# A curated collection of human vaccination response signatures

## Authors
Kenneth C. Smith[1]*, Daniel G. Chawla[2]*, Bhavjinder K. Dhillon[3]*, Zhou Ji[1], Randi Vita[4], Eva C. van der Leest[3], Jing Yi (Jessica) Weng[3], Ernest Tang[3], Amani Abid[3], The Human Immunology Project Consortium (HIPC)[&], Bjoern Peters[4,6], Robert E.W. Hancock[3,**], Aris Floratos[1,7,**], Steven H. Kleinstein[2,5,8,**]

*co-first authors
**co-senior authors

## Affiliations
[1]Department of Systems Biology, Columbia University, New York, New York, USA
[2]Interdepartmental Program in Computational Biology and Bioinformatics, Yale University, New Haven, CT, USA
[3]Department of Microbiology & Immunology, University of British Columbia, Vancouver, Canada
[4]Center for Infectious Disease and Vaccine Research, La Jolla Institute for Immunology, La Jolla, CA, USA
[5]Department of Pathology, Yale School of Medicine, New Haven, CT, USA
[6]Department of Medicine, Division of Infectious Diseases and Global Public Health, University of California, San Diego, La Jolla, CA, USA
[7]Department of Biomedical Informatics, Columbia University, New York, New York, USA
[8]Department of Immunobiology, Yale School of Medicine, New Haven, CT, USA

[&] Members of the HIPC Steering Committee: Raphael Gottardo (Fred Hutchinson Cancer Research Center, Seattle, WA, USA); Elias K. Haddad (Drexel University, Philadelphia, PA, USA); David A. Hafler (Yale School of Medicine, New Haven, CT, USA); Eva Harris (University of California, Berkeley, Berkeley, CA, USA); Donna Farber (Columbia University Medical Center, New York, NY, USA); Steven H. Kleinstein (Yale School of Medicine, New Haven, CT, USA); Ofer Levy (Boston Children's Hospital, Boston, MA, USA); Julie McElrath (Fred Hutchinson Cancer Research Center, Seattle, WA, USA); Ruth R. Montgomery (Yale School of Medicine, New Haven, CT, USA); Bjoern Peters (La Jolla Institute for Immunology, La Jolla, CA, USA.); Bali Pulendran (Stanford University School of Medicine, Stanford University, Stanford, CA, USA); Adeeb Rahman (Icahn School of Medicine at Mount Sinai, New York, New York, USA); Elaine F. Reed (David Geffen School of Medicine at University of California, Los Angeles, CA, USA); Nadine Rouphael (Emory University School of Medicine Decatur, Atlanta, GA, USA); Minnie Sarwal (University of California, San Francisco, San Francisco, CA, USA); Rafick Sekaly (Emory University School of Medicine Decatur, Atlanta, GA, USA); Ana Fernandez-Sesma (Icahn School of Medicine at Mount Sinai, New York, New York, USA); Alessandro Sette (La Jolla Institute for Immunology, La Jolla, CA, USA); Ken Stuart (Seattle Children's Research Institute, Seattle, WA, USA); John S Tsang (NIAID and Center for Human Immunology (CHI), NIH, Bethesda, MD, USA)

corresponding authors: Aris Floratos (af2202@cumc.columbia.edu), Steven Kleinstein (steven.kleinstein@yale.edu)

## Abstract
Recent advances in high-throughput experiments and systems biology approaches have resulted in hundreds of publications identifying "immune signatures". Unfortunately, these are often described within text, figures, or tables in a format not amenable to computational processing, thus severely hampering our ability to fully

exploit this information. Here we present a data model to represent immune signatures, along with the Human Immunology Project Consortium (HIPC) Dashboard (www.hipc-dashboard.org), a web-enabled application to facilitate signature access and querying. The data model captures the biological response components (e.g., genes, proteins, cell types or metabolites) and metadata describing the context under which the signature was identified using standardized terms from established resources (e.g., HGNC, Protein Ontology, Cell Ontology). We have manually curated a collection of >600 immune signatures from >60 published studies profiling human vaccination responses for the current release. The system will aid in building a broader understanding of the human immune response to stimuli by enabling researchers to easily access and interrogate published immune signatures.

## Introduction

Systems-level profiling of the human immune system has generated important insights into the mechanisms by which humans respond to exposures such as vaccination. These studies, including many conducted through the Human Immune Project Consortium (HIPC), have generated hundreds of publications. While repositories exist to promote re-use of primary experimental immunology data generated from these efforts, such as the Gene Expression Omnibus[1] (GEO) and the NIAID Division of Allergy, Immunology, and Transplantation (DAIT)-sponsored Immunology Database and Analysis Portal[2] (ImmPort), there is no centralized framework to aggregate and organize the published findings resulting from the analysis of this data, and particularly the coherent sets of biomarkers, termed here "signatures". Additionally, such signatures are not published in a consistent format between publications and may be presented as text, tables, or images. This heterogeneity presents a barrier to comparative analyses since identifying published signatures, for example of a vaccine response, requires extensive manual curation of the literature that must be repeated by investigators each time they wish to interpret a set of results. Here, we propose a model to standardize the representation of these published findings and present the Human Immunology Project Consortium (HIPC) Dashboard—a searchable interface to query curated signatures from the corpus of human immunology literature.

We define a 'signature' as the information required to specify a published result. This includes alterations in the levels of a set of one or more response component(s), i.e., biological entities such as genes or cell types, that are defined by a particular comparison in the context of an immune exposure. The signature also includes contextual information (termed metadata) such as the conditions and circumstances under which the signature was identified, the tissues or cells that were assayed, as well as clinical data such as demographic information about the groups that were included in the analysis. As a motivating example, a study by Bucasas et al.[3] identified a set of genes that are up-regulated in individuals with higher antibody responses (comparison) after vaccination with the 2008-2009 trivalent influenza vaccine (exposure) in an adult cohort. The expression of genes STAT1, IRF9, SPI1, CD74, HLA-E, and TNFSF13B one day after influenza vaccination was predictive of greater antibody responses. In the paper, these results were represented in a table, though similar findings often appear as text or within figures. Without standardization, such findings are not easily accessible to the wider scientific community for further analysis.

