



Published in final edited form as:

Nat Immunol. 2021 February ; 22(2): 216–228. doi:10.1038/s41590-020-00836-7.

Gut CD4⁺ T cell phenotypes are a continuum molded by microbes, not by T_H archetypes.

Evgeny Kiner^{1,2,9}, Elijah Willie³, Brinda Vijaykumar^{1,2}, Kaitavjeet Chowdhary^{1,2}, Hugo Schmutz^{1,2}, Jodie Chandler⁴, Alexandra Schnell², Pratiksha I. Thakore⁵, Graham LeGros⁴, Sara Mostafavi^{6,7,8}, Diane Mathis^{1,2}, Christophe Benoist^{1,2}, Immunological Genome Project Consortium

¹Department of Immunology, Harvard Medical School ²Evergrande Center for Immunologic Diseases, Harvard Medical School and Brigham and Women's Hospital, Boston, USA ³Bioinformatics Program, University of British Columbia, Vancouver, Canada ⁴Malaghan Institute of Medical Research, Wellington, New Zealand ⁵Klarman Cell Observatory, Broad Institute of MIT and Harvard, Cambridge, USA ⁶Departments of Statistics and Medical Genetics, University of British Columbia, Vancouver, Canada ⁷Canadian Institute for Advanced Research, Toronto, Canada ⁸Vector Institute, Toronto, Canada ⁹Present Address: Immunai, New York, NY, USA

Abstract

CD4⁺ effector lymphocytes (Teff) are traditionally classified by the cytokines they produce. To determine the states that Teff actually adopt in frontline tissues *in vivo*, we applied single-cell transcriptome and chromatin analysis on colonic Teff cells, in germ-free or conventional mice, or after challenge with a range of phenotypically biasing microbes. Subsets were marked by expression of interferon-signature or myeloid-specific transcripts, but transcriptome or chromatin structure could not resolve discrete clusters fitting classic T_H subsets. At baseline or at different times of infection, transcripts encoding cytokines or proteins commonly used as T_H markers distributed in a polarized continuum, which was also functionally validated. Clones derived from single progenitors gave rise to both IFN- γ and IL17-producing cells. Most transcriptional variance was tied to the infecting agent, independent of the cytokines produced, and chromatin variance primarily reflected activity of AP1 and IRF transcription factor families, not the canonical subset master regulators T-bet, GATA3, ROR γ .

Users may view, print, copy, and download text and data-mine the content in such documents, for the purposes of academic research, subject always to the full Conditions of use:http://www.nature.com/authors/editorial_policies/license.html#terms

*Address correspondence to: Diane Mathis and Christophe Benoist, Department of Immunology, Harvard Medical School, 77 Avenue Louis Pasteur, Boston, MA 02115, cbdm@hms.harvard.edu, Phone: (617) 432-7741, Fax: (617) 432-7744.

CONTRIBUTIONS

E.K. and E.W. performed experiments. E.K., B.V., K.C., H.S., S.M. and C.B. analysed and interpreted data. A.S., P.I.T., J.C. and G.L. provided data or reagents. E.K., S.M., D.M. and C.B. designed the study and wrote the manuscript.

COMPETING INTERESTS

The authors declare no competing interests.

INTRODUCTION

Effector CD4⁺ T lymphocytes (Teff) are key drivers of both humoral and cellular immune responses, orchestrating adaptive (antibodies, cytotoxic cells) and innate (macrophage, granulocytes) immune responses. This range of abilities has long raised the issue of functional diversity, which was documented by functional assays even before the molecular identification of MHC and TCR molecules, the central axis of T cell activation and differentiation^{1,2}. A key advance was the demonstration that functional phenotypes of different T cell clones were keyed to the cytokines they produce^{3,4}, coining the T_H1/T_H2 nomenclature. T_H1s secrete IFN- γ and mainly support inflammatory and cytotoxic responses; T_H2s produce IL4, IL5 or IL13 and principally help B cells produce antibodies. This division has since been revised several times to add more subsets (IL17-secreting T_H17s, IL9-secreting T_H9s, follicular helpers (T_{FH})⁵⁻⁷), but the core notion that Teff belong to discrete and largely stable states defined by the cytokines they produce has endured^{8,9}. Different types of infectious or allergic challenges elicit different Teff “flavors” (T_H1 generally associated with intracellular pathogens; T_H2 with helminth parasites, T_H17 with bacterial and fungal infections), and these T_H distinctions also have implications for immune-mediated diseases¹⁰. Indeed, the T_H paradigm has elicited parallel cosmologies in macrophages, $\gamma\delta$ T cells or innate lymphoid cells (ILC)¹¹.

However, this model was questioned almost since its inception^{12,13}. First, because its attractive simplicity could lead to shoe-horning of immune functions (e.g., 1990s’ publications erroneously tagged immune diseases to either T_H1 or T_H2). Second, many reports documented that secretion of IFN- γ , IL4 or IL17 is not always mutually exclusive¹⁴. Plasticity between T_H subtypes was demonstrated, suggesting that these cell states are not as stable and terminally differentiated as originally inferred from T_H lines grown in supra-physiological cytokine concentrations^{9,13}. Further, while some cell-surface markers were proposed as indicators of differentiated T_H types, they often proved non-exclusive. Thus, T_H subsets have most precisely been defined *in vitro*, their *in vivo* counterparts remaining elusive.

