

1 **APMAT analysis reveals the association between CD8 T cell receptors, cognate**
2 **antigen, and T cell phenotype and persistence**

3
4 Jingyi Xie^{1,2}, Daniel G. Chen^{1,3}, William Chour¹, Rachel H. Ng^{1,4}, Rongyu Zhang^{1,4}, Dan
5 Yuan^{1,4}, Jongchan Choi¹, Michaela McKasson¹, Pamela Troisch¹, Brett Smith¹, Lesley
6 Jones¹, Andrew Webster¹, Yusuf Rasheed¹, Sarah Li¹, Rick Edmark¹, Sunga Hong¹, Kim
7 M. Murray¹, Jennifer K. Logue⁵, Nicholas M. Franko⁵, Christopher G. Lausted¹, Brian
8 Piening⁶, Heather Algren^{6,7}, Julie Wallick^{6,7}, Andrew T. Magis¹, Kino Watanabe⁵, Phil
9 Mease⁷, Philip D. Greenberg^{3,5,8}, Helen Chu⁵, Jason D. Goldman^{3,5,6,7}, Yapeng Su^{1,3,9},
10 James R. Heath^{1,9,10*}

11
12 ¹ Institute of Systems Biology, Seattle, WA, 98109, USA
13 ² Molecular Engineering & Sciences Institute, University of Washington, Seattle, WA,
14 98195, USA
15 ³ Fred Hutchinson Cancer Center, Seattle, WA, 98109, USA
16 ⁴ Department of Bioengineering, University of Washington, Seattle, WA, 98195, USA
17 ⁵ Department of Medicine, Allergy and Infectious Diseases, University of Washington,
18 Seattle, WA 98109, USA
19 ⁶ Providence Health & Services, Seattle, WA, 99109, USA
20 ⁷ Providence Swedish Medical Center, Seattle, WA, 98122, USA
21 ⁸ Department of Immunology, University of Washington, Seattle, WA 98109, USA
22 ⁹ These authors jointly-supervised the work
23 ¹⁰ Corresponding author, Leading contact.
24 *Correspondence: jheath@isbscience.org
25

26 **Abstract**

27 Elucidating the relationships between a class I peptide antigen, a CD8 T cell receptor
28 (TCR) specific to that antigen, and the T cell phenotype that emerges following antigen
29 stimulation, remains a mostly unsolved problem, largely due to the lack of large data sets
30 that can be mined to resolve such relationships. Here, we describe Antigen-TCR Pairing
31 and Multiomic Analysis of T-cells (APMAT), an integrated experimental-computational
32 framework designed for the high-throughput capture and analysis of CD8 T cells, with
33 paired antigen, TCR sequence, and single-cell transcriptome. Starting with 951 putative
34 antigens representing a comprehensive survey of the SARS-CoV-2 viral proteome, we
35 utilize APMAT for the capture and single cell analysis of CD8 T cells from 62 HLA A*02:01
36 COVID-19 participants. We leverage this unique, comprehensive dataset to integrate with
37 peptide antigen properties, TCR CDR3 sequences, and T cell phenotypes to show that
38 distinct physicochemical features of the antigen-TCR pairs strongly associate with both T
39 cell phenotype and T cell persistence. This analysis suggests that CD8+ T cell phenotype
40 following antigen stimulation is at least partially deterministic, rather than the result of
41 stochastic biological properties.

42

43

44 Introduction

45 CD8 T cells play a major role in adaptive immunity against pathogens, exhibiting a
46 functional diversity that includes, as major phenotypes, naïve, memory, effector memory,
47 effector, and exhausted, with each of those phenotypes encompassing multiple sub-
48 phenotypes^{1,2}. Recent literature has suggested the existence of relationships between a
49 given antigen-specific T cell clonotype and the phenotypic trajectory that clonotype can
50 take following antigen stimulation. For example, tissue-resident, antigen-specific T cell
51 clonotypes in both tumor and chronic viral infection settings have been intimately
52 associated with specific phenotype differentiation trajectories³⁻⁵. We and others have
53 shown that, for at least certain immunogenic viral antigens, T cell clonotype is the
54 dominant factor in determining T cell phenotype⁶⁻⁸. These similar results, from tissue and
55 peripheral blood, and in very different disease settings, suggest that there may be rules
56 that underlie T cell clonotype-T cell phenotype relationships. Elucidating such
57 relationships requires experimental methods for collecting a large, longitudinal data set in
58 which the transcriptome-level phenotypes of antigen-specific T cells are co-measured
59 across a large number and diversity of antigens, coupled with computational methods for
60 elucidating what, if any, relationships between the T cell receptor (TCR) α/β genes,
61 cognate peptide antigen – major histocompatibility complex (pMHC), and T cell
62 phenotype exist.

63 Initial hints at what such a rule set might contain can be found in the literature. For
64 example, the importance of hydrophobicity at certain TCR residues is known to associate
65 with cytotoxicity and self-reactivity in CD4 T cells^{6,9-11}, and the hydrophobicity of certain
66 residues is known to associate with the immunogenicity of antigens¹². These observations
67 suggest that the biochemical nature of the TCR:pMHC interface may play an important
68 roles in determining T cell phenotype and cell fate trajectory. This question is not just of
69 fundamental importance to T cell immunology, but it is also highly relevant to the
70 engineering of T cells for use as therapeutic agents for treatments of both cancers¹³ and
71 autoimmune diseases^{14,15}. However, existing studies predominantly explore only one
72 aspect of the Antigen-TCR-phenotype interplay: either resolving the antigen specificity of
73 T cells from biological samples¹⁶⁻¹⁸, or conducting T cell receptor analysis on existing
74 dataset lacking antigen specificity (or containing only limited number of antigens)^{6,19,20}.
75 Recent advances highlight the need of high-throughput, high-dimensional approaches as
76 powerful tools for identifying and analyzing antigen-specific CD8 T cells^{17,20-22}.

77 We describe an integrated experimental-computational framework, APMAT for the
78 antigen-specific capture of T cells from many patient biospecimens in parallel, and is
79 integrated with single-cell multi-omics profiling designed to integrate paired antigen and
80 TCR sequence with T cell phenotype. For antigen-specific T cell capture, we utilize large
81 libraries of single-chain-trimers (SCTs)^{22,23}. We have recently reported on the feasibility
82 of these libraries to capture and characterize both viral antigen-specific and (tumor-
83 associated) neoantigen-specific T cell populations. Here we use SCTs to construct a 951-

84 element pMHC library that represents a comprehensive survey of putative Class I
85 antigens, presented by HLA A*02:01, from across the full SARS-CoV-2 genome. The
86 DNA-barcoded SCT-pMHC library (for n = 558 expressed SCT) is used to identify and,
87 through single cell (sc) RNA-seq analysis, characterize SARS-CoV-2 specific T cells from
88 62 HLA A*02:01 participants at the stages of acute and post-acute COVID-19. We identify
89 several viral antigen-specific T cell populations observed in previous work by ourselves
90 and others²⁴⁻²⁷ demonstrating the feasibility of this framework on a whole viral proteome
91 scale. Moreover, resultant data set allows for an in-depth exploration to reveal insights
92 into how the physicochemical properties of the TCR-pMHC interface associate with T cell
93 phenotype and T cell persistence. This study elucidates, for a whole viral proteome and
94 a single HLA allele, the physicochemical basis linking TCR:pMHC molecular interactions
95 to the phenotypic behavior of antigen-specific CD8+ T cells, and thus advances our
96 understanding of immunological mechanisms.

97

98 **Results**

99 **APMAT enables integrated multi-modal analysis of antigen-specific CD8 T cells**

100 We utilized APMAT to identify SARs-CoV-2 antigen-specific T cell responses to COVID-
101 19 infection in a previously described longitudinal cohort of 209 COVID-19
102 participants^{25,26}. From this cohort, we selected 62 HLA- A*02:01 participants representing
103 a range of COVID-19 disease severities (Fig. 1a) with longitudinal reference datasets
104 available at acute (diagnosis and approximately 1-week post-diagnosis), and post-acute
105 (2-3 months post-diagnosis) timepoints (Supplementary Table 1)^{25,26}.

106

107 For antigen-specific CD8 T cell capture, we utilized DNA-barcoded, pMHC-like SCT
108 multimers (SCT-dextramers). For library construction, we first utilized NetMHCpan²⁸ to
109 analyze the full SARS-CoV-2 viral genome (original strain) to resolve a list of 951 putative
110 9 - 10 mer antigens (Supplementary Table 2.1) for HLA-A*02:01. This peptide list was
111 converted into PCR-optimized DNA primers and used to build the plasmid library.
112 Plasmids were then transfected into Expi293 cells over four days to induce secretion of
113 the SCT protein product. 558 of the 951 putative antigens yielded usable SCT product
114 (Fig. 1b top). Each SCT was then purified, biotinylated, and assembled into a fluorophore-
115 labeled and DNA-barcoded dextramer, using previously reported protocols²³. To minimize
116 batch effects and enable integrated analysis, PBMCs from all 62 participants, collected
117 at the stage of acute disease, were pooled and profiled in parallel within one experiment
118 (Fig. 1b middle left). SCT-dextramer-positive CD8 T cells were sorted for scRNAseq to
119 profile gene expression, TCR α/β genes, and antigen specificity (Fig. 1b middle right). For
120 each cell captured, antigen-specificity was identified through the dominant pool hashtag
121 (Supplementary Fig. 1a) and dextramer (Supplementary Fig. 1b). Each cell was
122 associated with a specific patient through comparison of de novo computed single
123 nucleotide polymorphisms (SNPs) from scRNA-seq with germline whole genome

124 sequence (WGS) profiles of all patients (Supplementary Fig. 1c). Sex gene validation
125 further confirmed patient assignments (Supplementary Fig. 1d). Out of 19,103 cells from
126 scRNA-seq data, only 166 (0.87%) were not assignable to a donor (Supplementary Table
127 2.2). With such a multi-modular dataset, the characteristics of the peptide antigens and
128 the TCR CDR3 sequences can be integrated with T cell phenotype and longitudinal multi-
129 omics datasets (Fig. 1b bottom).

130

131 **APMAT enables high-throughput representation of whole SARS-CoV-2 genome**

132 We graph the distribution of putative antigens against the SARS-CoV-2 genome, and
133 highlighted the 558 (58.7%) constructs that were expressed and used for the T cell
134 capture experiment, and the 102 (18.3%) that captured at least 1 cell (Fig. 2a,
135 Supplementary Fig. 2a). While APMAT identified CD8 T cells against antigens from
136 across the SARS-CoV-2 genome, antigens from the Spike protein (S) exhibited the
137 highest rate (26.6%) of cell capture (SCTs that captured at least 1 cell) (Fig. 2b,
138 Supplementary Table 2.1). We identified immunodominant epitopes against SARS-CoV-
139 2 and common viruses across patients, confirming some previously reported epitopes
140 while discovering new ones (Supplementary Fig. 2b). We also found agreement for
141 dominant TCR gene usage for top epitopes in our dataset with reported literature
142 (Supplementary Table 2.3). For example, HLA-A*02:01/S269- YLQPRTFLL is
143 predominately recognized by TRAV12-1 containing TCRs in our dataset, consistent with
144 previous reports²⁹. Selected TCRs were validated via in-vitro Lenti-viral transduction and
145 tetramer assay to confirm antigen specificity (Supplementary Table 2.4, Supplementary
146 Fig. 3a, b).

147

148 We next investigated the antigen sequences associated with SCT expression and cell
149 capturing. In Fig 2c (left) we provide a sequence logo representation of the 9-mer SCTs
150 that were or were not expressed, and those expressed as SCTs that did or did not capture
151 cells (Fig. 2c right). We observed an enrichment of hydrophobic amino acids in non-
152 anchor residues for non-expressed SCT constructs (Fig 2c left bottom). For SCTs that
153 captured cells, polar and charged residues such as Threonine (T) and Arginine (R) were
154 enriched in non-anchor positions (Fig 2c right top). We next generated a matrix for each
155 peptide that contained the NetMHCpan prediction, the amino acid identities, and the
156 numeric physicochemical properties (such as hydrophobicity, polarity, etc) for peptide
157 anchor residues (Anchor) and residues exposed to the TCR (Exposed) (Fig. 2d,
158 Supplementary Table 3). As expected, expressed SCTs showed better NetMHCpan
159 prediction (lower prediction rank and lower predicted binding affinity) relative to non-
160 expressed SCT constructs (Fig. 2e left, Supplementary Fig. 4a). Notably, SCT expression
161 yield did not correlate with predicted affinity (Supplementary Fig. 4b). However, the
162 physicochemical properties analysis allowed for a quantitative validation of the
163 hydrophobic trends of non-expressed SCTs (Fig 2e right), and the polar/charged residue
164 trends of cell-capturing SCTs (Fig 2f right). We also show the agreement between SCT
165 expression and NetMHCPan prediction. We categorized peptides as weak-binders,
166 medium-binders, or strong-binders based on NetMHCPan prediction against the A*02:01

167 MHC molecule (See Methods). In fact, our observations that peptides that are expressed
168 as an SCT exhibit a lower average hydrophobicity closely mirror the NetMHCPan pan
169 predictions. The strong-binders indeed exhibit a relatively higher polarity and a lower
170 hydrophobicity relative to all attempted A*02:01 SCT constructs (Supplementary Fig.4c).
171 In addition, for those strong binders, the hydrophobicity and polarity of expressed vs non-
172 expressed SCTs are not significantly different (Supplementary Fig.4d). Finally, we found
173 that the likelihood that an SCT would be successfully expressed strongly agrees with
174 prediction. For weak-binders, SCT expression rate is 42.0%; For strong-binders, SCT
175 expression is elevated to 82.5% (Supplementary Fig.4e). Hence, the potential SCT-
176 expression biases match closely with the NetMHCPan predictions with respect to the
177 physicochemical properties of the putative epitopes, suggesting that little or no bias
178 originates from the SCT expression itself.