Several existing resources define pathway and gene module signatures through re-analysis of raw data, but few capture the original findings published with these data or are specifically geared towards human immunology research. Among these resources are the OMics Compendia Commons (OMiCC)[4], EnrichR[5,6], the integrative Library of Integrated Network-based Cellular Signatures (iLINCS)[7], the Molecular

111  Signatures Database (MSigDB)[8,9], and VaximmutorDB[10]. OMiCC crowdsources
112  annotations for gene expression data to be used in re-analysis and novel signature
113  generation. EnrichR and iLINCS offer biological annotations built from data re-
114  analyzed en masse, but similarly do not capture published findings. MSigDB does
115  include manually curated gene signatures along with those derived from data re-
116  analysis, albeit with fewer contextual details than captured for the HIPC Dashboard.
117  VaximmutorDB captures published gene expression and proteomic signatures but not
118  cell-type frequency signatures, and signatures from this database are not yet
119  downloadable in a machine-readable format.
120
121  To improve access and to promote reuse of published signatures, we designed a
122  data model that standardizes the content and context of published immune
123  signatures. Our initial curation efforts have focused on gene expression and cell-type
124  frequency/activation signatures of human vaccine responses, but this framework is
125  extensible to other domains such as response to infection. We captured what is
126  changing, (e.g. groups of genes), how that response component changed (e.g. up- or
127  down-regulation), where this change was observed (e.g. in sorted CD8+ T cells from
128  adults), and the comparison that was performed (e.g. individuals with high vs. low
129  antibody titers post-vaccination). We then manually curated signatures from
130  publications both within and outside of HIPC that described changes in gene
131  expression, cell-type frequencies, or cell activation state in response to vaccination.
132  To disseminate these immune signatures, we developed the HIPC Dashboard
133  (www.hipc-dashboard.org), a web-accessible, user-friendly interface to enable
134  signature searching and browsing, and to facilitate rapid comparative analyses. The
135  design of the HIPC Dashboard is based on a similar infrastructure we developed
136  previously for the Cancer Target Discovery and Development network, the CTD[2]
137  Dashboard[11], and leverages the same underlying ontological framework for the
138  standardized representation of research findings as well as the emphasis on the
139  consistent, curation-mediated use of controlled vocabularies for linking findings
140  reported in different publications.
141
142  **Results**
143  *A Data Model for Immune Signatures of Vaccination*
144
145  We developed a data model that captures, in a detailed and consistent format, the
146  essential information embedded in published immune signatures of vaccination for
147  dissemination through the HIPC Dashboard (**Table 1**). Key elements of this data
148  model (e.g., genes, vaccines, etc.) are specified using controlled vocabularies, thus
149  making immune signatures of vaccination amenable to data mining and promoting
150  compatibility with projects both within and outside of HIPC. A signature as defined in
151  this model encapsulates both a change in the behavior or abundance of a biological
152  response component as well as the metadata describing the context under which the
153  signature is identified, including (1) the tissue in which the signature was observed,
154  (2) the immune exposure and timing underlying the observed comparison, and (3)
155  clinical details of the cohort from which tissue samples were taken, including age
156  (**Figure 1**). The model accommodates many types of biological response
157  components (gene, protein, metabolite, pathway, and cell type (e.g. subsets of blood
158  cells). We focused on gene expression and cell type signatures of vaccination, but
159  the data model and HIPC Dashboard infrastructure are flexible and can be easily
160  expanded to accommodate arbitrary signature types.
161
162  To facilitate data mining and comparative analyses between different conditions
163  including vaccine types, and to afford consistency between this database and other
164  projects using the same controlled vocabulary terms, standardized terms and
165  ontology links were used for as many biological response components, immune

166  exposures, and demographic fields as possible. Gene and cell type response
167  components were standardized to the HGNC[12] (as provided through the NCBI) and
168  Cell Ontology[13,14] (CL) respectively to enable comparisons across publications using
169  different naming conventions. Cell types from the Cell Ontology were further
170  differentiated using protein marker terms drawn from Protein Ontology[15] (PRO),
171  where possible (e.g. IFNG+ T cells). This same naming convention is used to
172  describe the tissue in which the signature was observed[16].
173
174  To annotate the immune challenges driving each signature, we utilized the Immune
175  Exposure model[17], which provides a standardized description of a broad range of
176  potential and actual exposures to different immunological agents (e.g., vaccination,
177  laboratory confirmed infection, living in an endemic area, etc.). Immune exposures
178  are broken down into Exposure Process, Exposure Material, Disease Name, and
179  Disease Stage. Each of these components is modeled using standardized ontology
180  terminology. Within the data model for the HIPC Dashboard, Exposure Materials such
181  as vaccines are captured using terms in the Vaccine Ontology[18] (VO), which further
182  link to target pathogens and strains using the NCBI Taxonomy[19,20]. While these
183  ontology choices reflect our initial focus on vaccination, the data model can
184  accommodate other exposure processes beyond vaccination, with links to
185  appropriate ontologies. Integration of the Immune Exposure model in the HIPC
186  Dashboard data model promotes interoperability with other projects that have
187  adopted its use, both within and outside of HIPC, including data repositories such as
188  ImmPort and the Immune Epitope Database[21].
189
190  Cohort information that is important for interpreting signatures is also captured.
191  Cohort descriptors can vary widely between studies and can include, for example,
192  sex, antibody response titers, geographic location, health status, vaccination or
193  infection history, etc. This information is currently recorded as unstructured text to
194  maintain flexibility. Cohort age range is standardized separately by storing minimum
195  and maximum ages along with their units. Additional fields describe the particular
196  perturbations that drive the changes to the biological response components.  The
197  "comparison" field describes the cohort groups whose differential response under the
198  perturbation is measured. Examples of comparison groups include measurements
199  taken at two different time points (e.g. day1 vs. day 0), correlation with antibody
200  response, differing antibody response outcomes (e.g. high vs. low responders), or
201  comparisons across different demographic parameters such as age or sex (e.g.
202  younger vs. older, female vs. male). The "response_behavior" field captures the
203  directionality of the differential response (e.g. up or down, positively- or negatively-
204  correlated) under the specified comparison. Fields for which a formal controlled
205  vocabulary was not used, such as cohort descriptions and the comparison, are stored
206  as free text.
207
208  Finally, each signature is tagged with a Pubmed ID (PMID) and publication year field,
209  to connect observations to their literature sources.
210
211  *Manual Curation of Published Signatures*
212
213  A defined Pubmed search strategy was used (see Methods) to assemble an initial list
214  of publications comprising studies that involved a systems-level profiling of
215  measurable changes before and after vaccination in human subjects. The
216  publications were culled for signatures reporting statistically significant changes in
217  gene expression, cell-type frequency or cell activation state induced by an immune
218  exposure when comparing groups with different features (such as high vs low
219  responders as defined by antibody titers). We focused on components that also
220  included information about response behavior (e.g. up- or down-regulated). In total,

221 665 immune signatures were manually curated from 69 published studies. After
222 standardization and quality control (see Methods), these curated gene and cell type
223 signatures included 13,812 unique genes, 152 unique cell types (including protein
224 markers and additional type-modifiers), and 44 pathogens across 56 vaccines **(Table**
225 **2)**. Table 3 illustrates a typical gene-expression type signature after tissue, gene
226 symbol and pathogen standardization.
227
228 *The HIPC Signatures Dashboard*

229 The data model provides a means for representing immune signatures in a
230 structured, standardized, machine-readable manner while the curation process
231 enables the cross-referencing of signatures from different publications based on their
232 shared response components, by enforcing the consistent use of controlled
233 vocabularies for codifying these components (e.g., genes, cell types, tissues,
234 vaccines, and pathogen strains). These capabilities come together in the "HIPC
235 Dashboard" (http://hipc-dashboard.org), a web application developed to enable
236 dissemination of the curated set of immune response signatures. The HIPC
237 Dashboard allows signature browsing, as well as searching for one or multiple
238 response components (using the corresponding controlled vocabulary terms and their
239 synonyms), to retrieve all immune signatures involving the query response
240 component(s) across all curated publications.