Here, we aimed to assess the spectrum of phenotypic states that Teff can adopt *in vivo*, leveraging the unbiased potential of single-cell genomics¹⁵. In essence returning to the clonal analysis that founded the T_H paradigm^{3,4}, but now with the ability to evaluate the entirety of a cell’s transcriptome and chromatin structure, rather than only a few pre-selected cytokines or markers. We analysed T cells in the colonic lamina propria (LP), a front-line tissue under continuous and diverse challenge, comparing CD4⁺ T cells in mice under germ-free conditions, carrying normal commensal microbiota, or infected with agents that elicit diversely biased Teff responses. The results indicate that Teffs form a continuum in the transcriptional space, but highlight some novel phenotypes. The production of key cytokines did show skewed distributions, but these did not identify the discrete cell clusters that might have been expected from the T_H paradigm.

RESULTS

A continuum of effector phenotypes in colonic CD4⁺ Teff cells

To probe the transcriptional landscape of CD4⁺ Teff cells in an unbiased manner, we performed scRNAseq on total CD4⁺ T cells from the colonic LP starting with conventionally-housed (SPF) C57BL/6 mice. Two experiments were performed with droplet-based scRNAseq (Extended Data Fig. 1a,b; Supplementary Table 1; datasets analyzed individually, replicates serving to confirm conclusions). It was straightforward to parse, with standard clustering, CD4⁺ T cells into the four main groups expected from flow cytometry (Fig. 1a): regulatory T cells (Treg; *Foxp3*⁺), and their *Rorc*⁺ and *Ikzf2*⁺ (Helios) subsets, naive T conventional cells (naive Tconv; *Cd44*⁻ *Ccr7*⁺), and Teff cells (*Cd44*⁺ *Ccr7*⁻). To assess the influence of the commensal microbiota on this distribution, we generated scRNAseq datasets of colonic CD4⁺ T cells from SPF and germ-free mice (Fig. 1b), revealing similar clusters, with fewer ROR γ ⁺ Tregs and Teff cells in germ-free as expected.

To assess what phenotypic states gut Teff cells can adopt, we reclustered the Teff population from SPF mice. Here, with the exception of cycling cells, we could not observe any clear partitioning of cells, but rather a quasi-continuous cloud (Fig. 1c, left). To search for distinctions corresponding to the major recognized Teff types, we manually curated from published signatures short but robust and highly specific genesets, which included the defining cytokines, driving transcription factors, and a few correlated transcripts, but left out generic activation-associated or poorly specific transcripts (Supplementary Table 2). The T_H2 signature showed polarized expression, while cells expressing the T_H17- and especially T_H1-associated signatures were dispersed more widely across the continuum (Fig. 1c, right). To ensure that this continuum was not due to the high dropout rate of scRNAseq, we reanalyzed a published dataset from colonic Teff cells that included fewer cells, but sequenced to greater depths¹⁶. These also showed a continuous distribution and dispersion of the T_H signatures (Extended Data Fig. 1c). If cytokines do not represent the main axes of variance in colonic T cells, what does? To this end, we used a simple clustering strategy, which showed that the driving variance lies in the degree of activation of Teff cells, represented by typical activation transcripts like *Cd69* or *Nr4a1* (Fig. 1d). Teff cells with a lower degree of activation over-expressed *Klf2* and *S1pr1*, a combination shown to restrain CD4⁺ T cell differentiation¹⁷. Thus, the main heterogeneity of Teff in the colonic LP corresponds to a gradient of activation in response to commensal microbiota, but not predominantly to commitment to producing one cytokine or the other.

Different intestinal infections elicit divergent Teff phenotypes

It thus seemed difficult to identify discrete T_H1 or T_H17 cell populations in normal mice. We hypothesized that under baseline conditions Teffs were only partially polarized because of incomplete activation by commensals, with only “stubs” of more differentiated states that the cells could potentially reach. We thus polarized further the T cell pools by infecting mice with pathogens known to elicit biased immune responses: (i) AroA *Salmonella enterica* (serovar *Typhimurium*), a non-invasive mutant that elicits IFN- γ -dominated responses; (ii) *Citrobacter rodentium*, a strong IL-17 inducer; (iii) *Heligmosomoides polygyrus* and *Nippostrongylus brasiliensis*, two helminths that provoke prototypic Type-2 responses (Fig.

2a). Infection times (11-13 days) allowed responses to develop and achieve full bias. Flow cytometry confirmed the production of the expected cytokines, including some IFN- γ /IL17A double-producer cells, as expected (Fig. 2b). In a first experiment, Teffs from control or infected mice were tagged with DNA-coded antibodies¹⁸, “hash-tagged”) and comingled for sorting, microfluidic bead capture and library construction, making for a robust intra-batch comparison (Fig. 2a). As in uninfected mice, CD4⁺ T cells clustered into Tregs, naive T cells, and Teffs (Extended Data Fig. 2a).

