# Title

Reproducible single cell annotation of programs underlying T-cell subsets, activation states, and functions

# Authors

Dylan Kotliar[1,2,3,4,5,*], Michelle Curtis[1,2,3,4,*], Ryan Agnew[1,2,3,4], Kathryn Weinand[1,2,3,4,6], Aparna Nathan[1,2,3,4,6], Yuriy Baglaenko[1,7,8], Yu Zhao[1,2,3,4], Pardis C. Sabeti[4,9,10], Deepak A. Rao[2], Soumya Raychaudhuri[1,2,3,4,6,†]

# Affiliations

[1]Center for Data Sciences, Brigham and Women's Hospital and Harvard Medical School, Boston, MA 02115, USA.
[2]Division of Rheumatology, Inflammation, and Immunity, Department of Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, MA 02115, USA.
[3]Division of Genetics, Department of Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, MA 02115, USA.
[4]Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA.
[5]Harvard-MIT Division of Health Sciences and Technology, Harvard Medical School, Boston, MA 02115, USA.
[6]Department of Biomedical Informatics, Harvard Medical School, Boston, MA 02115, USA.
[7]Center for Autoimmune Genetics and Etiology and Division of Human Genetics, Cincinnati Children's Hospital Medical Center, Cincinnati, OH 45229, USA.
[8]Department of Pediatrics, University of Cincinnati, College of Medicine, Cincinnati, OH 45219, USA.
[9]Department of Organismic and Evolutionary Biology, FAS Center for Systems Biology, Harvard University, Cambridge, MA 02138, USA.
[10]Howard Hughes Medical Institute, Chevy Chase, MD 20815, USA.

[*]These authors contributed equally
[†]Address correspondence to:
Soumya Raychaudhuri
77 Avenue Louis Pasteur, Harvard New Research Building, Suite 250D
Boston, MA 02446, USA.
soumya@broadinstitute.org
617-525-4484 (tel); 617-525-4488 (fax)

# Abstract

37    T-cells recognize antigens and induce specialized gene expression programs (GEPs) enabling

38    functions including proliferation, cytotoxicity, and cytokine production. Traditionally, different

39    classes of helper T-cells express mutually exclusive responses – for example, Th1, Th2, and

40    Th17 programs. However, new single-cell RNA sequencing (scRNA-Seq) experiments have

41    revealed a continuum of T-cell states without discrete clusters corresponding to these subsets,

42    implying the need for new analytical frameworks. Here, we advance the characterization of T-

43    cells with T-CellAnnoTator (TCAT), a pipeline that simultaneously quantifies pre-defined GEPs

44    capturing activation states and cellular subsets. From 1,700,000 T-cells from 700 individuals

45    across 38 tissues and five diverse disease contexts, we discover 46 reproducible GEPs

46    reflecting the known core functions of T-cells including proliferation, cytotoxicity, exhaustion, and

47    T helper effector states. We experimentally characterize several novel activation programs and

48    apply TCAT to describe T-cell activation and exhaustion in Covid-19 and cancer, providing

49    insight into T-cell function in these diseases.

50

51

52

53

54

55

56

# Introduction

Canonically, T-cells are classified by membership in a hierarchy of discrete, mutually exclusive subsets associated with key transcription factors and surface markers. For example, expression of γδ or αβ T-cell receptors and CD4 or CD8 co-receptors divide T-cells into subsets recognizing different major histocompatibility complex (MHC) molecules. CD45 isoform and L-selectin expression subdivides naive and memory subsets. CD4 memory cells are further subcategorized into helper subsets, including Th1, Th2, and Th17, with distinct cytokine profiles upon activation[1].

Emerging evidence conflicts with this canonical model. T-cell states vary continuously[2], combine additively within a cell[3], and have plasticity in response to stimuli[4]. This may explain why single-cell RNA sequencing (scRNA-Seq) typically shows a continuum of T-cell states without well-delineated clusters corresponding to discrete subsets[5,6]. Even with incorporation of pre-defined surface protein markers based on cellular indexing of transcriptomes and epitopes by sequencing (CITE-seq)[7], unbiased clustering does not yield canonical discrete T-helper subsets[8]. Rather, scRNA-Seq has highlighted untraditional cell populations including cytotoxic CD4+ cells[9], CD8+ regulatory T-cells[10] and Th1/Th17 cells[11], consistent with the growing recognition of non discrete T-cell states.

While hard clustering is the predominant scRNA-Seq analysis technique, it has key limitations when cell states are not discrete or mutually exclusive. A cell's transcriptome reflects its complex identity through expression of multiple gene expression programs (GEPs) that reflect lineage, activation states, and lifecycle processes[12]. However, hard clustering forces cells into discrete groups that cannot easily reflect the multiplicity of GEPs they express. For example,

81    proliferating cells from multiple subsets may cluster together, obscuring information about their

82    subset. Hard clustering also cannot directly model continuous expression trajectories and

83    instead arbitrarily discretizes cells into distinct clusters.

84

85    Component-based models like non-negative matrix factorization (NMF), hierarchical Poisson

86    factorization, and SPECTRA can overcome some of these limitations of hard clustering[5,13–16].

87    These methods model GEPs as vectors of expression values for each gene, and cells as

88    weighted mixtures of GEPs. Unlike Principal Component Analysis (PCA), NMF components

89    have been shown to correspond to biologically distinct GEPs[14]. Thus, NMF can capture

90    instances where multiple GEPs reflecting cell-type and other functional states additively

91    contribute to a cell's transcriptome. Furthermore, unlike cluster assignments, GEP vectors may

92    be able to serve as a fixed coordinate system onto which new datasets can be projected,

93    enabling reproducible comparison of GEP activity across biological contexts. Previous analyses

94    of T-cells using component-based models have already recognized GEPs associated with T-cell

95    activation[5] and exhaustion[15].

96

97    We argue that scaling these approaches may further elucidate T-cell biology. First, most

98    previous analyses have only analyzed T-cells from a small number of individual donors in a

99    limited set of biological contexts. As a result, they have identified a modest number of GEPs.

100    Moreover, it is essential to demonstrate the possibility of transferring GEPs identified in one

101    dataset to new datasets. For example, it remains unclear whether reference GEPs learned in

102    one dataset can accurately infer cell subsets, T-cell receptor (TCR)-dependent activation, and

103    proliferation status for cells in a new dataset.

104

105    Here, we present CellAnnoTator (*CAT, pronounced starCAT), an approach to score cells

106    based on a fixed, multidataset catalog of GEPs from any tissues or cell-type (indicated by the

107 wildcard character "*"). We develop a catalog of GEPs reflecting the breadth of subsets,

108 activation states, and functions within T-cells by applying consensus NMF (cNMF)[14], a validated

109 implementation of NMF, to 7 scRNA-Seq datasets, spanning 1.7 million T-cells across 38

110 human tissues[6,8,11,17–20]. We observe striking concordance of many GEPs across contexts. After

111 combining analogous GEPs, we define a final catalog of 46 consensus GEPs (cGEPs) capturing

112 diverse features of T-cells (**Figure 1A**). We demonstrate *CAT by accurately inferring T-cell

113 subsets in query datasets and quantifying rates of TCR-dependent activation and exhaustion in

114 Covid-19 and cancer.

# Results

## 1. Annotating cells with pre-defined gene expression programs

117 We first augmented the published cNMF algorithm to enhance GEP discovery, which is the first

118 step of *CAT (**Figure 1A - top**). cNMF mitigates the randomness of individual NMF runs by

119 repeating NMF with multiple seeds and combining the results into robust estimates[14]. It outputs

120 GEP spectra, with gene weights for each GEP, and usages, reflecting the GEP's weighted

121 contribution to each cell. For our approach, it was essential to amalgamate the inferred GEP

122 spectra from multiple datasets. However, we found that dataset-specific batch effects could

123 hinder the identification of reproducible GEPs. Most batch correction methods are not

124 compatible with cNMF since they create many negative values or correct low-dimensional

125 embeddings rather than gene-level data. We therefore used Harmony[21], with modifications to

126 produce non-negative values for gene-level data rather than principal components. We also

127 adapted cNMF to incorporate surface proteins into the final spectra to aid in GEP interpretation

128 without impacting GEP discovery (**Methods**).

129

130    Next, we developed *CAT to enable GEPs learned in a reference dataset to be transferred to

131    previously unseen "query" datasets. Whereas cNMF simultaneously learns GEPs and scores

132    their usage in each cell's transcriptional profile, *CAT addresses the independent problem of

133    quantifying the usages of a fixed set of GEPs in a new dataset, using non-negative least

134    squares (NNLS) regression, similar to NMFproject[13]. The result is a vector of usages for each

135    cell representing the relative contribution of each GEP to the cell's profile (**Figure 1A - bottom**).

136

137    Using NNLS to refit GEPs as we do with *CAT provides significant advantages over direct

138    applications of cNMF or other matrix factorizations. First, *CAT uses a fixed set of GEPs from a

139    reference, instead of discovering GEPs *de novo* in the query. Thus, it provides a consistent

140    representation of cell states that can be compared across different datasets and biological

141    contexts. Second, *de novo* cNMF might miss GEPs that are active in small numbers of cells,

142    whereas *CAT can characterize activity in a query dataset with relatively few cells. Finally, *CAT

143    is significantly faster to run than cNMF.

144

145    We conducted simulations to benchmark *CAT in scenarios where the reference and query

146    datasets have only partially overlapping GEPs (**Methods**). We simulated two reference datasets

147    of 100,000 cells and a query dataset of 20,000 cells. Each cell could express up to eleven

148    GEPs, including one of ten mutually exclusive subset GEPs and up to ten non-subset GEPs.

149    One reference dataset included all 16 GEPs in the query data as well as four additional GEPs.

150    The other reference dataset was missing four GEPs present in the query (**Figure 1B**). We then

151    learned GEPs from each reference dataset with cNMF and fit them to the query using *CAT.

152    The reference and query datasets shared only 90% of genes in common, as datasets rarely

153    share all genes.

154

155 *CAT accurately inferred the usage of GEPs that overlapped between the reference and query

156 datasets (Pearson R>0.7) (**Figure 1C-D)**. *CAT had low predicted usage of the extra GEPs in

157 the reference panel that were not in the query dataset (**Figure S1A**). Surprisingly, *CAT

158 obtained better concordance with the simulated ground truth GEP usages than direct application

159 of cNMF to the query (**Figure 1E)**. This is striking because the reference GEPs had extra or

160 missing GEPs relative to the query, and were learned on different datasets, so could incorporate

161 dataset-specific noise. We hypothesized that *CAT's increased performance reflected the larger

162 reference datasets enabling more accurate GEP inference. We confirmed this by simulating

163 multiple query datasets with between 100 and 100,000 cells. While cNMFs performance

164 declined for small query datasets, *CATs remained constant, demonstrating that *CAT can out-

165 perform cNMF when the reference is larger than the query (**Figure 1F**).

166 ## 2. Gene expression programs for T-cell annotation

167 We next developed a catalog of GEPs to capture T-cell states; combining these GEPs with the

168 *CAT algorithm yields T-CellAnnoTator (TCAT). We analyzed T-cells from 7 diverse datasets

169 including blood and tissues from healthy individuals or individuals with Covid-19, cancer,

170 rheumatoid arthritis, or osteoarthritis (**Figure 1G).** After stringent quality control, there were 1.7

171 million cells from 905 samples from 695 individuals in our analysis. To preserve dataset-specific

172 GEPs, we applied cNMF to each batch-corrected dataset independently (**Supplementary item**

173 **1, Methods**).

174

175 We observed that GEPs were reproducible across the datasets. To quantify this, we clustered

176 highly correlated GEPs found in different datasets (**Methods**). Assuming that correlated

177 dataset-specific GEPs represented the same biological state, we defined a consensus gene

178 expression program (cGEP) as the average of a GEP cluster. Nine cGEPs derived from a

179 cluster of GEPs from all seven datasets (Average Pearson R=0.81) and 49 cGEPs derived from

180    a cluster of GEPs from two or more datasets (Average Pearson R=0.74) (**Figure 2A-B, S1B**).

181    Between 68.4% and 96.8% of GEPs identified in each of the seven reference datasets clustered

182    with at least one GEP from another reference, suggesting high reproducibility. By contrast, gene

183    expression principal components showed limited concordance between pairs of datasets,

184    suggesting they reflect more dataset-specific signals[14] (**Figure S1C**).

185

186    We curated a catalog of 46 cGEPs capturing diverse T-cell states, including 11 discovered only

187    in blood datasets, seven discovered only in tissue datasets, and 28 discovered in both (**Table**

188    **S1**, **Figure 2C**). This represents between 27 and 36 more programs than previous factorization

189    analyses of T-cells[13,15,16]. Of these cGEPs, 43 derived from multiple datasets, while three were

190    singletons found in a single dataset. We excluded 49 of the 52 initially identified singletons since

191    they likely reflect dataset-specific artifacts. The three retained singletons capture disease- or

192    tissue-specific GEPs with a biological justification. For example, the rheumatoid arthritis dataset

193    (referred to as AMP-RA), included a GEP highly enriched for T peripheral helper cells markers

194    (including PD-1 and CD4 protein, *LAG3*, and *CXCL13* RNA), which is characteristic of inflamed

195    rheumatoid arthritis synovium[22] (**Table S2**). Similarly, the pan-cancer dataset included a cancer-

196    specific exhaustion GEP (*HAVCR2*, *ENTPD1*, *LAG3*) which may be especially enriched in

197    cancer, and a GEP bearing markers for T follicular helper cells (PD-1 protein and *CXCR5*, *IL6R*,

198    and *CXCL13* RNA) which was distinct from a second Tfh-like GEP discovered in multiple non-

199    cancer tissue datasets. In addition to the main T-cell cGEPs, we identified six cGEPs

200    corresponding to non T-cell populations including erythrocytes (*HBA2*, *HBA1*, *HBB*) and

201    plasmablasts (*JCHAIN*, *IGKC*, *IGKV3-20*), potentially derived from doublets. We retained these

202    cGEPs to flag doublet-associated transcriptional signals.

203

204    To label cGEPs, we first examined their top weighted genes (**Figure 2D, Supplementary item**

205    **2, Table S1-2**). For example, the top 10 weighted genes in the Treg and Th2-Resting cGEPs

206      included the master regulators, *FOXP3* and *GATA3*, respectively. Similarly, top weighted genes

207      helped identify the Th2-Activated (*GATA3, IL4, IL5*) and Th17-Activated (*IL26, IL17A*, and

208      *RORC*) cGEPs. Many functional cGEPs could also be readily identified, such as Heatshock

209      (*HSPA1A, HSP90AA1, HSPA1B*), HLA (*HLA-DRA, HLA-DRB1, CD74*), Metallothionein (*MT1X*,

210      *MT2A, MT1E*), and Actin Cytoskeleton (*ACTB, ACTG1, PFN1*) (**Figure 2D**).

211

212      We also labeled cGEPs based on their ability to discriminate canonical T-cell subsets defined by

213      manual gating on surface markers. We gated PBMC-derived T-cells from the COMBAT CITE-

214      Seq reference dataset[18] and then used multivariate logistic regression to associate cGEPs with

215      subsets (**Figure S2A**, **Methods**). cGEPs labeled as regulatory T (Treg), gamma-delta T (gdT),

216      mucosal associated invariant T (MAIT), CD4 Naive, CD8 Naive, CD8 effector memory (CD8

217      EM), CD4 central memory (CD4 CM), and T Effector Memory-Expressing CD45RA (TEMRA)

218      were strongly associated with the expected manually gated populations (P-value<$1\times10^{-200}$,

219      Coefficient>0.35, **Figure S2B**). The CD4 effector memory gated population was most strongly

220      associated with cGEPs reflecting expected T-helper subsets labeled as Th17-Resting (*CCR6,*

221      *RORC, AQP3*) and Th1-like (*IFNG-AS1, CXCR3*, and CD195 protein) (P<$1\times10^{-200}$ and

222      P=$4.1\times10^{-190}$, coefficients 0.36 and 0.22, respectively, **Figure S2B**). Overall, this approach

223      enabled identification of 17 subset-associated cGEPs (**Figure 2C, Table S1**).

224

225      As a third strategy to label cGEPs, we used gene-set enrichment analysis with gene-sets from

226      the gene ontology database[23] and from T-cell polarization experiments[24] (**Methods, Table S3**).