241 The central viewable element of the HIPC Dashboard is the "Observation Summary",
242 a human-readable description of the information captured in an immune signature.
243 Observation summaries are constructed "on the fly", using template text devised as
244 part of the curation process. The template has placeholders for the various elements
245 of an immune signature, including the response component (gene or cell type) and
246 the response behavior type (up/down or correlation). When a specific signature is
247 selected in the process of browsing or searching the Dashboard, the observation
248 summary for that signature is instantiated by replacing the template placeholders with
249 the relevant values from that signature. For example, a joint search on the terms
250 "CD4" and "Zostavax" yields about 35 observation summaries. One of these is related
251 to a change in cell-type frequency:
252
253
254 In peripheral blood mononuclear cell, CD4-positive, alpha-beta memory T cell **&**
255 **CD38+, HLA-DR+, VZV tetramer+** frequency was **up** at **14 days** from **time of**
256 **vaccination** for the comparison **14d vs 0d** in cohort **50-75 yo** after exposure to
257 Zostavax targeting Human alphaherpesvirus 3 (**details »**)
258
259 Here, the "&" sign separates the Cell Ontology cell type from Protein Ontology
260 surface and other markers. In a second example, a joint search on the terms
261 "CXCL10" and "BCG" yields 6 observation summaries, one of which reports a
262 correlation of gene expression (at 1 day post-vaccination) to an ELISpot result at 28
263 days post-vaccination:
264
265 In blood, CXCL10 gene expression at **1 day** from **time of vaccination** was **positively**
266 **correlated** with **IFN-gamma ELISpot spot forming cell 28d** in cohort **4-6 mo**
267 **subgroup BCG-primed** after exposure to MVA85A targeting Mycobacterium
268 tuberculosis variant bovis BCG (**details »**)

269 In both cases, placeholders in the observation summary template have been
270 replaced by controlled terms for the response components and ontology-linked
271 metadata (blue, hyperlinked text) and by free-text metadata describing informative
272 experiment details (black, bold text). Following the hyperlink for a controlled term
273 leads to a dedicated page for the corresponding biological entity, providing additional

5

274 details (including links to relevant external annotation sources, e.g., Entrez,
275 GeneCards and UniProt for genes) as well as a listing of all the immune signatures
276 stored in the Dashboard that involve that entity **(Figure 2A)**. Further, the "details" link
277 at the end of each observation summary points to an "Observation" page **(Figure 2B)**
278 containing detailed information about the corresponding immune signature, including
279 a full listing of all its available metadata. This includes, for example, structured text for
280 values such as age group and days post-immunization, and links to download the full
281 signature source data (including all metadata) in tab-delimited form. Additionally,
282 each observation includes a link to a file containing the complete set of response
283 components from which it was derived, e.g. the full list of genes or cell types.

284

## Discussion

286

287 Users of the HIPC Dashboard can easily search and examine hundreds of immune
288 signatures related to human vaccination responses. Consolidating these publications,
289 standardizing their findings in a database, and disseminating them through the
290 Dashboard interface allows for rapid comparative analyses and re-use of published
291 findings. This is particularly important for identifying commonalities across studies
292 that may reflect shared mechanisms. The HIPC Dashboard can offer broad insights
293 into the mechanisms by which our immune systems respond to vaccination and will
294 be of great value to the vaccine research community. Although the HIPC Dashboard
295 is not designed as an analysis engine, all signatures are made available for download
296 so that users may perform more sophisticated and targeted downstream analyses.

297

298 Among data resources dedicated to the collection of vaccination signatures, the HIPC
299 Dashboard is nearly unique in its emphasis on manual curation of published
300 literature. To the best of our knowledge, only MSigDB and VaximmutorDB maintain
301 signatures curated from publications. MSigDB provides minimally redundant gene
302 sets for enrichment analyses, but unlike the HIPC Dashboard does not attempt to
303 capture the full biological context of published results. A reduced set of our curations
304 has recently been made available for gene set enrichment analyses through MSigDB
305 under the C7 VAX gene sets. VaximmutorDB provides access to a collection of
306 immune factors (genes/proteins) that change in response to vaccination against 46
307 pathogens. Compared to VaximmutorDB, the HIPC Dashboard offers several
308 advantages, including: (i) a wider breadth in the types of response components and
309 immune changes that are captured, (ii) improved browsing functionality that facilitates
310 comparisons of immune changes across studies and vaccines, and (iii) the ability to
311 download signature data.

312

313 As of the date of this publication, 152 unique cell types and 13,812 distinct genes
314 have been collected in the HIPC Dashboard; this large number of published results
315 allows users to quickly examine the role of particular biological response
316 components, such as individual genes or cell types, across studies. Most of the
317 currently curated gene signatures have fewer than 50 genes, with a range of 1 to
318 2,036 (**Figure 3A**), while most cell-type frequency or cell activation state signatures
319 have only a single cell type, with a range up to 9 (**Figure 3B**). These signatures
320 represent findings from 16 tissues and tissue extracts, including blood, PBMCs, T
321 cells, B cells, monocytes, and NK cells. Nearly 250 entries describe changes over
322 time, more than 75 capture antibody response-associated signatures, and several
323 others come from studies that report effects of age and T cell responses.

324

325 The frequency with which cell types and genes are reported in the Dashboard offers
326 insights about key players of the human immune response to vaccination. The most
327 commonly reported genes are STAT1, a key mediator of immune response activated