Teffs were then considered on their own, with dimensionality reduction on tSNE (Fig. 2c) or UMAP (Extended Data Fig. 2b) plots, which revealed a dominant partitioning according to the infectious agent. Outside the main “blob”, some Teffs did break out into discrete populations but we could not detect well-demarcated cell clusters that expressed the characteristic T_H genesets. These mapped to skewed but broad swaths (cells with high T_H2 geneset best demarcated, T_H17 geneset biased but dispersed, cells with high T_H1 geneset found almost throughout; Fig. 2c, Extended Data Fig. 2b). This lack of segregation was robust to the choice of genesets (if anything more diffuse using another curated signature set based on ref¹⁹ (Extended Data Fig. 2c). The expression of *Ifng* and *Il17a* transcripts also showed overlaps consistent with the double-producer cells detected by flow cytometry (Extended Data Fig. 2d). These conclusions were true also for a replicate set of colonic CD4⁺ Teff from mice infected with the same pathogens (Extended Data Fig. 2e). The dominant influence of the infectious microbe, over the T_H phenotype marked by cytokine production, was objectivized by comparing the overall Euclidean distance between all cells expressing *Il17a* and *Ifng* transcripts from the different conditions: Teff expressing *Ifng* or *Il17a* transcripts from any infection were much closer than their cytokine-sharing counterparts in the other mice (Fig. 2d).

We applied a panel of clustering and biclustering algorithms in attempts to break up the cell cloud into clusters that coincided with expression of the T_H signature sets, but none of the clusters thus generated were uniquely enriched for any one T_H signature or cytokine (Extended Data Fig. 3a–c). To objectively verify the continuity in the distribution of transcriptomes of Teff cells, we used Hartigan’s Dip Test of multi-modality²⁰, after applying a projection defined by minimum separation hyperplane²¹ to the expression of the most variable genes. The results showed that Teff, Treg and naive Tconv significantly segregate by Hartigan’s test (Fig. 2e, top), while there was no significant break in the distances within the Teff pools (Fig. 2e, middle and bottom). These results confirmed that Teffs occupy a continuum point cloud that is not easily separable into distinct clusters.

One explanation for this continuous Teff distribution was that they included different subsets of the canonical T_H1/2/17 archetypes. However, projection of differentiating genes reported for the “pathogenesis subsets” within T_H17 cells^{22,23} did not demarcate distinct subsets of IL17-producing cells, although showing skewed distribution more generally (Extended Data Fig. 3d). Similarly, a reported distinction between “homeostatic” and “inflammatory” T_H17 cells²⁴, the latter elicited by *C. rodentium* infection, may have mostly resulted from infection rather than from distinct T_H17 subsets, as the corresponding signature did not specifically demarcate IL-17-producing cells (Extended Data Fig. 3e).

Deep machine learning tools can efficiently discover combinatorial and non-linear patterns that are difficult to discern conventionally. In another attempt to identify patterns that would uniquely identify IL-17 or IFN- γ producing cells, we optimized and trained a Deep Neural Network (DNN) to classify cells into IL-17 vs IFN- γ producing groups based on their single-cell transcriptomes. As a positive control, this architecture could be trained to recognize Teff and Treg cells from the held-back test set (Methods). The DNN did partially identify *Ifng* or *Il17a* positive cells in the test set (Extended Data Fig. 4a,b; 90.2% and 60.7 % accuracy on *Ifng* and *Il17a*-positive cells, respectively). However, using the Integrated Gradients method to measure the importance of the transcripts used by the model to support this identification showed little reproducibility in independent training runs (Extended Data Fig. 4c). Beyond a few transcripts known to correlate with *Il17a* (*Tmem176a*, *Capg*); only *Il22* had a strong and reproducible influence, an internal control given its known co-regulation with *Il17*. Indeed, when *Il22* was left out, prediction efficacy dropped to 28.7%. Hence, even with a pliable artificial intelligence tool, it seemed difficult to identify robust T_H1 or T_H17 transcriptome patterns.

Finally, we assessed the distribution of surface markers that have been associated with T_H subsets and are commonly used for cell sorting (*Ccr5* and *Cxcr3* for T_H1; *Ccr6* and *Il1r2* for T_H17). *Ccr6* and *Il1r2* proved to be mutually exclusive, with only partial overlap with cells transcribing *Il17a* (Fig. 2f). *Ccr5* and *Cxcr3* transcripts were widely distributed across the cloud, with only partial overlap with the T_H1 signature. Flow cytometric analysis of LP cells after *Salmonella* infection confirmed these results (Fig. 2f). Thus, not only were classic T_H subsets not clearly identifiable in the transcriptional data, but the flow cytometry markers used to identify them had limited congruence in this context.