227      We found that the Th2-Resting and Th2-Activated cGEPs were the most significantly enriched

228      for genes upregulated following 16 hour stimulations of naive T-cells with Th2 polarizing

229      cytokines (Fisher Exact Test OR=22.7, 16.2,  P=$4.9\times10^{-5}$, $1.7\times10^{-4}$, respectively). Gene set

230      analysis also helped annotate 5 cGEPs corresponding to non-T-cell specific cellular functions

231      including early and late cell cycle S-phase (P=$3\times10^{-56}$ for DNA_REPLICATION and P=$2\times10^{-55}$

232    for MITOTIC_CELL_CYCLE), G2M-phase (P=9x10[-74] CELL DIVISION), interferon stimulated

233    genes (P=1x10[-59] for RESPONSE TO VIRUS), and translation (P=4x10[-163] for

234    GOCC_CYTOSOLIC_RIBOSOME).

235

236    Next, we identified technical artifact-associated cGEPs that correlate with low-quality cell

237    features (**Table S4**). A cGEP we label Mitochondria contains top markers that are exclusively

238    mitochondrially transcribed genes, which are frequently used to identify low-quality cells[25,26]; as

239    expected, this cGEP had a high correlation with the percentage of mitochondrial reads per cell

240    (average R=0.81 across datasets). We labeled another cGEP Poor-Quality based on its top

241    marker gene *MALAT1*, a long non-coding RNA linked to poor cell viability[27]; this cGEP also

242    correlated with the percentage of mitochondrial transcripts per cell (R=0.25 averaged across

243    datasets, **Figure S2C**) and was inversely correlated with the percentage of protein-coding

244    transcripts per cell (**Figure S2D,** average R=-0.50 across datasets). For the AMP-RA dataset,

245    we had access to raw sequence alignment files so we could quantify the percentage of reads

246    aligned to intergenic regions of the genome; the Poor-Quality cGEP was by far the most

247    correlated with the percentage of intergenic reads per cell (R=0.74, **Figure S2E**). Its usage may

248    be driven by higher levels of contaminating DNA or nascent RNA.

249

250    Finally, we label three correlated cGEPs as immediate early gene programs (IEG1, IEG2, IEG3

251    , pairwise R of 0.45-0.70). The top genes include canonical IEGs including *FOS*, *JUN*, and

252    *ZFP36,* and these cGEPs were all enriched for a published IEG gene set[28] (Fisher Exact Test

253    P<1x10[-53]). We suspect that IEG1 represents the core pathway as it was found in 6 out of 7

254    datasets (**Figure 2C**) whereas IEG2 and IEG3 represent mixtures with delayed immediate and

255    secondary response genes. We hypothesize that these cGEPs reflect sample processing

256    artifacts in scRNA-Seq, since IEGs are induced in as few as 30 minutes[29] in response to

257    mitogens or cell stress[30], and following processing steps like tissue dissociation[31,32]. As

258    evidence of the potential technical nature of these cGEPs, we calculated their mean usage per

259    sample in T-cells, B-cells, NK-cells and monocytes/DCs in the 3 PBMC references. We found

260    that their average usage in T-cells correlates with their usage in other cell-types (R=0.46-0.99,

261    average 0.77, **Figure S2F, Supplementary item 3**), suggesting that they are a sample-intrinsic

262    property, which would be expected of a sample-processing effect. However, in certain contexts,

263    these cGEPs may be biologically important.

264    ## 3. Benchmarking TCAT on an independent query dataset

265    Next, we benchmarked TCAT on predicting T-cell subsets in an independent CITE-seq dataset.

266    We analyzed 336,739 T-cells from PBMCs of 24 Covid-19-recovered and 17 healthy individuals

267    after flu vaccination[33] (**Figure 3A**). As ground truth, we assigned cells to one of ten subsets

268    through manual gating of surface proteins (**Figure S3A)**. We then predicted each subset by

269    thresholding the corresponding subset-associated cGEP (**Methods**). For all 10 subsets,

270    thresholding the single most-associated cGEP was comparable to RNA-based hard clustering,

271    across nine different clustering resolutions. Averaged across subsets, the accuracy difference

272    between TCAT and clustering ranged from 0.064 to -0.007 depending on the clustering

273    resolution (**Figure 3B-D**).

274

275    Since subsets can contain heterogeneity not captured in univariate analysis (e.g. multiple

276    polarized populations within CD4 effector memory), we performed multivariate analysis using all

277    cGEPs for simultaneous multi-label prediction (**Methods**). We trained the classifier on the

278    COMBAT dataset and evaluated its performance on the Flu-Vaccine dataset. The classifier was

279    more accurate than RNA clustering across all nine clustering resolutions tested, with average

280    accuracy differences ranging from 0.10 to 0.033 (**Figure 3B-C, Figure S3B-C**). Thus, for

281    PBMC-derived T-cells, TCAT can be combined with a multilabel classifier to predict subsets

282    without requiring manual annotation.

283

284 We also compared TCAT's subset classification accuracy against NMFproject[13] and gene-sets

285 derived from a recent NMF analysis of tumor-infiltrating T-cells[16] (**Methods**). TCAT single cGEP

286 and multi-label classification yielded higher area under the curve (AUC) for all lineage

287 predictions than these other approaches (**Figure S3B-C**).

288

289 Next, we validated TCAT's prediction of functional cGEPs relative to common continuous

290 metrics. Usage of the mitochondrial cGEP was highly correlated with percentage of

291 mitochondrial reads (R = 0.88, **Figure 3D**). In addition, predicted cell cycle cGEP usages

292 corresponding to the S and G2M phase were highly correlated with cell cycle scores calculated

293 from corresponding published gene sets[34,35] (R=0.75-0.81, **Figure 3D**).

294

295 Finally, we validated prediction of T-cell polarization against expression of canonical markers.

296 We discretized cells based on their expression of the Th1-Like, Th2-Resting, and Th17-Resting

297 cGEPs (usage>0.1) and computed per-sample pseudobulk profiles of high and low usage cells.

298 Th2-Resting-high samples expressed significantly more *GATA3*, *CCR4*, and *PTGDR2* than Th2-

299 Resting-low samples (P<1x10$^{-35}$ all, paired T-test) (**Figure 3E**). Th17-Resting-high samples also

300 had increased expression of Th17 markers including *CCR6*, *RORC*, and *AQP3* (P<1x10$^{-55}$ all).

301 The Th1-Like-high samples had increased expression of the Th1 markers *CXCR3*, *IFNG-AS1*,

302 and CD195 protein (P<1x10$^{-35}$ all). However, the Th1 markers *IFNG* and *TBX21* were also

303 expressed in Th1-Like-low samples (**Figure S3D**). We suspected this was due to the known

304 expression of these genes in cytotoxic T-cells[36,37]. When we excluded cells high in the cytotoxic

305 cGEP (usage>0.1) prior to pseudobulking, *IFNG* and *TBX21* were significantly higher in Th1-

306 Like-high samples (P=8.2x10$^{-13}$, P=9.6x10$^{-47}$, **Figure 3E**, **S3D**). Thus, TCAT can predict T-cell

307 polarization in query datasets.

308  ## 4. cGEPs capture multi-program identities of T-cells in scRNA-Seq

309  Next, we illustrate how TCAT can reveal cellular heterogeneity not visible with clustering. Using

310  the COMBAT dataset as an example, we analyzed cell cycle, a common signature that

311  frequently obscures other aspects of proliferating cells[38]. In the initial publication, two clusters

312  were annotated as proliferating CD4s and CD8s with subclusters that didn't clearly correspond

313  to subsets (e.g. CD4.TEFF.prolif.1, CD4.TEFF.prolif.GZMB.1). One sub-cluster labeled

314  CD4.TEFF.prolif.MKI67lo was enriched for the myeloid doublet cGEP (**Figure 4A-B**) and

315  expressed myeloid marker genes (e.g. *CD14*, *MNDA*, **Supplementary item 4**), illustrating how

316  cell cycle can drive cells with distinct cell lineages to cluster together. By contrast, TCAT readily

317  identified distinct proliferating subsets based on co-expression of cell cycle and subset cGEPs,

318  including CD8 EMs, TEMRAs, and Treg (**Figure 4C-D**).

319

320  Disentangling cell cycle and subset enabled us to quantify the percentage of proliferating cells

321  per subset and disease status. We assigned cells to subsets based on their most highly used

322  subset cGEP. This revealed increased expression of cell cycle cGEPs across many T-cell

323  subsets in Covid-19 compared to healthy cells, in both Covid-19 datasets (**Figure S4A**). The

324  most proliferative subsets in both Covid-19 and control samples expressed the T peripheral

325  helper cGEP, reflecting an inflammatory population that was recently identified in Covid-19[39].

326

327  We identified other functional cGEPs that obscured T-cell subsets, akin to proliferation. Many

328  CD4 memory subclusters in the original study were most strongly enriched for functional cGEPs

329  such as ISG, Cytotoxicity, and Poor-Quality, rather than subset cGEPs (**Figure 4B-D,**

330  **Supplementary item 4**). The CD4.Th.mitohi and CD4.Tem.mitohi.1 clusters were driven by

331  high usage of the Poor-Quality cGEP and contained cells expressing multiple subset cGEPs.

332  The CD4.TEM.IFN.resp and CD4.Th.IFN.resp clusters were both predominantly driven by the

333    interferon stimulated gene (ISG) cGEP. The CD4.TEM.IFN.resp cluster had high usage of the

334    Cytotoxicity and TEMRA cGEPs while the CD4.Th.IFN.resp cluster contained cells expressing

335    many subset cGEPs including CD4-Naive (**Figure 4B, S4B**). Cells with high usage of the CD4-

336    Naive cGEP expressed CD4 naive markers including CD45RA protein and *SELL* RNA,

337    confirming that clustering had misclassified them as memory T-cells (**Supplementary item 4**).

338

339    Clustering also obscured the subset of CD4 T-cells expressing the Cytotoxicity cGEP. We

340    visualized the per-cell usage of all cGEPs in cells from the CD4 memory sub-clusters that had

341    high Cytotoxicity cGEP usage (average cluster usage>0.1, **Figure 4B**). Intriguingly, these

342    clusters contained heterogeneous cells with high usage of many subset cGEPs including CD8-

343    EM, Th1-Like, TEMRA, and gdT (**Figure S4C**). Pseudobulk analyses showed that cells co-

344    expressing these cGEPs (usage>0.1 for both) co-expressed the expected cytotoxicity and

345    subset marker genes (**Figure S4D**). Thus, TCAT can reveal subset heterogeneity within

346    cytotoxic T-cells.

347

348    TCAT could readily annotate polarization status based on usage of the Th1-Like, Th2-Resting,

349    and Th17-Resting cGEPs (**Figure 4C**). By contrast, the published clustering did not identify a

350    Th2 cluster, and clusters annotated as Th1 and Th17 were only identified with a high clustering

351    resolution resulting in 243 clusters, likely due to other conflating signals. As expected, there was

352    significant enrichment between cells annotated as Th1 by clustering and high Th1-Like cGEP

353    usage, as well as Th17 clustering and high Th17-Resting cGEP usage (P<$1 \times 10^{-100}$ for both,

354    fisher exact test).

355

356    However, TCAT additionally identified expression of polarization cGEPs outside of the CD4

357    memory compartment (**Figure 4E**). We annotated polarization across manually gated T-cell

358    subsets with a usage threshold>0.1. As a control, we confirmed that the Treg cGEP was highly

359    enriched in the Treg gate, with an average of 88.1% of gated Tregs expressing the cGEP,

360    compared to 5.3% for the next highest population. Similarly, the Th17-Resting cGEP was most

361    enriched in the expected CD4 EM (22.1%) and CD4 CM (10.7%) populations compared to only

362    3.5% for MAITs, the next highest. Surprisingly, the Th2-Resting cGEP was most commonly

363    assigned within the CD8 CM (19.8%), CD4 CM (12.8%), and CD4/CD8 Double Positive (12.8%)

364    populations. The Th1-Like cGEP was also used by CD8 T-cells; it was most prevalent within the

365    CD8 CM (15.7%), CD4 EM (14.7%), CD8 EM (14.4%), and MAIT populations (12.3%). The

366    calculated subset polarization proportions were highly correlated between the COMBAT and

367    Flu-Vaccine datasets, the two datasets with the best quality manual gating (R>0.9, P<5.5x10$^{-5}$

368    for all three, **Figure S4E**). Furthermore, cells assigned to each polarization had high usage of

369    the expected marker genes for that polarization, irrespective of whether they were CD4+ or

370    CD8+ (**Figure S4F**). These findings support the emerging recognition of polarized CD8 T-cell

371    populations[40] and illustrate how these populations are easily revealed by TCAT.

## 372    5. cGEPs associated with TCR-dependent activation

373
374    Next we identified cGEPs induced following antigen recognition by the TCR. To do so, we

375    developed AIM-Seq (Activation-Induced Marker (AIM) assay followed by scRNA-Seq), an assay

376    to profile T-cells after antigen stimulus (**Figure 5A-D**). We collected PBMCs from 5 genome-

377    wide genotyped healthy donors and stimulated them for 24 hours using a pool of 176 peptide

378    antigens from common pathogens (CEFX, JPT)[41] and anti-CD28/CD49d co-stimulation. Using

379    flow cytometry, we separated T-cells expressing activation-induced markers (OX40 and PD-L1

380    for CD4s[42], CD137 for CD8s[43], AIM-positive) from unactivated cells (negative for these markers,

381    AIM-negative). As a negative control, we activated cells non-specifically with anti-CD28/CD49d

382    costimulation without peptides (Mock). We labeled cells from these conditions with hashtag

383 antibodies and pooled them for single-cell RNA, CITE, and TCR repertoire sequencing

384 (**Methods**).

385

386 As expected, CEFX stimulated CD4 and CD3+CD4- (hereafter labeled CD8) T-cells contained

387 higher proportions of AIM-positive cells than mock (**Figure 5B, S5A**). 4.21% of CD4 T-cells and

388 2.45% of CD8s were AIM-positive, compared to 0.049% and 0.54% of mock-stimulated CD4

389 and CD8 T-cells, respectively.

390

391 The CITE-Seq data showed that AIM-positive cells expressed additional surface activation

392 markers including CD54, CD25, CD71, and CD69 beyond the sorting markers (T-test $P<1x10^{-200}$,

393 **Figure S5C-E**). Moreover, AIM-positive cells were significantly depleted of naive T-cells (P=

394 0.027 and $P=8.6x10^{-4}$, for CD4 and CD8, respectively) and enriched for Tregs, CD4 central and

395 effector memory populations (P =0.00064, 0.0044 and 0.054, respectively, **Figure S5F**). This is

396 unsurprising as the peptide pool is derived from common pathogens and prior memory is

397 expected. However, 11.8% of the AIM-positive cells were CD4 naive and 1.4% were CD8 naive,

398 indicating we could detect both memory and naive cell responses.

399

400 Next, we identified cGEPs associated with antigen-specific activation in this assay. We used

401 pseudobulk sample-level regression to identify cGEPs upregulated in AIM-positive cells relative

402 to AIM-negatives (**Methods**). This identified 24 significant positively associated cGEPs (false

403 discovery rate (FDR) corrected P < 0.05), including two that are milieu regulated (I.e. non TCR-

404 dependent), five representing enriched subsets, and 17 functional cGEPs (**Figure 5E, S5G**).

405

406 The two milieu mediated cGEPs, Interferon Stimulated Gene (ISG) and Metallothionein, were

407 significantly upregulated in both AIM-negative and AIM-positive cells relative to mock (ISG: AIM-

408 negative - $P=8.9x10^{-7}$, AIM-positive - $P=3.1x10^{-5}$; Metallothionein: AIM-negative - $P=1.5x10^{-3}$,

409  AIM-positive - P=$3.3\times10^{-9}$). Interferon is a secreted cytokine that can activate nearby cells

410  independent of TCR-activation to induce the ISG cGEP. Shifting extracellular cytokine or ion

411  concentrations may similarly induce TCR-independent upregulation of the metallothionein

412  cGEP[44].

413

414  Five subset-associated cGEPs were increased in AIM-positive cells relative to AIM-negatives

415  (Th17-Resting, Treg, Tph, Th22, and Tfh-2) and 3 were increased in AIM-negatives (CD8-

416  Naive, CD4-Naive, and Th1-like) (**Table S5**). These associations likely reflect differential

417  abundance of cell populations rather than upregulation of the cGEPs, consistent with the

418  manual gating results (**Figure S5E**).