328    by cytokines and interferons; GBP1, an interferon induced gene involved in innate
329    immunity; IFI44L, a paralog of Interferon Stimulated Protein 44, and SERPING1, a
330    complement cascade protein **(Figure 3C)**. The most common cell types across
331    pathogens and comparisons are NK cells, CD4+ T cells, and CD8+ T cells. **(Figure**
332    **3D)**. We searched the HIPC Dashboard data for genes with vaccination signatures
333    across six or more pathogens and found a set of 36 associated genes across 12
334    target pathogens **(Figure 3E)**. Many are interferon stimulated genes, Toll-like
335    receptors, or members of the complement cascade, potentially reflecting a common
336    transcriptional program in response to many different vaccinations.
337
338    By design, our current implementation captures vaccination signatures as they are
339    reported in the literature, but it does not include related methodological or statistical
340    information regarding signature discovery (e.g. p-value cut-off) or provide analytical
341    tools that can be applied to the curated signatures. Test statistics are not usually
342    comparable across study designs, and we believe this information may give users a
343    false sense that some signatures are more statistically reliable than others. We
344    instead defer to the judgement of each study's authors and their peer reviewers, and
345    capture signatures as they were reported in each publication. We also caution that
346    bias regarding the number of times particular genes or cell types were investigated
347    might skew relationships in the Dashboard, thus precluding certain types of analyses.
348    As a result, high level analytical tools have not been integrated into the Dashboard,
349    although all of the signatures with full metadata can be downloaded to enable the use
350    of third-party tools. Despite this, we believe the signatures available in the HIPC
351    Dashboard will allow the research community to quickly query the literature and
352    provide valuable comparisons and context for their own experimental results.
353
354    The number of genes and cell types captured reflects publications curated through
355    January 2021, but we anticipate the HIPC Dashboard will undergo regular updates to
356    accommodate new findings and additional domains of interest. The current
357    implementation includes gene and cell type response components, as these
358    represent the most commonly published signature types, but it will be valuable to also
359    curate other response components, such as pathways, proteins, and metabolites.
360    Additionally, we recognize that researchers may wish to compare a vaccine response
361    against a particular pathogen to its corresponding disease response; it is easy to see
362    how future iterations could expand the existing vaccine signature framework to
363    capture signatures of infection. Based on our experience, we expanded the data
364    model to include figure numbers or supplementary file annotations within publications
365    as this can greatly simplify quality control during manual curation. We have provided
366    links in the Dashboard to original sources wherever possible. We are also keen on
367    exploring advancements in text-mining and artificial intelligence (AI), to assess how
368    they can assist in automating signature identification and coding. To that end, the
369    immune signatures in the HIPC Dashboard can be used as a data source for
370    training/testing such AI solutions in the future.
371
372    In summary, we present the HIPC Dashboard (hipc-dashboard.org) to provide the
373    vaccine research community with easy access to hundreds of published human
374    systems vaccinology signatures. This resource will allow researchers to rapidly
375    compare their own experimental results against existing findings that may otherwise
376    be difficult to locate in the literature. This resource encourages the re-use of
377    published results for advancing our understanding of human vaccine responses and
378    provides a framework that can be extended to capture signatures from other types of
379    immune exposures.
380

## Methods
381

382    *Manual Curation*

383
384  The initial list of publications to curate into the HIPC Dashboard were derived from a
385  PubMed search of papers matching the terms "Vaccine [AND] Signatures" or "Vaccine
386  [AND] Gene expression". Publications were further filtered to meet a set of inclusion
387  criteria: (i) study involved human subjects, (ii) provided a comparison of a measurable
388  change or correlation before and after vaccination (or challenge), and (iii) were
389  reported as statistically significant. Signatures were excluded if they were missing
390  directionality, or if they were derived from datasets external to the publication, to
391  avoid redundancy. Two data curators manually collected a standard set of
392  information from each study according to the designed data model (see **Figure 1**)
393  and recorded it into a spreadsheet. Each signature was entered by one curator, and
394  subsequently double-checked by the second curator. Table 3 shows a representative
395  portion of the standardized information captured for each signature (12 data fields out
396  of a total of 25).
397
398  Assays in the curated publications included gene expression analysis and measured
399  changes in cell-type frequency and cell activation state. Each publication could give
400  rise to any number and type of individual signatures. The signature content is
401  centered on a list of biological response components (genes or cell types) that had a
402  statistically significant change in the assay. These were designated as "response
403  components" to capture different types of entities in a single standardized template
404  column. For example, for gene expression assays, this was often a list of differentially
405  expressed genes.
406
407  For genes, an initial manual curation process was applied to make a first pass at
408  symbol standardization and detect any mistakes in copying.  Gene names, symbols
409  and or IDs (which may include HGNC symbols, Entrez IDs, Ensembl IDs etc.) were
410  searched in turn against Panther, (www.pantherdb.org)[22], using the  "Functional
411  classification viewed in gene list" search, followed if needed by searches against
412  UniProt (www.uniprot.org)[23] and NCBI (www.ncbi.nlm.nih.gov/search) until either a
413  match or updated symbol was found; in the case of no match, the original
414  representation was left unchanged. Any IDs that matched entries that were
415  deprecated or defined as pseudogenes were removed from the curation.
416
417  *Data Standardization*
418
419  The manually curated data required further steps to match terms encountered to their
420  appropriate ontology representations. A number of the translations described below
421  were orchestrated using an R script, which generated files ready for loading into the
422  HIPC Dashboard.
423
424  *Gene symbols* - Outdated gene symbols and known aliases were translated to their
425  current NCBI representation, which is the HGNC symbol in all but one case. The first
426  pass of conversion used the function alias2SymbolUsingNCBI() from the
427  Bioconductor limma package[24] with the most recent available gene annotation file.
428  This function returns either an exactly matching official symbol, or if none, the alias
429  with the lowest EntrezID.  We followed this by a second R package, HGNChelper[25],
430  which was able to resolve additional unmatched gene symbols to valid NCBI
431  symbols. Genbank accession numbers were converted to gene symbols where
432  possible using the org.Hs.egACCNUM2EG translation table which is part of the
433  Bioconductor org.Hs.eg.db package[26].  Selected symbols still not matching NCBI
434  names were investigated and corrected manually where possible after checking the
435  original publications for context or for errors in transcription. Symbols for which no
436  valid NCBI gene symbol was found, e.g. some pseudogenes, antisense, or
437  uncharacterized genes, are not included in the HIPC Dashboard proper, since a

438 requirement of the Dashboard framework is that all gene symbols must appear in the
439 controlled vocabulary (NCBI/HGNC). However, these symbols are included in
440 downloadable complete gene lists for each signature.
441

442 *Cell types* - Cell types as response components were first curated from the
443 publications as published using a combination of cell type terms and additional
444 descriptive terms, such as protein marker expression. This information was then
445 mapped to a combination of Cell Ontology and Protein Ontology terms, according to
446 a published model[16]. Note that cell types can appear in two different contexts, either
447 as response components themselves, or as the cell type isolated for gene expression
448 experiments. In some cases, additional information was provided which could not be
449 mapped to an ontology term[16]. This type of information related to a wide variety of cell
450 identification techniques and included the use of additional stains such as viability
451 dyes or tetramer staining. For each set of terms, an entry was created in a lookup
452 table by assigning (1) a parent cell type from the cell ontology, (2) mapping additional
453 protein marker terms to the protein ontology, and then (3) separately retaining as free
454 text descriptors not mapped to an ontology entry, such as tetramer specificity. Thus,
455 the original entry was mapped to up to three descriptor columns, which can be
456 combined as needed for display purposes. For a full, translated cell type, the
457 displayed format is the cell ontology name followed by, if there are additional terms,
458 the "&" symbol, followed by any PRO terms and then any free text.

