the expression profiles between the pre- (G1) and 24-hour-post-vaccination (G2) time points, but most of these differences disappear by 14 days (G3).³⁷ Consistent with expectations, in this binary-classification setting, there was wide variability in the performance of the methods to specifically discriminate the G3 time point from the G1 and G2 time points. While PFR and PLS performed poorly, CausER, ER, and LASSO had significantly better performances, with CausER being the best-performing method ($p < 0.01$; Figures 3D and 3E). In terms of correctly classifying just the true G3 profiles as G3, PLS and PFR had poor performances while CausER had the best performance, significantly better than other methods ($p < 0.01$, Figure 3F).

Next, we focused on the CausER hits, i.e., the significant latent factors from ER in the Markov blanket of the outcome variable (Figure 3G). Genes comprising these latent factors were seen to be differentially expressed between the G1 and G3 samples (Figures 3H and 3I). Our results suggest that beyond the initial divergence of immunoproteasome genes, there is a sustained divergence (2 weeks post-vaccination) of genes involved in immune-metabolic processes. These results complement recent findings that suggest that targeting immunometabolism is a promising direction in modulating trained immunity.⁴⁰ While a vaccine induces a rapid initial divergence in inflammatory signatures reflecting the activation of innate immune cells and their engagement with adaptive B and T cells, it may also induce alterations in the innate immune compartment that are discernible at later time points and contribute to a distinct form of immune memory.⁴⁰

Elucidating markers of latent and active tuberculosis (Tb)

To explore whether ER and CausER can predict clinically important outcomes, we applied these approaches to a dataset of high-dimensional antibody profiles for patients with latent and active Tb⁴¹ (Figure 4A). The high-dimensional antibody-omic dataset used a modern antibody-omic platform^{5,6,8,41} to quantify

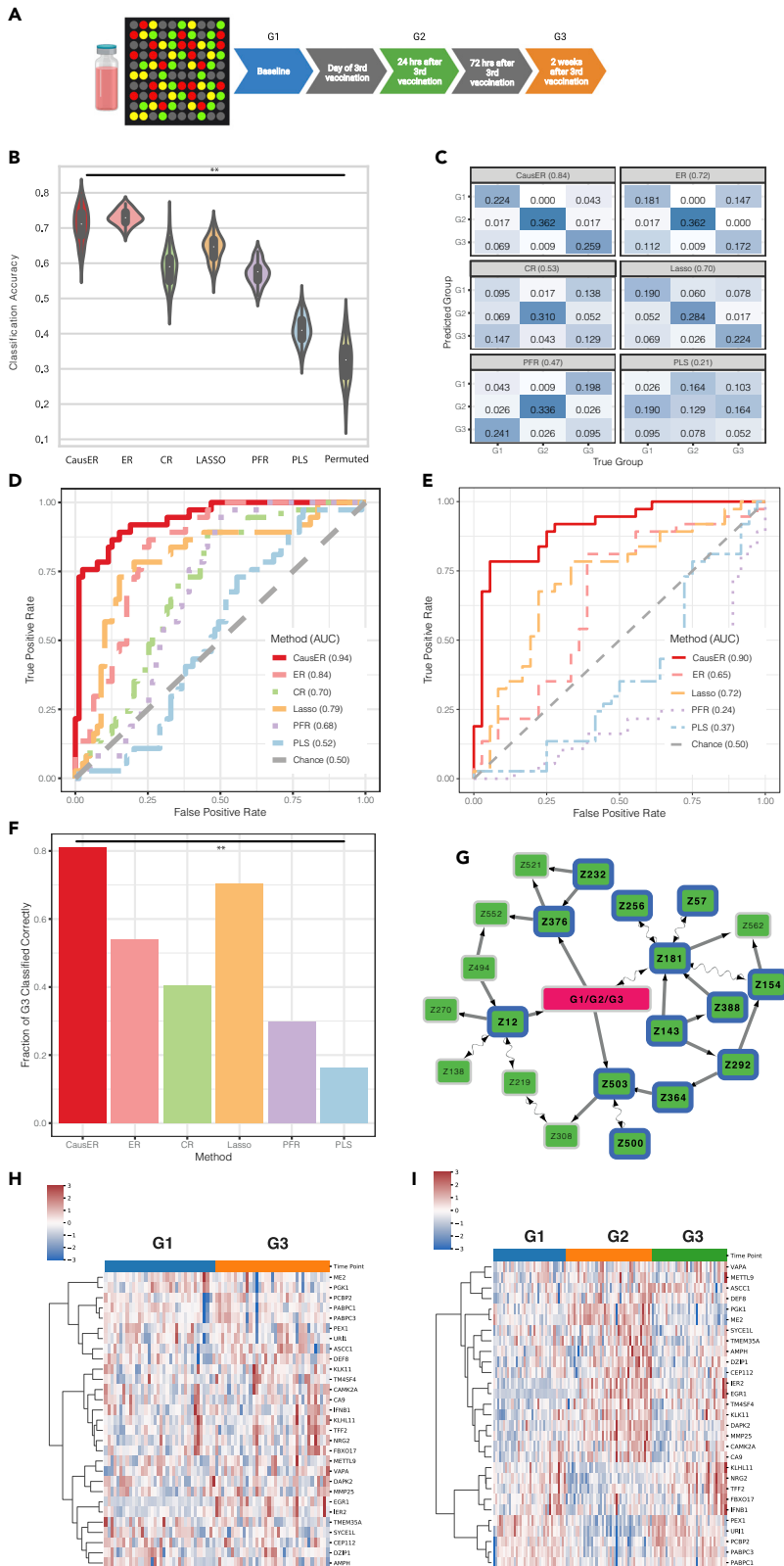


Figure 3. Identifying differences in vaccine-induced transcriptomic profiles over time

(A) Schematic summarizing the input data and the problem of interest.

(B) Ternary classification accuracy of the different methods at discriminating among G1, G2, and G3 in a replicated k -fold cross-validation framework.

(C) Confusion matrix summarizing the performance of the different methods at discriminating among G1, G2, and G3 in an LOOCV framework.

(D) ROC curves for the different methods at discriminating between G3 and G1 and G2 combined in an LOOCV framework.

(E) ROC curves for the different methods at discriminating between G3 and G1 in an LOOCV framework.

(F) Fraction of true G3 correctly classified as G3 (as measured in an LOOCV framework).