419

420  The remaining 17 AIM-associated programs are functional cGEPs including many with well-

421  known links to TCR-stimulation. Six of these are not T-cell specific, namely the three cell cycle

422  cGEPs[45] (P<$3.6\times10^{-4}$), actin cytoskeleton[46] (P=$3.3\times10^{-8}$), heatshock[47,48] (P=$1.7\times10^{-7}$), and MHC

423  class II[49] (P=0.012).

424

425  Excluding these leaves 11 functional AIM-associated cGEPs that may be specific to T-cell

426  activation. These include CTLA4/CD38 (P=$9.7\times10^{-9}$), ICOS/CD38 (P=$1.5\times10^{-6}$), NME1/FABP5

427  (P=$2.0\times10^{-6}$), OX40/EBI3 (P=$2.6\times10^{-5}$), Multi-cytokine (P=$5.4\times10^{-5}$), Exhaustion (P=$9.3\times10^{-5}$),

428  TIMD4/TIM3 (P=$5.0\times10^{-4}$), Th2-Activated (P=$5.9\times10^{-4}$), Th17-Activated (P=$2.1\times10^{-3}$) and

429  BCL2/FAM13A (P=$4.3\times10^{-3}$). We highlight 4 of these cGEPs here. CTLA4/CD38 showed the

430  most upregulation in Tregs and CD4 memory cells (**Figure 5F**) and is characterized by CD278

431  and CD38 protein levels as well as the anti-inflammatory genes *CTLA4* and *IL10*. ICOS/CD38

432  has similar top markers including CD278, CD71, and CD38 but shows broad upregulation

433  across naive T-cells and CD4 memory cells. The OX40/EBI3 cGEP includes many of the

434  activation-induced markers used to define AIM positivity in the first place including *TNFRSF4*

435    which encodes OX40 and *IL2RA* which encodes CD25. TIMD4/TIM3 is most expressed in

436    MAIT, gdT, and CD8 memory T-cells and is characterized by expression of activation markers

437    (CD38 protein and RNA) and cytotoxicity genes (*GZMB*, *GZMA*, *GNLY*), and likely represents a

438    cytotoxic activation response.

439

440    We hypothesized that AIM-associated cGEPs would be enriched in proliferating cells *in vivo*

441    since proliferation is a core response to TCR activation. To test this, we performed pseudobulk

442    sample-level association tests to identify cGEPs with higher usage in proliferating cells (sum of

443    cell cycle cGEPs>0.1) than non-proliferating cells (sum<0.1, **Methods**). The results were highly

444    concordant across datasets (**Table S6, Supplementary item 5**). 15 cGEPs were significantly

445    upregulated with proliferation in at least four out of six datasets. Meta-analysis across datasets

446    identified 12 functional cGEPs (including the three cell cycle cGEPs) and two subset cGEPs

447    (Th17-Activated and Tph) that were significantly associated with proliferation (**Figure S5H**).

448    Consistent with our hypothesis, 14 of 15 proliferation-associated cGEPs (including the 3 cell

449    cycle cGEPs) were upregulated with AIM positivity (Fisher exact test $P=2.1\times10^{-5}$). Thus, the

450    AIM-associated cGEPs are associated with proliferation *in vivo*, consistent with a role

451    downstream of TCR activation.

452    6. Annotating antigen-dependent activation *in vivo*

453    Next, we developed a per-cell antigen-specific activation (ASA) score to identify and

454    characterize TCR-activated T-cells in disease. We used forward stepwise selection to select

455    AIM-associated cGEPs that predicted co-expression of the activation markers CD71 and CD95

456    in the COMBAT and Flu-Vaccine datasets (**Methods**). These markers show sustained

457    upregulation within less than 24 hours of TCR activation[50–53], were upregulated in the AIM-

458    positive cells (**Figure S5D-F**), and had high quality across subsets in both datasets (**Figure**

459    **S6A**). Stepwise optimization defined ASA as the sum of four cGEPs – TIMD4/TIM3,

460    ICOS/CD38, CTLA4/CD38, and OX40/EBI3 (**Figure S6B**, **Methods**).

461

462    ASA accurately classified T-cells with CD71/CD95 co-expression suggestive of TCR-activation,

463    yielding AUCs of 0.920 and 0.818 in the COMBAT and Flu-Vaccine datasets (**Figure S6C-D**). It

464    also predicted AIM positivity with an AUC of 0.828 in the AIM-Seq assay (**Figure S6E**) and was

465    correlated with other surface markers of activation (e.g. R=0.43 (CD69) and 0.52 (CD25),

466    $P<1x10^{-100}$, **Supplementary item 6**). For cases where a discrete label is preferable to a

467    continuous score, we picked an ASA threshold of 0.0625 based on the trade-off between

468    sensitivity and specificity (**Figure S6C-E**). With this threshold, ASA annotated 76.7% of

469    CD71+CD95+ and 5.2% of non-CD71/CD95 double positive T-cells in the COMBAT dataset

470    (**Figure 6A**). In the AIM-Seq dataset, ASA annotated 60.6%, 7.0%, and 3.2% of stimulated AIM-

471    positive, stimulated AIM-negative, and mock stimulated cells, respectively (**Figure 6B**).

472

473    As proliferation is a core response to activation, we found high ASA in proliferating T-cell

474    clusters (**Figure 6E-F**) and significant overlap of ASA-high and proliferating cells (specifically,

475    cells with summed cell cycle usage > 0.1, Fisher Exact OR 2.8-58.8, $P<1x10^{-100}$, **Figure S6F -**

476    **left**). However, across reference datasets, substantially more cells were annotated as ASA-high

477    than proliferating ($P=8.8x10^{-189}$, paired T-test, **Figure 6H**). Consistent with this, correlation

478    between summed cell cycle cGEP usage and ASA was relatively low (mean=0.15) (**Figure S6F**

479    **- right**). Thus, while proliferation and antigen-specific activation overlap to some extent, ASA

480    offers greater sensitivity for classifying TCR-activation.

481

482    As clonal expansion often follows TCR activation, we tested whether high clonality was

483    associated with ASA in Covid-19 patients. ASA-high cells were more likely to be clonal, i.e. have

484    a TCR found in multiple cells from the same sample (Fisher Exact Test: COMBAT OR=2.50,

485    UK-Covid OR=2.28, P< $1x10^{-100}$ for both). Binarized ASA and cell cycle status were

486    independently associated with clonality in a multivariate logistic regression (ASA Beta = 0.45,

487    0.50; Cell cycle Beta = 0.66, 0.52 in COMBAT and UK-Covid respectively, P<$1x10^{-22}$,

488    **Methods**). Furthermore, the absolute number of cells sharing a TCR sequence in a sample was

489    significantly higher in ASA-high than ASA-low cells (Mann Whitney U test P<$1x10^{-100}$, both

490    datasets, **Figure 6C, S6G**).

491

492    Next, we evaluated how ASA varied between Covid-19 and healthy samples across T-cell

493    subsets. The percentage of activated (I.e. ASA positive) conventional T-cells varied widely

494    across samples, between 2.7%-41.2% (mean 10.3%) and 4.9%-44.7% (mean 22.1%), in the

495    COMBAT and UK-Covid datasets, respectively (**Figure 6D**). Activation rates were significantly

496    higher in conventional T-cells in Covid-19 samples than in healthy controls (COMBAT P=$1.9x10^{-7}$,

497    UK-Covid P=$1.5x10^{-6}$), even in CD4+ and CD8+ T-cells separately (**Figure S6H-J**). Activation

498    rates were similar between CD4s and CD8s (median activation of 8.3%, 21.8% for CD4s and

499    7.8%, 21.7% for CD8s in COMBAT and UK-Covid). By contrast, there was greater Treg

500    activation in both healthy and Covid-19 samples, with a median of 33.6 and 35.3% of cells

501    activated in COMBAT and UK-Covid (**Figure S6J**). This coincided with substantial overlap of

502    ASA with the Treg cluster (**Figure 6E-F**). Tregs were the most ASA-enriched subset in healthy

503    control samples in the COMBAT (OR=11.4, P<$1x10^{-100}$) and Flu-Vaccine datasets (OR=4.1,

504    P<$1x10^{-10}$) (**Figure 6G**). Outside of acute infection, we would expect Tregs to be actively

505    suppressing inappropriate activation. By contrast, in acute Covid-19 samples, we saw less

506    enrichment for Tregs (OR=4.8 down from 11.4) and more for CD8 central memory (OR=4.8),

507    CD8 effector memory (OR=2.8), and double negative populations (OR=3.1), reflecting the

508    antiviral response (all P<$1x10^{-10}$).

509

510    Next, we quantified levels of T-cell exhaustion and activation per sample and subset within the

511    pan-cancer dataset. CD4 conventional T-cell (CD4 Conv) activation rates varied widely across

512    and between tumor types (**Figure 6I**). The highest rates of activation were in esophageal cancer

513    (ESCA - median 48.0%) and the lowest were in bladder cancer (BC - median 5.4%, **Figure 6I -**

514    **left**). As expected, there was minimal exhaustion usage by CD4 Convs across cancer types[54]

515    but highly variable levels of CD8 conventional T-cell (CD8 Conv) exhaustion (**Figure 6I -**

516    **middle**). The percentage of activated CD4 Convs and CD8 Convs was correlated (R=0.70,

517    P=2.6x10$^{-9}$). In addition, CD4 conv activation was somewhat correlated with CD8 Conv

518    exhaustion (R=0.38, P=4.0x10$^{-3}$, **Figure S6K**). CD4 Treg activation levels were higher in healthy

519    tissues and tumors than CD4 and CD8 Conv T-cells (**Figure S6L**). In addition, Treg activation

520    was significantly higher in thyroid cancer (P=3.0x10-6) and esophageal cancer (P=0.0045)

521    relative to matched normal tissues.

522

523    Observing that many tumor-infiltrating T-cells had both low ASA and exhaustion usage, we

524    defined bystanders as cells with summed ASA and exhaustion usage below 0.0625. The

525    percentage of CD4 bystanders varied widely by cancer from 42.0% (esophageal) to 91.2%

526    (bladder) and CD8 bystanders varied similarly from 35.5% (endometrial) to 90.1% (bladder).

527

528    Within tumor samples, we tested which T-cell subset cGEPs were enriched for bystanders

529    (**Figure 6J**). The most bystander-enriched subsets were CD4-Naive (OR=15.9), Th2-Resting

530    (OR=10.6), Th1-like (OR=7.3), MAIT (OR=4.42), and CD8-Naive (OR=4.03) (Fisher Exact Test

531    P< 1x10$^{-100}$ for all comparisons). The subsets most depleted of bystanders were also those most

532    enriched for activation, namely Tph (OR=0.19), Treg (OR=0.23), and CD8-Trm (OR=0.61)

533    (P<1x10$^{-21}$, all comparisons). These analyses illustrate how TCAT and ASA scoring can

534    facilitate exploration of disease.

## 7. Identifying disease-associated cGEPs

Next, we associated cGEPs with sample-level disease phenotypes in infection, autoimmunity, and cancer (**Table S7**). First, we tested cGEP associations with Covid-19 (**Methods**). We applied ordinary least squares using psuedobulk sample-level features to two PBMC-derived T-cell datasets: UK-Covid (80 Covid-19, 21 healthy donors, **Figure 7A**) and COMBAT (77 Covid-19, 10 healthy donors, **Figure 7B**). We observed overall concordant cGEP associations (Pearson R=0.64, P=$2.8\times10^{-7}$, **Figure 7C**). Consistent with the key role of interferon in viral infections[17,18], ISG was the most positively upregulated cGEP in both datasets (FDR-corrected P, denoted as Q<0.05). AIM-associated functional cGEPs were up-regulated in acute Covid-19, consistent with viral activation of T-cells. These included exhaustion, cell cycle, TIMD4/TIM3, OX40/EBI3, NME1/FABP5, and CTLA4/CD38 (Q<0.05 for both datasets). We also found increased Tph cGEP usage in Covid-19 relative to controls (Q<$1\times10^{-8}$ for both datasets), consistent with recent demonstration of increased abundance of this subset in infection[39]. An intriguing novel finding is that the Th1-like cGEP was significantly negatively associated with Covid-19 in both datasets (Q<$1\times10^{-4}$). This negative association was seen within manually gated CD4 memory (Q=$1.1\times10^{-4}$) and CD4 effector memory subsets (Q=$4.5\times10^{-6}$), suggesting it is not due to differential abundance of circulating memory CD4 T-cells. Consistent with this, pseudobulk expression of the Th1 markers *CXCR3* RNA and protein levels were significantly lower in Covid-19 samples relative to controls (P=$8.1\times10^{-7}$ and 0.010 respectively, COMBAT). Immediate early gene cGEPs (IEG1, IEG2, IEG3) were also significantly associated with Covid-19 in the COMBAT dataset (FDR-corrected P<$1\times10^{-5}$) but not in the UK-Covid dataset (P>0.5), perhaps related to sample processing differences (see section 2).

Next, we identified cGEPs associated with inflamed synovial tissue in rheumatoid arthritis (RA) using the AMP-RA dataset, which includes synovial biopsies from 70 RA and 8 osteoarthritis

560 (OA) patients (**Figure 7D**)[20]. Ten out of the eleven significantly associated cGEPs were AIM-

561 associated, including the metallothionein (Q=2.9x10[-5]), ISG (Q=0.0020), Tph (Q=0.0020), HLA

562 (Q=4.9x10[-5]), ICOS/CD38 (Q=0.00010), Exhaustion (Q=0.041), and cell cycle (Q<.05 for all

563 three). Of note, Metallothionein was shown to be increased in the plasma of RA patients and

564 within the synovia of mouse models of RA[55]. The Tph association is consistent with prior

565 observations by us and others of Tph enrichment within RA synovia[22]. The Th22 cGEP was also

566 associated with RA (Q=0.0027), confirming a prior observation of increased Th22 cell

567 abundance in RA synovia, where they may stimulate osteoclasts[56].

568

569 Lastly, we identified cGEPs associated with T-cells in tumors relative to matched healthy tissues

570 (**Figure 7E**). We utilized a pan-cancer dataset containing 89 tumor and 47 matched normal

571 samples from 13 cancer types. First, we analyzed all samples together, controlling for tumor

572 type and sequencing technology as fixed effects. The Treg cGEP was the most strongly

573 associated, consistent with the known importance of Tregs in tumors (Q=7.4x10[-12])[57]. The

574 exhaustion and ISG cGEPs were also strongly associated with cancer, as expected (Q=8.5x10[-6]

575 and 9.3x10[-6], respectively)[58,59]. There was also substantial upregulation of AIM-associated

576 functional cGEPs, including CTLA4/CD38 (Q=1.3x10[-9]), TIMD4/TIM3 (Q=1.3x10[-9]), and

577 OX40/EBI3 (Q=4.9x10[-9]). Overall, 17 of the 21 significantly upregulated cGEPs in tumor-

578 infiltrating T-cells were AIM-associated (Fisher exact test P=7.4x10[-6]).

579

580 We also separately tested for cGEP association in each of the six cancer types with at least two

581 normal and two tumor samples (**Methods**). The results were highly concordant across cancers

582 (P<.05, sign test, for 14 out of 15 pairs of tumor types, **Figure 7F**). For example, the Treg,

583 Exhaustion, and CTLA4/CD38 cGEPs were significantly upregulated in all six tumor types

584 tested (P<.05). However, some signals were more specific. The Th17-Activated cGEP was only

585 significant in thyroid and hepatocellular carcinoma (P=5.3x10[-6] and P=0.013), while the Th2-

586　Activated cGEP was upregulated in esophageal, uterine, thyroid and hepatocellular carcinoma

587　(P=0.023, P=0.023, P=0.00057, P=0.0019).

588

589　Surprisingly, the Tfh-2 and Tph cGEPs were both upregulated in cancer ($Q=3.6 \times 10^{-4}$, $Q=3.3 \times 10^{-10}$

590　$^{-10}$). T follicular helper (Tfh) and T peripheral helpers (Tph) are *CXCL13*-producing CD4 subsets

591　that recruit B-cells and aid in antibody production. Tfhs are found primarily in lymphoid organs

592　and Tphs are predominantly in inflamed tissues[60], including likely within tumors[61].