459

460 *Vaccines* - Vaccine names collected from the literature were manually mapped to the
461 most specific vaccine ontology term available. If a specific vaccine could not be found
462 in VO, new terms were requested. Some examples of terms we included were
463 VO_0004899 (2012-2013 seasonal trivalent inactivated influenza vaccine),
464 VO_0003961 (ChAd63-KH vaccine, *Leishmania donovani*), VO_0004890 (gH1-Qbeta
465 vaccine, novel pandemic-influenza), and VO_0004891 (CN54gp140□+□GLA, HIV-1).
466

467 *Pathogens* - The viral, bacterial or protozoan pathogens targeted by each vaccine are
468 represented with terms from the NCBI Taxonomy. For the case of influenza
469 vaccines, a table was created mapping vaccines by year of administration and type
470 (e.g. trivalent or quadrivalent) to their seasonal viral components, unless otherwise
471 indicated in the publication (e.g. for monovalent or specialized vaccines such as
472 *Pandemrix*). For the few cases where the exact viral strain was not present in the
473 NCBI Taxonomy, the closest more general term in the hierarchy was used. This
474 mapping table was used to substitute in the actual viral pathogen names.
475

## Data Availability

477 Curated signatures are available on the HIPC Dashboard website (http://hipc-
478 dashboard.org) and Github (see Code Availability for details). Files listing all
479 response components for a signature can be downloaded from within individual
480 observations in the Dashboard. The complete set of signature data can be
481 downloaded from the GitHub repository at https://github.com/floratos-lab/hipc-
482 dashboard-pipeline. This repository contains copies of (1) the original curated data
483 sheets, (2) the response components in individual files, one per signature, (3) the
484 response components in the Broad GMT format[8], and (4) the actual tab-delimited
485 Dashboard load files, in which the complete signature data is fully denormalized into
486 an easy-to-parse format. Further details about the file formats are available on the
487 GitHub project page. R session information for the Dashboard signature pre-
488 processing pipeline is available on GitHub.

489    Code Availability

490    Source curated data and mapping files (cell types, vaccine components, the NCBI
491    gene file used, etc.), as well as the R code for the processing pipeline used to create
492    the    Dashboard    submission    files,    are    available    on    GitHub    at
493    https://github.com/floratos-lab/hipc-dashboard-pipeline.    The    data    in    this    paper
494    corresponds to pipeline and data version 1.2.1 in the pipeline GitHub repository.
495    Code    for    the    HIPC    Dashboard    web    interface    is    available    on    GitHub    at
496    https://github.com/floratos-lab/hipc-signature.

497    *Supported web browsers* - The HIPC Dashboard has been tested on recent versions
498    of

499    • Chrome (Version  93.0.4577.63 (Official Build) (64-bit))
500    • Firefox (Version  92.0 (64-bit))
501    • Edge (Version 93.0.961.38 (Official build) (64-bit))

502

## Acknowledgements

## Author contributions

510    Conceptualization: AF, SHK; Methodology: KCS, DC, BKD, ZJ, RV, BP, REWH, AF,
511    SHK; Software: KCS, ZJ, AF; Investigation: KCS, DC, BKD, EvdL, JW, ET, AA;
512    Writing - Original Draft: KCS, DC, BKD, AF, SHK; Writing - Review and Editing: all
513    authors
514

## Competing interests

516    S.H.K. receives consulting fees from Northrop Grumman and Peraton.
517

## Figure Legends

519    **Figure 1.** Overview of the manual curation process for extracting immune signatures
520    from relevant publications into the HIPC signatures database and a web-accessible
521    HIPC Dashboard. The middle panel highlights the various fields that are captured for
522    a  given  immune  signature,  with  examples  provided  in  red  font.  Key  fields  are
523    standardized  using  existing  ontologies  or  pre-defined  criteria  in  order  to  capture  a
524    wide array of signatures.
525

526    **Figure 2.**  HIPC  Dashboard  web  interface.  **A**.  Subject  page  for  cell  type  "CD4-
527    positive, alpha-beta T cell" showing a link-out to the Cell Ontology, the filtering box to
528    further  narrow  the  displayed  observations,  and  the  first  two  observation  summaries
529    ("Related  observations").  **B**.  Partial  view  of  a  details  page  for  a  CD4-positive,  alpha-
530    beta T cell observation.  For each controlled term, its name, plus its class, role, and
531    description are shown. Linked pages list details from the relevant ontology and list all
532    observations  containing  that  term.  The  class  equates  to  its  controlled  vocabulary
533    type; values are cell subset, gene, pathogen, and vaccine. Roles are used to further
534    differentiate  how  each  term,  whether  controlled  or  standardized,  is  being  used.