179
180 Overall, APMAT enabled a direct and comprehensive analysis of putative epitopes,
181 supporting prediction algorithms for the common HLA A*02.01 allele. The data further
182 suggests that an analysis of the physicochemical properties of the antigens (and possibly
183 their cognate TCRs) may provide insights for interpreting the large multimodal data set
184 generated through APMAT.

185
186 **Three peptide groups distinguished by sequence physicochemical properties**
187 We probed for potential relationships between the physicochemical properties of putative
188 epitopes and the antigenicity of those epitopes. We encoded the 951 peptides by residue-
189 level descriptors of their physicochemical properties (Fig. 3a, Methods). Unsupervised
190 clustering based on peptide amino acid identity and properties resulted in a two-
191 dimensional peptide Uniform Manifold Approximation and Projection (Pep-UMAP) (Fig.
192 3b, Supplementary Table 4). The upper right wing of this UMAP exhibits higher
193 hydrophobicity and bulkiness, and is dominated by non-expressed SCTs. Conversely, the
194 lower left wing displays greater polarity, and is enriched for expressed SCTs (Fig. 3c,
195 Supplementary Fig. 5a). However, SCTs that successfully captured T cells are uniformly
196 distributed across the UMAP, indicating that additional factors beyond hydrophobicity
197 modulate antigenicity.

198
199 We next utilized unsupervised clustering to resolve whether combinations of
200 physicochemical properties could be used to further classify the peptide antigens. Such
201 analysis clearly distinguished three peptide groups (Pep-Groups), PG1-3 (Fig. 3d).
202 Specifically, PG1 exhibited higher hydrophobicity yet lower charge and polarity (Fig. 3d
203 left top). PG2 and PG3 both showed higher polarity, but differ in anchor and secondary
204 anchor (Position 6) properties (Fig. 3d left bottom). When analyzing only the cell-capturing
205 SCTs, the hydrophobicity of PG1 becomes even more prominent while PG2-3 distinctions
206 diminish (Fig. 3d right). Accordingly, Pep-Groups occupied different regions on the Pep-
207 UMAP (Fig. 3e, Supplementary Fig. 5b). Comparisons of PG1-3 for hydrophobicity, SCT
208 expression rate, and cell-capture rate (Fig. 3f) reveal the most hydrophobic group (PG1)

209 is the most challenging to express, but also exhibits the highest rate for cell capture. Thus,
210 we categorized all putative antigen peptides into unique Pep-Groups for downstream
211 analysis.

212

213 **Pep-groups associate with different T cell phenotypes**

214 We next investigated whether the PG1-3 peptide groups associate with distinct CD8 T
215 cell phenotypes during the acute COVID-19 response. The SARS-CoV-2 SCT-dextramer-
216 positive CD8 T cells were filtered to only include those with an assigned patient ID,
217 antigen specificity, and paired TCR α/β sequences. These cells were then projected onto
218 a gene expression UMAP (GEX-UMAP) based on the scRNAseq data (Fig. 4a). Unbiased
219 clustering and differential gene expression analysis defined canonical CD8 T cell
220 phenotypes including naïve, central memory (CM), effector memory (EM), hybrid, and
221 cytotoxic phenotypes (Fig. 4b). For instance, CCR7, LEF1, TCF7, and SELL were up-
222 regulated in naïve and central memory (CM) cells, while memory markers such as IL7R,
223 GZMK were elevated in CM and effector memory (EM) phenotypes (Fig. 4b,
224 Supplementary Fig. 6a-b). Phenotype assignment based on transcriptomics was further
225 validated by surface protein expression measured by scCITE-seq (see Methods)
226 (Supplementary Fig. 6c). As expected, cytotoxic and EM phenotypes dominate the SARS-
227 Cov-2-specific CD8 T cell response during acute disease^{26,30,31}, with elevated effector
228 markers such as GZMB, GZMA and PRF1. In addition to gene expression, in Fig. 4c we
229 projected the polarity of exposed residues on the antigen recognized by the individual T
230 cells (left), the distribution of T cells captured by a given antigen across multiple
231 participants (middle), and all captured SARS-CoV-2 specific T cells for a given participant.

232

233 The projection of antigen polarity on the GEX-UMAP (Fig. 4c) showed a strong skewing
234 towards cytotoxic CD8 T cell phenotypes, suggesting that antigens with particular
235 physicochemical properties may associate with specific T cell phenotypes. We explored
236 this concept by projecting each cell's antigen specificity (Pep-Groups) onto the GEX-
237 UMAP (Fig. 4d top). PG1-specific cells were dominated by naïve-like phenotypes and
238 exhibited less clonal expansion relative to the other two groups (Supplementary Fig. 6d).
239 PG2-specific cells were enriched with EM and CM phenotypes, and PG3-specific cells
240 were mainly cytotoxic (Fig. 4d bottom). Furthermore, differentially expressed genes
241 (DEGs) in PG3-captured cells have enriched pathway signatures related to immune
242 synapse formation, PD-1 signaling, CD28 co-stimulation, which further highlighted the
243 effector state of PG3-specific cells relative to PG2 (Fig. 4e, Supplementary Table 5-6).
244 This analysis suggests strong associations exist between the physicochemical properties
245 of the peptide antigens, and the phenotypic characteristics of the T cells specific to those
246 antigens.

247

248 **TCR hydrophobicity is an important factor for effector function**

249 Intrigued by the link between antigenic peptides and T cell phenotypes, we investigated
250 whether a similar connection exists for the physicochemical features of each antigen-
251 associated TCR-CDR3 sequence by overlaying those features on the GEX-UMAP (Fig
252 5a, Supplementary Fig. 7a). Specifically, we categorized CDR3 β residues into V, J, and
253 CDR3 β mer (central) regions: The V and J regions comprise the highly conserved n- and
254 c-terminal motifs respectively; while the central CDR3 β mer region, which primarily
255 contacts the antigen, is the most diverse in length and amino acid usage. As expected,
256 cytotoxic cells showed maximal clonal expansion followed by effector memory (EM) cells
257 (Supplementary Fig. 7b). We compared the CDR3 β features of effector cells (cytotoxic
258 and EM) relative to other cell types (Supplementary Table 7-8). TCR feature differences
259 were not significant for the V and J regions as expected (Fig. 5c middle). However, for
260 the CDR3 β mer region, a preference of sequence features was observed. Effector T cells
261 were marked by higher CDR3 β mer hydrophobicity and bulkiness, lower polarity and
262 charge, and shorter length (Fig. 5c top and bottom). Note that CDR3 β hydrophobicity was
263 defined both by the percentage of hydrophobic residues, and by the average
264 hydrophobicity across central residues (See method). We validated the above trends by
265 plotting how the selected CDR3 β physicochemical properties varied across all cell
266 phenotypes. The CDR3 β sequences displayed the trend of increased hydrophobicity and
267 decrease in charge and length from naïve, to EM and cytotoxic phenotypes (Fig. 5d).

268
269 The hydrophobicity of the CDR3 β exhibited the strongest significant association with T
270 cell phenotype. This prompted us to define a binary classifier (HPhobic-High and
271 HPhobic-low) based on the percentage of hydrophobic residues in CDR3 β mer (cutoff =
272 25%) (Fig. 5e top, see Methods). Cells expressing HPhobic-High TCRs were more
273 clonally expanded than HPhobic-Low ones (Fig. 5e bottom). We first validated that TCR-
274 Groups still preserve the physicochemical features observed earlier: Indeed, HPhobic-
275 High TCRs exhibited higher hydrophobicity, shorter CDR3 β length, and lower charge
276 (Supplementary Fig. 7c). Density mapping validated that HPhobic-High TCRs were more
277 prevalent in cytotoxic cells than in memory and naïve subsets (Fig. 5f), with elevated
278 exhaustion markers such as LAG3 and TIGIT (Supplementary Fig. 7d). Gene set
279 enrichment analysis linked HPhobic-High clonotypes to TCR activation (e.g. CD3 and
280 TCR Zeta-chain phosphorylation), inflammation, and apoptosis pathways (Fig. 5g,
281 Supplementary Table 9-10). To test whether dominant clonotypes affect our result,
282 additional analysis was performed by removing dominant TCRs, although dominant TCRs
283 were only found for a few of the epitopes in our dataset (Supplementary Fig. 8a, Details
284 in Method). We observed consistent results with or without the dominant clones. Cells
285 with hydrophobic-high CDR3 β s consistently showed elevated PRF1 and GZMB gene
286 expression, as well as higher cytotoxic scores relative to cells with hydrophobic-low
287 CDR3 β s after removal of large clones (Supplementary Fig. 8b top). On contrast,
288 hydrophobic-low CDR3 β s associated with higher expression of naïve and memory related
289 genes including CCR7 and IL7R (Supplementary Fig. 8b bottom). These results indicate
290 that our original conclusions are not biased by specific dominant clones.

291

292 Overall, we demonstrated that, the physicochemical properties of both the peptide antigen
293 and the TCR CDR3 β exhibit strong associations with T cell phenotype for SARS-CoV-2
294 specific CD8 T cells.

295

296 **Integrated analysis of both peptide and TCR physicochemical features orchestrate** 297 **phenotypes of SARS-CoV-2-specific CD8 T cells in acute disease**

298 To elucidate how interplay between antigen and TCR interactions influence T cell function,
299 we systematically linked peptide and TCR features for the subset of SARS-CoV-2 CD8 T
300 cells with fully paired antigen-TCR information (See Method). This encompassed 87
301 unique antigenic peptides (SCT-pMHCs) and 440 paired TCR clones (Fig. 6a). Notably,
302 distinct TCR Groups (HPhobic-Low and HPhobic-High) were discerned within each
303 Peptide Group, prompting a refined categorization into PG-TCR Groups based on
304 combined antigen peptide and TCR features (Fig. 6b). For example, PG3-High denotes
305 cells captured by PG3 peptides with HPhobic-High TCRs.

306

307 We then further evaluated the combinatorial effect of peptide and TCR features for each
308 PG-TCR group. We depicted key physicochemical properties identified earlier on radar
309 plots, revealing a significant shift in overall characteristics between PG-TCR groups (Fig.
310 6c, Supplementary Fig. 9a). For example, PG3:High exhibits high peptide charge and low
311 hydrophobicity on exposed residues, combined with high CDR3 β mer hydrophobicity and
312 bulkiness (Fig. 6c left). In contrast, PG2:Low (Fig. 6c middle) emphasizes CDR3 β mer
313 length and charge, while PG1:Low (Fig. 6c right) emphasizes distinct characteristics
314 compared to PG3:High.

315

316 Building on our findings, we explored how PG-TCR groups associate with T cell
317 phenotypes. PG3:High cells displayed the strongest cytotoxicity – marked by the highest
318 percentage of effector cells and elevated cytotoxic cytokines (e.g. GZMB, PRF1, GZMH)
319 (Fig. 6d left, Supplementary Fig. 9b-c). PG1:Low describes cells with TCR:pMHC
320 interfacial properties that are opposite that of PG3:High, and those cells similarly exhibit
321 phenotypes that contrast with PG3:High – marked by highest frequency of naïve cells and
322 elevated naïve-associated markers (e.g. CCR7, LEF1, TCF7) (Fig. 6d right). Notably,
323 PG1 consistently exhibits a Naïve-like phenotype, regardless of TCR-group. With
324 intermediate peptide-TCR properties, PG2:Low and PG3:Low represent transitional
325 phenotypes. These associations between PG-TCR groups and gene expression were
326 validated using protein markers by integrating existing scCITE-seq data (See Methods)
327 (Supplementary Fig. 9d). We performed further analysis by directly linking the key
328 physicochemical properties to phenotypes independent of the PG-TCR groups. This
329 heatmap highlights the balance between the properties of the antigen and the TCR
330 CDR3 β sequence (Fig. 6e), and illustrates the near opposite relationships between
331 naïve/CM cells relative to cytotoxic cells at the stage of acute disease. Hence, by

332 integrating peptide and TCR sequence features, APMAT reveals fundamental rules in the
333 peptide-TCR-phenotype synergy.

334

335 **Distinct peptide-TCR groups associate with distinct longitudinal fates of SARS-** 336 **CoV-2-Specific CD8 T cells**

337 T cell activation strength through pMHC-TCR interaction can influence cell fate over
338 time^{32–34}. Certain T cell phenotypes exhibit long-term *in vivo* persistence following acute
339 illness, while highly cytotoxic phenotypes can undergo activation-induced cell death
340 (AICD) and contract after clearance of pathogen^{35–37}. We hypothesized that distinct Pep-
341 TCR groups may associate with T cell fate decisions across longitudinal disease
342 trajectories as well. Leveraging the longitudinal scRNA-seq and scCITE-seq data
343 generated from the same COVID-19 cohort²⁵, we tracked SARS-CoV-2-specific CD8 T
344 cells matched on patient ID and TCR sequences from acute to post-acute timepoints (Fig.
345 7a, Supplementary Table 11). This analysis further revealed longitudinal differences
346 between PG-TCR groups. As expected, the overall percentage of SARS-CoV-2 specific
347 CD8 T cells identified from the reference dataset decreased from 3% to 1.3% at the post-
348 acute timepoint. Specifically, PG3:High cells were short-lived, showing the greatest drop
349 in abundance at the post-acute timepoint. By contrast, PG1 cells (including PG1:High and
350 PG1:Low) were the most persistent, showing an increased abundance at the later
351 timepoint (Fig. 7b, Supplementary Fig. 9e).