593

594　Consistent with functional Tph activity, the expression of the B-cell chemoattractant *CXCL13*

595　was highly correlated with average Tph cGEP usage across samples ($R=0.67$, $P=1.2 \times 10^{-30}$,

596　**Figure S7A**). This correlation was stronger in tumor ($R=0.69$, $P=1.2 \times 10^{-13}$) than normal samples

597　($R=0.34$, $P=0.021$). We hypothesized that average Tph usage would correlate with plasma cell

598　abundance in tumors. To test this, we re-analyzed a published pan-cancer dataset containing

599　other cell-types besides T-cells from 148 primary tumors, 53 matched adjacent tissues, and 25

600　healthy donor samples[62]. Tph usage and *CXCL13* expression remained correlated in this

601　dataset ($R=0.67$, $P=1.2 \times 10^{-30}$, **Figure S7B**). Average Tph, Tfh-1, and Tfh-2 cGEP usage were

602　significantly correlated with plasma cell percentage within the tumors (Spearman $\rho=0.23$, 0.34,

603　0.28, respectively, $P<1 \times 10^{-2}$, **Figure S7C**). In a multivariate regression across all samples, Tfh-1

604　and Tph usage were independently associated with plasma cell abundance (P=0.042, P=0.051

605　respectively). Subsetting to non-tumor samples, Tfh-1 and Tfh-2 remained statistically

606　significant (P=0.017, P=0.027, respectively), but Tph was no longer significant (P=0.351). These

607　findings suggest that Tph cells are functional within tumors and are associated with increased

608　abundance of plasma cells.

## Discussion

609

610    Here, we introduced CellAnnoTator (abbreviated *CAT) for annotating scRNA-Seq data with

611    predefined GEPs. *CAT exploits the observation that functionally informative GEPs learned by

612    cNMF are reproducible across different datasets and contexts (**Figure 2**). This enables GEPs

613    identified across multiple reference datasets to aid in interpreting new datasets.  We

614    demonstrated *CAT with a GEP catalog derived from T-cells across diverse tissues and

615    diseases, yielding T-Cell AnnoTator (TCAT). We meta-analyzed a range of reference datasets,

616    obtaining the most comprehensive T-cell GEP catalog to date, including 16 subset-associated,

617    five technical artifact, and 25 functional programs.

618

619    TCAT demonstrated key advantages over clustering of T-cells. First, it simultaneously

620    annotated functional and subset GEPs within the same cells, disentangling signals that

621    clustering conflated (**Figure 4**). Second, TCAT out-performed RNA-based clustering for

622    annotation of T-cell subsets without requiring manual curation of the cluster labels (**Figure 3**).

623    Third, TCAT cGEP activity could be assessed across diverse disease states (**Figure 7**). TCAT

624    also improved upon prior matrix factorizations of T-cells by yielding a more comprehensive

625    catalog of T-cell GEPs. It was faster than running *de novo* matrix factorization, avoided the need

626    to manually re-label GEPs, and increased accuracy for smaller datasets (**Figure 1C-F**).

627

628    TCAT explained why traditional T-cell subsets have been challenging to identify in scRNA-Seq.

629    T-cell transcriptional clusters were heavily influenced by many non-subset GEPs, including

630    technical artifacts, cell cycle, interferon response, and cytotoxicity (**Figure 4**). TCAT overcame

631    this by annotating subset-associated cGEPs in parallel with functional cGEPs. In addition, TCAT

632    revealed how cGEPs can be expressed in different contexts. For example, the cytotoxic cGEP

633    was expressed in multiple subsets, and polarization cGEPs were expressed in both CD4 and

634    CD8 T-cells (**Figure 4E, S4**). There has recently been increased recognition of polarized CD8

635    populations such as Tc2 which can secrete cytokines typically associated with Th2-polarized

636    CD4 memory T-cells[40]. TCAT helped reveal these overlooked populations in scRNA-Seq data.

637

638    TCAT also highlighted the growing recognition of T peripheral helper (Tph) cells in disease. The

639    Tph cGEP was significantly associated with Rheumatoid Arthritis (RA), Covid-19, and Cancer

640    (**Figure 6**). While the association with RA was expected since Tph cells were discovered there,

641    and recent data has identified Tph cells in Covid-19[39], the association with cancer is less well

642    established[63]. Tph usage was associated with expression of *CXCL13* and plasma cell

643    abundance in tumors, suggesting Tph cells may drive lymphoid aggregation.

644

645    We also demonstrated that many cGEPs were induced following a TCR-dependent activation

646    stimulus using the novel AIM-Seq assay (**Figure 5**). AIM-Seq produces TCR and CITE-Seq

647    profiles for T-cells that are labeled based on their response to activation-induced marker

648    assays. This identified 24 cGEPs associated with TCR-dependent activation, including 11 that

649    may reflect context-dependent activation responses such as Th17-activated in Th17-polarized

650    cells and CTLA4/CD38 in Tregs. Many of the AIM-associated GEPs were strongly associated

651    with Covid-19, rheumatoid arthritis, and cancer, consistent with the importance of TCR-

652    dependent activation in these diseases (**Figure 6**).

653

654    We aggregated several AIM-associated cGEPs into an antigen-specific activation (ASA) score

655    to compare activation rates across diseases and cell subsets. This revealed impressive

656    variability in the percentage of activated and exhausted CD4 and CD8 T-cells within and

657    between different tumor types (**Figure 7**). In all tumor types, many T-cells lacked activation or

658    exhaustion signatures and were labeled as bystanders. Bystanders were enriched for naive and

659    unconventional T-cell subsets, whereas activated cells were enriched for Treg, Tph, and

660    resident memory subsets. This approach shows how TCAT can aid in characterizing activation

661    and exhaustion *in vivo*.

662

663    We highlight some current limitations of TCAT. First, TCAT's output can be non-sparse, leading

664    to non-zero usage of cGEPs contributing little biological function. This necessitates the use of

665    thresholds balancing sensitivity and specificity to decide if a cGEP is active in a cell. For

666    example, annotating TCR-activation or polarization currently relies on score thresholds. This

667    limitation can be mitigated by algorithmic improvements that increase TCAT's sparsity. Second,

668    several cGEPs lack a clear interpretation, or may be redundant with other cGEPs in the catalog.

669    For example, three cGEPs labeled IEG1-IEG3 are strongly enriched for immediate early genes.

670    We used reproducibility of spectra across multiple datasets to enrich for biologically meaningful

671    GEPs. As more datasets get incorporated, we anticipate increasing robustness of the catalog.

672    Furthermore, new experimental perturbation datasets can facilitate linkage of cGEPs with

673    upstream regulators to aid in interpretation.

674

675    We demonstrated application of *CAT to T-cells, but it is equally applicable to other cell types or

676    tissues. We make the *CAT software publicly available and have created a repository to host

677    cGEP catalogs, enabling easy application to new datasets. Furthermore, users studying other

678    tissues and cell-types can contribute their own catalogs to the repository. We envision this as a

679    resource akin to the molecular signatures database (MSigDB)[64,65] , but hosting GEPs for

680    annotation of scRNA-Seq data rather than gene-sets for enrichment testing. We hope it will aid

681    in comprehensive identification of GEPs underlying cell behavior across tissues and diseases.

682 # Methods

683 **Materials and reagents**

684
685

| Reagent or Resource | Source | Identifier |
|---|---|---|
| XVIVO15 culture media | Lonza | Catalog #: 02-060Q |
| RPMI 1640 Medium | ThermoFisher | Catalog #: 11875093 |
| Benzonase Nuclease | Sigma Aldrich | CAS #: 9025-65-4 |
| Anti-CD28 antibody | Biolegend | Catalog #: 302933<br>RRID: AB_11150591 |
| Anti-CD49d antibody | Biolegend | Catalog #: 304339<br>RRID: AB_2810443 |
| Human TruStain FcX™ (Fc Receptor Blocking Solution) | Biolegend | Catalog #: 422302<br>RRID: AB_2818986 |
| Zombie Yellow™ Fixable Viability Kit | Biolegend | Catalog #: 423104 |
| TotalSeq™-C Human Universal Cocktail, V1.0 | Biolegend | Catalog #: 399905 |
| Human TOTAL-SeqC Repertoire (5') Hashing Antibodies | BioLegend | Catalog #: 394661, 394663, 394665 |
| Anti-CD3-BV421 (SK7) | Biolegend | Catalog #: 344833<br>RRID: AB_2565674 |
| Anti-CD134-PE (Ber-ACT35) | Biolegend | Catalog #: 350003<br>RRID: AB_10641708 |
| Anti-CD274-BV785 (29E.2A3) | Biolegend | Catalog #: 329735<br>RRID: AB_2629581 |
| Anti-CD137-APC (4-B4-1) | Biolegend | Catalog #: 309809<br>RRID: AB_830671 |
| Anti-CD4-FITC (RPA-T4) | Biolegend | Catalog #: 300505<br>RRID: AB_314073 |

| Chromium Next GEM Single Cell 5' Kit v2, 16 rxns | 10X | Catalog #:1000263 |
|---|---|---|
| Dual Index Kit TN Set A, 96 rxn | 10X | Catalog #: 1000250 |
| Chromium Next GEM Chip K Single Cell Kit, 48 rxns | 10X | Catalog #: 1000286 |
| Chromium Single Cell Human TCR Amplification Kit, 16 rxns | 10X | Catalog #: 1000252 |
| Library Construction Kit, 16 rxns | 10X | Catalog #: 1000190 |
| 5' Feature Barcode Kit, 16 rxns | 10X | Catalog #: 1000256 |

686
687
688

**CellAnnoTator (\*CAT) Algorithm**

690

691 Whereas cNMF learns both GEPs and their usage in cells, \*CAT has the simpler problem of

692 fitting the usage for a fixed set of GEPs. Specifically cNMF runs NMF multiple times, each time

693 solving the following optimization:

694

695 $$ArgMin_{G,U} \, | X - UG |_F \text{ where } U \geq 0, G \geq 0$$

696

697 where $X$ is a NxH matrix of N cells by the top H overdispersed genes, $U$ is a learned NxK matrix

698 of the usages of K GEPs in each cell, and $G$ is a learned KxH matrix where each row encodes

699 the relative contribution of each highly variable gene in a GEP. H is usually a parameter set to

700 ~2000 overdispersed genes. $| \,|_F$ denotes the Frobenius norm. $X$ includes variance-normalized

701 overdispersed genes to ensure biologically informative genes are included and contribute

702    similar amounts of information even when they may be expressed on different scales. For

703    cNMF, the optimization is solved multiple times and the resulting $G$ matrices are concatenated,

704    filtered, and clustered to determine a final average estimate of $G$. Ultimately cNMF refits the

705    GEP spectra into two separate representations, one reflecting the average expression of the

706    GEP and units of transcripts per million $G^{tpm}$ and on in Z-scored units used to define marker

707    genes $G^{scores}$ (see Kotliar et, al., 2019[14] for details).

708

709    Analogously, *CAT takes a fixed catalog of GEPs as input, denoted as $G^*$, and a new query

710    dataset $X^{query}$ and solves the optimization:

711

712    $$ArgMin_U \mid X^{query} - UG^* \mid_F \text{ where } U \geq 0$$

713    The columns of $X^{query}$ and $G^*$ correspond to a pre-specified set of overdispersed genes.

714    Analogous to cNMF, we use gene-wise standard-deviation-normalized counts for $X^{query}$. See

715    below for how $G^*$ is calculated for T-CellAnnoTator. We solve for $U$ with non-negative least

716    squares using the NMF package in scikit-learn version 1.1.3[66] with $G^*$ fixed. We use the

717    Frobenius error, the multiplicative update ("mu") solver, tolerance of $1 \times 10^{-4}$, and max iterations

718    of 1000. We then row-normalize the $U$ matrix so that each cell's aggregate usage across all K

719    GEPs sums to 1.

720

721    **Dataset pre-processing and batch-effect correction**

722

723    To generate the input matrix for cNMF for each dataset, we first filtered genes detected in fewer

724    than 10 cells and cells with fewer than 500 unique molecular identifiers (UMIs). We also

725    excluded antibody-derived tags (ADTs) and genes containing a period in their gene name. We

726    subsequently subsetted the data to the top 2000 most overdispersed genes, identified by the

727    "seurat_v3" algorithm as implemented in Scanpy[35]. Next, we scaled each gene to unit variance.

728    To avoid outliers with excessively high values, we calculated the 99.99th percentile value across

729    all cells and genes and set this as a ceiling. We denote this matrix as $X^{raw}$.

730

731    We used an adapted version of harmonypy to correct batch effect and other technical variables

732    from $X^{raw}$ prior to cNMF[21]. For this, we computed Harmony's maximum diversity clustering

733    matrix from principal components calculated from a normalized version of $X$ which we label

734    $X^{norm}$. Specifically, to compute $X^{norm}$, we started from the same initial gene list described

735    above but first normalized the rows of the matrix so that each cell's counts sum to 10,000

736    (TP10K normalization). We then subsetted to the top 2000 overdispersed genes, and scaled

737    each column (gene) to unit variance, resulting in $X^{norm}$. We then performed principal

738    component analysis (PCA) on $X^{norm}$ and supplied those principal components to the

739    run_harmony function of harmonypy. We then used the mixture of experts model correction,

740    implemented in harmonypy with the computed maximum diversity clustering matrix, but instead

741    of correcting the PCs using this model, as standard Harmony does, we corrected $X^{raw}$. This

742    creates a small amount of variability around 0 for the smallest values in $X^{raw}$. We therefore set

743    a floor of 0, resulting in the corrected matrix $X^c$ used as the count matrix for cNMF.

744

745    **Consensus non-negative matrix factorization (cNMF)**

746

747    We ran cNMF on the batch-corrected $X^c$ matrix which only includes the top 2000 overdispersed

748    RNA genes. Spectra for the resulting GEPs were then refit by cNMF including all genes that

749    passed the initial set of filters including ADTs. Specifically, RNA counts were normalized to sum

750    to 10,000, and ADT counts were separately normalized to sum to 10,000 and the combined

751    matrix was passed as the –tpm argument for cNMF. Thus the GEP spectra output by cNMF

752    incorporate ADTs and genes not included in the 2000 overdispersed genes.

753

754    cNMF was run for each dataset with the number of components (K) varying between 15 and 55

755    and with 20 iterations. The final number of NMF components used for each dataset, K*, was

756    chosen by visualizing the trade-off between reconstruction error and stability for these runs

757    (**Supplementary item 1**). Once K* was selected, we ran cNMF a final time with only this value

758    for K and with 200 iterations to generate the final GEP spectra estimates.

759

760    **Constructing a catalog of consensus GEPs (cGEPs)**

761

762    Next, we identified consensus GEP spectra – I.e. the average of correlated GEP spectra

763    identified by cNMF in different datasets. Normalized input GEP vectors, denoted as $g_i$, were

764    computed by starting from the spectra_tpm output from cNMF, renormalizing each vector to

765    sum to $10^6$, and then dividing each element by the standard deviation of the corresponding gene

766    in the –tpm input to cNMF.  Then, we created an undirected graph where the 267 GEPs

767    identified across all reference datasets were represented as nodes $g_1$ … $g_{267}$. We drew edges,

768    denoted as $E_{i,j}$ connecting a pair of GEPs $g_i$ and $g_j$ if the following criteria were met:

769

770    1.  $g_i$ and $g_j$ were from different datasets

771    2.  $R_{ij} > 0.5$ where $R_{ij}$ denotes the Pearson correlation between $g_i$ and $g_j$. For computing $R_{ij}$ ,

772        $g_i$ and $g_j$ were subset to the union of the overdispersed genes for each dataset.

773    3.  $g_i$ was among the top seven most correlated GEPs with $g_j$, and $g_j$ was among the top

774        seven most correlated GEPs with $g_i$ with correlation defined as in 2.

775

776    Next, we initialized a set for each GEP: $x_1 = \{g_1\}$ … $x_{267} = \{g_{267}\}$. We then iterated through all

777    edges $E_{i,j}$ in the graph in order of decreasing $R_{ij}$ and merged the sets $x_i$ and $x_j$ into a new set $x_{i,j} =$

778    $\{g_i , g_j\}$. If either $g_i$ or $g_j$ were already members of a merged set from previous merges, we

779    merged their containing sets only if at least two thirds of the GEP pairs in the resulting

780    consensus set were connected by edges. For example, if there is an edge $E_{4,9}$ and $g_4$ is already

781    merged into a set {$g_1$, $g_2$, $g_4$}, then we only merged {$g_1$, $g_2$, $g_4$} and {$g_9$} if there were also

782    edges $E_{1,9}$ and $E_{2,9}$. This resulted in 52 merged sets and 52 unmerged "singleton" sets. We

783    filtered 49 of the 52 singletons and retained 3 that had a biological explanation for being

784    identified in only one dataset.