535   Among the classes in the HIPC Dashboard, only the class "cell subset" has more
536   than one role, these being "tissue" and "cell_biomarker". Full metadata, not shown
537   here, is contained in the table labeled "Evidence" at the bottom.
538
539   **Figure 3.**  Summarization of HIPC Dashboard contents. Signature size distributions
540   showing the number of response components across **A.** gene and **B.** cell type
541   signatures. **C.** Word cloud of the top 50 most frequent gene symbols and **D.** top 10
542   most frequent cell types, where size corresponds to the total number of observations
543   in the Dashboard. **E.** Heatmap of recurring genes across vaccines targeting different
544   pathogens. Temporally associated genes in adult whole blood or PBMCs were
545   filtered to those with signatures for six or more pathogens. Color indicates up (red) or
546   down (blue) regulation.   Genes with opposing directions in multiple studies were
547   marked 'trends up' or 'trends down' according to the most common direction (or
548   marked 'no consensus' for perfect ties).
549
550   **Tables**
551
552   **Table 1.** Data model for capturing immune signatures. Genes and cell types are
553   captured as response components, with terms standardized against NCBI/HGNC or
554   CL+PRO, respectively. Exposures are captured and standardized against VO terms
555   and NCBI Taxonomy IDs. Metadata includes observed tissue, study timing, cohort
556   descriptors, and age characteristics. *fields in the Immune Exposure model[17]
557
558   **Table 2.** Dashboard summary statistics for gene and cell-type signatures.  "Joint"
559   refers to the union of the two signature types, as they overlap in the various
560   components. "Total Response Components" lists the number of genes or parent cell
561   types from the Cell Ontology (CL) across all signatures.  When additional cell-type
562   markers are included, e.g. from the Protein Ontology (PRO), there are 152 unique
563   cell types represented among the signatures.   "Response Components per
564   Signature" shows the range of the number of response components found in
565   individual signatures.
566
567   **Table 3.** Key fields in the immune signature data model for a gene expression
568   signature.  This signature reports positive correlation between gene expression and a
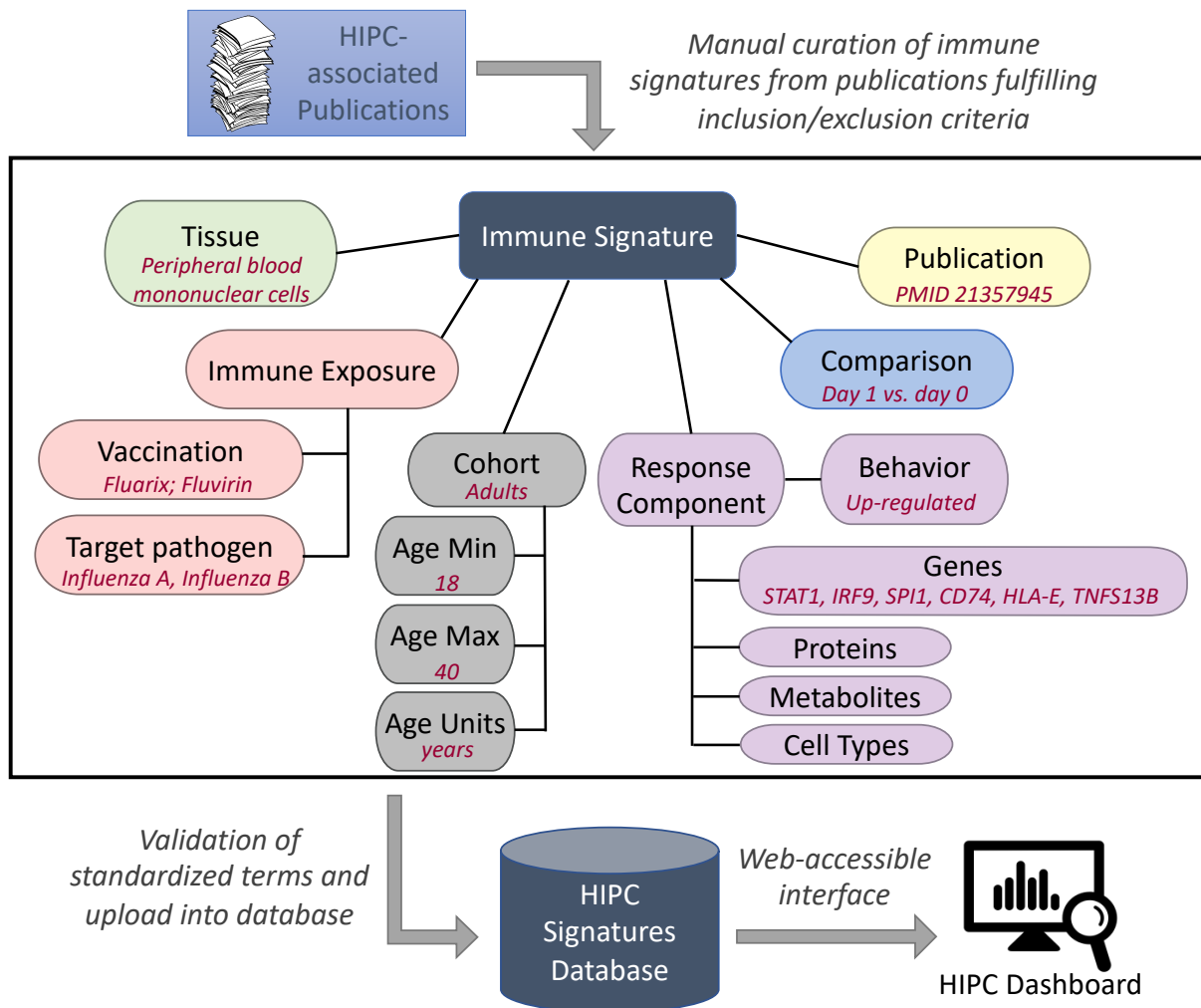569   computed titer response index (TRI).
570
571

11

572    1.  Edgar, R., Domrachev, M. & Lash, A. E. Gene Expression Omnibus: NCBI gene

573        expression and hybridization array data repository. *Nucleic Acids Res* **30**, 207–

574        210 (2002).

575    2.  Bhattacharya, S. *et al.* ImmPort, toward repurposing of open access

576        immunological assay data for translational and clinical research. *Sci Data* **5**,

577        180015 (2018).

578    3.  Bucasas, K. L. *et al.* Early patterns of gene expression correlate with the humoral

579        immune response to influenza vaccination in humans. *The Journal of infectious*

580        *diseases* **203**, 921–929 (2011).

581    4.  Shah, N. *et al.* A crowdsourcing approach for reusing and meta-analyzing gene

582        expression data. *Nat. Biotechnol.* **34**, 803–806 (2016).

583    5.  Chen, E. Y. *et al.* Enrichr: interactive and collaborative HTML5 gene list

584        enrichment analysis tool. *BMC Bioinformatics* **14**, 128 (2013).

585    6.  Kuleshov, M. V. *et al.* Enrichr: a comprehensive gene set enrichment analysis

586        web server 2016 update. *Nucleic Acids Res.* **44**, W90-97 (2016).

587    7.  Pilarczyk, M. *et al. Connecting omics signatures of diseases, drugs, and*

588        *mechanisms of actions with iLINCS.* 826271

589        https://www.biorxiv.org/content/10.1101/826271v2 (2020) doi:10.1101/826271.

590    8.  Subramanian, A. *et al.* Gene set enrichment analysis: A knowledge-based

591        approach for interpreting genome-wide expression profiles. *PNAS* **102**, 15545–

592        15550 (2005).

593    9.  Liberzon, A. *et al.* The Molecular Signatures Database (MSigDB) hallmark gene

594        set collection. *Cell Syst* **1**, 417–425 (2015).

595    10. Berke, K. *et al.* VaximmutorDB: A Web-Based Vaccine Immune Factor Database

596        and Its Application for Understanding Vaccine-Induced Immune Mechanisms.

597        *Frontiers in Immunology* **12**, 645 (2021).

598   11. Aksoy, B. A. *et al.* CTD2 Dashboard: a searchable web interface to connect

599        validated results from the Cancer Target Discovery and Development Network.

600        *Database (Oxford)* **2017**, (2017).

601   12. Yates, B. *et al.* Genenames.org: the HGNC and VGNC resources in 2017.

602        *Nucleic Acids Res.* **45**, D619–D625 (2017).

603   13. Bard, J., Rhee, S. Y. & Ashburner, M. An ontology for cell types. *Genome Biology*

604        **6**, R21 (2005).

605   14. Diehl, A. D. *et al.* The Cell Ontology 2016: enhanced content, modularization, and

606        ontology interoperability. *J Biomed Semant* **7**, 44 (2016).

607   15. Natale, D. A. *et al.* Protein Ontology (PRO): enhancing and scaling up the

608        representation of protein entities. *Nucleic Acids Res.* **45**, D339–D346 (2017).

609   16. Overton, J. A. *et al.* Reporting and connecting cell type names and gating

610        definitions through ontologies. *BMC Bioinformatics* **20**, 182 (2019).