352

353 In addition to abundance changes over time, we examined gene expression changes for
354 SARS-CoV-2 specific CD8 T cells. Combining all PG-TCR groups, we observed that
355 persisting cells at the post-acute timepoint showed higher expression of genes that were
356 associated with long-lived memory signatures (such as CCR7, IL7R, HLA-DR, MKI67, etc)
357 (Supplementary Fig. 9f), in agreement with previous studies^{16,18,20}. In addition, our
358 analysis further suggested a few other trends. Specifically, PG1 cells identified at the late
359 timepoint showed a relatively lower naïve signature (CCR7, TCF7, etc) and slightly higher
360 effector functions (GZMB, GZMH, etc) than those cells at the earlier timepoint (Fig. 7c
361 right, Supplementary Fig. 9g right). By contrast, the few remaining post-acute PG3:High
362 cells evolve to express lower cytotoxicity (GZMB, PRF1, etc) and higher CCR7, indicating
363 a shift towards less effector or central memory phenotypes^{16,38,39} (Fig. 7c left,
364 Supplementary Fig. 9g left). In summary, APMAT revealed that distinct combinations of
365 peptide and TCR physicochemical properties exhibit clear associations with not only
366 cellular phenotypes at acute disease, but also divergent cell fates over time (Fig. 7d).

367

368 **Discussion**

369 The conceptual advance we report as APMAT lies in the combination of the large and
370 comprehensive experimental data set and the biophysical analysis of the TCR-pMHC
371 interface. APMAT provides a framework for the integrated analysis of antigen-specific

372 CD8 T cells paired with phenotypic data on those T cells. We applied APMAT to
373 investigate circulating CD8 T cells collected from 62 participants at the acute and
374 convalescent stages of COVID-19. These cells exhibited specificity to SARS-CoV-2
375 antigens presented by HLA A*02:01 from across the full viral proteome. Our analysis
376 uncovered relationships between the physicochemical characteristics of the pMHC:TCR
377 interface and the corresponding T cell phenotypes and T cell phenotypic evolution. Our
378 analysis is, in several aspects, aligned with existing literature. We identified T cell
379 populations targeting previously reported immunogenic antigens^{17,24,40,41}, and we find that
380 the prediction rankings from the commonly-used NetMHCpan algorithm effectively assist
381 in designing pMHC multimer constructs that will most likely capture T cells. In addition,
382 the longitudinal behavior of the antigen-specific clonotypes is consistent with current
383 literature. Highly-cytotoxic effector clonotypes, which are expanded at acute disease,
384 contract significantly at the time of convalescence, potentially from antigen-induced cell
385 death. In contrast, those persisting clonotypes with mild effector properties evolve
386 towards central or effector memory phenotypes⁴². Further, clonotypes that exhibit
387 memory and progenitor-like phenotypes at acute disease persist or expand at
388 convalescence^{16,43,44}.

389 APMAT analysis uncovers novel relationships between viral antigen-specific T cells, TCR
390 clonotype, and T cell phenotype for the common HLA allele A*02.01. Take, for example,
391 the above-described case of effector T cells that expand during an acute infection and
392 contract at convalescence⁴². Our analysis further suggests that T cells possessing TCRs
393 characterized by a hydrophobic CDR3 β , which recognize antigens featuring hydrophobic
394 anchor residues alongside charged or polar exposed residues, are statistically more
395 inclined towards effector phenotypes at the stage of acute COVID-19, and contract during
396 convalescence. An analogous description of the pMHC:TCR interface (but with near
397 opposite characteristics) can be used to identify those clonotypes that exhibit naïve or
398 central memory phenotypes at acute disease and remain at convalescence. We focus on
399 CDR3 β , but CDR3 α chains can be similarly analyzed to identify specific physicochemical
400 properties that significantly associate with effector (Cytotoxic, EM, Hybrid) or non-effector
401 T cells. Those TCR α properties are distinct from those for TCR β , suggesting that the
402 influences from the α and β chains might be complementary rather than independent
403 (Supplementary Fig. 10a).

404 An additional significant aspect of this study is the detailed characterization of
405 paired peptide antigen with TCR. Notably, our analysis reveals that the anchor and
406 exposed residues of the antigen exert different influences on T cell phenotype. Unlike
407 many prior studies that relied on a rough annotation of the antigen identity, our approach
408 takes into account the position of each individual residue, as well as the classification of
409 antigens that may be comprised of distinct sequences, and yet exhibit biochemical
410 similarities at the TCR:pMHC interface.

411 Whether or not these physicochemical determinants are general to other disease
412 contexts or other HLA Class I alleles is not resolved here, although the consistency of the
413 picture painted here with what has been observed in tissue settings in murine models of
414 chronic viral infection,⁴ as well as in tumors,^{5,6} does suggest some level of generality.

415 Further, an analysis of a separate study on SARS-CoV-2 reactive T cells (Supplementary
416 Fig. 11a)⁴⁵, as well as our recently reported phenotypic analysis of T cell clonotypes
417 specific to three highly immunogenic viral antigens, including two from influenza and
418 CMV⁸, (Supplementary Fig. 11b) also suggest a degree of generality. Datasets of similar
419 breadth and depth for a range of diseases and for antigens presented by different HLA
420 alleles are needed to more fully resolve this picture, and such work represents an exciting
421 future direction.

422 We hypothesize that an analogous study as the one reported here, but directed at
423 antigens presented by different HLA alleles, might yield a different set of rules that are
424 dependent upon HLA-specific docking geometries. While Class I HLA A, B, and C alleles
425 are highly polymorphic, over 90% of the world's population carries at least one of the
426 dozen most common. An APMAT analysis centered around those most common alleles
427 should offer valuable insights for assessing the therapeutic potential of tumor-targeting
428 TCRs and T cell vaccines.

429

430 **Methods**

431

432 **Lead Contact**

433 Further information and requests for resources and reagents should be directed to and
434 will be fulfilled by the Lead Contact, Dr. James R. Heath (jim.heath@isbscience.org).

435

436 **Participants and sample collection**

437 The study cohort is a subset of the INCOV cohort published previously^{25,26}. Procedures
438 for the INCOV study were approved by the Institutional Review Board (IRB) at Providence
439 St. Joseph Health with IRB study number STUDY2020000175 and the Western
440 Institutional Review Board with IRB study number 20170658. This research complies with
441 all relevant ethical regulations. Potential participants were identified at five hospitals of
442 Swedish Medical Center and affiliated clinics located in the Puget Sound region near
443 Seattle, WA. All enrolled participants provided written in-person informed consent and
444 samples were de-identified prior to analysis. 62 HLA A*02:01 individuals from the INCOV
445 cohort were selected for this study. PBMCs collected at Acute timepoint, including
446 enrollment close to diagnosis (Acute-1) and 1 week (Acute-2), were used for antigen-TCR
447 paring assay in this study.

448

449 **Large-scale preparation of peptide-HLA complex libraries**

450 Single chain trimer (SCT) peptide-MHC (pMHC) libraries of the virus antigens were
451 generated as described previously²³ 1/8/25 4:32:00 PM. Briefly, potential HLA A*02:01
452 binding epitopes (9-11 mer peptides) were generated from the complete SARS-CoV-2
453 genome (Wuhan-Hu-1 strain, GenBank ID: MN908947.3) and filtered by NetMHCpan-4.1
454 prediction²⁸. We categorized peptides as weak-binders (column BindLevel = 0), medium-

455 binders (column BindLevel = WB), or strong-binders (column BindLevel = SB) based on
456 NetMHCPan prediction. Note that even the weak-binders are relatively good candidates
457 compared to peptides that were not attempted for SCT expression. A plasmid library of
458 pcDNA3.1 vectors encoding covalently linked peptide antigen, β 2M, HLA was built and
459 verified by SANGER sequencing. Plasmids were transfected into Exp293 cells
460 (ExpiFectamine™, Thermo Fisher) following manufacturer protocol. SCT expression yield
461 was measured and normalized. Expressed pMHC-like SCTs were biotinylated (BirA
462 ligase Kit, Avidity) and Histag purified (IMAC PhyTip columns, PhyNexus) in 96-well
463 format. Individual SCT concentration was measured by protein absorbance at 290nm.
464

465 **SCT-dextramer generation and cell staining**

466 SCT dextramers were individually DNA barcoded using dCODE Klickmers (dCODE
467 Klickmer, Immudex). Briefly, SCT pMHC monomer was mixed with barcoded dCODE-
468 PE-dextramer at a ratio of 20 ligands per dextran and incubated for at least 1h on ice
469 before adding biotin (100 μ M) to block free binding sites. Dextramer cocktails were
470 prepared by mixing 31-65 unique SARS-CoV-2 SCT dextramers and CMV (NLVPMVATV)
471 SCT dextramers freshly before cell staining. PBMCs from 62 participants were thawed
472 for CD8 T cell enrichment (Human CD8 T cell Isolation Kit, Miltenyi Biotec) according to
473 the manufacturer's protocol then incubated with Human TruStain FcX blocking reagent
474 (422302, BioLegend) for 10 min at 4 °C before wash. Cells were then divided into tubes
475 and processed simultaneously on ice. Each tube of CD8 T cells were stained with a
476 cocktail of dextramers for 25 min on ice in the presence of herring sperm DNA according
477 to the manufacturer's instructions, with individual dextramer concentration at 1.1 nM.
478 Cells were washed three times before surface antibody staining. BV421 anti-human CD8
479 Antibody (BioLegend, 344748, clone SK1) and Apotracker™ Green viability dye
480 (Biolegend, 427403) was added into each tube, in addition to one unique TotalSeq-C anti-
481 human hashtag antibody (BioLegend) to identify each tube. Samples were incubated for
482 30 min on ice and washed 3 times before sorting.
483

484 **10X genomics single cell sequencing**

485 Single, live, CD8, dextramer-positive T cells were sorted into FACS buffer (PBS, 2%FBS,
486 2mM EDTA and 10mM HEPES) using a BD FACSArial cell sorter. Sorted cells were
487 immediately pelleted, resuspended and loaded into a 10X Chromium reaction for
488 single cell RNA sequencing (scRNA-seq). GEX, VDJ and Surface Protein libraries were
489 generated using Chromium Next GEM Single-Cell 5' kits v2 (10X Genomics) according
490 to the manufacturer's protocol. Libraries were sequenced on an Illumina NovaSeq at a
491 read length of 26x90 bp.
492

493 **Whole genome sequencing**

494 DNA extraction from whole blood was performed via bead-based enrichment on an
495 automated extraction platform (Qiagen Qiasymphony and/or Promega Maxwell). The
496 resultant extracts were quantified by Nanodrop. WGS library preparation was performed

497 using Illumina DNA Prep kits and the final barcoded libraries were quantified by
498 fluorometer. Libraries were multiplexed and loaded onto an Illumina flow cell for
499 sequencing at 30x or higher coverage on a NovaSeq 6000 instrument. Raw sequencing
500 data was analyzed for sequence variants using the Illumina DRAGEN field-programmable
501 gate array (FPGA) platform. Briefly this platform performs the following automated steps:
502 conversion of raw sequencing image data into demultiplexed fastq files, alignment to the
503 reference human genome (hg19), analysis of single nucleotide variants, indels and copy
504 number/structural variants using variant calling algorithms as well as assessment of
505 sequencing data quality. Analyses with hg38 were computed after a liftover was done
506 using the UCSC browser⁴⁶. WGS information was de-identified.

507

508 **Single-cell sequencing data processing**

509 Transcriptome, TCR, surface protein levels and antigen specificity were simultaneously
510 analyzed for each cell. Raw data were processed via Cell Ranger Single-Cell Software
511 Suite (v3.1.0, 10X Genomics) using GRCh38 as a reference. Cells that fit any of the
512 following filters were excluded due to low quality: n-counts <1000 or >10,000, n-genes
513 <250 or >2500, mitochondrial percentage >10%. Gene counts for each cell were
514 normalized by total expression, multiplied by a scale factor of 10,000 and transformed to
515 log scale.

516

517 **Single CD8 T cell phenotype assignment**

518 Single cells were assigned phenotypes by clusters determined through the leiden
519 algorithm. Phenotype associated transcripts were acquired from literature as
520 follows^{16,26,39,47}. Naïve/memory: LEF1, TCF7, CCR7. Memory: IL7R. Effector Memory:
521 GZMK. Cytotoxic: PRF1, GZMB. Exhaustion: TIGIT, PDCD1. Proliferation: MKI67.

522

523 **Transcriptomics data validation via scCITE-seq**

524 To validate the transcriptomics data, we leverage our previously reported scCITE-seq
525 data that simultaneously measured transcriptomic and surface protein levels and TCR α/β
526 from the same single cell. We extracted each cell's scCITE-seq data by TCR-based
527 inquiry (illustrated in manuscript Fig. 7a, and methods). This leads to single cells with
528 matching transcriptomics and surface protein data (which is the data for plotting main Fig.
529 7a-c, Supplementary Fig.6c, Supplementary Fig.9d). Our phenotype assignment based
530 on transcriptomics (as the reviewer mentioned for Main Figs. 4 and 5) can be validated
531 by surface protein expression (measured by scCITE-seq). These pairs including CCR7
532 and CD197-CCR7-CITEseq, IL7R and CD127-IL-7RA-CITEseq. In addition, well-
533 established surface protein markers were supported by our transcriptomics data, such as
534 CD45RO for memory phenotype, CD45RA for naïve phenotype, CD25 for central
535 memory⁴⁸.