785

786    Lastly, we subset each GEP to the union of overdispersed genes across all 7 reference

787    datasets that were present in all dataset and obtained the final consensus GEPs by taking the

788    element-wise average GEPs in each merged set. This matrix was used as the reference for

789    TCAT. For marker gene analyses (e.g. **Figure 2B, D, Supplementary Item 2**), we element-wise

790    averaged the Z-score representation of GEPs output by cNMF for GEPs in a consensus set.

791

792    **Simulation analysis**

793

794    We adapted the scsim simulation framework described in the cNMF publication[14] and based on

795    Splatter[67] into a new iteration, scsim2. Like with scsim, we distinguished between subset GEPs

796    which are mutually exclusive and non-subset or "activity" GEPs which are not. For the original

797    scsim framework, cells used one of multiple subset GEPs and potentially used a single activity

798    GEP. We adapted scsim to allow cells to use anywhere from none to all of the activity GEPs in

799    addition to their single subset GEP. We kept the Splatter parameters used in the cNMF

800    publication to describe the distribution of gene expression data: mean_rate=7.68,

801    mean_shape=0.34, libloc=7.64, libscale=0.78, expoutprob=0.00286, expoutloc=6.15,

802    expoutscale=0.49, diffexpprob=.025, diffexpdownprob=.025, diffexploc=1.0, diffexpscale=1.0,

803    bcv_dispersion=0.448, bcv_dof=22.087.

804

805     For figure 1, we simulated 10 subset GEPs and 10 activity GEPs based on 10,000 total genes.

806     The extra-GEP reference included all 20, the missing-GEP reference included 6 of the subset

807     GEPs and 6 of the non-subset GEPs, and the query dataset included 8 subset GEPs and 8 non-

808     subset GEPs. Each dataset consisted of 9000 genes, randomly sampled from the 10,000. Each

809     cell was randomly assigned a subset GEP with uniform probability (shown in the UMAP in figure

810     3B), and each cell randomly selected whether it expressed each activity GEP with probability of

811     0.3. The degree of usage of each activity GEP was sampled uniformly between 0.1 and 0.7. If

812     the sum of the activity GEPs exceeded 0.8 for a cell, they were renormalized to sum to 0.8.

813     Thus each cell's usage of its subset GEP always exceeded 0.2. We simulated 100,000 cells

814     each for the extra-GEP and missing GEP references. We simulated multiple query datasets

815     containing 100, 500, 1000, 5000, 10,000, 20,000, 50,000, or 100,000 cells.

816

817     We subsequently ran cNMF using 1000 overdispersed genes, 20 iterations,

818     local_neighborhood_size=0.3 and density_threshold=0.15. We used K=20, K=12, and K=16 for

819     the extra-GEP reference, missing-GEP reference, and query datasets respectively. We then

820     used *CAT to fit the usage of the reference GEPs on the query dataset. To evaluate the

821     performance of *CAT and cNMF, we calculated the Pearson correlation of the inferred GEP

822     usage with the simulated ground truth usage.

823

824     **Gene-set enrichment analysis**

825

826     We used Fisher Exact Test in Python's Scipy library to associate cGEPs with gene sets. For the

827     T-cell polarization dataset[24] we defined polarization gene sets as genes that had FDR-corrected

828     P-value < 0.05 and fold change > 2 with the stimulation condition. We excluded genes with

829     FDR-corrected P-value between 0.05 and 0.2 and fold-change>1, as many of these are up-

830     regulated by the stimulation but just did not reach FDR significance. We also obtained literature

831    gene sets corresponding to immediate early genes[28] and gene ontologies[23,68]. We tested these

832    literature gene-sets for enrichments with gene sets derived from the Z-score representation of

833    cGEPs based on a score threshold of 0.015, which corresponded to the 99th percentile across

834    all genes and cGEPs. We then tested for association using Fisher's Exact Test as implemented

835    in scipy.stats in Python.

836

837    **Manual subset gating analysis**

838

839    We library-size normalized antibody derived tag (ADT) protein measurements to sum to $10^4$

840    (TP10K) and applied the centered log ratio (CLR) transformation. We then scaled each protein

841    to unit variance, and truncated at 15 to remove excessively high outliers. Next, we performed

842    principal component analysis (PCA) and ran batch correction using harmonypy with the same

843    batch features as for cNMF. We then computed the K-nearest neighbor graph with K=5

844    neighbors, using the Harmony-corrected principal components. We then smoothed the

845    normalized protein estimates using MAGIC[69] using the K-nearest neighbor graph computed

846    above and the diffusion operator powered to t=3.

847

848    We gated canonical T-cell subsets using the smoothed normalized ADTs. First, we gated

849    gamma-delta (γδ) T-cells using expression of Vδ2 TCR. Then, we separated MAIT cells using

850    expression of CD161 and TCR Vα 7.2. We then used CD4 and CD8 to separate CD4

851    (CD4+CD8-), CD8 (CD4-CD8+), double positive (DP) (CD4+CD8-), and double negative (DN)

852    (CD4-CD8-) T-cells. We then subset to CD4 T-cells and gated regulatory T-cells (Tregs) using

853    expression of CD25 and CD39. Of the remaining CD4 T-cells, we used CD62L and CD45RA to

854    define CD4 Naive (CD62L+CD45RA+), CD4 Central Memory (CD62L+CD45RA-), CD4 Effector

855    Memory (CD62L-CD45RA-), and CD4 TEMRA (CD62L-CD45RA+) populations. For the CD8 T-

856    cells, we similarly used CD62L and CD45RA to define CD8 Naive (CD62L+CD45RA+), CD8

857 Central Memory (CD62L+CD45RA-), CD8 Effector Memory (CD62L-CD45RA-), and CD8

858 TEMRA (CD62L-CD45RA+) populations.

859

860 **T-cell subset classification benchmarking analyses**

861 We used T-cell subsets defined by manual gating of ADTs in the Flu-Vaccine dataset as ground

862 truth for prediction. For single cGEP prediction, we ran TCAT to predict cGEP usage, and

863 identified the cGEP that best predicted the lineage based on area under the curve (AUC).

864

865 We also used all of the cGEP simultaneously to perform simultaneous multi-label prediction. We

866 scaled the normalized usages for all cGEPs to zero mean and unit variance. Using COMBAT as

867 a training dataset, we trained a multinomial logistic regression using scikit-learn[66] version 1.0.2

868 with lbfgs solver to predict gated subset from usages. Model weights were adjusted by the

869 inverse of subset size using class_weight="balanced", allowing subsets with different cell counts

870 to contribute to the model equally. We excluded CD4 TEMRA, double negative, and double

871 positive subsets from this analysis due to low cell counts in both the training and testing

872 datasets. We evaluated this model in the independent Flu-Vaccine query dataset.

873

874 Analogous comparisons were made using GEPs from Yasumizu et. al, 2024 fit to the data using

875 the NMFproject software[13]. We also obtained gene sets derived from NMF analyses of T-cell in

876 a pan-cancer dataset[16]. To assess the ability of these gene sets to predict gated subsets, we

877 used the score_genes function in Scanpy[35] on data normalized following the standard pipeline

878 (library size normalizing to TP10K, log transformation, scaling each each gene to unit variance).

879 We then assigning each subset to the gene set that yielded the maximal AUC.

880

881 To evaluate clustering, we first normalized the data as above, and subset to highly variable

882 genes using the highly_variable_genes function in Scanpy with default parameters. We then ran

883 principal component analysis (PCA) and Harmony batch correction of the PCs[21]. We then

884 computed the K nearest neighbor graph using 31 harmony-corrected PCs and 30 nearest

885 neighbors. We then performed Leiden clustering[70] with resolution parameters ranging from 0.25

886 to 2.25 increasing by 0.25. For each clustering resolution, we performed a greedy search to

887 assign clusters to manually gated subsets based on maximization of the balanced accuracy (I.e.

888 the average recall across all subsets). In each iteration, we considered all unassigned clusters

889 and possible gated subset assignments, and selected the cluster and assignment that most

890 increased the overall balanced accuracy. When no remaining cluster assignments would

891 increase the balanced accuracy, we assigned the cluster to a subset that least decreased the

892 balanced accuracy. We continued this process until each cluster was assigned to a subset.

893

894 **Activation Induced Marker assay followed by scRNA-Seq (AIM-Seq)**

895

896 PBMCs were quickly thawed and placed in pre-warmed xVIVO15 cell culture medium (Lonza)

897 supplemented with 5% heat-inactivated FBS. To reduce cell clumping, PBMCs were incubated

898 in xVIVO15 containing 50 U/mL of benzonase nuclease (Sigma-Aldrich) for 15 minutes at 37

899 degrees and filtered using a 70 μm cell strainer. Washed and nuclease treated cells were

900 seeded in a 96 well cell culture plate at a concentration of $2.5 \times 10^6$/mL. Peptide stimulations

901 were performed using the CEFX Ultra SuperStim Pool (JPT Peptide Technologies, Product

902 Code: PM-CEFX-1) at a final concentration of 1.25 μg/mL per peptide for 22 hours at 37

903 degrees and 5% $CO_2$. Recombinant anti-CD28 and anti-CD49d antibodies (BioLegend) were

904 added at a final concentration of 5 μg/mL and 0.625 μg/mL, respectively, to provide co-

905 stimulation for peptide reactive T-cells. Separately mock-stimulated cells were treated with anti-

906 CD28 and anti-CD49d antibodies at the same concentration.

907

908 Peptide responsive T-cells were detected by the expression of the surface activation markers

909 PD-L1, OX40, and CD137 via flow cytometry. Following the stimulation, peptide treated and

910 mock-stimulated cells were washed in cell staining buffer (PBS + 2mM EDTA + 2% FBS) to end

911 the stimulation. Fc receptor blocking was performed using a 1:50 dilution of Human TruStain

912 FcX (Biolegend) in cell staining buffer for 10 minutes at 4 degrees. Cell viability staining was

913 performed using a 1:500 dilution of Zombie Yellow Fixable Viability Dye (BioLegend) prepared

914 in PBS for 30 minutes at 4 degrees. Surface staining was performed using 1:100 dilutions of

915 BV421 conjugated anti-CD3, FITC conjugated anti-CD4, BV786 conjugated anti-PD-L1, PE

916 conjugated anti-OX40, and APC conjugated anti-CD137 (BioLegend) for 25 minutes at 4

917 degrees in cell staining buffer. Following cell staining, antigen reactive and non-reactive T-cells

918 were identified using a BD FACSAria II cell sorter and collected in cRPMI medium (100 U/mL

919 penicillin-streptomycin + 2 mM L-glutamine + 10 mM HEPES + 0.1 mM non-essential amino

920 acids + 1 mM sodium pyruvate + .05 mM 2-Mercaptoethanol) supplemented with 20% FBS.

921 Sorted T-cell populations were then labeled with 75 uL of TotalSeq oligo conjugated hashing

922 antibody mix, incubated for 30 minutes at 4 degrees with gentle mixing after 15 minutes, and

923 pooled in equal quantities. Staining with the TotalSeq-C Human Universal Cocktail (BioLegend)

924 was then performed according to the manufacturer's instructions. The cells were then

925 resuspended in PBS supplemented with .04% FBS at a final concentration of 500 cells/μL and

926 submitted for single-cell profiling on the Chromium Next GEM instrument. Library preparation

927 was completed for the hashtag oligos, single-cell rna-seq, cite-seq, and TCR-repertoire

928 sequencing following the manufacturer's instructions.

929

930 We collected AIM-Seq data from two separate 10X runs. In the first experiment, PBMCs from

931 three donors were processed independently as described above and were pooled together after

932 fluorescence activated cell sorting (FACS). In the second run, PBMCs from four donors, two of

933 which overlapped with the first run, were stimulated separately and pooled prior to FACS.

934

**Preprocessing the AIM-Seq dataset**

936     The AIM-Seq data was processed using Cell Ranger version 6.1.1 with default parameters and

937     alignment to hg38 reference genome. The donor of origin for each cell was determined using

938     Demuxlet version 1.0 with doublet-prior of 0.1[71]. Cells with null or ambiguous demuxlet result,

939     fewer than 10 counts of the hashtag oligos, or fewer than 50 total RNA counts were filtered. To

940     account for staining differences between the hashtag oligos and different sequencing depths of

941     the two 10X runs, the counts for each hashtag oligo in each 10X run were scaled to have the

942     same median value. Next we added a pseudocount to the hashtag oligo counts and log10

943     transformed this data. Then we ran Gaussian Mixture models separately for each hashtag oligo

944     with K=2 clusters. Each cell was assigned to a single condition if it was in the high cluster for

945     one oligo and the low clusters for all others, a doublet if it was in the high cluster for more than

946     one oligo, or an empty droplet if it was in the low cluster for all oligos. Empty droplets or

947     doublets based on the hashtag oligo clustering were filtered, as were doublets based on

948     demuxlet. Genes detected in fewer than 10 cells were filtered prior to running TCAT.

949

**cGEP associations with AIM-positivity, proliferation, and disease**

951

952     To associate cGEPs with the AIM-Seq stimulus, we first ran TCAT to fit the usages of the

953     cGEPs in the AIM-Seq dataset. We then computed the average usage of each cGEP in cells

954     from each sort condition in each donor. We created two dummy variables, the first indicating

955     whether a sample was treated with CEFX or Mock, and the second indicating whether a CEFX-

956     treated sample was AIM-positive or not. We fit these two variables and an intercept to average

957     cGEP usage in the sample. cGEPs associated with the CEFX-or-Mock dummy variable were

958     labeled milieu-associated while cGEPs positively associated with the AIM-positive dummy were

959     labeled AIM-associated.

960

961    To associate cGEPs with proliferation, we defined cells as proliferating or non-proliferating in

962    each dataset by setting a threshold of 0.1 on the sum of the three cell cycle cGEPs, S-phase,

963    late S-phase, and G2M-phase. We then computed the mean usage of each cGEP per sample

964    separately in high cell-cycle (sum usage > 0.1) and low cell-cycle (sum usage < 0.1) cells. We

965    filtered samples that did not have at least 10 high cell-cycle cells and 100 low cell-cycle cells.

966    Then, for each cGEP, we performed a two-sample T-test paired by individual (ttest_rel in Scipy,

967    default parameters) between average cGEP usage for high and low cell-cycle cells. We meta-

968    analyzed P-values across datasets using Fisher's Method (combine_pvalues in Scipy).

969

970    To associate cGEPs with sample-level disease phenotypes, we calculated the average usage of

971    each cGEP in each sample for a given dataset. We then used ordinary least squares regression

972    to find cGEPs with higher average usage in disease samples than controls, controlling for

973    sample-level batch variables as covariates. For all datasets, disease status was modeled as a

974    binary dummy variable, and an intercept was included. For UK-Covid, the processing site was

975    included as dummy variable covariates. For COMBAT, sequencing pool, and processing

976    institute were included as dummy variable covariates. For the Pan-cancer dataset, all cancer

977    types were initially included in the analysis and dummy variable covariates were included for

978    tissue of origin. In addition, sequencing technology was included as a dummy variables. When

979    there were multiple tumor samples or matched normal samples from the same donor, we

980    retained only the duplicate sample with the most cells prior to the regression.

981

982    For all association tests, we performed FDR-correction of the P-values using the Benjamini

983    Hochberg method (fdrcorrection in Statsmodels with method='indep').

984

985 **Defining the antigen-specific activation (ASA) score**

986

987 We used CD71+CD95+ surface protein co-expression in the COMBAT and Flu-Vaccine

988 datasets as an *in vivo* correlate of TCR activation to help prioritize AIM-associated cGEPs for

989 predicting TCR-activated cells. First we preprocessed the ADT surface proteins in these

990 datasets as described in the manual subset gating section. We then subsetted cells by their

991 manual gating-defined broad cell types (CD4 Conv, CD4 Treg, CD8 Conv, other) and gated

992 CD71+CD95+ cells separately for each cell type as the response feature to be predicted by

993 AIM-associated cGEPs.

994

995 We then performed forward stepwise selection, evaluating how well the summation of usages of

996 different combinations of AIM-associated cGEPs would predict CD71+CD95+ gating. At each

997 stage, the per-cell ASA score was computed as the sum of normalized usages of cGEPs in the

998 predictive set. At each forward step, we determined which cGEP should be added to the

999 predictive set based on which would most improve the average AUC across the Flu-Vaccine

1000 and COMBAT datasets. We used a reduction in AUC in both datasets as the stopping criterion

1001 for adding cGEPs. We considered all AIM-associated cGEPs identified in section 6 as

1002 candidates for this, excluding those known to have a broader function outside of T-cell activation

1003 (e.g. cytoskeleton, metallothionein, cell cycle) and those reflecting activation-associated T-cell

1004 subsets (Tph and Th17-Activated). We also excluded Exhaustion from the ASA score as it

1005 reflects a distinct inhibitory response to antigen-stimulation that users may wish to annotate

1006 separately.