611   17. Vita, R. *et al.* A structured model for immune exposures. *Database (Oxford)* **2020**,

612        (2020).

613   18. He, Y. *et al.* VO: Vaccine Ontology. *Nature Precedings* 1–1 (2009)

614        doi:10.1038/npre.2009.3552.1.

615   19. Sayers, E. W. *et al.* GenBank. *Nucleic Acids Res* **47**, D94–D99 (2019).

616   20. Schoch, C. L. *et al.* NCBI Taxonomy: a comprehensive update on curation,

617        resources and tools. *Database (Oxford)* **2020**, (2020).

618   21. Vita, R. *et al.* The Immune Epitope Database (IEDB): 2018 update. *Nucleic Acids*

619        *Res.* **47**, D339–D343 (2019).

620   22. Mi, H., Poudel, S., Muruganujan, A., Casagrande, J. T. & Thomas, P. D.

621        PANTHER version 10: expanded protein families and functions, and analysis

622        tools. *Nucleic Acids Res* **44**, D336–D342 (2016).

623   23. The UniProt Consortium. UniProt: a worldwide hub of protein knowledge. *Nucleic*

624        *Acids Research* **47**, D506–D515 (2019).

625   24. Ritchie, M. E. *et al.* limma powers differential expression analyses for RNA-
626       sequencing and microarray studies. *Nucleic Acids Research* **43**, e47 (2015).
627   25. Oh, S. *et al.* HGNChelper: identification and correction of invalid gene symbols for
628       human and mouse. *F1000Res* **9**, 1493 (2020).
629   26. Carlson, M. org.Hs.eg.db: Genome wide annotation for Human. R package
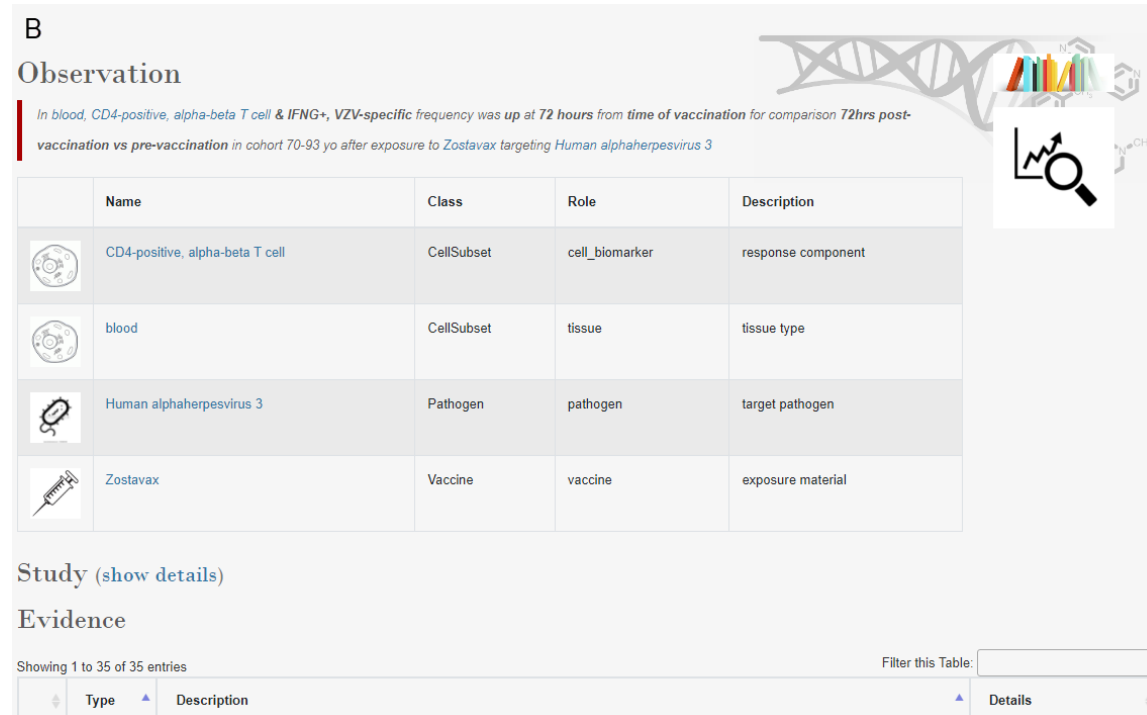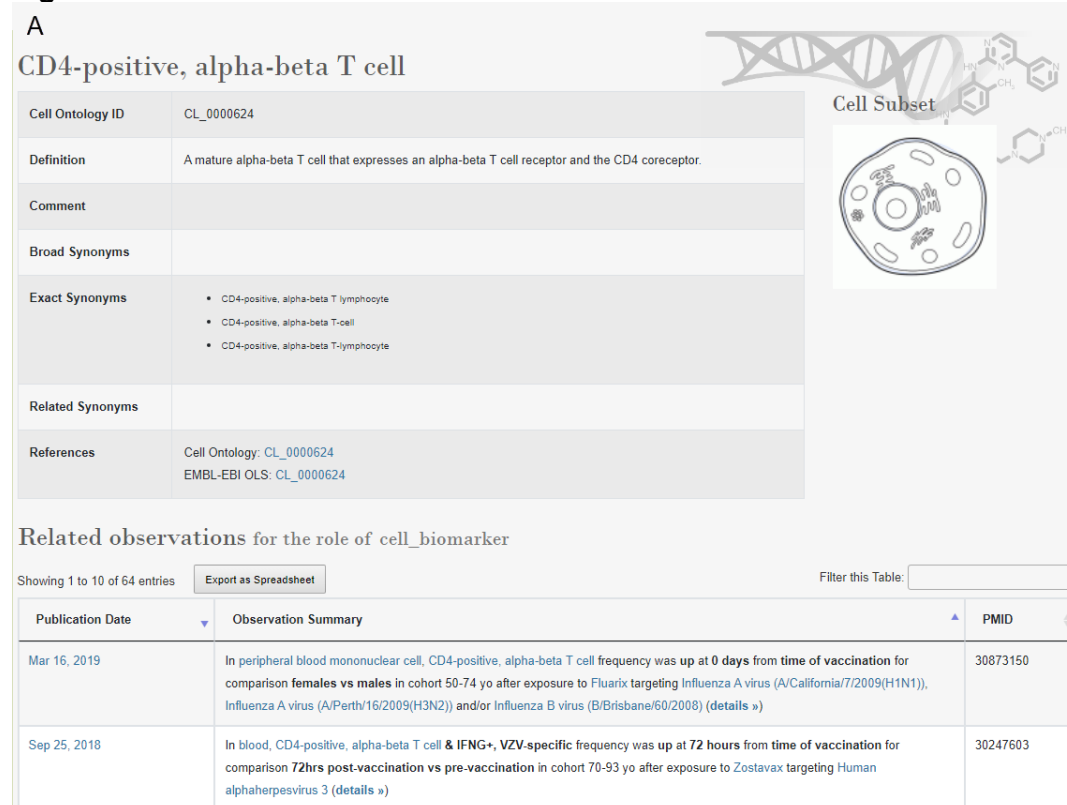630       version 3.10.0.
631