Teff phenotypes distinguished by infecting agents, not by T_H type

Colonic Teff cells clustered according to the type of infection, rather than by the cytokine they expressed. Accordingly, analysis of variable transcripts present in *Il17a*- and *Ifng*-expressing cells revealed divergent patterns, with blocs of co-expressed transcripts that largely aligned with the infection (Fig. 3a). To validate this result and exclude technical pitfalls in scRNAseq, we used an *Il17a-gfp* reporter mouse line and performed population RNA-seq on colonic GFP-positive cells, at baseline or after infection with *Salmonella* or *Citrobacter* (Fig. 3b). Echoing the single-cell data, Principal Components Analysis (PCA) showed that IL-17A⁺ cells from each condition clustered separately from each other (Fig. 3c). The direct comparison of IL17A⁺ cells from *Salmonella* or *Citrobacter*-infected mice brought forth 277 differential transcripts (at FC>2, FDR <0.05; Fig. 3d). Among this set, transcripts with differential representation in the single-cell data showed similar biases. Thus, the majority of changes imparted by infection were unrelated to *Il17* or *Ifng* expression, or membership in a T_H class.

The primary determinants of Teff variability

Their expression patterns within the projection plots of Figs 1c/2c indicated that prototypic T_H1 or T_H17 signature sets did not mark discrete sets of cells. To turn the question to a gene-centric perspective, we asked what co-regulated modules of transcripts existed among these CD4⁺ Teff cells, and whether these might track with cytokine production. First, a PCA

showed that the genesets in the PCs with most variance contained few Th-associated signature genes (Extended Data Fig. 5a). Next we analyzed gene-gene correlation, leveraging co-expression across thousands of individual cells²⁵. Some cytokines did show significant positive co-expression (*Il4* or *Il5* and *Il13*; *Il17a* and *Il17f*; Extended Data Fig. 5b). We separated coregulated gene modules (Affinity Propagation) that defined independent transcriptional programs (Extended Data Fig. 5c). GeneOntology analysis showed that most related to generic functions (Supplementary Table 3), a few small modules (M7, M11, M13) including some elements of the prototypical T_H signatures, e.g. cytokines and transcription factors (*Ifng*, *Il13*, *Tbx21*, *Gata3*). But when projected across the cell space, most showed broadly differential representation as a gradient across all cells, cutting across cells expressing *Ifng* or *Il17a* transcripts (Extended Data Fig. 5d, with the exception of cell cycle genes in M1/2, the MHC-II module in M9, or the T_H2-like cluster in M7). Thus, the major components of variability among Teff highlighted a continuous cloud of phenotypic variance, rather than discrete cell sets.

Teff phenotypes across time

A possible explanation for the lack of discrete T_H1 and T_H17 identities was that the 13-day timepoint chosen for analysis might be off, and that polarized cells might have faded, or have yet appear. To test this possibility, we analyzed LP CD4⁺ cells at different times after *Salmonella* infection, again hash-tagged in a single batch. The CD4 response, denoted by total CD4⁺ T proportions and the effector/naïve ratio, was highest in the d10–17 window (Extended Data Fig. 6a, 6b). A marked shift in the overall Teff transcriptomes occurred from d10 onwards (Fig. 4a). Transcripts that distinguished these two super-clusters included many of the *Salmonella*-specific transcripts identified above, but no prototypical T_H signature transcripts (Extended Data Fig. 6c). In these samples, IL17⁺ cells were better demarcated than in earlier experiments, IFN- γ ⁺ cells again broadly spread out, with no indication of a time-dependent convergence (Fig. 4b). Both types of cytokine-producing cells were shifted during the “d10 transition”, again implicating the infectious agent as the dominant driver of Teff phenotypes at the height of infection.

Next, we asked whether one could identify distinct lineages of *Il17* and *Ifng* expressing cells within CD4⁺ Teffs at different infection times, using the sequence of rearranged TCR α and TCR β to lineage-trace cells originating from the same progenitor. A total of 579 repeated clonotypes were observed (defined by shared nucleotide sequences for both chains, and P/N nucleotide addition that ensured true clonal amplification; examples in Fig. 4c). These repeated clonotypes expanded with time in Teff, but not in naive Tconv or in Treg cells, consistent with infection-driven expansion (Extended Data Fig. 6d). Importantly, expanded Teff clones were not restricted to expression of one cytokine: most *Il17*-producing cells within a clonotype had cousins which produced *Ifng*, or both (Fig. 4d). That expanded clonotypes do not appear committed to produce a single cytokine could be explained by parallel differentiation of the initial precursor. However, the median Euclidean distance between members of a clonotype did not increase with time, if anything contracting beyond d10, whether computed from the T_H signature genesets (Fig. 4e) or the most variable genes (Extended Data Fig. 6e), indicating that the cells were not diversifying. Thus, this lineage tracing revealed no parallel tracks of differentiation for *Ifng* and *Il17*-production, the

Salmonella-driven dominance of IFN- γ production extending across all amplified clonotypes.