1007

1008 **Code availability**

1009 The code for CellAnnotator (starCAT) is available at

1010 https://github.com/immunogenomics/starCAT. The analysis scripts used in this paper are

1011 available at https://github.com/immunogenomics/TCAT_analysis.

# 1012 Acknowledgements

# 1017 Author contribution

1018 Conceptualization (D.K., M.C., S.R.), Software (D.K., M.C.), Formal Analysis (D.K., M.C.),

1019 Methodology (D.K., M.C., R.A., S.R.), Writing (D.K., M.C., K.W., A.N., Y.B., Y.Z., P.C.S., D.A.R.,

1020 S.R.), Funding acquisition (S.R.)

# 1021 Declaration of interest

1022 S.R. is a founder for Mestag, Inc, on advisory boards for Pfizer, Janssen and Sonoma, and a

1023 consultant for Abbvie, Biogen, Nimbus and Magnet. D.A.R. is a co-inventor on a patent using

1024 Tph cells as a biomarker in autoimmune diseases. P.C.S. is a co-founder of, shareholder in, and

1025 consultant to Sherlock Biosciences, Inc. and Delve Bio, as well as a Board member of and

1026 shareholder in Danaher Corporation. Other authors declare no competing interests.

1027

1028


1029    # Tables / Legends


1030    **Table S1. cGEP Summary.** Summary of cGEPs including their full name, abbreviated name,

1031    assigned class, top 3 most strongly associated genes, and which datasets it was derived from.

1032

1033    **Table S2. Marker genes.** Top 200 marker genes associated with each cGEP, colored by their

1034    strength of association with the cGEP, based on the average gene score.

1035

1036    **Table S3. Gene-set enrichment.** The "GO_Enrichment" tab includes the top 10 associated

1037    gene sets for each cGEP including the GEP name, gene-set name, fisher exact test odds ratio,

1038    and P-value. The subsequent tabs include the same information but for enrichment tests for

1039    gene sets defined from a dataset that polarized T-cells for either 16 hours (16h) or 5 days (5d)

1040    starting from either naive (TN) or memory T-cells (TM)[24]. The tab name indicates the stimulation

1041    conditions.

1042

1043    **Table S4. Correlation with cell quality features.** Each tab includes the Pearson correlation of

1044    each cGEP's usage (rows) with different per-cell quality features (tab names) for each dataset

1045    (columns). MitoFrac denotes the % of unique molecular identifiers from MT- genes.

1046    RNA_Detected denotes the number of unique genes detected per cell. RNA_Count denotes the

1047    number of unique molecular identifiers per cell. PCFrac denotes the percentage of unique

1048    molecular identifiers that are assigned to a protein coding gene in Gencode version 44.

1049

1050    **Table S5. AIM-Seq association.** Provides regression coefficients and P-values for the

1051    association between cGEP usage and binary variables reflecting CEFX vs. mock stimulation or

1052    AIM-positive vs. AIM-negative. Coef. represents the regression coefficient, P represents the P-

1053    value, and Q represents the FDR-corrected P-value.

1054

1055    **Table S6. Association with proliferation.** T-statistics, P-values, and $\log_2$ odds ratios for the

1056    paired T-test of proliferating and non-proliferating T-cells in each dataset (tabs). For the meta-

1057    analysis across datasets it provides the Fisher's method combined P-value and the average $\log_2$

1058    odds ratio.

1059

1060    **Table S7. Association with disease.** ordinary least squares regression coefficients (Beta), P-

1061    values (P), FDR-corrected Q-values (Q), and average fold changes (FC) for phenotype

1062    associations shown in **Figure 7**. Each tab represents a different phenotype.

1063

1064

1065 # Figures / Legends



1066

1067

1068 **Figure 1. Overview of CellAnnoTator (*CAT).** (A) Schematic of the *CellAnnoTator (*CAT)

1069 pipeline. (B) Schematic of simulation strategy (left) with resulting Uniform Manifold

1070 Approximation and Projection (UMAP) plot (right). Cells are colored by lineage gene expression

1071 program (GEP). (C-E) Pearson correlation of ground truth simulated usages of each GEP

1072 (columns) vs inferred usages (rows) for *CAT with the 20 GEP reference (C), *CAT with the 12

1073 GEP reference (D) or cNMF of the query with 16 inferred components (E). (F) Pearson

1074 correlation of ground truth and inferred usages by *CAT and cNMF for different query dataset

1075 sizes. Marker represents mean and error bars represent range. (G) Summary of reference

1076 datasets including number of individual donors (x-axis), number of cells (y-axis), and tissue

1077 source (dot color). Phenotypes are listed below the dataset names.

1078

1079

1080

1081

1082

1083

1084

1086 **Figure 2. Cataloging consensus gene expression programs (GEPs) across datasets.** (A)

1087 Pairwise correlations of GEPs discovered across reference datasets with insets for consensus

1088 GEPs derived from all seven references. Inset row and column orders are the same for all

1089 cGEPs. (B) Scatter plots of selected correlated GEP pairs. X and Y axis labels indicate the

1090 datasets the GEP was found in (P<1x10$^{-100}$ for all correlations). (C) Heatmap of cGEPs (rows)

1091 and which datasets the comprising GEPs were found in (columns). Green boxes indicate a GEP

1092 was found in a dataset. Colorbar indicates the cGEP's assigned class. cGEPs corresponding to

1093 non T-cell lineages were excluded. (D) Marker genes for selected example cGEPs in cNMF

1094 gene score units.

1095

1097 **Figure 3. Benchmarking T-CellAnnoTator on a query dataset.** (A) UMAP of the Flu-Vaccine

1098 dataset colored by the manually gating shown in **Figure S4A**. (B) Same UMAP as (A) but

1099 demonstrating prediction of manual gating of Treg and CD8 EM populations with the most

1100 associated individual cGEP (usage > 0.025), the multilabel classifier based on multiple cGEPs,

1101 or Ledien clustering with resolution 1.0. (C) Comparison of balanced accuracy for prediction of

1102 manually gated subsets, including clustering with multiple Leiden resolution parameters. (D)

1103 Usage of the mitochondria cGEP against the percentage of mitochondrial reads per cell (left).

1104 Usage of the CellCycle-S (middle) and CellCycle-G2M (right) cGEPs against the S and G2M

1105 scores output by Scanpy's score_genes_cell_cycle function with published proliferation gene

1106 sets[34]. (E) Heatmap of pseudobulk expression in cGEP-high and low cells, per sample. Samples

1107 are normalized by library size and expression is z-scored across rows.

1108
1109
1110

1111

1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124

1125
1126
1127

1128

1129

1130  **Figure 4. Comparing TCAT to clustering in the COMBAT dataset.**

1131  (A) UMAP of T-cells showing published sub-clusters of clusters annotated as CD4 memory with

1132  other clusters shown in gray. (B) Average usage of selected cGEPs across CD4 memory

1133  subclusters. (C) Same UMAP as (A) but colored by usage of selected subset, functional, and

1134  artifact cGEPs usage. Intensities are averaged over 20 nearest neighbors to reduce

1135  overplotting. (D) Usage of selected cGEPs in cells with high or low usage of cell cycle GEPs.

1136  Cells are grouped by their most highly used subset GEPs.(E) Percentage of cells within each

1137  manual gate assigned to each polarization (usage > 0.1). Bar represents the average and

1138  whiskers represents the 95% confidence interval, across samples.

1139

1140

1141

1142

1143

1144

1145
1146
1147
1148
1149

1150

1151    **Figure 5. Identifying cGEPs associated with TCR-dependent activation.** (A) Schematic of

1152    AIM-Seq. (B) FACS experiment from an AIM-Seq run showing surface activation markers in

1153    CD3+CD4+ and CD3+CD4- gated populations with the gates used for AIM-positive (+), AIM-

1154    negative (-) and Mock (M)  populations. (C-D) UMAP of AIM-Seq dataset colored by sorting

1155    condition (C) or manually gated population (D).  (E) cGEP association with AIM-positive

1156    samples. X-axis shows the mean $Log_2$ ratio of average usages. Y-axis shows the -$Log_{10}$ P-

1157    value. cGEPs are labeled by assigned category. (F) Average usage of selected Aim-associated

1158    cGEPs in +, -, and U cells from different gated subsets. Boxes represent interquartile range.

1159    Error bars represent 95th percentiles.

1160

1161

1162
1163
1164

1166 **Figure 6. Annotating antigen-specific activation (ASA) *in vivo*.** (A) Box plot of ASA score for

1167 cells stratified as activated (CD71+CD95+) or not activated. (B) Same as (A) but for AIM-Seq

1168 with cells stratified by sort condition. (C) Clonality in manually gated conventional CD4 and CD8

1169 T-cells annotated as activated (ASA>0.065) or not activated (ASA<0.065). Clonality is defined

1170 as the number of cells in the same sample with an identical alpha and beta CDR3 amino acid

1171 sequence. (D) Percentage of activated CD4 and CD8 convs (ASA>0.065) in Covid-19 and

1172 healthy control samples, by cohort. (E-F) UMAP of the COMBAT dataset colored by ASA score

1173 or low-resolution published clustering. (G) $Log_{10}$ odds ratio for 2x2 association of ASA positivity

1174 and manual gating subset assignment. * indicates P-value<0.05. (H) Percentage of activated

1175 (ASA>0.065) or proliferating (sum of cell cycle cGEPs>0.1) cells per sample across datasets.

1176 Boxes represent the interquartile range and whiskers represent 95% quantile range. (I)

1177 Percentage of activated, exhausted (exhaustion cGEP usage>0.065), or bystander (ASA +

1178 exhaustion usage<0.065) T-cells in CD4 and CD8 Convs, per sample stratified by tumor type

1179 and corresponding healthy tissues. (J) $Log_2$ odds ratio for enrichment of bystander T-cells by

1180 subset cGEP assignment. Error bars represent 95% confidence intervals.

1181

1182

1183

1184

1185   **Figure 7. cGEPs association with disease.** (A-B)  Associations of cGEP usage with Covid-19

1186   status for UK-Covid and COMBAT datasets. X-axis shows the regression coefficient. Y-axis

1187   shows the -Log10 FDR-corrected Q-value. (C) Scatter plot of regression coefficients from (A)

1188   and (B). (D-E) Same as (A) but comparing synovial T-cells from patients with Rheumatoid

1189   Arthritis and Osteoarthritis, or from tumors and healthy adjacent tumors. (F) Regression

1190   coefficients for tumor vs. normal samples for each tissue of origin. * denotes P<.05 for the

1191   corresponding coefficient. Cancer type abbreviations are: bladder cancer (BC), esophageal

1192   cancer (ESCA), hepatocellular carcinoma (HCC), renal cell carcinoma (RC), thyroid carcinoma

1193   (THCA), and endometrial cancer (UCEC).

1194

1195

1196

1197

1198

1199

1200      # Supplemental Figures / Legends



1201

1202 **Figure S1. Characterizing \*CAT.** (A) \*CAT predicted GEP usage for cells that use a GEPs with

1203 ground-truth usage>0.2, 0.1-0.2, or 0. Also shows the predicted usage for GEPs present in the

1204 reference data that are not present in the query (labeled unused GEP). (B) Number of GEPs

1205 identified in each dataset. The color indicates whether each GEP clustered with one or more

1206 GEPs from another dataset as part of a consensus GEP (purple, red, or green), did not cluster

1207 with a GEP from another dataset but was kept in the catalog as a dataset-specific GEP

1208 (orange), or did not cluster with a GEP from another dataset and was filtered (blue). (C)

1209 Absolute value of Pearson correlation of spectra learned by cNMF (top) or PCA (bottom)

1210 between different pairs of datasets. PCs are learned on the same matrices of batch-corrected

1211 matrices used for cNMF. Mean correlation refers to the mean value along the matrix diagonal,

1212 which corresponds to pairs of components with highest correlation across the two datasets.

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1224    **Figure S2. Annotating cGEPs.** (A) Manual gating of COMBAT dataset using smoothed surface

1225    protein antibody-derived tag (ADTs). (B) Multivariate logistic regression coefficients of cGEPs

1226    (columns) against manually gated populations (rows). For visualization, the minimum value is

1227    thresholded to 0 and the maximum is threshold to 1.25. Seven selected non-subset cGEPs are

1228    shown on the right as examples. (C) Pearson correlation of cGEPs with percentage of

1229    mitochondrial transcript per cell, for each dataset. All cGEPs excluding Mito and Poor-Quality

1230    are included in the "Other" column. P-values are from a Ranksum test of the selected cGEP

1231    against the Other cGEPs. (D) Same as (C) but showing correlation with the percentage of UMIs

1232    assigned to protein coding genes. (E) Scatter plot of the proportion of UMIs mapping to

1233    intergenic regions in the genome against Poor-Quality cGEP usage for cells in the AMP-RA

1234    dataset. (F) Correlation of per-sample average cGEP usage in T-cells with that in B-cells, NK-

1235    cells for the 3 immediate early gene cGEPs, in the COMBAT, UK-Covid, and HIV-Vaccine

1236    datasets.

1237

1239 **Figure S3. Benchmarking CellAnnoTator on simulated and real datasets.** (A) Manual

1240 gating for the Flu-Vaccine dataset analogous to **Figure S2A**. (B) Receiver operator curves

1241 (ROCs) for prediction of manually gated subset based on a single most associated subset (dark

1242 blue), TCAT multilabel prediction (light blue), analogous predictions using the single most

1243 associated NMF component published in Yasumizu et al., 2024[13], or using gene sets from NMF

1244 components in Gavish et al., 2023[16]. Individual points show accuracies of discrete predictions

1245 based on cGEP multilabel regression, or clustering with the leiden resolution specified in the

1246 legend. (C) Areas under the curve (AUC) from receiver operator curves in (B). (D) Heatmap of

1247 pseudobulk expression in Th1-Like-high and low cells, per sample. Cytotoxic-high cells are

1248 included (left) and filtered (right). Sample expression is normalized by library size and z-scored

1249 across rows, separately for the two filtering conditions.

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1264 **Figure S4. Comparing TCAT with COMBAT dataset clustering.** (A) Fraction of proliferating

1265 cells (cell cycle usage>0.1) assigned to each subset based on the most highly used subset-

1266 associated GEPs, for cells from Covid-19 or healthy donors in the two Covid-19 datasets. Error

1267 bars represent 95% bootstrap confidence intervals. (B) Usage of selected cGEPs (columns) in

1268 cells (rows) grouped by maximum subset cGEP. Cells are drawn from subclusters with high

1269 usage of the ISG cGEP, indicated in the colorbar. (C) Same as (B) but only showing cells from

1270 subclusters with high cytotoxicity cGEP usage. (D) Heatmap of pseudobulk expression of

1271 marker genes in cytotoxic-high and low cells and subset cGEP high and low cells, per sample.

1272 Expression is normalized by library size and z-scored across rows. (E) Average fraction of

1273 polarized cells (usage>0.1) per gated subset, across samples, within COMBAT and Flu-Vaccine

1274 datasets. (F) Heatmap of pseudobulk expression of marker genes in polarization-high and low

1275 cells, separately for gated CD4 and CD8s T-cells, per sample. Sample expression is normalized

1276 by library size and z-scored across rows, for each polarization.

1277
1278

1280    **Figure S5. Identifying activation associated cGEPs with AIM-Seq.** (A-B) Flow cytometry

1281    data of CD3+CD4+ and CD3+CD4- gated populations for 3 donor samples for CEFX and mock

1282    conditions. (C-E) Activation-induced marker (AIM) surface protein expression based on CITE-

1283    Seq for CD4+, CD8+, and Treg subsets, stratified by sort condition. Boxes represent

1284    interquartile range and whiskers represent 95% percentiles. (F) Percentage of each sample

1285    assigned to each subset based on manual gating, colored by stimulation condition. * indicates t-

1286    test P<.05 comparing + and U. (G) Average cGEP usage in each donor and condition, for AIM-

1287    associated cGEPs. (H) Paired t-test of pseudobulk cGEP usage in high and low cell cycle usage

1288    cells (threshold 0.1) from each sample. X-axis shows the mean $Log_2$ ratio of average usages

1289    across datasets. Y-axis shows the -$Log_{10}$ P-value. Statistically significant and positively

1290    associated cGEPs are indicated in red.