**Figures**
**Figure 1**

**Figure 2**

A

## CD4-positive, alpha-beta T cell

| | |
|---|---|
| Cell Ontology ID | CL_0000624 |
| Definition | A mature alpha-beta T cell that expresses an alpha-beta T cell receptor and the CD4 coreceptor. |
| Comment | |
| Broad Synonyms | |
| Exact Synonyms | • CD4-positive, alpha-beta T lymphocyte<br>• CD4-positive, alpha-beta T-cell<br>• CD4-positive, alpha-beta T-lymphocyte |
| Related Synonyms | |
| References | Cell Ontology: CL_0000624<br>EMBL-EBI OLS: CL_0000624 |

Cell Subset

**Related observations** for the role of cell_biomarker

Showing 1 to 10 of 64 entries    Export as Spreadsheet                    Filter this Table: [            ]

| Publication Date ▼ | Observation Summary ▲ | PMID ⬍ |
|---|---|---|
| Mar 16, 2019 | In peripheral blood mononuclear cell, CD4-positive, alpha-beta T cell frequency was **up** at **0 days** from **time of vaccination** for comparison **females vs males** in cohort 50-74 yo after exposure to Fluarix targeting Influenza A virus (A/California/7/2009(H1N1)), Influenza A virus (A/Perth/16/2009(H3N2)) and/or Influenza B virus (B/Brisbane/60/2008) (**details »**) | 30873150 |
| Sep 25, 2018 | In blood, CD4-positive, alpha-beta T cell **& IFNG+, VZV-specific** frequency was **up** at **72 hours** from **time of vaccination** for comparison **72hrs post-vaccination vs pre-vaccination** in cohort 70-93 yo after exposure to Zostavax targeting Human alphaherpesvirus 3 (**details »**) | 30247603 |

B

## Observation

*In blood, CD4-positive, alpha-beta T cell* **& IFNG+, VZV-specific** *frequency was* **up** *at* **72 hours** *from* **time of vaccination** *for comparison* **72hrs post-vaccination vs pre-vaccination** *in cohort 70-93 yo after exposure to Zostavax targeting Human alphaherpesvirus 3*

| | Name | Class | Role | Description |
|---|---|---|---|---|
| | CD4-positive, alpha-beta T cell | CellSubset | cell_biomarker | response component |
| | blood | CellSubset | tissue | tissue type |
| | Human alphaherpesvirus 3 | Pathogen | pathogen | target pathogen |
| | Zostavax | Vaccine | vaccine | exposure material |

## Study (show details)

## Evidence

Showing 1 to 35 of 35 entries                          Filter this Table: [            ]

| ⬍ | Type ▲ | Description ▲ | Details ⬍ |
|---|---|---|---|

**Figure 3**

**Table 1**

| Elements | Content Definition | Field value Example | Ontology | Data Type | Content Format |
|---|---|---|---|---|---|
| Response Component | The biological entity being observed: either a specific gene or cell subset | CKS1B | NCBI/HGNC for genes, Cell Ontology for cell types, PRO for proteins) | String | Controlled |
| Response Component Type | Type of response agent, e.g. gene, cell subset | gene | | String | Fixed List |
| Tissue Type | The type of cells analyzed, e.g. whole blood, gated cells | whole blood (UBERON:0000178) | Blood + PRO/CL | String | Controlled |
| Exposure Process Type * | Category of immune exposure | Vaccination | | String | Fixed List |
| Exposure Material * | Eg, the type of vaccine administered | VO:0001176 | VO | String | Controlled |
| Disease Name * | The condition being observed | None | DO | String | Controlled |
| Disease Stage * | Stage of infection (e.g. acute, chronic, post, etc.) | None | OGMS | String | Controlled |
| Target Pathogen | The pathogenic organism (e.g. virus, bacterium, prion, fungus) being studied | Influenza A, Influenza B | NCBI Taxonomy ID | String | Controlled |
| Vaccine Year | Year for seasonal vaccines | 2012 | | Numeric | Free text |
| Adjuvant | A substance added to vaccines to increase the body's immune response to the vaccine | | VO | String | Controlled |
| Route | Eg oral, nasal spray, intradermal injection, etc. | Intramuscular, intradermal, transcutaneous | VO | String | Controlled |
| Scheduling | The number of times a substance is administered within a specific time period. | e.g. prime / boost scheme | | String | Free text |
| Time Point | | 1 | | Numeric | |
| Time Point Units | | day | | String | Fixed List |
| Baseline Time Event | The starting point against which other events are compared | time of vaccination (first dose) | | String | Free Text |
| Cohort | Features that describe study cohorts | e.g. antibody responders | | String | Free Text |
| Age Min | | 18 | | Numeric | |
| Age Max | | 45 | | Numeric | |
| Age Units | {days, weeks, months, years} | years | | String | Fixed List |
| Publication Reference | PubMed unique identifier of an article. | 30843873 | | Numeric | Controlled |
| Publication Year | The year in which the study was published | 2019 | | Numeric | |
| Publication Reference URL | Link to article in PubMed | https://www.ncbi.nlm.nih.gov/pubmed/30843873 | | String | URL |
| Comparison | The contrast used for deriving the signature | 1d-0d | | String | Free Text |
| Response Behavior | Observed change in the response agent under the comparison | Up | | String | Free Text |

**Table 2**

| Signature Type | Vaccines | Target Pathogens | Publications | Signatures | Total Response Components | Response Components per Signature |
|---|---|---|---|---|---|---|
| Gene | 52 | 38 | 62 | 480 | 13,812 genes | 1 to 2,036 genes |
| Cell type | 28 | 26 | 31 | 185 | 47 cell types | 1 to 9 cell types |
| Joint | 56 | 44 | 69 | 665 | | |

**Table 3**

| Column name | Ontology | Values |
|---|---|---|
| response_component | HGNC | STAT1, IRF9, SPI1, CD74, HLA-E, TNFSF13B |
| tissue_type | Cell Ontology | peripheral blood mononuclear cell |
| exposure_material | Vaccine Ontology | VO:0000045; VO:0000046 (Fluarix; Fluvirin) |
| target_pathogen | NCBI Taxonomy | Influenza A virus (A/Brisbane/59/2007(H1N1)); Influenza A virus (A/Brisbane/10/2007(H3N2)); Influenza B virus (B/Florida/4/2006) |
| vaccine_year | | 2008 |
| time_point | | 1 |
| time_point_units | | days |
| baseline_time_event | | time of vaccination (first dose) |
| cohort | | 18-40 yo, subgroup high responders |
| publication_reference (PMID) | | 21357945 |
| comparison | | correlated with titer response index (TRI) |
| response_behavior | | positive |