Teff phenotypes at the chromatin level

Accessibility of enhancer elements in chromatin is a more proximal readout of a cell's differentiated state than mRNA levels, which are affected by post-transcriptional events. To explore the relationship between *Ill7*- and *Ifng*-expressing cells at the chromatin level, we performed single cell ATACseq²⁶ on 4,671 LP CD4⁺ T cells from colons of *Salmonella*-infected mice. As with scRNAseq, 3 distinct clusters could be distinguished by clustering and identification based on accessibility of typical indicator genes (Tregs, naive Tconv and a cloud of Teffs; Fig. 5a). We leveraged a framework of pan-immune OCRs (Open Chromatin Regions) and charts of those most likely to be associated with activity of a given gene²⁷, extracting OCRs that best predict the expression of *Tbx21* and *Rorc* and averaging their accessibility to calculate per-cell *Tbx21* and *Rorc* "chromatin scores". We validated these scores by showing that they clearly distinguished ATACseq profiles from *in vitro*-derived TH1 and TH17 cells (Fig. 5b, top). However, when projected on the UMAP plot of *ex vivo* Teff cells, the *Tbx21* and *Rorc* scores were broadly distributed, with diffuse local maxima but no cell cluster displaying either exclusively (Fig. 5b, bottom). We also examined chromatin profiles across the *Rorc* and *Tbx21* loci themselves, by collapsing the reads from cells selected as having high or low signals at *Rorc*- or *Tbx21*-controlling OCRs, and asking whether one would be anti-correlated to the other. Clearly, chromatin openness at one locus was independent of the state at the other (Fig. 5c). Thus, chromatin opening at master regulator loci did not split identifiable TH1 and TH17 subsets.

As an alternative to analyzing the *Rorc* and *Tbx21* loci, we computationally mapped the differential activity of OCRs enriched in DNA motifs recognized by these transcription factors (TF) relative to background OCRs²⁸. T-Bet and GATA3 motif scores were broadly distributed (Fig. 5d), with a more concentrated over-representation of ROR γ motif scores (acknowledging the caveat that these motifs may be recognized by the related EOMES and ROR α , respectively).

If ROR γ and T-bet are not the main discriminators of chromatin accessibility of Teff cells, then what is? We broadened the analysis to all TF motifs in the Jaspar database, ranking them by their overall variability (null distribution from randomized data; Fig. 5e). This ranking was dominated by motifs for several factors, foremost those for the AP1 (Fos, Jun, etc.) and the IRF (Irf4, Irf2, Irf9) families, or for other factors related to T cell activation (BACH2), while the T-Box and Nuclear Receptor families (T-Bet, EOMES and RORC, RORA) figured less prominently. Correspondingly, scores for Fos and Irf4 segregated most distinctly (Fig. 5f). Thus, in line with mRNA which showed generic activation as the main driver of Teff diversity, activation drivers (AP1, Irf4, Bach2) seemed to have a more important contribution in parsing Teffs than the classic master regulators.

A functional continuum of CD4⁺ Teff cells

A continuum in which different functions are distributed along poles and gradients is more challenging to address experimentally than demarcated groups of cells. To validate the

notion of a continuum of Teff phenotypic states, we followed a strategy similar to one described recently^{29–31} wherein cell sorting was not steered to well-defined cell populations, but by integrating information in a multidimensional marker space (Fig. 6a). We first identified in the scRNAseq data transcripts that showed different gradients of expression through the Teff continuum and encoded cell-surface molecules detectable by flow cytometry (*Klrg1*, *Cxcr6*, *Icos*, *Cd69*, *Ly6a* (encodes Sca-1); Fig. 6b). Colon LP cells were resolved by flow cytometry with antibodies against these markers, combining results a multiparameter tSNE projection (Fig. 6c). In this proteomic space, no specific clusters of cells were identified by any one marker (perhaps with the exception of KLRG1), all distributing as quantitative gradients as the mRNA had. We then determined empirically gates useable to pilot a cell sorter to purify cells belonging to specific areas of the cell cloud (Fig. 6d), yielding three distinct cell populations. Such cells were sorted from colon LP of *Salmonella*-infected mice for phenotypic and functional testing.

Conventional RNA-seq on these sorted population showed a differential transcript representation, with enrichments that corresponded well to the signatures predicted from the scRNAseq data (Fig. 6e). Differentially expressed genes included, in population B, transcripts associated with a more resting state (*Ccr7*, *Sell* and *Tcf7*), while *Il17*-associated transcripts (*Rorc*, *Il23r*, *Il17re*) were over-represented in Population A. For a test of function, we stimulated these sorted cells, and measured cytokine secretion by ELISA (Fig. 6f). Distinctive patterns were observed, although as expected no single pool was associated with exclusive secretion of any one cytokine. Populations A and C secreted significantly more IFN- γ than did population B, whereas population A secreted more IL-17A and IL-22. But both encompassed all potentialities, only in quantitatively different amounts, confirming that the scRNAseq data captured true continuous Teff heterogeneity.

Novel Teff populations

As presented above, the Teff pool in SPF or infected mice included, beyond the main “cloud”, a few well distinguished populations (Fig. 2c).

(1) A small Teff population (ISG-T) was peculiar because it expressed high levels of IFN-induced signature transcripts (ISGs) (Fig. 7a, left), and over-represented after *Salmonella* or *Citrobacter* infection. Comparison with profiles induced in T cells by Type-1 or Type-2 interferons (IFN) indicated that ISG-T likely respond to Type-1 IFNs (Fig. 7a, right). Their existence suggested either a small subset uniquely responsive to IFN, or normal Teffs that happened to reside in a small anatomical compartment where IFN was particularly abundant. Similar subsets have been described in CD4⁺ T cells from House-dust-mite infected lungs, and kidneys from lupus nephritis patients^{32,33}.