536

537 **Demultiplexing using genetic variants**

538 We wrote a Snakemake workflow (<https://github.com/racng/snakemake-merge-wgs>) for
539 processing GVCF files from WGS to generate a multi-sample VCF file of exon variants.

540 GVCFs are combined for specified genomics region (autosomes and sex
541 chromosomes) using GATK (v4.1.9.0) GenomicsDBImport to generate a GenomicsDB
542 datastore, which is then used by GenotypeGVCFs for joint calling of variants. After
543 removing indels using vcftools (v0.1.16)⁴⁹ and excluding intron variants via bcftools
544 (v1.8)⁵⁰ the remaining exon variants were lifted to GRCh38 (hg38) using CrossMap
545 (v0.5.2)⁵¹ and filtered again to remove KI27 contigs and duplicated variants. For each
546 10x library, BAM alignment files from cellranger were filtered for reads from autosomes
547 and sex chromosomes. Using the processed VCF file (929,678 SNPs), single cell
548 variants were extracted from the filtered BAM files via 10x Genomics VarTrix
549 (<https://github.com/10XGenomics/vartrix>) with the coverage scoring method. To reduce
550 memory usage, single cell variants were kept if both ALT and REF alleles are detected
551 in the dataset. We then used vireo (v0.5.0)⁵² to assign donor identity and doublets
552 based on the processed single-cell and WGS variants. Doublets were removed. Cells
553 that unassigned with donor identity were removed.

554

555 **Antigen assignment based on dextramers and hashtags**

556 Raw reads for each Hashtags and Dextramers were normalized. Cells were assigned to
557 their maximally expressed Hashtag. Cells that expressed multiple Hashtags at a high
558 level were removed as potential doublets. Hashtag identities were then used to identify
559 cells' SCT-dextramer cocktail. For each cell, we calculated the number of unique
560 molecular identifiers (UMIs) for each dextramer, and the percentage of each dextramer.
561 We assigned each cell an antigen only if their UMI count was >25 and the UMIs specific
562 for that dextramer occupied >25% of that cells' dextramer reads. Antigens were then
563 assigned by the maximally mapped dextramer for each cell. Ambiguous cells that didn't
564 assigned with any dextramer were removed.

565

566 **Peptide physicochemical property assignment and Pep-UMAP**

567 We first transform peptide sequences into a numerical peptide matrix, where each row
568 represents a residue position, and each column represents a feature characterizing amino
569 acids, including amino acid identity, charge, hydrophobicity, weight, bulkiness, polarity,
570 sulfur presence, aromaticity (Supplementary Table 3: Amino acid property scales used
571 for peptide and TCR residues)¹². The numeric values were scaled to ensure consistency
572 range of values. For quantitative comparison of peptides in Figure 2, we calculated the
573 average value for each property for anchor residues and TCR-exposing residues.
574 Specifically, position 2 and 9 (and occasionally 6) tend to serve as anchors for HLA-
575 A*0201 binding, while other exposed residues potentially contact the TCR (CITE). Logo
576 plot in Figure 2 was generated by Seq2Logo - 2.0⁵³ using default settings. In figure 3, to
577 visualize peptide features and similarities, we applied UMAP followed by an autoencoder
578 for dimensional reduction. In detail, since peptides have varying lengths, the residue
579 positions are mapped to a normalized scale of 0 to 100. Each amino acid's features are

580 replicated across the corresponding positions in the matrix. We then implemented an
581 autoencoder using an MLPRegressor from scikit-learn to reduce the dimensionality of the
582 peptide matrix. Finally, we computed a two dimensional peptide UMAP (Pep-UMAP)
583 using scanpy.tl.umap to visualize peptide features and similarities.

584

585 **TCR physicochemical property assignment**

586 TCR-related physicochemical properties were computed for each cell based on its TCR
587 CDR3 sequence of the beta chain. These properties, including charge, hydrophobicity,
588 weight, bulkiness, polarity, sulfur presence, aromaticity, were calculated based on
589 Supplementary Table 3 (Amino acid property scales used for peptide and TCR residues).
590 Charge is absolute value unless specifically indicated. Specifically, we categorized
591 CDR3 β residues into V, J, and CDR3bmer regions. The V/J region comprise the first four,
592 and last five amino acids, respectively, while the central CDR3bmer contains the amino
593 acids in between. We evaluated TCR CDR3 hydrophobicity numerically by the average
594 hydrophobicity of the CDR3bmer region. Additionally, we introduced a categorical score
595 called HPhobic%. HPhobic% represents the percentage of strongly hydrophobic residues
596 (A, V, L, I, F, M) in the CDR3 β , excluding the first and last four amino acids.

597

598 **TCR physicochemical property analysis**

599 We calculated 45 properties for each cell's TCR sequences, including charge,
600 hydrophobicity, weight, bulkiness, polarity, for three regions of CDR3 β (V, J, CDR3 β mer),
601 as well as full CDR3 α chain. The Mann-Whitney U test is applied to compare the
602 distribution of each property between effector (Cytotoxic, Effector Memory, Hybrid) and
603 non-effector (Naïve and Central Memory) cell phenotypes. Log2 fold change (log2fc) for
604 each property was calculated between the mean values of effector cell types and non-
605 effector cell types. Top selected properties are based on the criteria of log2fc absolute
606 value > 0.05, and p value <0.05. (Supplementary Table 7 and 8). For Supplementary Fig.
607 8b, Dominant TCRs were removed to test whether our result is biased by large clones. In
608 detail, we performed additional analysis by removing T cell clonotypes and then evaluated
609 the relationship between TCR Groups and phenotypes without these dominant clones.
610 Dominant clones were defined as clone size \geq 50 based on 10X VDJ library readout
611 (which counted all cells in the pre-filtered dataset). We also defined phenotype scores
612 as: Cytotoxic score (GZMB, GZMA, GNLY, PRF1) and Naïve-Memory score (TCF7,
613 CCR7, SELL, LEF1, GZMK, IL7R). These were calculated by Scanpy.tl.score to represent
614 the average expression of the given set of genes.

615

616 **Peptide and TCR groups density analysis**

617 Densities for peptide groups were projected onto GEX-UMAP by matching the antigen
618 specificity of each single CD8 T cell to the peptide group that antigen peptide belongs to.
619 Embedding density was first calculated via sc.tl.embedding_density then a 5 n-neighbor
620 kNN graph was used to diffuse the values via five iterations to create a whole UMAP

621 score for the density scores, as reported previously^{8,54}. TCR group density calculations
622 were implemented this same methodology. The UMAP densities in the original Fig. 4d
623 and Fig. 5f were calculated as an odds ratio. In the area of low cell frequency of certain
624 PG/TCR groups, the density has a low value – but not necessarily zero.

625

626 **Differential gene expression and signature analysis**

627 Differentially expressed genes were called via `scanpy.tl.rank_genes_groups` through the
628 Scanpy package using the Wilcoxon method which implements the Mann-Whitney U
629 test⁵⁵. Differentially expressed genes (DEGs) between peptide groups (PG2 vs PG3)
630 were filtered for p values < 0.05 (Supplementary Table 5). Enriched pathways through
631 filtered DEGs were computed via Enrichr⁵⁶ and provided in Supplementary Table 6. Top
632 enriched pathways in PG3 than PG2 were reported in Fig. 4e. DEGs between TCR groups
633 (HPhobic-High vs HPhobic-Low) were called similarly to generate list of DEGs
634 (Supplementary Table 9) and enriched pathways (Supplementary Table 10). Top
635 enriched pathways in HPhobic-High than HPhobic-Low were reported in Fig. 5g.

636

637 **Longitudinal T cell inquiry by GLIPH2 analysis for SARS-CoV-2 specific TCRs**

638 We utilized a reference dataset - our previously reported longitudinal dataset on 209
639 COVID-19 participants contained both scRNA-seq and scTCR-seq data from the same
640 single cell along with assigned donors²⁵. We perform GLIPH2 analysis on SARS-CoV-2
641 specific TCRs identified in this study that assigned with antigen specificity, and TCRs from
642 the longitudinal dataset with unknown antigen specificity. GLIPH2 was run on
643 <http://50.255.35.37:8080/> using the GLIPH2 algorithm⁵⁷, version 1.0 reference for CD8,
644 with `all_aa_interchangeable` set to YES. Antigen specific CD8 T cells were identified as
645 those that belonged to the same GLIPH group of the SARS-CoV-2 specific TCRs
646 identified in this study. Longitudinal SARS-CoV-2 specific CD8 T cells were then subset
647 from the reference dataset for analysis in Figure 7. Cell count for each PG-TCR groups
648 at Acute and Post-Acute timepoints were provided in Supplementary Table 11.

649

650 **Statistical analysis**

651 Microarray like datasets were analyzed using SCANpy and statistical comparisons were
652 generated using `scanpy.tl.rank_genes_groups` using the wilcoxon method. All
653 correlations were calculated using Pearson correlation. All p values were calculated using
654 Mann-Whitney U test unless otherwise specified. Bar charts were provided with error bars
655 when multiple values were present, and these bars represented standard errors. Bar level
656 represented the mean variable value.

657

658 **Data Availability**

659 The single cell sequencing data generated in this study have been deposited in the
660 ArrayExpress database under accession number: E-MTAB-14002, or using the URL:

661 <https://www.ebi.ac.uk/biostudies/arrayexpress/studies/E-MTAB-14002?key=3ae4f467->
662 [fe01-4d03-b2fe-8a6a522e1cabA](https://www.ebi.ac.uk/biostudies/arrayexpress/studies/E-MTAB-14002?key=3ae4f467-fe01-4d03-b2fe-8a6a522e1cabA).

663 Source data in this study are provided in the Supplementary Information/Source Data file.

664 Any additional information required to reanalyze the data reported in this work paper is

665 available from the Lead Contact upon request.

666 **References**

- 667 1. Chang, J. T., Wherry, E. J. & Goldrath, A. W. Molecular regulation of effector and memory T
668 cell differentiation. *Nat Immunol* **15**, 1104–1115 (2014).
- 669 2. Zhang, N. & Bevan, M. J. CD8+ T Cells: Foot Soldiers of the Immune System. *Immunity* **35**,
670 161–168 (2011).
- 671 3. Connolly, K. A. *et al.* A reservoir of stem-like CD8+ T cells in the tumor-draining lymph node
672 preserves the ongoing antitumor immune response. *Science Immunology* **6**, eabg7836
673 (2021).
- 674 4. Daniel, B. *et al.* Divergent clonal differentiation trajectories of T cell exhaustion. *Nature*
675 *Immunology* **23**, 1614–1627 (2022).
- 676 5. Pai, J. A. *et al.* Lineage tracing reveals clonal progenitors and long-term persistence of
677 tumor-specific T cells during immune checkpoint blockade. *Cancer Cell* **41**, 776-790.e7
678 (2023).
- 679 6. Lagattuta, K. A. *et al.* Repertoire analyses reveal T cell antigen receptor sequence features
680 that influence T cell fate. *Nature Immunology* **23**, 446–457 (2022).
- 681 7. Pettmann, J. *et al.* The discriminatory power of the T cell receptor. *eLife* **10**, e67092 (2021).
- 682 8. Chen, D. G., Xie, J., Su, Y. & Heath, J. R. T cell receptor sequences are the dominant factor
683 contributing to the phenotype of CD8+ T cells with specificities against immunogenic viral
684 antigens. *Cell Reports* **42**, 113279 (2023).
- 685 9. Daley, S. R. *et al.* Cysteine and hydrophobic residues in CDR3 serve as distinct T-cell self-
686 reactivity indices. *Journal of Allergy and Clinical Immunology* **144**, 333–336 (2019).
- 687 10. Schmidt, J. *et al.* Prediction of neo-epitope immunogenicity reveals TCR recognition
688 determinants and provides insight into immunoediting. *Cell Reports Medicine* **2**, 100194
689 (2021).
- 690 11. Textor, J. *et al.* Machine learning analysis of the T cell receptor repertoire identifies
691 sequence features of self-reactivity. *Cell Systems* **14**, 1059-1073.e5 (2023).

- 692 12. Chowell, D. *et al.* TCR contact residue hydrophobicity is a hallmark of immunogenic CD8+ T
693 cell epitopes. *Proceedings of the National Academy of Sciences of the United States of*
694 *America* **112**, E1754–E1762 (2015).
- 695 13. Leidner, R. *et al.* Neoantigen T-Cell Receptor Gene Therapy in Pancreatic Cancer. *New*
696 *England Journal of Medicine* **386**, 2112–2119 (2022).
- 697 14. Pauken, K. E. *et al.* TCR-sequencing in cancer and autoimmunity: barcodes and beyond.
698 *Trends in Immunology* **43**, 180–194 (2022).
- 699 15. Shah, K., Al-Haidari, A., Sun, J. & Kazi, J. U. T cell receptor (TCR) signaling in health and
700 disease. *Signal Transduction and Targeted Therapy* **6**, (2021).
- 701 16. Adamo, S. *et al.* Signature of long-lived memory CD8+ T cells in acute SARS-CoV-2
702 infection. *Nature* **602**, 148–155 (2022).
- 703 17. Saini, S. K. *et al.* SARS-CoV-2 genome-wide T cell epitope mapping reveals
704 immunodominance and substantial CD8⁺ T cell activation in COVID-19 patients. *Sci.*
705 *Immunol.* **6**, eabf7550 (2021).
- 706 18. Mold, J. E. *et al.* Divergent clonal differentiation trajectories establish CD8+ memory T cell
707 heterogeneity during acute viral infections in humans. *Cell Reports* **35**, 109174 (2021).
- 708 19. Stadinski, B. D. *et al.* Hydrophobic CDR3 residues promote the development of self-reactive
709 T cells. *Nat Immunol* **17**, 946–955 (2016).
- 710 20. Minervina, A. A. *et al.* SARS-CoV-2 antigen exposure history shapes phenotypes and
711 specificity of memory CD8+ T cells. *Nat Immunol* (2022) doi:10.1038/s41590-022-01184-4.
- 712 21. Ma, K.-Y. *et al.* High-throughput and high-dimensional single-cell analysis of antigen-specific
713 CD8+ T cells. *Nat Immunol* **22**, 1590–1598 (2021).
- 714 22. Puig-Saus, C. *et al.* Neoantigen-targeted CD8+ T cell responses with PD-1 blockade
715 therapy. *Nature* **615**, (2023).
- 716 23. Chour, W. *et al.* Large libraries of single-chain trimer peptide-MHCs enable antigen-specific
717 CD8+ T cell discovery and analysis. *Communications Biology* **6**, 1–13 (2023).