(G) CausER graph i.e., CausMGM on the significant Zs from ER. The Markov blanket is highlighted with a blue border and bolder fonts. A directed edge $X \rightarrow Y$ indicates X is a cause of Y , while a bidirected edge $X \leftrightarrow Y$ indicates the presence of a latent confounder that is a common cause of X and Y . A partially oriented edge $X \circ \rightarrow Y$ indicates that Y is not a cause of X but that either X or a latent confounder causes Y . Unoriented edge indicates directionality could not be inferred for that edge.

(H) Heatmap of genes in CausER hits (significant Zs in the Markov blanket) for G1 and G3 samples.

(I) Heatmap of genes in CausER hits (significant Zs in the Markov blanket) for G1, G2, and G3 samples.

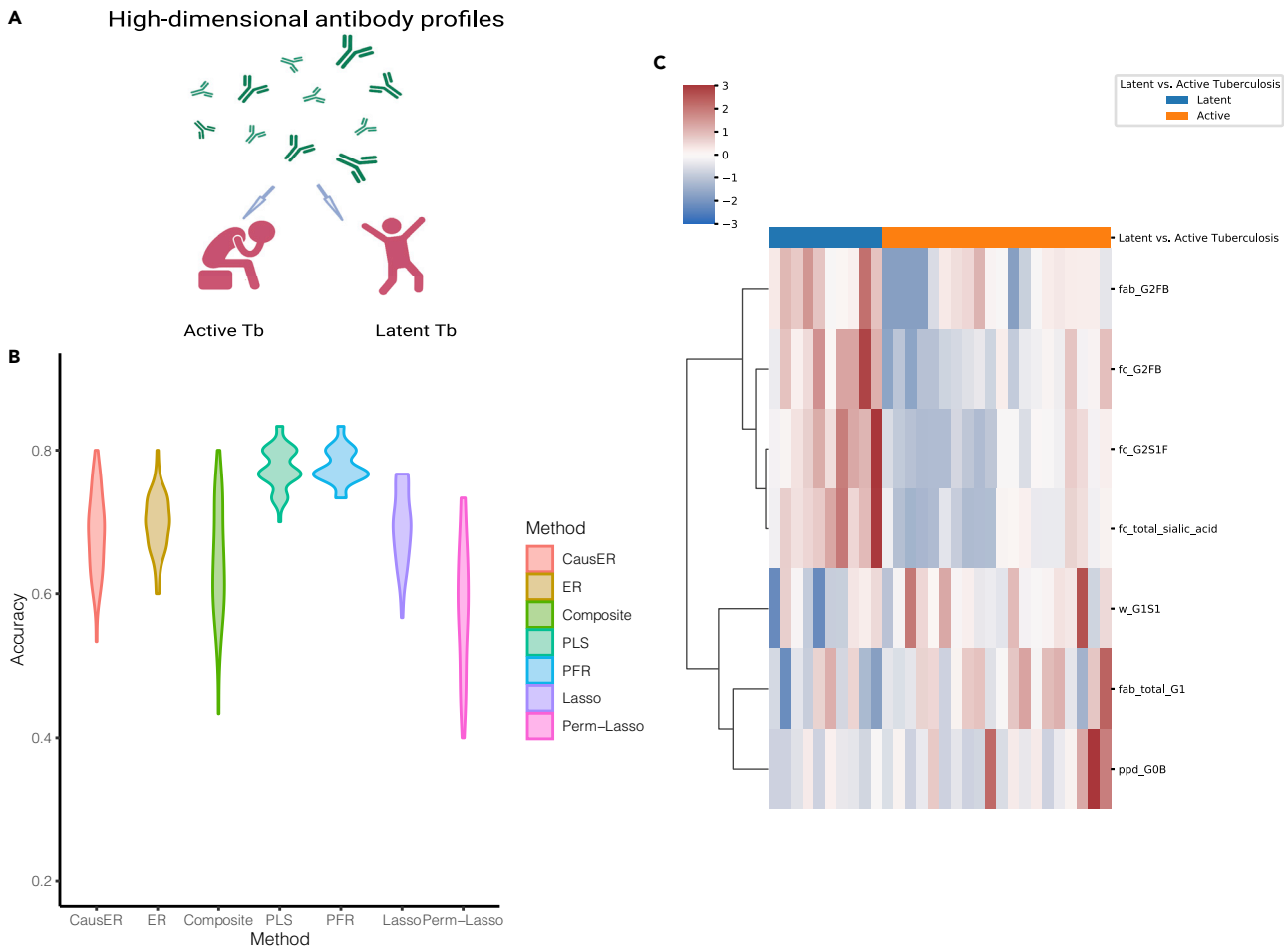


Figure 4. Elucidating markers of latent and active tuberculosis (Tb)

(A) Schematic summarizing the input data and the problem of interest.

(B) Classification accuracy of the different methods at discriminating between latent and active Tb, measured in a replicated k -fold cross-validation framework.

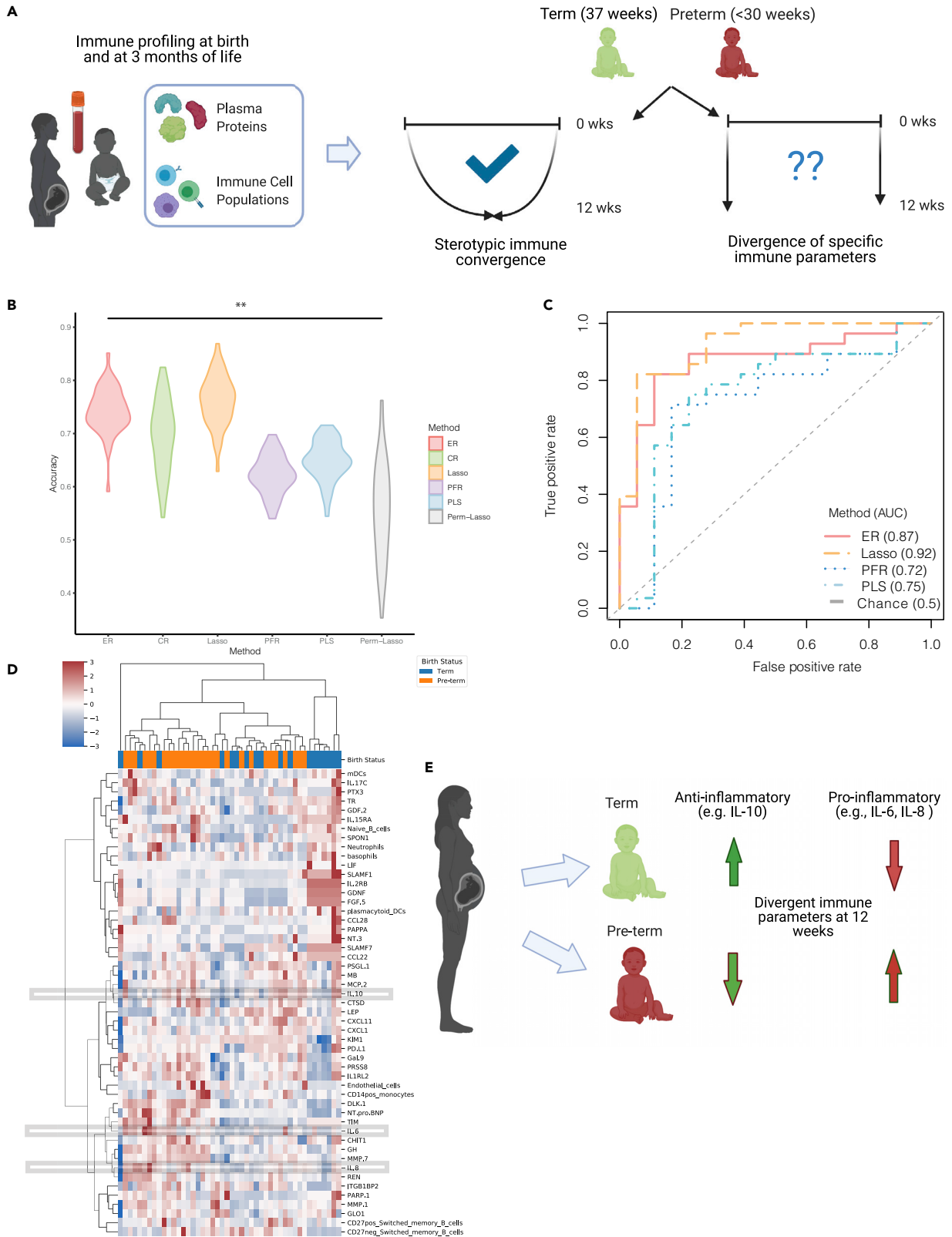
(C) Heatmap of features in the single CausER hit.