(2) Another population expressed high levels of the surface markers *Cd160*, *Crtam* and *Lag3*, the neural gene *Nrgn* and several chemokines (Fig. 7b). We sorted this CD4⁺CRTAM⁺ population for bulk RNA-Seq, confirming the particular signature (Fig. 7c). Pathway analysis showed enrichment of STAT3, prolactin and neuregulin signaling pathways, hinting at a possible origin.

(3) The most intriguing population was MyT, which unexpectedly showed many myeloid cell transcripts, such as *ApoE*, *Lyz2* or *C1qa*, and several MHC-II (Fig. 7d). This expression of myeloid transcripts was not wholesale: only a fraction of genes with strong T vs myeloid differential expression was represented in MyTs (Fig. 7e), several corresponding to innate antimicrobial receptors or defense mechanisms (*Lyz2*, *C1qa*, *Cfp*, *Tyrobp*). Correspondingly, a small MHC-II⁺ subset of TCRβ⁺CD4⁺CD44^{hi} Teffs was detected by cytometry (Extended Data Fig. 7a), whose RNAseq transcriptome confirmed the single-cell data (Extended Data Fig. 7b). We applied “CITEseq” for protein detection with DNA-barcoded antibody³⁴, revealing a good correspondence between mRNA and of surface H2-A^b and CD14 proteins (Fig. 7f) in MyTs, at protein levels that were only somewhat lower than those seen in true myeloid cells. In the experiments of Fig. 4, αβTCRs detected in MyTs were shared with other Teffs from the same mice, suggesting that the MyT phenotype is not acquired during thymic differentiation, but late in the periphery after antigen encounter. MyTs may correspond to the unusual CD3⁺CD14⁺ cells in human blood recently reported to also increase upon infection³⁵, which was attributed to stable doublets. Several arguments suggested that MyT cells are not doublets, however (very few myeloid cells in the sorted CD4⁺TCRβ⁺ datasets, normal UMI/cell counts, partial myeloid gene representation). Other than doublet formation or cell fusions, explanations for MyTs included exosomal transfer of transcripts from myeloid to Teff cells, or the activation of unusual transcriptional modules. To formally resolve this issue, we created bone-marrow chimeras with a 50/50 mix of congenically marked stem cells from wild-type (WT) and MHC-II-deficient donors with an inactivating *neo* insertion in *H2-Ab1*, a mutation that results in altered *Ab1* transcripts³⁶ (Fig. 7g)). After 10 weeks’ reconstitution, we sorted TCRβ⁺CD4⁺ cells of both donor origins for RNAseq, analyzing the sequence reads at *H2-Ab1*. Should MyT cells result from doublets or mRNA transfer, TCRβ⁺CD4⁺ of knockout origin would have picked up WT *H2-Ab1* transcripts. This was not the case (Fig. 7h), as those T cells expressed transcripts from their own *H2-Ab1*. Thus, MyT cells are bona fide TCRαβ⁺ T cells that activate a segment of the myeloid transcriptome. Their origin and significance remain to be established. However, there may be a precedent in the myeloid-like T cells that constitute the high-risk “Mixed Phenotype Acute Leukemia”³⁷.

DISCUSSION

Our study set out to map the landscape of phenotypes that Teff cells in the gut can adopt when stressed by microbial infection, relating to the long-running question of Teff cell heterogeneity. Whether evaluated at the transcriptome or the chromatin levels, our results show that Teff cells are molded by the infections in a profound and specific manner, one that does not readily conform itself to T_H stereotypes, and also bring out other intriguing new cell states.

From the realization over 40 years ago that distinct functions of helper T cells reside in different cells¹, the field has strived to subdivide Teff cells into discrete subsets. Since the seminal discoveries of Mossman and Coffman³ and the coining of the T_H1/2 semantic, these distinctions have been anchored by cytokine production, an anchor which has persisted despite repeated demonstrations of dual-expressing cells, T_H sub-subsets^{22,23,38}, and plasticity between T_H states^{9,13,14}. Our results suggest that Teff transcriptional identities

form a “polarized continuity” and cannot be parsed out into discrete T_H cell-types, even in the context of infections expected to drive focused differentiation. Nor does progressing infection result in phenotypic divergence between clearly distinct states. This model does not imply homogeneity, however, as the different poles of the phenotypic cloud do show strong preference for production of one cytokine or the other (most marked for IL-4 or IL-5).

This view of Teff cell heterogeneity differs from previously proposed concepts of cell plasticity, in which cells of defined phenotypes can switch between states that are otherwise coherent and reproducible^{9,13,14}. The plasticity concept implies that discrete states do exist, but are not irrevocable. We find that there are no defined states to inter-convert between. This view also diverges from the notion of sub-subsets [e.g. pathogenic T_H17^{22,23,39,40}], which also implied discrete cell sets that could be further subdivided. Such sub-subsets also seemed absent, and in hindsight may represent the spread of IL17-producing cells across different regions of the phenotypic cloud.