1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310

1312 **Figure S6. Annotating antigen-specific activation *in vivo*.** (A) Definition of activation used for

1313 training the antigen-specific activation (ASA) score in the COMBAT dataset for manually gated

1314 subsets. (B) AUC estimates averaged for predicting CD71/CD95 co-expression based on

1315 summation of cGEPs sequentially added to the score from left to right. (C-D) Receiver operator

1316 curve (ROC) for ASA prediction of CD71/CD95-based activation labels, with various thresholds

1317 denoted as colored points. (E) ROC for ASA prediction of AIM-positivity in the AIM-Seq dataset.

1318 (F) Left - Odds ratio of enrichment between proliferation (aggregate cell cycle cGEP usage>0.1)

1319 and activation (ASA>0.065) for each dataset. Error bars denote 95% confidence intervals. Right

1320 - Pearson correlation between ASA and aggregate cell cycle cGEP usage with colors mapping

1321 to dataset. (G) Clonality in manually gated conventional CD4 and CD8 T-cells annotated as

1322 activated (ASA>0.065) or not activated (ASA<0.065). Clonality is defined as the number of cells

1323 in the same sample with an identical alpha and beta CDR3 amino acid sequence. (H-J)

1324 Percentage of activated CD4 convs, CD8 convs, and Tregs based on ASA>0.065 in Covid-19

1325 and healthy control samples from COMBAT and UK-Covid datasets. (K) Percentage of activated

1326 conventional CD4 T-cells (ASA>0.065) versus percentage of activated or exhausted

1327 (exhaustion usage>0.065) conventional CD8 T-cells across tumor samples. (L) Percentage of

1328 activated, exhausted, or bystander (ASA + exhaustion usage<0.065) Tregs in tumors and match

1329 normal samples.
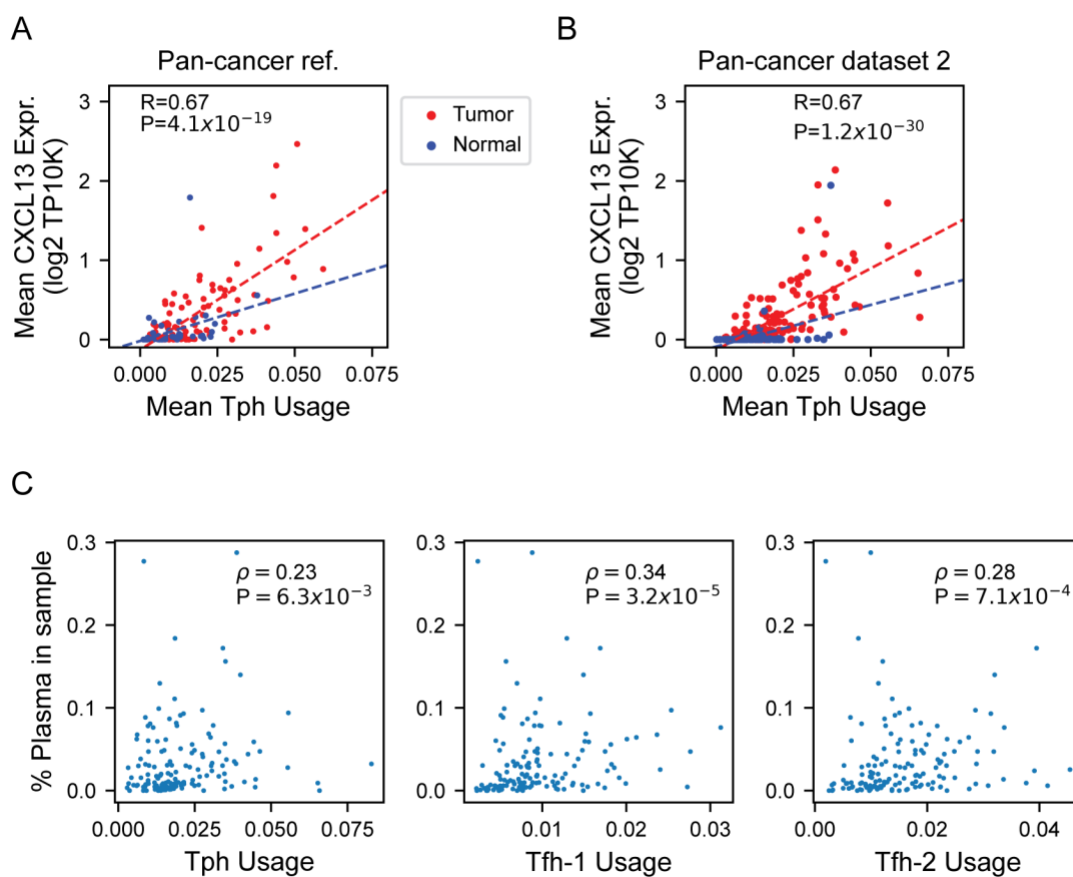
1330
1331
1332
1333
1334
1335
1336
1337
1338

**Figure S7. Identifying cGEPs associated with disease phenotypes.** (A-B) Average usage of the T peripheral helper (Tph) cGEP compared to average *CXCL13* expression from T-cells within tumors and matched normal tissue samples in Pan-cancer reference and Luo et al., 2022[62]. Trend lines show the regression coefficients fit for tumors and normal samples separately (D-F) Percentage of cells annotated as plasma cells against the average Tph, Tfh-1, or Tfh-2 usage within T-cells from tumor samples in Luo et al., 2022[62].

1350

# Supplementary Item Figures / Legends

1351

1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370

1372 **Supplementary item 1. K selection plots for consensus NMF runs on reference datasets.**

1373 Vertical line denotes the selected number of components.

cGEP Marker Genes

1375

1376

1377 **Supplementary item 2. Example marker genes for all cGEPs.** Color indicates average cNMF

1378 gene score units which denotes how much 1 additional count of usage of the cGEP would be

1379 expected to increase expression of the gene in Z-scored units.

1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393

1394
1395 **Supplementary item 3. Immediate early gene usage across circulating blood cell types**.

1396 Average per-sample usage of each IEG cGEP in T-cells versus monocytes and dendritic cells,

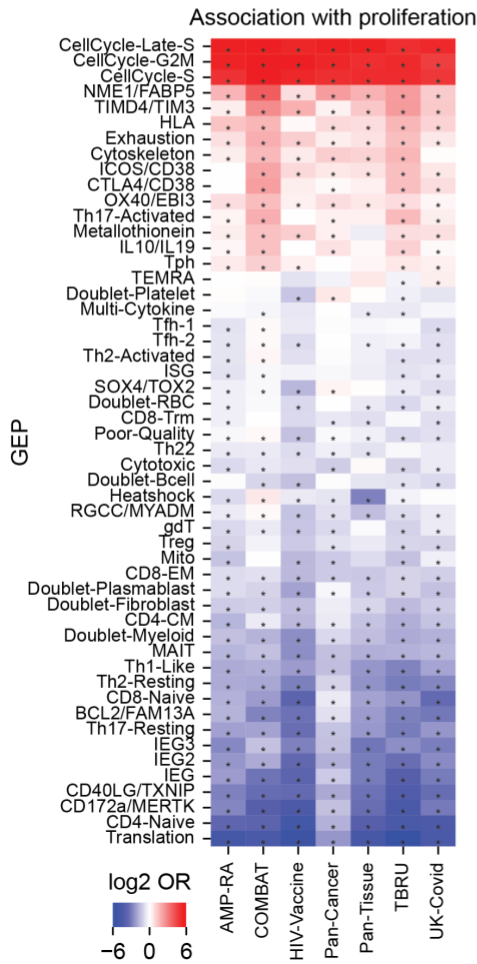1397 NK cells, or B-cells, in the three reference PBMC datasets.

1398
1399
1400
1401
1402
1403

1404
1405
1406 **Supplementary item 4. Characterization of COMBAT dataset clustering.** (A) Violin plot for

1407 myeloid cell marker genes in cells originally annotated as CD4 memory T-cells broken out by

1408 the CD4.TEFF.prolif.MKI67lo subcluster, or all other subclusters combined. (B) Usage of the

1409 ISG, Cytotoxic, and Poor-quality cGEPs in cells stratified by their CD4 memory subcluster. (C)

1410 Expression of CD4 naive marker genes in cells initially clustered as CD4 memories (blue and

1411 orange boxes) or CD4 naives (green cluster). Cells initially clustered as CD4 memory are

1412 stratified by their usage of the CD4 naive cGEP with a threshold of 0.1.
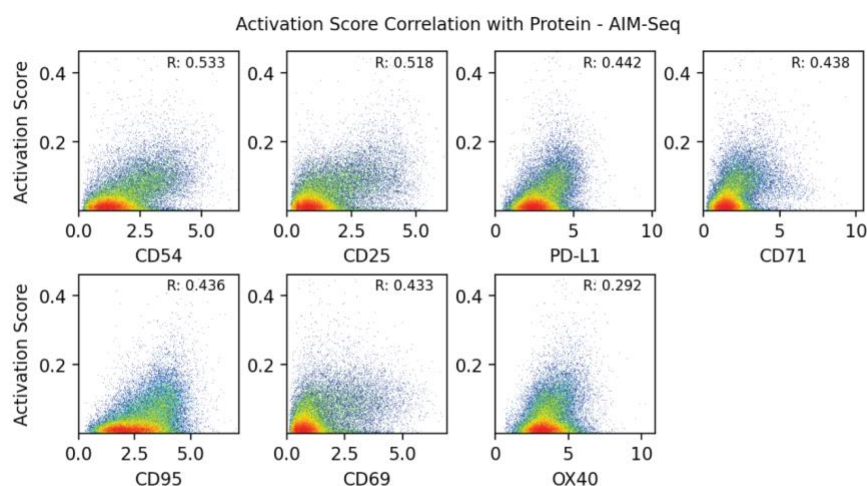
1413
1414

1415

1416

1418 **Supplementary item 5. cGEP associations with proliferation across datasets.** (A) Gating

1419 strategy to identify CD3+ CD4+ and CD3+ CD4- populations in the AIM-Seq experiment. (B)

1420 Heatmap of the average Log2 ratio of mean usage in proliferating cells (usage>0.1 of

1421 proliferation GEPs) and non-proliferating cells (usage<0.1) for all GEPs (rows) and datasets

1422 (columns). An absolute value ceiling of 6 is used to aid visualization. * indicates paired t-test

1423 P<.05.

1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443

Activation Score Correlation with Protein - AIM-Seq

1444
1445    **Supplementary item 6.** Antigen-specific activation (ASA) score correlation with surface protein

1446    activation markers in the AIM-Seq dataset.

1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466

# References

1.  Raphael, I., Nalawade, S., Eagar, T.N., and Forsthuber, T.G. (2015). T cell subsets and their signature cytokines in autoimmune and inflammatory diseases. Cytokine *74*, 5–17.

2.  Kiner, E., Willie, E., Vijaykumar, B., Chowdhary, K., Schmutz, H., Chandler, J., Schnell, A., Thakore, P.I., LeGros, G., Mostafavi, S., et al. (2021). Gut CD4+ T cell phenotypes are a continuum molded by microbes, not by TH archetypes. Nat. Immunol. *22*, 216–228.

3.  Eizenberg-Magar, I., Rimer, J., Zaretsky, I., Lara-Astiaso, D., Reich-Zeliger, S., and Friedman, N. (2017). Diverse continuum of CD4+ T-cell states is determined by hierarchical additive integration of cytokine signals. Proc. Natl. Acad. Sci. U. S. A. *114*, E6447–E6456.

4.  DuPage, M., and Bluestone, J.A. (2016). Harnessing the plasticity of CD4(+) T cells to treat immune-mediated disease. Nat. Rev. Immunol. *16*, 149–163.

5.  Szabo, P.A., Levitin, H.M., Miron, M., Snyder, M.E., Senda, T., Yuan, J., Cheng, Y.L., Bush, E.C., Dogra, P., Thapa, P., et al. (2019). Single-cell transcriptomics of human T cells reveals tissue and activation signatures in health and disease. Nat. Commun. *10*, 4706.

6.  Zheng, L., Qin, S., Si, W., Wang, A., Xing, B., Gao, R., Ren, X., Wang, L., Wu, X., Zhang, J., et al. (2021). Pan-cancer single-cell landscape of tumor-infiltrating T cells. Science *374*, abe6474.

7.  Stoeckius, M., Hafemeister, C., Stephenson, W., Houck-Loomis, B., Chattopadhyay, P.K., Swerdlow, H., Satija, R., and Smibert, P. (2017). Simultaneous epitope and transcriptome measurement in single cells. Nat. Methods *14*, 865–868.

8.  Integrated analysis of multimodal single-cell data (2021). Cell *184*, 3573–3587.e29.

9.  Takeuchi, A., and Saito, T. (2017). CD4 CTL, a Cytotoxic Subset of CD4+ T Cells, Their Differentiation and Function. Front. Immunol. *8*, 194.

10. Li, J., Zaslavsky, M., Su, Y., Guo, J., Sikora, M.J., van Unen, V., Christophersen, A., Chiou, S.-H., Chen, L., Li, J., et al. (2022). KIR+CD8+ T cells suppress pathogenic T cells and are active in autoimmune diseases and COVID-19. Science *376*, eabi9591.

11. Nathan, A., Beynor, J.I., Baglaenko, Y., Suliman, S., Ishigaki, K., Asgari, S., Huang, C.-C., Luo, Y., Zhang, Z., Lopez, K., et al. (2021). Multimodally profiling memory T cells from a tuberculosis cohort identifies cell state associations with demographics, environment and disease. Nat. Immunol. *22*, 781–793.

12. Wagner, A., Regev, A., and Yosef, N. (2016). Revealing the vectors of cellular identity with single-cell genomics. Nat. Biotechnol. *34*, 1145–1160.

13. Yasumizu, Y., Takeuchi, D., Morimoto, R., Takeshima, Y., Okuno, T., Kinoshita, M., Morita, T., Kato, Y., Wang, M., Motooka, D., et al. (2024). Single-cell transcriptome landscape of circulating CD4+ T cell populations in autoimmune diseases. Cell Genomics, 100473.

14. Kotliar, D., Veres, A., Nagy, M.A., Tabrizi, S., Hodis, E., Melton, D.A., and Sabeti, P.C.

1503   (2019). Identifying gene expression programs of cell-type identity and cellular activity with
1504   single-cell RNA-Seq. Elife *8*. 10.7554/eLife.43803.

1505   15. Kunes, R.Z., Walle, T., Land, M., Nawy, T., and Pe'er, D. (2023). Supervised discovery of
1506   interpretable gene programs from single-cell data. Nat. Biotechnol. 10.1038/s41587-023-
1507   01940-3.

1508   16. Gavish, A., Tyler, M., Greenwald, A.C., Hoefflin, R., Simkin, D., Tschernichovsky, R., Galili
1509   Darnell, N., Somech, E., Barbolin, C., Antman, T., et al. (2023). Hallmarks of transcriptional
1510   intratumour heterogeneity across a thousand tumours. Nature *618*, 598–606.

1511   17. Stephenson, E., Reynolds, G., Botting, R.A., Calero-Nieto, F.J., Morgan, M.D., Tuong, Z.K.,
1512   Bach, K., Sungnak, W., Worlock, K.B., Yoshida, M., et al. (2021). Single-cell multi-omics
1513   analysis of the immune response in COVID-19. Nat. Med. *27*, 904–916.

1514   18. COvid-19 Multi-omics Blood ATlas (COMBAT) Consortium. Electronic address:
1515   julian.knight@well.ox.ac.uk, and COvid-19 Multi-omics Blood ATlas (COMBAT) Consortium
1516   (2022). A blood atlas of COVID-19 defines hallmarks of disease severity and specificity.
1517   Cell *185*, 916–938.e58.

1518   19. Domínguez Conde, C., Xu, C., Jarvis, L.B., Rainbow, D.B., Wells, S.B., Gomes, T., Howlett,
1519   S.K., Suchanek, O., Polanski, K., King, H.W., et al. (2022). Cross-tissue immune cell
1520   analysis reveals tissue-specific features in humans. Science *376*, eabl5197.