functional and biophysical properties of a polyclonal pool of antigen-specific antibodies. Each of these properties has its own inherent distribution so this was an appropriate test of the ability of ER and CausER to integrate multi-modal datasets for clinical outcome prediction. CausER and ER along with PLS, PFR, and LASSO were able to accurately discriminate between latent and active Tb patients using the antibody-omic profiles (Figure 4B). Notably, only one latent factor was identified as significant by ER and in the Markov blanket of outcome, i.e., this latent factor was the sole CausER hit. It consisted of specific glycosylation profiles (Figure 4C), and the majority of these glycosylation-based biomarkers were in perfect agreement with our previous study.⁴¹ These analyses demonstrate that ER and CausER are able to accurately predict clinically important outcomes.

Uncovering latent factors that distinguish immune-system states of term and pre-term infants

Finally, we focused on a multi-omic longitudinal cohort that analyzed immune-cell populations and plasma proteins in 100 newborn children during their first 3 months of life⁴² (Figure 5A).

Striking differences were observed in immune parameters between pre-term and term children at birth. However, the immune trajectories appeared to achieve a stereotypic convergence within the first 3 months of life⁴² (Figure 5A). We hypothesized that ER might be able to uncover latent factors that distinguish immune-system states of term and pre-term infants after 3 months of life and therefore reveal features that could impact later life (Figure 5A). As expected, based on the striking differences at birth between term and pre-term children, all methods (LASSO, PLS, PFR, ER, CR, and CausER) were able to discriminate between these 2 groups using immune parameters measured in the first week of life (Figure S2). All model performances were measured in a rigorous cross-validation framework (experimental procedures). However, given the stereotypic convergence in the first 3 months (12 weeks) of life,⁴² we found that PLS and PFR were unable to accurately discriminate between term and pre-term children using immune parameters measured at 12 weeks of life (Figures 5B and 5C). However, LASSO was able to accurately distinguish between term and pre-term births using the 12-week profiles (Figures 5B and 5C),



(legend on next page)

suggesting that despite broad convergence, a small subset of immune parameters still remain different in term and pre-term infants at 3 months of life. More importantly, ER and CR were able to accurately discriminate between term and pre-term births using immune profiles at 3 months of life, significantly better than other methods (Figures 5B and 5C, $p < 0.01$). ER identified only 2 significant latent factors, and based on CausalMGM analyses, one of these 2 significant latent factors was in the Markov blanket, i.e., for this dataset, this single latent factor was the sole CausER hit (Figure 5D).

We visualized the immune-cell populations and plasma proteins in this hit (Figure 5D). These profiles had clearly remained divergent even at 3 months of life (Figure 5D) despite the broad stereotypic convergence of most other immune parameters. At 3 months of life, term infants had an anti-inflammatory milieu including high interleukin (IL)-10 while pre-term infants had a pro-inflammatory milieu including elevated IL-6 and IL-8 (Figure 5D). These findings agree with a previous study that IL-10 is highly expressed in the uterus and placenta and has a key role in controlling inflammation-induced pre-term labor in a murine model.⁴³ Furthermore, regulatory B cells are a key source of IL-10 and appear to be important in sustaining pregnancy until term.^{44–46} It is also known that modulation of pro- versus anti-inflammatory environments by relevant cytokines and chemokines at the maternal-fetal interface (decidua) is a critical component of the bifurcation between term and pre-term births.⁴⁴ Thus, our analyses of immune-system states of term and pre-term infants at 3 months of life revealed that pre-term infants had a pro-inflammatory state while term infants had an anti-inflammatory state (Figure 5E). These findings could have long-term implications for the health of pre-term infants.

DISCUSSION

Over the last two decades, while there have been rapid advances in high-throughput experimental technologies to generate deep molecular profiles, computational analyses of these high-dimensional datasets have primarily focused on biomarker discovery.⁴⁷ This is because rigorous statistical approaches for analyzing high-dimensional datasets, such as regularized regression and bootstrap-aggregated classification, are focused on uncovering predictive biomarkers, which may simply be correlative surrogates of outcome or system-wide property but are unrelated to the underlying causal factors. Incorrect extrapolation of insights derived from biomarker-based approaches can lead to perturbation experiments with low success. Alternatively, efforts to move beyond biomarkers to mechanistic insights often use biological priors, which may be incomplete or suffer from sampling/study biases.⁴⁸ Furthermore, while there have been advances in causal modeling,⁴⁹ existing approaches are difficult to apply to high-dimensional datasets due to the

computational intractability of applying these approaches on²² and the multi-collinearity of the data. The methods presented in this article address this fundamental limitation in systems biology. ER is a first-in-class machine-learning method that can both handle high-dimensional multi-omic datasets with collinear variables and prioritize cause-effect relationships between the input features and the outcome of interest. Our framework is also complementary to modern approaches that combine multi-omic datasets with prior knowledge to uncover causal relationships.¹⁴ ER generates mechanistic hypotheses solely based on latent factors identified from multi-omic data without the incorporation of any prior knowledge. It is thus applicable in contexts where prior knowledge is weak or unavailable and is not limited by the nature and quality of available prior knowledge. ER is compatible with all existing batch correction/normalization approaches as it makes no assumptions regarding data-generating mechanisms. However, data need to be appropriately normalized/batch corrected before being used as inputs to ER. Further, ER is also able to handle complex replicate structures. Biological and technical replicates may be pre-processed using a suitable context-specific approach; ER does not impose any restrictions on/is robust to how replicates are handled (which depends entirely on the underlying biological context/question). ER works downstream of these methods to integrate appropriately pre-processed/normalized multi-omic datasets and uncover causal latent factors underlying groups/outcomes of interest.