One might argue that the polarized continuity represents transient intermediates between cell states. But, then, most cells would be intermediates. Velocity testing of differentiation within the Teff continuum⁴¹ gave no indication of directional progression, and the time course study showed no particular convergence towards more distinct Teff phenotypes, overall or for amplified progeny of the same precursor. Importantly, chromatin analysis revealed that key controlling loci, *Rorc* and *Tbx21*, opened largely independently of each other.

Several studies are also consistent with this view of “polarized continuity” within Teff cells that is dominantly molded by microbes. Cloned human memory CD4⁺ T cells showed phenotypic divergence related to the initiating microbe⁴². Proteomic analysis by mass cytometry revealed a wide phenotypic range in CD4⁺ Teff unleashed by *Ctla4* deficiency⁴³. In tumor infiltrating cells, scRNAseq studies also found gradients of transcriptional phenotypes^{44,45}, as in other broad “landscape” studies where T cells have been notoriously difficult to parse finely^{46,47}. A recent analysis of airway-resident T cells also reported a continuous disposition of Teffs in house-dust-mite infection³², showing that our results are not gut-specific. A continuous phenotypic spectrum has been described for ILCs⁴⁸, contrasting with commonly used categorization¹¹. Rather, ILC phenotypes can be described by a series of “topics”⁴⁸, conceptually similar to, and partially overlapping with, the modules reported here. While this work was under review, Cano-Gamez et al also proposed a model of human T cell activation *in vitro* dominated by “continuous effectorness”⁴⁹.

In conclusion, this study sheds light on the T cell response to infectious challenges: broad responses that adapt to each microbe, dominant co-regulated gene modules that are not anchored by cytokines, different leading transcriptional drivers, and intriguing new cell subsets.

METHODS

Mice

Male C57BL/6 mice were purchased from Jackson Laboratories. IL17A^{GFP/+} mice (Jax:C57BL/6-Il17atm1Bcgen/J) were a gift from Dr. Jun Huh. OT-II TCR-transgenic mice

were obtained from Jackson Laboratories (B6.Cg-Tg(Tcr α Tcr β)425Cbn/J). *H2-Ab1*-deficient mice were previously described³⁶. To construct bone marrow chimeras, bone marrow cells were harvested from both femurs and tibias, and treated with ACK buffer (Lonza) to remove red blood cells. CD45.1⁺/CD45.2⁺ mice were irradiated (10 Gy) and reconstituted with equal proportions (~5 million cells each) of CD45.1 and CD45.2 (Ab KO) bone marrow cells. All mice were bred and maintained in our specific-pathogen-free facilities at Harvard Medical School (IACUC protocols IS1257, IS187-3, IS2221).

Flow cytometry

Cells from colon LP were prepared as previously described⁵¹. Briefly, Intestinal tissues were treated with RPMI containing 1 mM DTT, 20 mM EDTA and 2% FBS at 37°C for 15 min to remove epithelial cells, minced and dissociated in collagenase solution (1.5mg/ml collagenase II (Gibco), 0.5mg/ml Dispase (Gibco) and 1% FBS in RPMI) with constant stirring at 37°C for 45min. Single cell suspensions were then filtered and washed with 4% RPMI solution. For cytokine analysis, cells were treated with 10ng/ml phorbol 12-myristate 13-acetate (Sigma), 1 μ M Ionomycin (Sigma) and 1X protein transport inhibitor cocktail (eBioscience cat. #00-4980-03) for 3.5 hours in 10% FBS RPMI. For intracellular staining of cytokines and transcription factors, cells were stained for surface markers and fixed in eBioscience Foxp3 buffer overnight, followed by permeabilization in eBioscience (both cat. #00-5523-00) buffer for 45 min in the presence of antibodies. Fluorescence profiles were acquired on a BD Symphony and analysis was performed with FlowJo (Tree Star) software.

Antibodies used in the study included anti-mCD45 (30-F11), anti-mCD19 (6D5), anti-mCD4 (RM4-5), anti-mTCR β (H57-597), anti-mCD44 (IM7), anti-mCD25 (PC61), anti-mFOXP3 (FJK-16s), anti-mIFN- γ (XMG1.2), anti-mIL-17A (TC-11-18H10.1), anti-mIL-5 (TRFK5), anti-mIL-13 (W17010B), anti-mCCR6 (29-2L17), anti-mIL1R2 (4E2), anti-mKLRG1 (2F1/KLRG1), anti-mICOS (C398.4A), anti-mCXCR6 (SA051D1), anti-mCD69 (H1.2F3), anti-SCA1 (D7), anti-mCRTAM (11-5/CRTAM), anti-I/A I-E(M5/114.15.2), , anti-mCD45.1 (A20), and anti-mCD45.2 (104). All antibodies were diluted 1:100 with the exception of CD25 (1:50).