- 718 24. Grifoni, A. *et al.* SARS-CoV-2 human T cell epitopes: Adaptive immune response against
719 COVID-19. *Cell Host & Microbe* **29**, 1076–1092 (2021).
- 720 25. Su, Y. *et al.* Multiple early factors anticipate post-acute COVID-19 sequelae. *Cell* **185**, 881-
721 895.e20 (2022).
- 722 26. Su, Y. *et al.* Multi-Omics Resolves a Sharp Disease-State Shift between Mild and Moderate
723 COVID-19. *Cell* **183**, 1479-1495.e20 (2020).
- 724 27. Lee, J. W. *et al.* Integrated analysis of plasma and single immune cells uncovers metabolic
725 changes in individuals with COVID-19. *Nature Biotechnology* (2021) doi:10.1038/s41587-
726 021-01020-4.
- 727 28. Reynisson, B., Alvarez, B., Paul, S., Peters, B. & Nielsen, M. NetMHCpan-4.1 and
728 NetMHCIIpan-4.0: Improved predictions of MHC antigen presentation by concurrent motif
729 deconvolution and integration of MS MHC eluted ligand data. *Nucleic Acids Research* **48**,
730 W449–W454 (2021).
- 731 29. Chaurasia, P. *et al.* Structural basis of biased T cell receptor recognition of an
732 immunodominant HLA-A2 epitope of the SARS-CoV-2 spike protein. *Journal of Biological*
733 *Chemistry* **297**, 101065 (2021).
- 734 30. Zheng, H. *et al.* Multi-cohort analysis of host immune response identifies conserved
735 protective and detrimental modules associated with severity across viruses. *Immunity* **54**,
736 753-768.e5 (2021).
- 737 31. Zhang, J.-Y. *et al.* Single-cell landscape of immunological responses in patients with
738 COVID-19. *Nature Immunology* **21**, 1107–1118 (2020).
- 739 32. Mckeithan, T. W. Kinetic proofreading in T-cell receptor signal transduction. *Proceedings of*
740 *the National Academy of Sciences of the United States of America* **92**, 5042–5046 (1995).
- 741 33. Irving, M. *et al.* Interplay between T cell receptor binding kinetics and the level of cognate
742 peptide presented by major histocompatibility complexes governs CD8+ T cell
743 responsiveness. *Journal of Biological Chemistry* **287**, 23068–23078 (2012).

- 744 34. Chin, S. S. *et al.* T cell receptor and IL-2 signaling strength control memory CD8+ T cell
745 functional fitness via chromatin remodeling. *Nature Communications* **13**, (2022).
- 746 35. Kaech, S. M. & Cui, W. Transcriptional control of effector and memory CD8+ T cell
747 differentiation. *Nature Reviews Immunology* **12**, 749–761 (2012).
- 748 36. Joshi, N. S. *et al.* Inflammation Directs Memory Precursor and Short-Lived Effector CD8+ T
749 Cell Fates via the Graded Expression of T-bet Transcription Factor. *Immunity* **27**, 281–295
750 (2007).
- 751 37. Kaech, S. M. & Wherry, E. J. Heterogeneity and Cell-Fate Decisions in Effector and Memory
752 CD8+ T Cell Differentiation during Viral Infection. *Immunity* **27**, 393–405 (2007).
- 753 38. Zhou, X. *et al.* Differentiation and Persistence of Memory CD8+ T Cells Depend on T Cell
754 Factor 1. *Immunity* **33**, 229–240 (2010).
- 755 39. Campbell, J. J. *et al.* CCR7 Expression and Memory T Cell Diversity in Humans. *The*
756 *Journal of Immunology* **166**, 877 LP – 884 (2001).
- 757 40. Dolton, G. *et al.* Emergence of immune escape at dominant SARS-CoV-2 killer T cell
758 epitope. *Cell* **185**, 2936-2951.e19 (2022).
- 759 41. Wu, D. *et al.* Structural assessment of HLA-A2-restricted SARS-CoV-2 spike epitopes
760 recognized by public and private T-cell receptors. *Nature Communications* **13**, 1–14 (2022).
- 761 42. Wherry, E. J. & Ahmed, R. Memory CD8 T-Cell Differentiation during Viral Infection. *Journal*
762 *of Virology* **78**, 5535–5545 (2004).
- 763 43. Giles, J. R. *et al.* *Shared and Distinct Biological Circuits in Effector, Memory and Exhausted*
764 *CD8+ T Cells Revealed by Temporal Single-Cell Transcriptomics and Epigenetics*. *Nature*
765 *Immunology* vol. 23 (Springer US, 2022).
- 766 44. Harty, J. T. & Badovinac, V. P. Shaping and reshaping CD8+ T-cell memory. *Nature Reviews*
767 *Immunology* **8**, 107–119 (2008).
- 768 45. Fischer, D. S. *et al.* Single-cell RNA sequencing reveals ex vivo signatures of SARS-CoV-2-
769 reactive T cells through ‘reverse phenotyping’. *Nature Communications* **12**, (2021).

- 770 46. Kent, W. J. *et al.* The Human Genome Browser at UCSC. *Genome Research* **12**, 996–1006
771 (2002).
- 772 47. Szabo, P. A. *et al.* Single-cell transcriptomics of human T cells reveals tissue and activation
773 signatures in health and disease. *Nature Communications* **10**, 4706 (2019).
- 774 48. Herndler-Brandstetter, D. *et al.* CD25-Expressing CD8+ T Cells Are Potent Memory Cells in
775 Old Age. *The Journal of Immunology* **175**, 1566–1574 (2005).
- 776 49. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158
777 (2011).
- 778 50. Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *GigaScience* **10**, 1–4 (2021).
- 779 51. Zhao, H. *et al.* CrossMap: A versatile tool for coordinate conversion between genome
780 assemblies. *Bioinformatics* **30**, 1006–1007 (2014).
- 781 52. Huang, Y., McCarthy, D. J. & Stegle, O. Vireo: Bayesian demultiplexing of pooled single-cell
782 RNA-seq data without genotype reference. *Genome Biology* **20**, 1–12 (2019).
- 783 53. Thomsen, M. C. F. & Nielsen, M. Seq2Logo: A method for construction and visualization of
784 amino acid binding motifs and sequence profiles including sequence weighting, pseudo
785 counts and two-sided representation of amino acid enrichment and depletion. *Nucleic Acids*
786 *Research* **40**, 281–287 (2012).
- 787 54. Chen, D. G. *et al.* Integrative systems biology reveals NKG2A-biased immune responses
788 correlate with protection in infectious disease, autoimmune disease, and cancer. *Cell*
789 *Reports* **43**, 113872 (2024).
- 790 55. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression data
791 analysis. *Genome Biology* **19**, 15 (2018).
- 792 56. Xie, Z. *et al.* Gene Set Knowledge Discovery with Enrichr. *Current Protocols* **1**, 1–51 (2021).
- 793 57. Huang, H., Wang, C., Rubelt, F., Scriba, T. J. & Davis, M. M. Analyzing the Mycobacterium
794 tuberculosis immune response by T-cell receptor clustering with GLIPH2 and genome-wide
795 antigen screening. *Nat Biotechnol* **38**, 1194–1202 (2020).

796 58. Fischer, D. S. *et al.* Single-cell RNA sequencing reveals ex vivo signatures of SARS-CoV-2-
797 reactive T cells through ‘reverse phenotyping’. *Nat Commun* **12**, 4515 (2021).

798

799

800 **Acknowledgements**

801 We are grateful to all participants in this study and to the medical teams at
802 Providence Swedish Medical Center and UW for their support. We thank the Core Facility
803 at Institute for Systems Biology, the Northwest Genomic Center and Fred Hutchinson
804 Cancer Research Center for help with sequencing services, and the ISB-Swedish COVID-
805 19 Biobanking Unit. We acknowledge funding support from the Parker Institute for Cancer
806 Immunotherapy (J.R.H., P.D.G.), Merck, the Biomedical Advanced Research and
807 Development Authority (HHSO10201600031C to J.R.H.), the NIH (1 R01 CA264090-01
808 to J.R.H. and P.D.G, and CTSI UL1TR001881 to A.M.X). Y.S. was supported by the
809 Damon Runyon Quantitative Biology Fellowship from the Damon Runyon Cancer
810 Research Foundation (DRQ-13-22), the Mahan Fellowship at Herbold Computational
811 Biology Program of Fred Hutchinson Cancer Research Center, the Translational Data
812 Science Integrated Research Center New Collaboration Award and Integrated Research
813 Center at Fred Hutchinson Cancer Research Center, Pilot Award of Immunotherapy
814 Integrated Research Center, and in part through the NIH/NCI Cancer Center Support
815 Grant P30 CA015704. We wish to thank John Heath and Dr. Alphonsus Ng, and all the
816 colleagues from the Institute for Systems Biology for the technical support. Schematic
817 figures in Figure 1, 2d, 3a, 4a, 5a, 6a, 7a were created with BioRender.com, released
818 under BioRender Academic Publication Licenses.

819
820

821 **Author Contributions Statement**

822 Conceptualization, J.X., Y.S., and J.R.H.; Resources, H.C., J.D.G., and J.R.H.;
823 Methodology, J.X., D.G.C., and J.R.H.; Investigation, J.X., D.G.C., W.C., R.N., R.Z., D.Y.,
824 J.C., M.K., P.T., B.S., L.J., A.W., Y.R., S.L., R.E., S.H., K.M.M, J.K.L, N.M.F, C.G.L, B.P,
825 H.A., J.W., A.T.M., K.W., P.M., P.D.G., H.C., J.D.G., Y.S., and J.R.H.; Formal analysis,
826 J.X., and J.R.H.; writing – original draft preparation, J.X., and J.R.H.; writing – review &
827 editing, J.X., D.G.C., W.C., R.N., R.Z., D.Y., J.C., M.K., P.T., B.S., L.J., A.W., Y.R., S.L.,
828 R.E., S.H., K.M.M, J.K.L, N.M.F, C.G.L, B.P, H.A., J.W., A.M., P.M., P.D.G., H.C., J.D.G.,
829 Y.S., and J.R.H.

830
831

832 **Competing Interests Statement**

833 J.R.H. is a consultant to Regeneron, and has received research support from Gilead and
834 Merck. J.D.G. declared contracted research with Gilead, Lilly, and Regeneron. P.D.G. is
835 on the Scientific Advisory Board of Affini-T, Catalio, Earli, Elpiscience, Fibrogen,
836 Immunoscapes, Metagenomi, Nextech, and Rapt; was a scientific founder of Juno
837 Therapeutics and Affini-T; and receives research support from Lonza and Affini-T
838 Therapeutics. H.C. reported consulting with Ellume, Pfizer, and the Bill and Melinda Gates
839 Foundation and has served on advisory boards for Vir, Merck and Abbvie. H.C. has
840 conducted CME teaching with Medscape, Vindico, Cataylst CME, and Clinical Care
841 Options. H.C. has received research funding from Gates Ventures, and support and

842 reagents from Ellume and Cepheid outside of the submitted work. The remaining authors
843 declare no competing interests.
844
845

846 **Figure Legends**

847

848 **Main Figures**

849 Fig. 1: Schematic overview of antigen-TCR pairing and multiomic analysis of T-cells
850 (APMAT) in COVID-19 participants

851 a. Participants in the study, broken down by COVID-19 severity. b. Experimental and
852 computational flow of APMAT. i. Construction of the large SCT-pMHC library representing
853 HLA-A*02:01 peptide-MHC complexes covering antigens from across the full SARS-CoV-
854 2 proteome. ii. Antigen-specific T cells are captured from pooled patient PBMCs with
855 barcoded SCT multimers for scRNAseq analysis. iii. Patient i.d.'s are assigned to each
856 single cell by matching SNP analysis from whole genome with scRNAseq data. iv-vii. For
857 each T cell, the physicochemical properties of the peptide antigen and the CDR3 β domain
858 of the TCRs are analyzed to identify statistical associations between the CDR3 β -peptide
859 antigen interface with T cell phenotype. viii. T cell clonotype persistence from acute
860 disease to convalescence is similarly associated with the physicochemical properties of
861 the TCR:antigen.

862

863 Fig. 2: APMAT enables high-throughput representation of whole SARS-CoV-2 genome

864 a. Graphic relating the SARS-CoV-2 proteome to putative HLA A*02:01 restricted
865 antigens (top row), to those SCT constructs that were expressed in usable yield (middle
866 row), and those that captured antigen-specific CD8 T cells from COVID-19 participants
867 (bottom row). For the SCT Expression row, darker red lines correspond to higher SCT
868 expression, while grey means low/no expression. For the Cell Capture row, darker red
869 lines mean more cells captured.