Importantly, ER is fundamentally different from classical factor regression models used exclusively for prediction. In those models, one first seeks a low-dimensional factor $Z = XV$ constructed via some projection matrix V . Although, Z can then be used to regress to Y , this framework can only be used for prediction and not inference as Z is not uniquely identifiable and this makes inference on the regression coefficients impossible. However, in the ER framework, the latent factors (Z s) and the corresponding linear coefficients (between Y and Z s) are uniquely identifiable, making the inference problem well-posed. Thus, our framework addresses a key limitation of classical factor regression models where the recovered factors have ambiguous meaning. However, the unique identifiability of the latent factors in the unsupervised step of ER makes inference meaningful. Thus, we cannot simply replace it with other modern clustering approaches with no guarantees regarding identifiability. The identifiability criterion tied to the guarantees regarding inference make ER a first-in-class interpretable latent-factor regression framework for high-dimensional multi-omic datasets.

Our framework pushes the envelope on multiple key challenges in systems biology. First, it establishes a rigorous framework with provable statistical guarantees that explores a large space of higher-order relationships from high-dimensional features and uncovers latent factors tied to the outcome variable via directed cause-effect relationships. Second, unlike existing

Figure 5. Uncovering specific immune parameters from term and pre-term infants that do not achieve stereotypic convergence

(A) Schematic summarizing the input data and the problem of interest.

(B) Classification accuracy of the different methods at discriminating between term and pre-term births using immune profiles at 3 months after birth, measured in a replicated k -fold cross-validation framework.

(C) ROC curves for the different methods at discriminating between term and pre-term births as measured in an LOOCV framework.

(D) Heatmap of features (plasma proteins and immune cells) in the single hit (significant Z identified by ER in the Markov blanket of outcome).

(E) Mechanistic insights obtained from ER.

causal-reasoning approaches that are constrained by the size of the input data, ER can be applied to modern high-dimensional datasets. The time complexities of the different steps are essentially quadratic and not exponential like some other causal-reasoning approaches. Third, ER makes no assumptions regarding data-generating mechanisms, and ER can integrate multi-omic datasets to capture the interplay across a plethora of biological processes at multiple scales of organization of the system. Fourth, ER is able to home in on one or a few causal latent factors of outcome comprising a small number of observable features from thousands of input features, many of which are completely uninformative. Finally, ER converges on these causal latent factors without the use of any prior knowledge; however, we find that the uncovered factors include both previously elucidated and novel mechanistic bases. The ability of ER to converge on meaningful biological insights without any prior knowledge makes it applicable in the broadest sense even in contexts where there are weak or no priors.

An important elaboration of our framework is the sequential use of two orthogonal methods for statistical inference, ER and causal graphical modeling. These methods have different theoretical bases and assumptions, and yet the ER hits are validated by CausalMGM, underscoring the robustness of our approach. The order is critical, with ER being the key first step offering a two-stage dimensionality reduction: first from observables to latent factors, and then the identification of significant latent factors. Without these two steps, the application of causal graphical models on the initial set of observables would be computationally intractable due to the high dimensionality of the dataset. Thus, ER solves a long-standing limitation with causal graphical-modeling approaches and enables, for the first time, causal inference on high-dimensional data. ER also has polynomial time complexity that makes it efficient and scalable for extremely large datasets. For all the datasets analyzed in this study, it resulted in runtimes of core minutes for each cross-validation replicate, which translates into tens to hundreds of core hours after accounting for cross-validation replicates. The datasets included tens to hundreds of samples and up to 10^3 – 10^4 features/sample. So, this all suggests that ER is extremely efficient with modern multi-omic datasets.

ER has a number of limitations. One relates to the constraint each latent factor is anchored by at least 2 pure variables (i.e., variables that belong to only that and no other latent factor). However, this is a reasonable assumption as most observables are allowed to be mixed, i.e., they can belong to one or more latent factors, and each latent factor only requires 2 pure variables to anchor it. Also, in some instances, the linearity assumption between Y and Z could be restrictive. For example, when the number of latent factors is small, this restrictive assumption could be overcome by including high-order terms of Z to predict Y . It is also possible to extend the current framework to a more general setting where Y and Z follow generalized linear models with any appropriate link function, such as the logistic and probit functions. While it is relatively straightforward to incorporate suitable link functions in the setting of prediction, achieving theoretical guarantees for the inference of the coefficients of Z needs more careful theoretical analyses.

The coupling of causal graphical models to ER and the inference of causality from observational data also has some

assumptions. First, it is assumed that the structure of the cause-effect relationships of all variables in the dataset form a DAG. Next, the causal Markov assumption requires that the Markov condition for DAGs holds for the causal graph. Finally, the causal faithfulness assumption states that the conditional-independence relationships in the dataset are faithful to the causal graph. A distribution is faithful to its causal DAG when there are no additional conditional-independence relationships that are not entailed by the Markov condition of the DAG. Importantly, while some algorithms also require the assumption of causal sufficiency, which states that there are no unobserved confounders of the variables in the dataset, the fast causal inference (FCI) algorithm used here does not have this constraint. Further, for full identification of the causal graph, the assumptions of the conditional-independence test, in this case linearity, must be met, and the sample size must be asymptotically large. Thus, the inference of true causality is constrained by these assumptions, which may not always hold. However, importantly, these assumptions are tied to the causal graphical-modeling framework. ER (without coupling to CausalMGM) can be used to identify significant latent factors, with only very minimal assumptions, as described above. Thus, while true causality may, in some instances, be difficult to infer from observational data, the significant latent factors identified by ER provide inference into generative processes beyond just prediction.