1521   20. Zhang, F., Jonsson, A.H., Nathan, A., Millard, N., Curtis, M., Xiao, Q., Gutierrez-Arcelus,
1522   M., Apruzzese, W., Watts, G.F.M., Weisenfeld, D., et al. (2023). Deconstruction of
1523   rheumatoid arthritis synovium defines inflammatory subtypes. Nature *623*, 616–624.

1524   21. Korsunsky, I., Millard, N., Fan, J., Slowikowski, K., Zhang, F., Wei, K., Baglaenko, Y.,
1525   Brenner, M., Loh, P.-R., and Raychaudhuri, S. (2019). Fast, sensitive and accurate
1526   integration of single-cell data with Harmony. Nat. Methods *16*, 1289–1296.

1527   22. Rao, D.A., Gurish, M.F., Marshall, J.L., Slowikowski, K., Fonseka, C.Y., Liu, Y., Donlin,
1528   L.T., Henderson, L.A., Wei, K., Mizoguchi, F., et al. (2017). Pathologically expanded
1529   peripheral T helper cell subset drives B cells in rheumatoid arthritis. Nature *542*, 110–114.

1530   23. Gene Ontology Consortium, Aleksander, S.A., Balhoff, J., Carbon, S., Cherry, J.M.,
1531   Drabkin, H.J., Ebert, D., Feuermann, M., Gaudet, P., Harris, N.L., et al. (2023). The Gene
1532   Ontology knowledgebase in 2023. Genetics *224*. 10.1093/genetics/iyad031.

1533   24. Cano-Gamez, E., Soskic, B., Roumeliotis, T.I., So, E., Smyth, D.J., Baldrighi, M., Willé, D.,
1534   Nakic, N., Esparza-Gordillo, J., Larminie, C.G.C., et al. (2020). Single-cell transcriptomics
1535   identifies an effectorness gradient shaping the response of CD4+ T cells to cytokines. Nat.
1536   Commun. *11*, 1801.

1537   25. Ilicic, T., Kim, J.K., Kolodziejczyk, A.A., Bagger, F.O., McCarthy, D.J., Marioni, J.C., and
1538   Teichmann, S.A. (2016). Classification of low quality cells from single-cell RNA-seq data.
1539   Genome Biol. *17*, 29.

1540   26. Osorio, D., and Cai, J.J. (2021). Systematic determination of the mitochondrial proportion in
1541   human and mice tissues for single-cell RNA-sequencing data quality control. Bioinformatics
1542   *37*, 963–967.

1543 27. Wong, W.K., Jiang, G., Sørensen, A.E., Chew, Y.V., Lee-Maynard, C., Liuwantara, D., Williams, L., O'Connell, P.J., Dalgaard, L.T., Ma, R.C., et al. (2019). The long noncoding RNA MALAT1 predicts human pancreatic islet isolation quality. JCI Insight *5*. 10.1172/jci.insight.129299.

1547 28. Arner, E., Daub, C.O., Vitting-Seerup, K., Andersson, R., Lilje, B., Drabløs, F., Lennartsson, A., Rönnerblad, M., Hrydziuszko, O., Vitezic, M., et al. (2015). Transcribed enhancers lead waves of coordinated transcription in transitioning mammalian cells. Science *347*, 1010–1014.

1551 29. Vacca, A., Itoh, M., Kawaji, H., Arner, E., Lassmann, T., Daub, C.O., Carninci, P., Forrest, A.R.R., Hayashizaki, Y., FANTOM Consortium, et al. (2018). Conserved temporal ordering of promoter activation implicates common mechanisms governing the immediate early response across cell types and stimuli. Open Biol. *8*. 10.1098/rsob.180011.

1555 30. Bahrami, S., and Drabløs, F. (2016). Gene regulation in the immediate-early response process. Adv. Biol. Regul. *62*, 37–49.

1557 31. van den Brink, S.C., Sage, F., Vértesy, Á., Spanjaard, B., Peterson-Maduro, J., Baron, C.S., Robin, C., and van Oudenaarden, A. (2017). Single-cell sequencing reveals dissociation-induced gene expression in tissue subpopulations. Nat. Methods *14*, 935–936.

1560 32. Lacar, B., Linker, S.B., Jaeger, B.N., Krishnaswami, S.R., Barron, J.J., Kelder, M.J.E., Parylak, S.L., Paquola, A.C.M., Venepally, P., Novotny, M., et al. (2016). Nuclear RNA-seq of single neurons reveals molecular signatures of activation. Nat. Commun. *7*, 11022.

1563 33. Sparks, R., Lau, W.W., Liu, C., Han, K.L., Vrindten, K.L., Sun, G., Cox, M., Andrews, S.F., Bansal, N., Failla, L.E., et al. (2023). Influenza vaccination reveals sex dimorphic imprints of prior mild COVID-19. Nature *614*, 752–761.

1566 34. Tirosh, I., Izar, B., Prakadan, S.M., Wadsworth, M.H., 2nd, Treacy, D., Trombetta, J.J., Rotem, A., Rodman, C., Lian, C., Murphy, G., et al. (2016). Dissecting the multicellular ecosystem of metastatic melanoma by single-cell RNA-seq. Science *352*, 189–196.

1569 35. Wolf, F.A., Angerer, P., and Theis, F.J. (2018). SCANPY: large-scale single-cell gene expression data analysis. Genome Biol. *19*, 15.

1571 36. Szabo, S.J., Sullivan, B.M., Stemmann, C., Satoskar, A.R., Sleckman, B.P., and Glimcher, L.H. (2002). Distinct effects of T-bet in TH1 lineage commitment and IFN-gamma production in CD4 and CD8 T cells. Science *295*, 338–342.

1574 37. McLane, L.M., Banerjee, P.P., Cosma, G.L., Makedonas, G., Wherry, E.J., Orange, J.S., and Betts, M.R. (2013). Differential localization of T-bet and Eomes in CD8 T cell memory populations. J. Immunol. *190*, 3207–3215.

1577 38. Buettner, F., Natarajan, K.N., Casale, F.P., Proserpio, V., Scialdone, A., Theis, F.J., Teichmann, S.A., Marioni, J.C., and Stegle, O. (2015). Computational analysis of cell-to-cell heterogeneity in single-cell RNA-sequencing data reveals hidden subpopulations of cells. Nat. Biotechnol. *33*, 155–160.

1581 39. Asashima, H., Mohanty, S., Comi, M., Ruff, W.E., Hoehn, K.B., Wong, P., Klein, J., Lucas, C., Cohen, I., Coffey, S., et al. (2023). PD-1highCXCR5-CD4+ peripheral helper T cells

1583        promote CXCR3+ plasmablasts in human acute viral infection. Cell Rep. *42*, 111895.

1584   40. Koh, C.-H., Lee, S., Kwak, M., Kim, B.-S., and Chung, Y. (2023). CD8 T-cell subsets:
1585        heterogeneity, functions, and therapeutic potential. Exp. Mol. Med. *55*, 2287–2299.

1586   41. Lehmann, A.A., Reche, P.A., Zhang, T., Suwansaard, M., and Lehmann, P.V. (2021).
1587        CERI, CEFX, and CPI: Largely Improved Positive Controls for Testing Antigen-Specific T
1588        Cell Function in PBMC Compared to CEF. Cells *10*. 10.3390/cells10020248.

1589   42. Reiss, S., Baxter, A.E., Cirelli, K.M., Dan, J.M., Morou, A., Daigneault, A., Brassard, N.,
1590        Silvestri, G., Routy, J.-P., Havenar-Daughton, C., et al. (2017). Comparative analysis of
1591        activation induced marker (AIM) assays for sensitive identification of antigen-specific CD4 T
1592        cells. PLoS One *12*, e0186998.

1593   43. Wolfl, M., Kuball, J., Ho, W.Y., Nguyen, H., Manley, T.J., Bleakley, M., and Greenberg, P.D.
1594        (2007). Activation-induced expression of CD137 permits detection, isolation, and expansion
1595        of the full repertoire of CD8+ T cells responding to antigen without requiring knowledge of
1596        epitope specificities. Blood *110*, 201–210.

1597   44. Subramanian Vignesh, K., and Deepe, G.S., Jr (2017). Metallothioneins: Emerging
1598        Modulators in Immunity and Infection. Int. J. Mol. Sci. *18*. 10.3390/ijms18102197.

1599   45. Boyman, O., and Sprent, J. (2012). The role of interleukin-2 during homeostasis and
1600        activation of the immune system. Nat. Rev. Immunol. *12*, 180–190.

1601   46. Billadeau, D.D., Nolz, J.C., and Gomez, T.S. (2007). Regulation of T-cell activation by the
1602        cytoskeleton. Nat. Rev. Immunol. *7*, 131–143.

1603   47. Scarneo, S.A., Smith, A.P., Favret, J., O'Connell, R., Pickeral, J., Yang, K.W., Ferrari, G.,
1604        Loiselle, D.R., Hughes, P.F., Kulkarni, M.M., et al. (2022). Expression of membrane Hsp90
1605        is a molecular signature of T cell activation. Sci. Rep. *12*, 18091.

1606   48. Di Conza, G., Ho, P.-C., Cubillos-Ruiz, J.R., and Huang, S.C.-C. (2023). Control of immune
1607        cell function by the unfolded protein response. Nat. Rev. Immunol. *23*, 546–562.

1608   49. Holling, T.M., van der Stoep, N., Quinten, E., and van den Elsen, P.J. (2002). Activated
1609        human T cells accomplish MHC class II expression through T cell-specific occupation of
1610        class II transactivator promoter III. J. Immunol. *168*, 763–770.

1611   50. Motamedi, M., Xu, L., and Elahi, S. (2016). Correlation of transferrin receptor (CD71) with
1612        Ki67 expression on stimulated human and mouse T cells: The kinetics of expression of T
1613        cell activation markers. J. Immunol. Methods *437*, 43–52.

1614   51. Paulsen, M., and Janssen, O. (2011). Pro- and anti-apoptotic CD95 signaling in T cells. Cell
1615        Commun. Signal. *9*, 7.

1616   52. Meyer Zu Horste, G., Przybylski, D., Schramm, M.A., Wang, C., Schnell, A., Lee, Y., Sobel,
1617        R., Regev, A., and Kuchroo, V.K. (2018). Fas Promotes T Helper 17 Cell Differentiation and
1618        Inhibits T Helper 1 Cell Development by Binding and Sequestering Transcription Factor
1619        STAT1. Immunity *48*, 556–569.e7.

1620   53. Flores-Mendoza, G., Rodríguez-Rodríguez, N., Rubio, R.M., Madera-Salcedo, I.K., Rosetti,
1621        F., and Crispín, J.C. (2021). Fas/FasL Signaling Regulates CD8 Expression During

1622        Exposure to Self-Antigens. Front. Immunol. *12*, 635862.

1623   54.  Miggelbrink, A.M., Jackson, J.D., Lorrey, S.J., Srinivasan, E.S., Waibl-Polania, J.,
1624        Wilkinson, D.S., and Fecci, P.E. (2021). CD4 T-Cell Exhaustion: Does It Exist and What
1625        Are Its Roles in Cancer? Clin. Cancer Res. *27*, 5742–5752.

1626   55.  Sun, J., Li, L., Li, L., Ding, L., Liu, X., Chen, X., Zhang, J., Qi, X., Du, J., and Huang, Z.
1627        (2018). Metallothionein-1 suppresses rheumatoid arthritis pathogenesis by shifting the
1628        Th17/Treg balance. Eur. J. Immunol. *48*, 1550–1562.

1629   56.  Miyazaki, Y., Nakayamada, S., Kubo, S., Nakano, K., Iwata, S., Miyagawa, I., Ma, X.,
1630        Trimova, G., Sakata, K., and Tanaka, Y. (2018). Th22 Cells Promote Osteoclast
1631        Differentiation via Production of IL-22 in Rheumatoid Arthritis. Front. Immunol. *9*, 2901.

1632   57.  Nishikawa, H., and Sakaguchi, S. (2014). Regulatory T cells in cancer immunotherapy.
1633        Curr. Opin. Immunol. *27*, 1–7.

1634   58.  Thommen, D.S., and Schumacher, T.N. (2018). T Cell Dysfunction in Cancer. Cancer Cell
1635        *33*, 547–562.

1636   59.  Jorgovanovic, D., Song, M., Wang, L., and Zhang, Y. (2020). Roles of IFN-γ in tumor
1637        progression and regression: a review. Biomark Res *8*, 49.

1638   60.  Yoshitomi, H., and Ueno, H. (2021). Shared and distinct roles of T peripheral helper and T
1639        follicular helper cells in human diseases. Cell. Mol. Immunol. *18*, 523–527.

1640   61.  Gu-Trantien, C., Migliori, E., Buisseret, L., de Wind, A., Brohée, S., Garaud, S., Noël, G.,
1641        Dang Chi, V.L., Lodewyckx, J.-N., Naveaux, C., et al. (2017). CXCL13-producing TFH cells
1642        link immune suppression and adaptive memory in human breast cancer. JCI Insight *2*.
1643        10.1172/jci.insight.91487.

1644   62.  Luo, H., Xia, X., Huang, L.-B., An, H., Cao, M., Kim, G.D., Chen, H.-N., Zhang, W.-H., Shu,
1645        Y., Kong, X., et al. (2022). Pan-cancer single-cell analysis reveals the heterogeneity and
1646        plasticity of cancer-associated fibroblasts in the tumor microenvironment. Nat. Commun.
1647        *13*, 6619.

1648   63.  Garaud, S., Buisseret, L., Solinas, C., Gu-Trantien, C., de Wind, A., Van den Eynden, G.,
1649        Naveaux, C., Lodewyckx, J.-N., Boisson, A., Duvillier, H., et al. (2019). Tumor infiltrating B-
1650        cells signal functional humoral immune responses in breast cancer. JCI Insight *5*.
1651        10.1172/jci.insight.129641.

1652   64.  Subramanian, A., Tamayo, P., Mootha, V.K., Mukherjee, S., Ebert, B.L., Gillette, M.A.,
1653        Paulovich, A., Pomeroy, S.L., Golub, T.R., Lander, E.S., et al. (2005). Gene set enrichment
1654        analysis: a knowledge-based approach for interpreting genome-wide expression profiles.
1655        Proc. Natl. Acad. Sci. U. S. A. *102*, 15545–15550.

1656   65.  Liberzon, A., Subramanian, A., Pinchback, R., Thorvaldsdóttir, H., Tamayo, P., and
1657        Mesirov, J.P. (2011). Molecular signatures database (MSigDB) 3.0. Bioinformatics *27*,
1658        1739–1740.

1659   66.  Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,
1660        Prettenhofer, P., Weiss, R., Dubourg, V., et al. (2011). Scikit-learn: Machine Learning in

1661        Python. J. Mach. Learn. Res. *12*, 2825–2830.

1662    67. Zappia, L., Phipson, B., and Oshlack, A. (2017). Splatter: simulation of single-cell RNA
1663        sequencing data. Genome Biol. *18*, 174.

1664    68. Ashburner, M., Ball, C.A., Blake, J.A., Botstein, D., Butler, H., Cherry, J.M., Davis, A.P.,
1665        Dolinski, K., Dwight, S.S., Eppig, J.T., et al. (2000). Gene ontology: tool for the unification of
1666        biology. The Gene Ontology Consortium. Nat. Genet. *25*, 25–29.

1667    69. van Dijk, D., Sharma, R., Nainys, J., Yim, K., Kathail, P., Carr, A.J., Burdziak, C., Moon,
1668        K.R., Chaffer, C.L., Pattabiraman, D., et al. (2018). Recovering Gene Interactions from
1669        Single-Cell Data Using Data Diffusion. Cell *174*, 716–729.e27.

1670    70. Traag, V.A., Waltman, L., and van Eck, N.J. (2019). From Louvain to Leiden: guaranteeing
1671        well-connected communities. Sci. Rep. *9*, 1–12.

1672    71. Kang, H.M., Subramaniam, M., Targ, S., Nguyen, M., Maliskova, L., McCarthy, E., Wan, E.,
1673        Wong, S., Byrnes, L., Lanata, C.M., et al. (2018). Multiplexed droplet single-cell RNA-
1674        sequencing using natural genetic variation. Nat. Biotechnol. *36*, 89–94.