870 b. Bar plot showing the distribution of expressed SCTs for each SARS-COV-2 protein.
871 SCTs that successfully captured CD8 T cells are shown in colored bars. The color code
872 of the captured cells is that used in SARS-CoV-2 proteome of panel a.

873 c. Antigen sequence motif of expressed (top left) vs non-expressed (bottom left) SCT
874 constructs; and those expressed SCTs that did (top right) or did not capture CD8 T cells
875 (bottom right).

876 d. For each peptide, conventional anchor positions (R2, R9) and non-anchor residues are
877 assigned. The physicochemical properties are tabulated for each residue.

878 e. Bar plots comparing expressed (N = 560) vs non-expressed (N = 391) SCT constructs.
879 X-axis: SCT expression status. Note that two of the expressed SCT constructs were not
880 included for 10X experiment due to low sample volume. Y-axis: NetMHC prediction rank
881 (left) and average hydrophobicity of exposed residues (right).

882 f. Bar plots comparing expressed SCTs that captured CD8 T cells (N = 102) vs those that
883 did not (N = 456). X-axis: Cell capture status. Y-axis: NetMHC prediction rank (left) and

884 average charge of exposed residues (right). Charge values represent the average of
885 positive and negative charges rather than the absolute value.

886 For bar plots, data are presented as mean values +/- SEM, with corresponding individual
887 data points overlaid as hollow dots when possible. Dots outside of the range of y-axis
888 are not shown. The Statistical significance was determined using the two-sided Mann-
889 Whitney U test, and p values are annotated on all relevant plots with exact p-values
890 provided unless $p < 0.0001$.

891

892 Fig. 3: Three peptide groups distinguished by sequence physicochemical properties

893 a. All 951 putative SARS-CoV-2 peptides for HLA-A*02:01 are encoded by the
894 physicochemical properties of amino acids at each position for unsupervised clustering.

895 b. UMAP embedding of all peptides (Pep-UMAP) based on their physicochemical
896 properties with peptide clusters colored (legend on the right), each dot represents a
897 unique SARS-CoV-2 peptide.

898 c. Pep-UMAP colored with selected physicochemical properties, including average (Avg)
899 hydrophobicity (HPhobic), average polarity, whether the SCT expressed, and whether the
900 SCT captured CD8 T cells. Average values were calculated for all residues, including
901 anchor and exposed residues.

902 d. Left: Clustermap of the 951 peptide antigens by their normalized physicochemical
903 properties, revealing 3 major peptide groups (Pep-Groups). The key signatures that
904 distinguish the individual Pep-Groups are highlighted in red boxes. Right: Clustermap of
905 the Pep-Groups including only those peptides that captured T cells.

906 e. Pep-UMAP with densities of PG1-3 depicted, legend on the bottom.

907 f. Left: Violin plot of peptide hydrophobicity for Pep-Groups, sorted by mean value. Middle:
908 SCT protein expression efficiency for the Pep-Groups. Right: SCT cell capture efficiency
909 for Pep-Groups. Mean values +/- SEM are utilized for violin plots. The Statistical
910 significance was determined using the two-sided Mann-Whitney U test, and p values are
911 annotated on all relevant plots with exact p-values provided unless $p < 0.0001$.

912

913 Fig. 4: Pep-groups associate with different T cell phenotypes.

914 a. Illustration of the mapping of peptide physicochemical properties on to T cell gene
915 expression profiles. CD8 T cells were captured by SCT-dextramers from 62 HLA-matched
916 COVID-19 participants for scRNAseq and TCR sequencing. Gene expression UMAP
917 (GEX-UMAP) was generated based on scRNAseq. We used the SCT identity to connect
918 Pep-Groups with gene expression.

919 b. GEX-UMAP of Sars-CoV-2-specific CD8 T cells with different phenotypes color-
920 encoded (legend on the top), and then color-coded by expression levels of selected
921 mRNA transcripts.

922 c. GEX-UMAP, color-coded by (left) the polarity of exposed peptide antigen residues;
923 (middle) all T cells specific to a given antigen; and (right) SARS-CoV-2 specific T cells
924 from two study participants.

925 d. Top: GEX-UMAP color encoded with T cell densities specific for antigens from each
926 Pep-Group. Bottom: Bar plot of the relative abundance of phenotypes for T cells
927 associated with each Pep-Group. Note that the UMAP densities in the original were
928 calculated as an odds ratio.

929 e. Top enriched biological pathways of genes significantly elevated in cells captured by
930 PG3 relative to PG2. Adjusted p-values are generated by EnrichR

931

932 Fig. 5: TCR hydrophobicity is an important factor for effector function.

933 a. Each TCR beta chain was split into V, CDR3 β mer, and J regions, and then encoded
934 by the physicochemical properties of amino acids at each residue position for overlay on
935 the GEX-UMAP.

936 b. GEX-UMAP, color-coded by phenotypes (top) and overlaid with the percentage of
937 hydrophobic residues within CDR3 β mer (HPhobic %) (bottom).

938 c. The top differential TCR physicochemical properties for effector (Cytotoxic, EM and
939 Hybrid) phenotypes, relative to non-effector (CM and naïve) phenotypes. The Statistical
940 significance was determined using the two-sided Mann-Whitney U test, p-values < 0.0001
941 were marked with ****. See color key inset. Source data are provided as a Source Data
942 file.

943 d. Bar plots showing the variation of selected TCR physicochemical properties across T
944 cell phenotypes. Number of cells: Naïve (n = 68), CM (n = 133), Hybrid (n = 43), EM (n =
945 93), Cytotoxic (n = 384). For bar plots, data are presented as mean values +/- SEM. p
946 values are annotated on all relevant plots with exact p-values provided unless p < 0.0001.
947 Source data are provided as a Source Data file.

948 e. Top: Binary TCR-Groups (HPhobic-Low and HPhobic-High) are defined based on the
949 median value of hydrophobic residue percentages for all SARS-CoV-2 cells. Bottom:
950 Stacked bar plot of clonal size distribution for each TCR-Groups with legend on the right.
951 f. GEX-UMAP with densities of the TCR-Groups projected. The relative abundance of
952 phenotypes in HPhobic-High cells is plotted at bottom. The UMAP densities were
953 calculated as an odds ratio.

954 g. Top enriched biological pathways of genes significantly elevated in cells with HPhobic-
955 High vs HPhobic-Low TCRs. Adjusted p-values are generated by EnrichR.

956

957 Fig. 6: Combination of peptide-TCR features associated with cell phenotypes

958 a. We investigated the combination of peptide and TCR properties for antigen-TCR paired
959 SARS-CoV-2 CD8 T cells. The plotted peptides are those that captured cells.

960 b. Cell distribution between peptide groups (Left) and TCR groups.

961 c. Radar graphs of physicochemical properties for representative PG-TCR groups. Blue
962 and grey shaded areas of the outer rings indicate Pep-Group, and TCR β properties,
963 respectively. Each axis displays the normalized average value for each property, with

964 lowest value in the center. The shaded polygons reflect the property space occupied by
965 the peptide-TCR β groupings. Legend on the bottom.
966 d. Top: Cell percentage of Cytotoxic, EM and Naïve phenotypes for each PG-TCR group.
967 Bottom: Heatmap showing selected mRNA levels for each PG-TCR group. Underlined
968 PG-TCR groups are those from panel c.
969 e. Heatmap showing average value of peptide and TCR physicochemical properties for
970 each phenotype.

971
972 Fig. 7: Longitudinal analysis for PG-TCR groups
973 a. Tracking SARS-CoV-2 CD8 T cells from acute to post-acute timepoint based on TCR-
974 GLIPH query from a previously reported longitudinal dataset.
975 b. Cell abundancy changes from post-acute to acute timepoint for each PG-TCR grouping.
976 Red or blue means higher abundancy at acute or post-acute timepoint, respectively.
977 c. Selected gene expression changes from the acute to post-acute timepoints for each
978 PG-TCR groups. Gene annotations on the bottom.
979 d. Summary: The association between distinct peptide-TCR properties and cell fates for
980 antigen-specific CD8 T cells.

981
982

983 Supplementary Figures

984 Supplementary Fig. 1: Technical validation and methods for antigen and patient
985 assignment

986 a. Heatmap where each row is a Hashtag used for each SCT-dextramer pool and each
987 column is a single cell. Legend at bottom.
988 b. Scatter plots of antigen assignment for representative CMV and SARS-CoV-2 antigens.
989 X-axis is the SCT intensity, defined by numbers of UMIs mapped to the SCT-dextramer.
990 Y-axis is the SCT dominance, defined by percent of the cell's SCT-dextramer associated
991 UMIs that mapped to this antigen's SCT-dextramer. Cells assigned with the antigen are
992 shown in red positive zone, and have SCT intensity >25, and SCT dominance >25%.
993 c. Patient assignment by comparison of the WGS SNPs and derived de novo SNPs
994 derived from scRNA-seq data, exemplified by Donor-58 and Donor-194. Each panel
995 represents the comparison for each chromosome. Allele frequency and nucleotide identity
996 of reference and alternate are shown in lower right legend.
997 d. Heatmap with rows as sex specific genes (RPS4Y1 for male, XIST for female), columns
998 are assigned participants. Upper row indicates patient sex from clinical records with blue
999 as male and pink as female. Legend on upper right.

1000

1001 Supplementary Fig. 2: SARS-CoV-2 SCT expression

1002 a. Left: Number of constructed SCT plasmids for putative peptides. Middle: Expressed
1003 SCT constructs with useable protein expression yield that were constructed as DNA-

1004 barcoded SCT dextramers during 10X Chromium experiment. Note that two of the
1005 expressed SCTs were not included due to low sample volume. Right: SCTs that captured
1006 SARS-CoV-2 CD8 T cells from COVID-19 participants.

1007 b. Immunodominant A*02:01 epitopes among individuals detected by SCT. Left: CD8+ T
1008 cell distribution for each antigen specificity (rows) identified from individual COVID-19
1009 infected participants (each color represents a different participant). Right: total cell counts
1010 for each antigen (log₁₀ scale). Only antigens assigned with more than 5 cells were plotted.

1011

1012 Supplementary Fig. 3: TCR Validation via in-vitro Lenti-virus transduction

1013 a. Flow cytometry gating strategy.

1014 b. Selected SCT pMHC tetramer assay on untransduced or TCR-transduced Jurkat cells.
1015 TCR#1, previously published by Chour et.al., provides a positive control. We further
1016 validate TCR#2 and #3, which are identified from this dataset.

1017

1018 Supplementary Fig. 4: Agreement between SCT expression and NetMHCpan prediction

1019 a. Bar plots comparing expressed vs non-expressed SCT constructs. X-axis: SCT
1020 expression status. Y-axis: Normalized SCT expression yield (left) and NetMHCpan
1021 predicted binding affinity between peptide and HLA-A*02:01 (right).

1022 b. Scatter plot for normalized SCT expression yield (Y-axis) and predicted binding affinity
1023 (X-axis) for all expressed SCTs, each dot represents a SCT construct.

1024 c. Bar plots comparing three groups of peptides based on NetMHCpan prediction output.
1025 Y-axis: average hydrophobicity and polarity of peptides within each prediction group.

1026 d. For predicted strong binders, a comparison of the average hydrophobicity and polarity
1027 for the ones that showed no/low SCT expression, and the ones showed successful SCT
1028 expression.

1029 e. Pie chart showing the SCT expression percentage for peptides in each prediction group.
1030 Total peptide number were shown in the top.

1031 For bar plots, data are presented as mean values +/- SEM. The Statistical significance
1032 was determined using the two-sided Mann-Whitney U test, and p values are annotated
1033 on all relevant plots with exact p-values provided unless $p < 0.0001$.

1034

1035 Supplementary Fig. 5: Characteristics of Pep-UMAP and Pep-Groups

1036 a. Pep-UMAP colored with selected physicochemical properties, including anchor
1037 hydrophobicity, hydrophobicity of the 6th residue, average hydrophobicity, bulkiness,
1038 polarity, and absolute charge of exposed (Exp) residues.

1039 b. UMAP embedding of Pep-Groups.

1040

1041 Supplementary Fig. 6: Characteristics of GEX-UMAP

- 1042 a. UMAP embedding of gene expression (GEX-UMAP) for SARS-CoV-2 specific CD8 T
1043 cells color-coded by expression levels of selected mRNA transcripts, and UMAP leiden
1044 groups.
1045 b. Dot plot showing normalized expression levels of selected marker genes in each T cell
1046 phenotype. The size and color of each dot represent the fraction of expressing cells and
1047 the mean of normalized expression levels in each phenotype.
1048 c. Cluster map of phenotype-related genes and surface protein measured via scCITE-seq
1049 from Su et. al (2022)²⁵. Values are row-normalized, legend on the right. Proteins were
1050 marked as red triangle next to surface marker's name while mRNAs were unlabeled.
1051 d. Stacked bar plot of clone size distribution for T cells captured by threes PG-group with
1052 legend on the right.