Here, we applied ER to diverse contexts. First, we applied ER on simulated datasets and demonstrated that it performed better than LASSO, PLS, and PFR across a range of parameter settings. Next, we utilized ER on two recent human systems-immunology studies that had generated high-dimensional multi-omic profiles. Using ER, we were able to address key questions that had not been the focus of the original studies, in part because of limitations of methods used. Such questions could now be addressed by the methodological advances of ER over state-of-the-art approaches. We demonstrated that ER significantly outperforms PFR and PLS across contexts and either outperforms or matches LASSO in terms of predictive performance. While we used three examples to illustrate the superior performance of ER, these methods come with broad theoretical guarantees to outperform PLS, PFR, and LASSO across contexts ([experimental procedures](#); [Figure 1](#)). Furthermore, while the existing methods simply identify correlates, without using any prior knowledge, ER provides mechanistic insights. Some of these outcomes are consistent with previous mechanistic experiments while others are novel. ER can also be used for noisier and older datasets not generated using state-of-the-art methods. Our findings have broad implications across domains in systems biology and are likely to transform both computational workflows used to analyze multi-omic datasets and downstream experiments designed based on the insights gleaned via these analyses.

EXPERIMENTAL PROCEDURES

Resource availability

Lead contact

Requests for data and code used for the study should be directed to and will be fulfilled by the lead contact, Jishnu Das (jishnu@pitt.edu).

Materials availability

This study did not generate new unique reagents.

Data and code availability

Detailed code, associated datasets, and documentation for ER, CR, and CausER are available at <https://github.com/jishnu-lab/ER>. A corresponding stable release can be accessed at <https://doi.org/10.5281/zenodo.6178063>.

Any queries regarding the code or data should be directed to the lead contact, Jishnu Das (jishnu@pitt.edu).

Theoretical underpinnings of ER

We provide brief descriptions of the methods, associated tuning parameters, cross-validation strategies, and data pre-processing in this section. Additional theoretical details are included in [Note S1](#).

Processing of systems-immunology datasets

For the dataset of multi-omic responses to the Zostavax vaccine, we included the following multi-scale measurements of immune state: IgG titers, blood transcriptional modules, metabolic clusters, CD4⁺ T cell populations, T follicular helper (TFH) cell populations, flow-cytometry cell populations, cytokine profiles, and IFN T cells. We used subject age as the response variable for $n = 72$ subjects. We excluded features that had missing values for more than a half of subjects. We also excluded 5 subjects that had no observed features. The remaining datasets were merged via the unique IDs of subjects. The final dataset contains $p = 1,721$ features of $n = 67$ subjects.

For the transcriptomic (microarray) dataset pre- and post-malaria vaccination, we had $n = 116$ samples with 22,277 probes. We filtered out ambiguous probes (i.e., those that could map to multiple genes) and then averaged technical replicates (multiple probes/gene) with the limma package in R. The final dataset comprised 116 samples and $p = 12,424$ genes. Y is a categorical variable with 3 levels corresponding to three time points.

For the dataset of high-dimensional antibody profiles, we had $n = 30$ subjects (20 latent Tb, 10 active Tb) with $p > 100$ features/subject. The features included titers, Fc effector functions and whole, Fab, and Fc glycan profiles (independent of antigen) as well PPD- and Ag85-specific titers and glycan profiles.

For the dataset of term and pre-term infants, we included all available immune parameters as features and only removed clinical metadata (such as “gender,” “mode of delivery,” “family,” etc.). The final dataset we used has $n = 183$ samples and $p = 282$ features with 56 samples from week 1 and 46 samples from week 12. The response is binary, either “control” (representing term) or “pre-term” (representing pre-term). We used the 5-NN to impute the missing values.

Cross validation

Two cross-validation techniques were used to assess the predictive performance of the different methods: (1) replicated 10-fold cross validation and 2) LOOCV. (1) To assess the accuracy of the classifiers for the term/pre-term immune profile, 50 replicates of nested 10-fold cross validation were performed. For each replicate, we independently ran each of the methods and assessed the predictive accuracy. For ER, the latent factors were learned on each fold and each replicate, and the regression and final latent-factor selection were repeated. For CausER, a causal model was learned over the latent factors selected as significant by ER for each fold and replicate. The average cross-validation accuracy across the 10-folds was calculated for each of the 50 replicates. (2) For the datasets, we also performed LOOCV to assess the accuracy of each method. In LOOCV, each sample in the dataset was held out as the predictive models were trained on the remaining samples, and then the held-out sample was predicted with the trained models. Assessment of model performance (AUC) was done with the set of predictions for the left-out values.

ER

The first step in ER is the estimation of all latent factors. The identification of latent factors is unsupervised. This is done based on the empirical sample covariance matrix using a three-step procedure. The first step involves the identification of latent-variable structure using the sample covariance matrix. A key component of this step is the identification of K (reduced dimensionality) from p (original dimensionality). The second step involves inference of the clusters: each cluster (latent factor) is anchored by at least 2 pure variables. Variables that are associated with multiple clusters are designated mixed

variables. The third step involves determination of the overall allocation matrix based on the cluster assignments in the earlier step. Formal descriptions of all 3 steps are provided in [Note S1](#), Section 2.

After the identification of Z_s , the regression coefficients linking the Z_s to Y are estimated using a theoretical framework we recently established for estimation in latent-factor regression models.⁵⁰ This is the supervised part of the ER algorithm. A detailed description of the estimation procedure is provided in [Note S1](#), Section 2.

CR

CR utilizes a 2-step procedure. First, it uses ER to identify significant Z_s as described above. Then, it uses LASSO on the X_s associated only with the significant Z_s to identify a sparse basis for the system-wide property/outcome of interest. For LASSO on the significant Z_s identified by ER, lambda is tuned using k -fold cross validation. The lambda tuning is specific to a given fold for a given replicate and utilizes only the fold-specific training data.

ER coupled to CausalMGM

We implemented CausalMGM as previously described²² on all Z_s for the Zostavax dataset and only the significant Z_s identified by ER for the term/pre-term, malaria, and Tb datasets. Briefly, when constructing the causal model, we first learned an undirected graphical model with MGM⁵¹/GLASSO.⁵² The optimal regularization parameters were selected based on graph stability using StEPS³³/StARS.⁵³ The resulting undirected graph was then used as an initial graph for performing causal inference with the FCI algorithm. To build a predictor of the outcome variable, the Markov blanket was used. The Markov blanket was defined as the set of variables that, when conditioned on, make the response variable independent of every other variable in the dataset according to the structure of the causal graph. For a DAG, this comprises the parents, children, and spouses (other parents of the children) of the response variable.