1053
1054 Supplementary Fig. 7: Characteristics of TCR Groups

- 1055 a. GEX-UMAP colored with selected physicochemical properties of the TCR β sequences.
1056 b. Stacked bar plot of clone size distribution for each of the T cell phenotypes with legend
1057 on the right.
1058 c. Distribution of CDR3 β mer hydrophobicity (top), CDR3 β length (middle), CDR3 β mer
1059 absolute charge (bottom) for each of the TCR-Groups. Bar plots for the respective
1060 property are on the right.
1061 d. Bar plots for the mRNA levels of LAG3 and TIGIT.
1062 For bar plots, data are presented as mean values +/- SEM. The Statistical significance
1063 was determined using the two-sided Mann-Whitney U test, and p values are annotated
1064 on all relevant plots with exact p-values provided unless $p < 0.0001$.

1065
1066 Supplementary Fig. 8: Validation after removal of dominant TCR clones

- 1067 a. CD8⁺ T Cell clonotype distribution for each antigen specificity (rows) (for each antigen,
1068 each color represents cells expressing a unique TCR). Only top 15 were plotted.
1069 b. Bar plots evaluating TCR-Groups after removal of large clones. Y-axis: average
1070 expression level of representative genes and scores related with Cytotoxicity (top) and
1071 Naïve-Memory (bottom). Red (HPhobic-High) and black (HPhobic-Low) bars represent
1072 cells within each TCR-Group. Legend on the right. For bar plots, data are presented as
1073 mean values +/- SEM. The Statistical significance was determined using the two-sided
1074 Mann-Whitney U test, and p values are annotated on all relevant plots with exact p-values
1075 provided unless $p < 0.0001$.

1076
1077
1078 Supplementary Fig. 9: Characteristics of PG-TCR groups

- 1079 a. Radar graphs of physicochemical properties for example PG-TCR groups. Blue and
1080 grey shaded areas of the outer rings indicate Pep-Group, and TCR β properties,
1081 respectively. Each axis displays the normalized average value for each property, with

1082 lowest value in the center. The shaded polygons reflect the property space occupied by
1083 the peptide-TCR β groupings. Legend on the bottom.
1084 b. GEX-UMAP color encoded with densities of PG-TCR groups based on each cell's
1085 antigen specificity and TCR sequence.
1086 c. Stacked bar plot of phenotype distribution for each of the PG-TCR groups with legend
1087 on the right.
1088 d. Surface protein validation with respect to findings in Main Fig. 6. Heatmap with X-axis
1089 as PG-TCR group assignment and Y-axis as level of a given protein normalized per row
1090 and column, see legend on right.
1091 e. Stacked bar plot of PG-TCR group distribution for acute and post-acute timepoint with
1092 legend on the right.
1093 f. Clustermap of cells by expression levels of selected mRNA transcripts for each
1094 timepoint. Legend on the right.
1095 g. Bar plots comparing expression levels of selected mRNA transcripts for cells at acute
1096 and post-acute timepoint. Number of cells for each group: PG3:High Acute (n=448), Post-
1097 Acute (n=35); PG1: Acute (n=116), Post-Acute (n=156). For bar plots, data are presented
1098 as mean values +/- SEM. The Statistical significance was determined using the two-sided
1099 Mann-Whitney U test, and p values are annotated on all relevant plots with exact p-values
1100 provided unless $p < 0.0001$.

1101
1102 Supplementary Fig. 10: TCR alpha chain analysis

1103 a. Analysis of which physicochemical characteristics of the TCR alpha chain exhibit
1104 significant associations with T cell phenotype. The top differential TCR CDR3 α
1105 physicochemical properties for effector (Cytotoxic, EM and Hybrid) phenotypes, relative
1106 to non-effector (CM and naïve) phenotypes. The Statistical significance was determined
1107 using the two-sided Mann-Whitney U test, and p values are annotated on all relevant plots
1108 with exact p-values provided unless $p < 0.0001$.

1109
1110 Supplementary Fig. 11: External dataset validation

1111 a. Bar plots comparing expression levels of selected mRNA for cells within each TCR-
1112 Groups. Number of cells in each group: HPhobic-High (n = 1841), HPhobic-Low (n =
1113 1988). Original dataset from Fischer et. al (2021)⁵⁸. For bar plots, data are presented as
1114 mean values +/- SEM. The Statistical significance was determined using the two-sided
1115 Mann-Whitney U test, and p values are annotated on all relevant plots with exact p-values
1116 provided unless $p < 0.0001$.

1117 b. Pearson correlation coefficients between selected mRNA levels and TCR CDR3 β
1118 properties for CMV-NLVP (left) and Influenza-GILG (right) specific CD8 T cells. Original
1119 dataset from Chen et. al (2023)⁸.

1120

1121 **Supplementary Tables**

1122 Table S1: Clinical characteristics, medical history of INCOV sub cohort in this study

1123 Table S2.1: List of putative SARS-CoV-2 peptide antigens for SCT-pMHC expression

1124 Table S2.2: Vireo output for SNP demultiplexing

1125 Table S2.3: TCR gene usage for top 5 SARS-CoV-2 epitopes

1126 Table S2.4: List of TCRs used for validation (via lenti-virus transduction and SCT-tetramer
1127 binding assay)

1128 Table S3: Amino acid property scales used for peptide and TCR residues

1129 Table S4: Pep-UMAP characteristics

1130 Table S5: Filtered DEGs in cells assigned to PG3 compared to PG2 antigen

1131 Table S6: Enriched pathways (Reactome 2022) for upregulated genes in PG3 vs PG2

1132 Table S7: Log2 fold change of TCR physicochemical properties, effector phenotypes vs
1133 non-effector phenotypes

1134 Table S8: Mann-Whitney U test of TCR physicochemical properties, effector phenotypes
1135 vs non-effector phenotypes

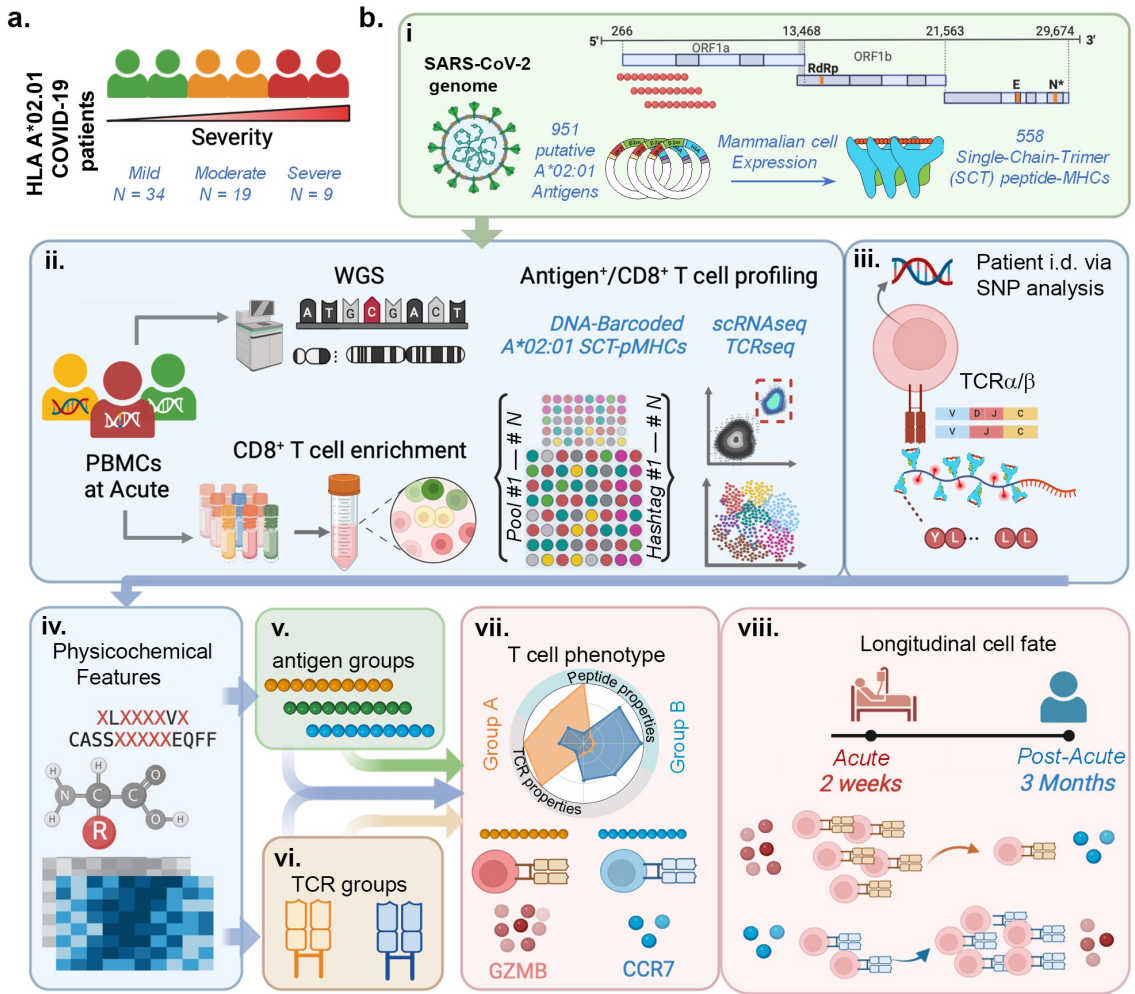
1136 Table S9: DEGs in cells in HPhobic-High compared to HPhobic-Low

1137 Table S10: Enriched pathways (Reactome 2022 and GO Biological Process 2023) for
1138 upregulated genes in HPhobic-High vs HPhobic-Low

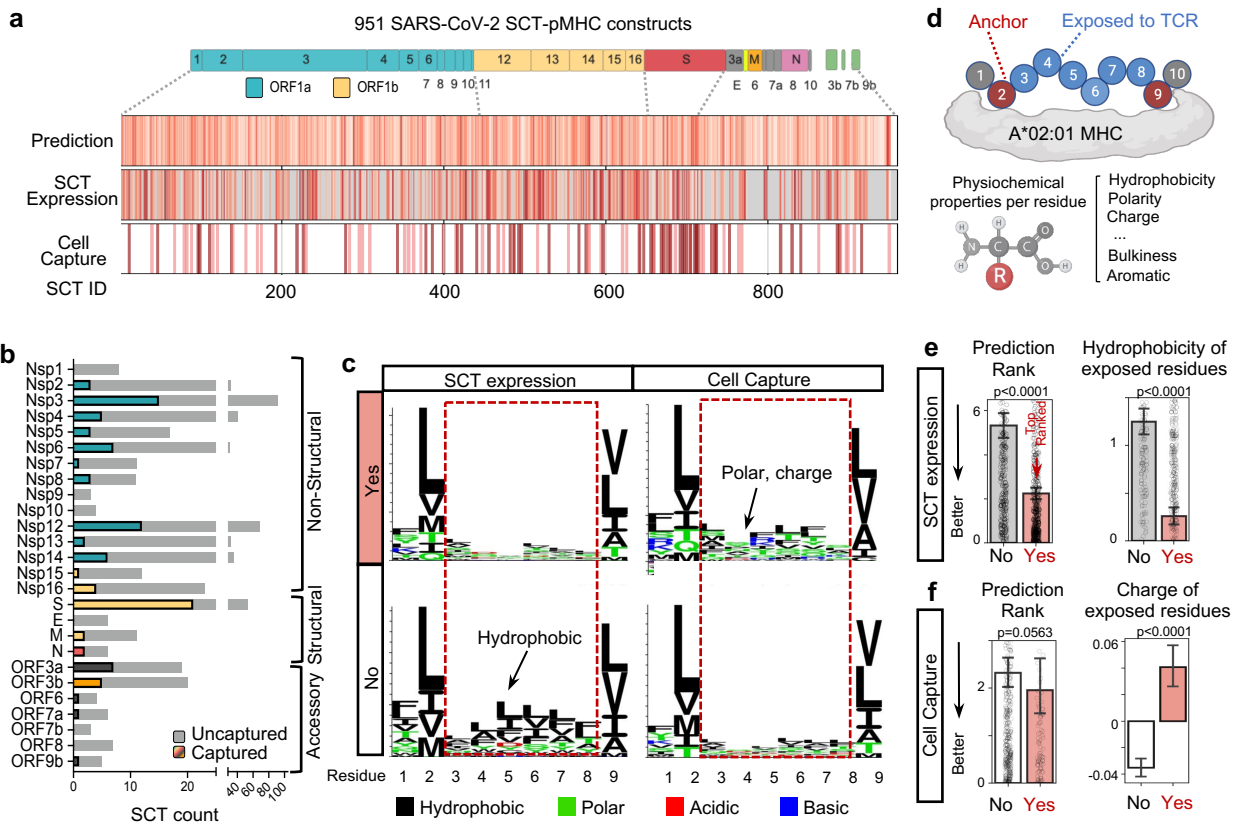
1139 Table S11: Cell count at Acute and Post-Acute time points

1140

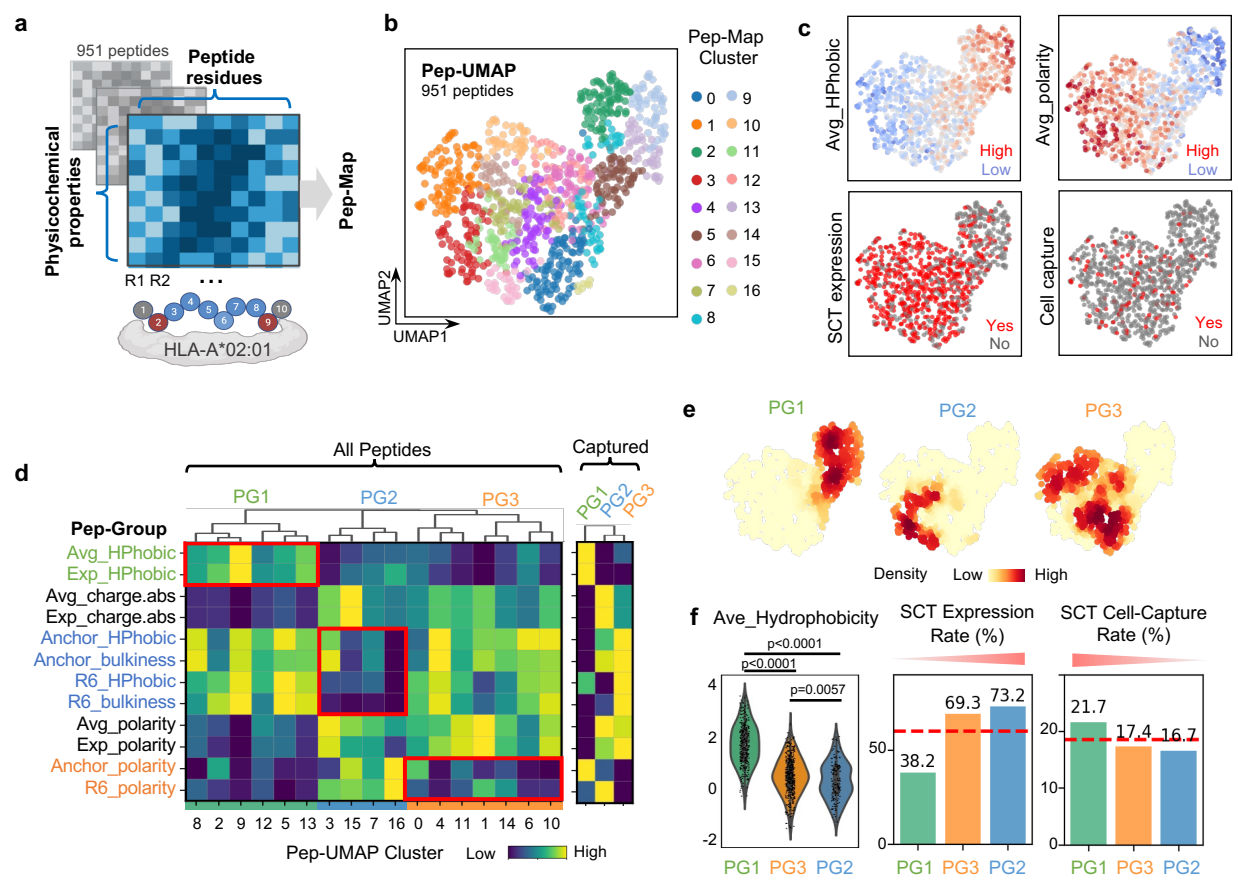
1141



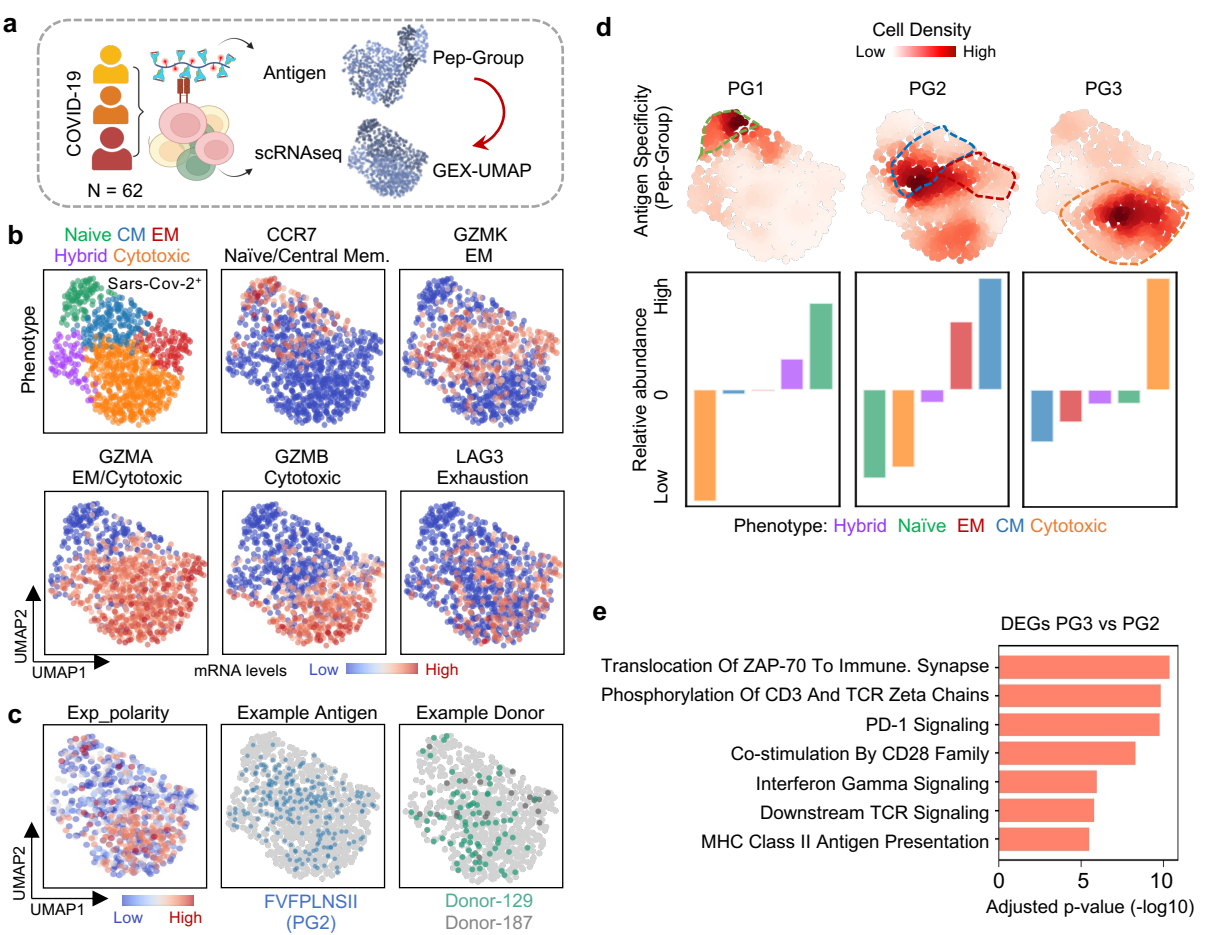
Main Fig. 1



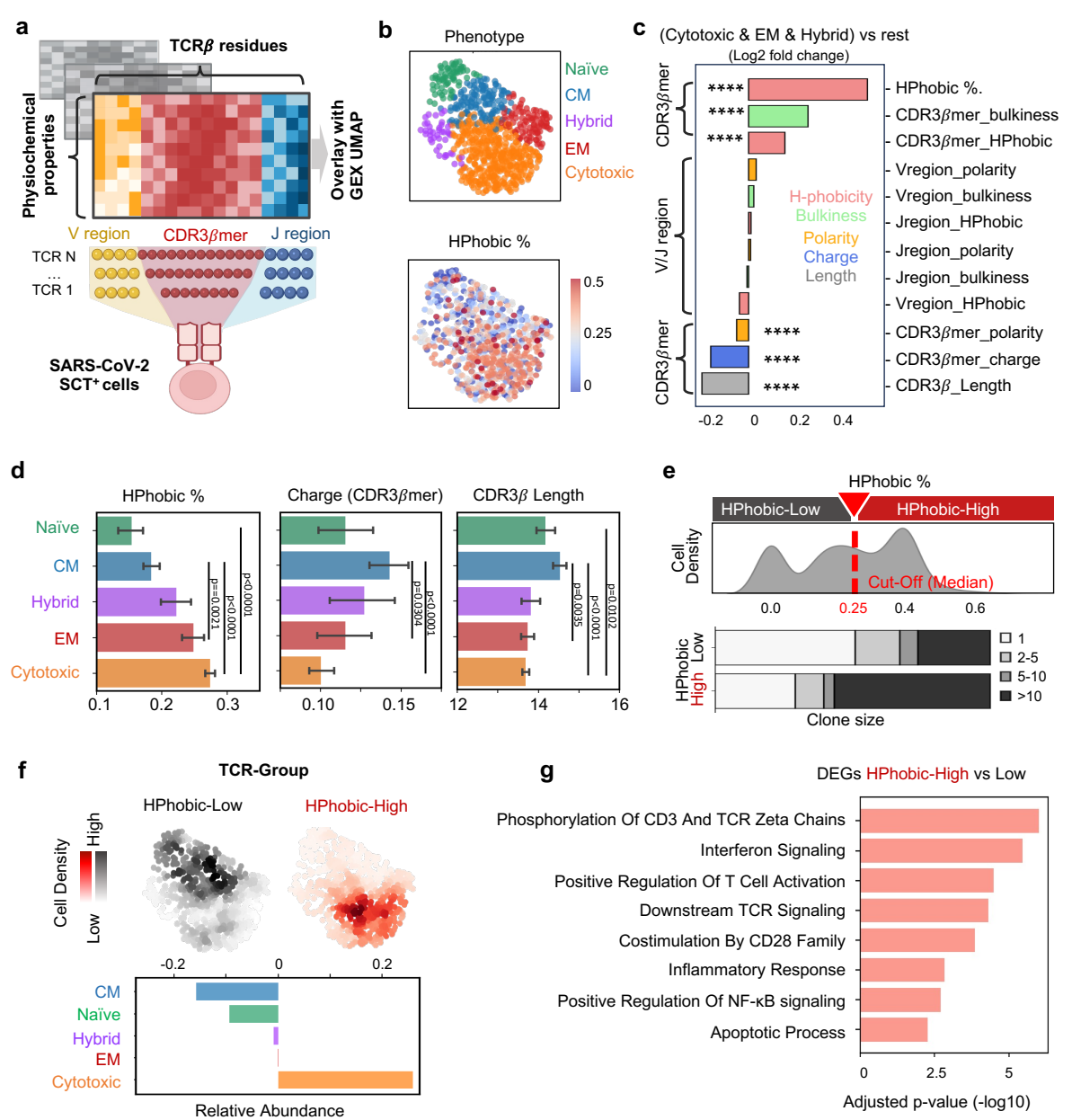
Main Fig. 2

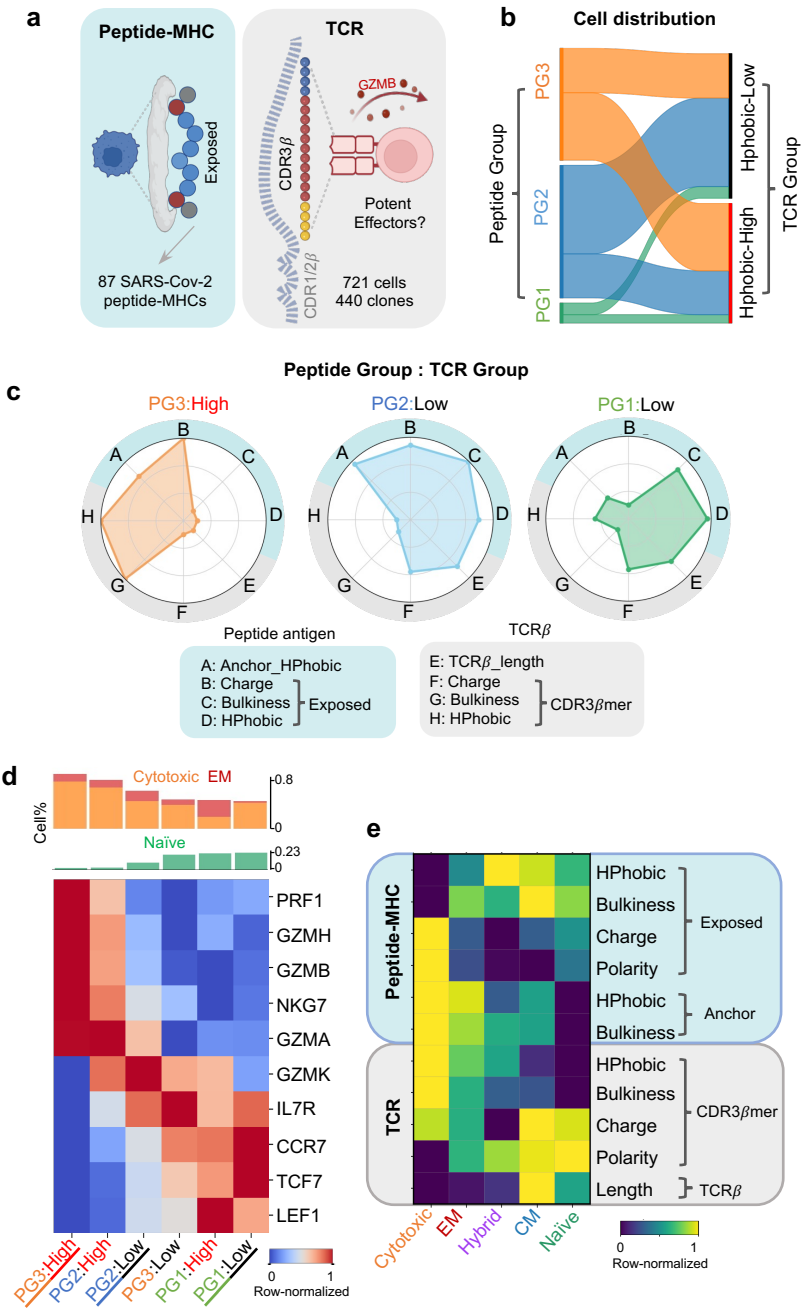


Main Fig. 3



Main Fig. 4





Main Fig. 6