Implementation of LASSO, PLS, and PFR

LASSO was implemented using glmnet in R with parameter tuning done in a manner analogous to that described above for ER and CR. If no feature was selected by LASSO in a specific fold for a given replicate, we randomly selected 5 features (only for that fold in that replicate) and used an ordinary-least-squares estimator. Thus, the feature selection in each case is specific to each fold for a given replicate; this is the most stringent and unbiased way to evaluate model performance. PLS was implemented using the pls function in R with the number of components selected by the default function selectNcomp. For PFR, which regresses Y on the first K principal components of X , the number of principal components K is selected based on the ratios of non-decreasing eigenvalues of X^*X/n using previously established criteria.⁵⁴

SUPPLEMENTAL INFORMATION

Supplemental information can be found online at <https://doi.org/10.1016/j.patter.2022.100473>.

ACKNOWLEDGMENTS

The authors would like to thank Mirko Paiardini and Maria Elena Bottazzi for their support. This study was partially supported by NIH grants DP2AI164325 to J.D., U01HL137159, R01HL140963, R01HL159805, and R01HL157879 to P.V.B., U01AI141990 to H.S., and F31LM013966 to T.L.; NSF grants DMS-1712709 and DMS-2015195 to F.B. and M.W.; and DoD grant W81XWH2110864 to J.D. H.S. also acknowledges support from the UPMC ITTC fund. S.P.K. acknowledges support from the Yerkes Pilot Research Pilot Program (part of the Yerkes NPRC Base Grant, P51-OD011132). Several images were created with BioRender.com.

AUTHOR CONTRIBUTIONS

J.D. designed the study and oversaw all aspects of it. X.B., F.B., and M.W. jointly conceived the theoretical basis of the ER framework. J.D. and X.B. jointly conceived the application of the ER and CR frameworks to real data.

J.D., P.V.B., and H.S. jointly conceived the CausER framework. X.B. and T.L. implemented the ER, CR, and CausER frameworks and carried out all computational analyses. S.P.K. provided data. J.D., P.V.B., and H.S. interpreted the results. J.D., H.S., and P.V.B. wrote the main text. X.B., T.L., F.B., and M.W. wrote the supplementary methods including formal proofs.

DECLARATION OF INTERESTS

The authors declare no competing interests.

Received: July 28, 2021

Revised: September 17, 2021

Accepted: March 1, 2022

Published: March 24, 2022

REFERENCES

- Hagan, T., and Pulendran, B. (2018). Will systems biology deliver its promise and contribute to the development of new or improved vaccines? From data to understanding through systems biology. *Cold Spring Harb. Perspect. Biol.* *10*, a028894. <https://doi.org/10.1101/cshperspect.a028894>.
- Pulendran, B., Li, S., and Nakaya, H.I. (2010). Systems vaccinology. *Immunity* *33*, 516–529. <https://doi.org/10.1016/j.immuni.2010.10.006>.
- Davis, M.M., Tato, C.M., and Furman, D. (2017). Systems immunology: just getting started. *Nat. Immunol.* *18*, 725–732. <https://doi.org/10.1038/ni.3768>.
- Villani, A.C., Sarkizova, S., and Hacohen, N. (2018). Systems immunology: learning the rules of the immune system. *Annu. Rev. Immunol.* *36*, 813–842. <https://doi.org/10.1146/annurev-immunol-042617-053035>.
- Suscovitch, T.J., Fallon, J.K., Das, J., Demas, A.R., Crain, J., Linde, C.H., Michell, A., Natarajan, H., Arevalo, C., Broge, T., et al. (2020). Mapping functional humoral correlates of protection against malaria challenge following RTS,S/AS01 vaccination. *Sci. Transl. Med.* *12*, eabb4757. <https://doi.org/10.1126/scitranslmed.abb4757>.
- Das, J., Devadhasan, A., Linde, C., Broge, T., Sassic, J., Mangano, M., O'Keefe, S., Suscovitch, T., Streeck, H., Irrinki, A., et al. (2020). Mining for humoral correlates of HIV control and latent reservoir size. *PLoS Pathog.* *16*, e1008868. <https://doi.org/10.1371/journal.ppat.1008868>.
- Goetghebuer, T., Smolen, K.K., Adler, C., Das, J., McBride, T., Smits, G., Lecomte, S., Haelterman, E., Barlow, P., Piedra, P.A., et al. (2019). Initiation of antiretroviral therapy before pregnancy reduces the risk of infection-related hospitalization in human immunodeficiency virus-exposed uninfected infants born in a high-income country. *Clin. Infect. Dis.* *68*, 1193–1203. <https://doi.org/10.1093/cid/ciy673>.
- Ackerman, M.E., Das, J., Pittala, S., Broge, T., Linde, C., Suscovitch, T.J., Brown, E.P., Bradley, T., Natarajan, H., Lin, S., et al. (2018). Route of immunization defines multiple mechanisms of vaccine-mediated protection against SIV. *Nat. Med.* *24*, 1590–1598. <https://doi.org/10.1038/s41591-018-0161-0>.
- Sadanand, S., Das, J., Chung, A.W., Schoen, M.K., Lane, S., Suscovitch, T.J., Streeck, H., Smith, D.M., Little, S.J., Lauffenburger, D.A., et al. (2018). Temporal variation in HIV-specific IgG subclass antibodies during acute infection differentiates spontaneous controllers from chronic progressors. *AIDS* *32*, 443–450. <https://doi.org/10.1097/QAD.0000000000001716>.
- Vafaee, F., Diakos, C., Kirschner, M.B., Reid, G., Michael, M.Z., Horvath, L.G., Alinejad-Rokny, H., Cheng, Z.J., Kuncic, Z., and Clarke, S. (2018). A data-driven, knowledge-based approach to biomarker discovery: application to circulating microRNA markers of colorectal cancer prognosis. *NPJ Syst. Biol. Appl.* *4*, 20. <https://doi.org/10.1038/s41540-018-0056-1>.
- Li, S., Roupael, N., Duraisingham, S., Romero-Steiner, S., Presnell, S., Davis, C., Schmidt, D.S., Johnson, S.E., Milton, A., Rajam, G., et al. (2014). Molecular signatures of antibody responses derived from a systems biology study of five human vaccines. *Nat. Immunol.* *15*, 195–204. <https://doi.org/10.1038/ni.2789>.
- Nakaya, H.I., Wrammert, J., Lee, E.K., Racioppi, L., Marie-Kunze, S., Haining, W.N., Means, A.R., Kasturi, S.P., Khan, N., Li, G.M., et al. (2011). Systems biology of vaccination for seasonal influenza in humans. *Nat. Immunol.* *12*, 786–795. <https://doi.org/10.1038/ni.2067>.
- Argelaguet, R., Arnol, D., Bredikhin, D., Deloro, Y., Velten, B., Marioni, J.C., and Stegle, O. (2020). MOFA+: a statistical framework for comprehensive integration of multi-modal single-cell data. *Genome Biol.* *21*, 111. <https://doi.org/10.1186/s13059-020-02015-1>.
- Dugourd, A., Kuppe, C., Sciacovelli, M., Gjerga, E., Gabor, A., Emdal, K.B., Vieira, V., Bekker-Jensen, D.B., Kranz, J., Bindels, E.M.J., et al. (2021). Causal integration of multi-omics data with prior knowledge to generate mechanistic hypotheses. *Mol. Syst. Biol.* *17*, e9730. <https://doi.org/10.15252/msb.20209730>.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. Ser. B (Methodological)* *58*, 267–288.
- Breiman, L. (2001). Random forests. *Mach. Learn.* *45*, 5–32. <https://doi.org/10.1023/A:1010933404324>.
- Bing, X., Bunea, F., Ning, Y., and Wegkamp, M. (2020). Adaptive estimation in structured factor models with applications to overlapping clustering. *Ann. Stat.* *48*, 2055–2081, 2027.
- Bing, X., Bunea, F., and Wegkamp, M. (2019). Inference in latent factor regression with clusterable features. Preprint at arXiv. 1905.12696. <https://doi.org/10.48550/arXiv.1905.12696>.
- Boulesteix, A.L., and Strimmer, K. (2007). Partial least squares: a versatile tool for the analysis of high-dimensional genomic data. *Brief Bioinform.* *8*, 32–44. <https://doi.org/10.1093/bib/bbl016>.
- Bair, E., Hastie, T., Paul, D., and Tibshirani, R. (2006). Prediction by supervised principal components. *J. Am. Stat. Assoc.* *101*, 119–137. <https://doi.org/10.1198/016214505000000628>.
- Li, S., Sullivan, N.L., Roupael, N., Yu, T., Banton, S., Maddur, M.S., McCausland, M., Chiu, C., Canniff, J., Dubey, S., et al. (2017). Metabolic phenotypes of response to vaccination in humans. *Cell* *169*, 862–877.e17. <https://doi.org/10.1016/j.cell.2017.04.026>.
- Ge, X., Raghu, V.K., Chrysanthis, P.K., and Benos, P.V. (2020). CausalMGM: an interactive web-based causal discovery tool. *Nucleic Acids Res.* *48*, W597–W602. <https://doi.org/10.1093/nar/gkaa350>.
- Sedgewick, A.J., Buschur, K., Shi, I., Ramsey, J.D., Raghu, V.K., Manatakis, D.V., Zhang, Y., Bon, J., Chandra, D., Karoleski, C., et al. (2019). Mixed graphical models for integrative causal analysis with application to chronic lung disease diagnosis and prognosis. *Bioinformatics* *35*, 1204–1212. <https://doi.org/10.1093/bioinformatics/bty769>.
- Schadt, E.E., Lamb, J., Yang, X., Zhu, J., Edwards, S., Guhathakurta, D., Sieberts, S.K., Monks, S., Reitman, M., Zhang, C., et al. (2005). An integrative genomics approach to infer causal associations between gene expression and disease. *Nat. Genet.* *37*, 710–717. <https://doi.org/10.1038/ng1589>.
- Sachs, K., Perez, O., Pe'er, D., Lauffenburger, D.A., and Nolan, G.P. (2005). Causal protein-signaling networks derived from multiparameter single-cell data. *Science* *308*, 523–529.
- Manatakis, D.V., Raghu, V.K., and Benos, P.V. (2018). piMGM: incorporating multi-source priors in mixed graphical models for learning disease networks. *Bioinformatics* *34*, i848–i856. <https://doi.org/10.1093/bioinformatics/bty591>.
- Kitsios, G.D., Fitch, A., Manatakis, D.V., Rapport, S.F., Li, K., Qin, S., Huwe, J., Zhang, Y., Doi, Y., Evankovich, J., et al. (2018). Respiratory microbiome profiling for etiologic diagnosis of pneumonia in mechanically ventilated patients. *Front. Microbiol.* *9*, 1413. <https://doi.org/10.3389/fmicb.2018.01413>.
- Abecassis, I., Sedgewick, A.J., Romkes, M., Buch, S., Nukui, T., Kapetanaki, M.G., Vogt, A., Kirkwood, J.M., Benos, P.V., and Tawbi, H. (2019). PARP1 rs1805407 increases sensitivity to PARP1 inhibitors in cancer cells suggesting an improved therapeutic strategy. *Sci. Rep.* *9*, 3309. <https://doi.org/10.1038/s41598-019-39542-2>.

29. Raghu, V.K., Zhao, W., Pu, J., Leader, J.K., Wang, R., Herman, J., Yuan, J.M., Benos, P.V., and Wilson, D.O. (2019). Feasibility of lung cancer prediction from low-dose CT scan and smoking factors using causal models. *Thorax* 74, 643–649. <https://doi.org/10.1136/thoraxjnl-2018-212638>.
30. Raghu, V.K., Beckwitt, C.H., Warita, K., Wells, A., Benos, P.V., and Oltvai, Z.N. (2018). Biomarker identification for statin sensitivity of cancer cell lines. *Biochem. Biophys. Res. Commun.* 495, 659–665. <https://doi.org/10.1016/j.bbrc.2017.11.065>.
31. Raghu, V.K., Ramsey, J.D., Morris, A., Manatakis, D.V., Sprites, P., Chrysanthis, P.K., Glymour, C., and Benos, P.V. (2018). Comparison of strategies for scalable causal discovery of latent variable models from mixed data. *Int. J. Data Sci. Anal.* 6, 33–45. <https://doi.org/10.1007/s41060-018-0104-3>.
32. Raghu, V.K., Poon, A., and Benos, P.V. (2018). Evaluation of causal structure learning methods on mixed data types. *Proceedings of 2018 ACM SIGKDD Workshop on Causal Discovery (PMLR)*.
33. Sedgewick, A.J., Shi, I., Donovan, R.M., and Benos, P.V. (2016). Learning mixed graphical models with separate sparsity parameters and stability-based model selection. *BMC Bioinformatics* 17 (Suppl 5), 175. <https://doi.org/10.1186/s12859-016-1039-0>.
34. Rydzynski, C., Daniels, K.A., Karnele, E.P., Brooks, T.R., Mahl, S.E., Moran, M.T., Li, C., Sutiwisesak, R., Welsh, R.M., and Waggoner, S.N. (2015). Generation of cellular immune memory and B-cell immunity is impaired by natural killer cells. *Nat. Commun.* 6, 6375. <https://doi.org/10.1038/ncomms7375>.
35. Rydzynski, C.E., Cranert, S.A., Zhou, J.Q., Xu, H., Kleinstein, S.H., Singh, H., and Waggoner, S.N. (2018). Affinity maturation is impaired by natural killer cell suppression of germinal centers. *Cell Rep.* 24, 3367–3373.e4. <https://doi.org/10.1016/j.celrep.2018.08.075>.
36. Pino, M., Abid, T., Pereira Ribeiro, S., Edara, V.V., Floyd, K., Smith, J.C., Latif, M.B., Pacheco-Sanchez, G., Dutta, D., Wang, S., et al. (2021). A yeast expressed RBD-based SARS-CoV-2 vaccine formulated with 3M-052-alum adjuvant promotes protective efficacy in non-human primates. *Sci. Immunol.* 6, eabh3634. <https://doi.org/10.1126/sciimmunol.abh3634>.
37. Vahey, M.T., Wang, Z., Kester, K.E., Cummings, J., Heppner, D.G., Jr., Nau, M.E., Ofori-Anyinam, O., Cohen, J., Coche, T., Ballou, W.R., and Ockenhouse, C.F. (2010). Expression of genes associated with immunoproteasome processing of major histocompatibility complex peptides is indicative of protection with adjuvanted RTS,S malaria vaccine. *J. Infect Dis.* 201, 580–589. <https://doi.org/10.1086/650310>.
38. Kazmin, D., Nakaya, H.I., Lee, E.K., Johnson, M.J., van der Most, R., van den Berg, R.A., Ballou, W.R., Jongert, E., Wille-Reece, U., Ockenhouse, C., et al. (2017). Systems analysis of protective immune responses to RTS,S malaria vaccination in humans. *Proc. Natl. Acad. Sci. U S A* 114, 2425–2430. <https://doi.org/10.1073/pnas.1621489114>.
39. Neafsey, D.E., Juraska, M., Bedford, T., Benkeser, D., Valim, C., Griggs, A., Lievens, M., Abdulla, S., Adjei, S., Agbenyega, T., et al. (2015). Genetic diversity and protective efficacy of the RTS,S/AS01 malaria vaccine. *N. Engl. J. Med.* 373, 2025–2037. <https://doi.org/10.1056/NEJMoa1505819>.
40. Arts, R.J., Joosten, L.A., and Netea, M.G. (2016). Immunometabolic circuits in trained immunity. *Semin. Immunol.* 28, 425–430. <https://doi.org/10.1016/j.smim.2016.09.002>.
41. Lu, L.L., Das, J., Grace, P.S., Fortune, S.M., Restrepo, B.I., and Alter, G. (2020). Antibody Fc glycosylation discriminates between latent and active tuberculosis. *J. Infect Dis.* 222, 2093–2102. <https://doi.org/10.1093/infdis/jiz643>.
42. Olin, A., Henckel, E., Chen, Y., Lakshmikanth, T., Pou, C., Mikes, J., Gustafsson, A., Bernhardsson, A.K., Zhang, C., Bohlin, K., and Brodin, P. (2018). Stereotypic immune system development in newborn children. *Cell* 174, 1277–1292.e14. <https://doi.org/10.1016/j.cell.2018.06.045>.
43. Robertson, S.A., Skinner, R.J., and Care, A.S. (2006). Essential role for IL-10 in resistance to lipopolysaccharide-induced preterm labor in mice. *J. Immunol.* 177, 4888–4896. <https://doi.org/10.4049/jimmunol.177.7.4888>.
44. Gomez-Lopez, N., StLouis, D., Lehr, M.A., Sanchez-Rodriguez, E.N., and Arenas-Hernandez, M. (2014). Immune cells in term and preterm labor. *Cell Mol. Immunol.* 11, 571–581. <https://doi.org/10.1038/cmi.2014.46>.
45. Rolle, L., Memarzadeh Tehran, M., Morell-Garcia, A., Raeva, Y., Schumacher, A., Hartig, R., Costa, S.D., Jensen, F., and Zenclussen, A.C. (2013). Cutting edge: IL-10-producing regulatory B cells in early human pregnancy. *Am. J. Reprod. Immunol.* 70, 448–453. <https://doi.org/10.1111/aji.12157>.
46. Jensen, F., Muzzio, D., Soldati, R., Fest, S., and Zenclussen, A.C. (2013). Regulatory B10 cells restore pregnancy tolerance in a mouse model. *Biol. Reprod.* 89, 90. <https://doi.org/10.1095/bioreprod.113.110791>.
47. Libbrecht, M.W., and Noble, W.S. (2015). Machine learning applications in genetics and genomics. *Nat. Rev. Genet.* 16, 321–332. <https://doi.org/10.1038/nrg3920>.
48. Cusick, M.E., Yu, H., Smolyar, A., Venkatesan, K., Carvunis, A.R., Simonis, N., Rual, J.F., Borick, H., Braun, P., Dreze, M., et al. (2009). Literature-curated protein interaction datasets. *Nat. Methods* 6, 39–46. <https://doi.org/10.1038/nmeth.1284>.
49. Pearl, J. (2010). An introduction to causal inference. *Int. J. Biostat.* 6, Article 7. <https://doi.org/10.2202/1557-4679.1203>.
50. Bing, X., Bunea, F., Strimas-Mackey, S., and Wegkamp, M. (2021). Prediction in latent factor regression: adaptive PCR, interpolating predictors and beyond. *J. Mach. Learn. Res.* 22, 1–50. <https://www.jmlr.org/papers/volume22/20-768/20-768.pdf>.
51. Lee, J.D., and Hastie, T.J. (2015). Learning the structure of mixed graphical models. *J. Comput. Graph Stat.* 24, 230–253. <https://doi.org/10.1080/10618600.2014.900500>.
52. Friedman, J., Hastie, T., and Tibshirani, R. (2008). Sparse inverse covariance estimation with the graphical lasso. *Biostatistics* 9, 432–441. <https://doi.org/10.1093/biostatistics/kxm045>.
53. Liu, H., Roeder, K., and Wasserman, L. (2010). Stability approach to regularization selection (StARS) for high dimensional graphical models. *Proceedings of the 23rd International Conference on Neural Information Processing Systems - Volume 2 (Curran Associates Inc.)*.
54. Lam, C., and Yao, Q. (2012). Factor modeling for high-dimensional time series: inference for the number of factors. *Ann. Stat.* 40, 694–726, 633.