For the “tSNE sort”, the goal was to sort cells defined combinatorially by a panel of markers, even if they were not readily identifiable as well demarcated populations on conventional 2-parameter flow cytometry profiles. Flow cytometry tSNE plots were generated in FlowJo v10 from gated CD4⁺TCR β ⁺CD44⁺CD25⁻ cells stained for markers found by manual inspection to have non-discrete and non-correlated expression in the scRNAseq data, and thus most appropriate to represent the phenotypic continuity within the Teff phenotypic cloud (KLRG1, ICOS, CD69, SCA1, and CXCR6). Selected regions that demarcated clusters on the tSNE continuum were then backgated onto normal 2-parameter plots, from which gating instructions interpretable by a cell-sorter were drawn (by a manual and iterative process). The gates were drawn such that more than 90% of the events in the selected region of would be within the sort gates. These combinatorial gates were then applied to sort colonic Teff cells.

Multiplex ELISA

1000 to 10,000 cells were sorted (single sort) into 100 μ l of T cell medium (RPMI 1640, 10% FBS, 20mM HEPES, 1mM sodium pyruvate, 0.05mM 2-mercaptoethanol, 2mM L-glutamine, 100 mg/ml streptomycin and 100 mg/ml penicillin). Cells were plated in round bottom 96-well plates with 1:1 ratio of anti-CD3/CD28 beads (Miltenyi) and were incubated at 37C for 24 hours. Supernatants were collected and analyzed with the LegendPlex™ T Helper Cytokine Version 2 and Proinflammatory Chemokine panel kits (Biolegend) per manufacturer's instructions. Samples were acquired via the BD Symphony and analyzed with LegendPlex™ software. Paired student t-test was used for significance assessment.

Infections

For *Salmonella* infection, mice were gavaged with 100 μ l of 200 mg/mL streptomycin in water, and 24 hours later, gavaged with 10^9 *Salmonella enterica* (Serovar Typhimurium) AroA⁵²; gift of D. Littman). For *Citrobacter* infection, mice were gavaged with 5×10^8 *Citrobacter rodentium*⁵³. Unless noted otherwise, mice were sacrificed at day 13 after infection. For helminth infections⁵⁴, mice were gavaged with 200 *Heligmosomoides polygyrus* L3 larvae in 200 μ l H₂O or subcutaneously injected with 500 L3 larvae of *Nippostrongylus brasiliensis* in 100 μ l PBS, and sacrificed 11 days later.

Low Input RNA-Seq

All cells were double-sorted. For the final sort, 1,000 cells were collected directly into 5 μ l lysis buffer (TCLBuffer (Qiagen) with 1% 2-Mercaptoethanol), and the lysates frozen after 5 min. Smart-seq2 libraries were prepared as previously described¹². Reads were aligned to the mouse genome (GENCODE GRCm38/mm10 primary assembly and gene annotations vM16; https://www.genecodegenes.org/mouse_releases/16.html) or human genome (GENCODE Human Release 27; Reference genome sequence: GRCh38/hg38; annotation: GENCODE v27) with STAR 2.5.4a. The ribosomal RNA gene annotations were removed from the GTF (General Transfer Format) file. The gene-level quantification was calculated by featureCounts (<http://subread.sourceforge.net/>). Raw reads counts tables were normalized by median of ratios method with DESeq2 package from Bioconductor, and then converted to Gene Pattern GCT and CLS format. Samples with less than 3 million uniquely mapped reads were automatically excluded from normalization to mitigate the effect of poor quality samples on normalized counts. Normalized read counts were filtered for robust expression (>10) to avoid confounders from low-level noise, and processed in the Multiplot suite and Morpheus (<https://software.broadinstitute.org/morpheus/>), or with custom code in R. PCA was done using the pcomp function in R on all genes that had expression higher than 0 in any sample.

Single-cell RNA-seq

Intestinal tissues were treated with RPMI containing 1 mM DTT, 20 mM EDTA and 2% FBS at 37°C for 15 min to remove epithelial cells, minced and dissociated in 1 mg/ml collagenase VIII (Sigma), 1 μ g/ml DNase and 1% FCS in RPMI with constant stirring at 37°C for 20min. Single cell suspensions were then filtered and washed with 4% FCS in RPMI medium. Single-cell suspensions were stained on ice for 30 min with antibodies

against CD4, TCR β , CD19 and CD45 (Biolegend) with 20 ng/ml Dapi (Biolegend) as a viability dye. T cells were then sorted on an Astrios MoFlo (Beckman Coulter), as Dapi⁻CD45⁺CD4⁺TCR β ⁺CD19⁻. For single sample processing, cells were sorted directly into PBS to 0.04% BSA final concentration. For cell hash-tagging, TotalSeq-A™ hashtag antibodies (SPF: Hashtag1, *C. rodentium*: Hashtag2, *S. Enterica*: Hashtag3, *N. brasiliensis*: Hashtag4, *H. polygyrus*: Hashtag5) were added to each sample individually at the same time as other antibodies: All samples were sorted together directly into RPMI with 2% FCS and subsequently spun down and reconstituted in 33 ul of PBS with 0.04% BSA. All samples were loaded on the 10x Chromium Controller (10x Genomics) within 30 minutes of sorting. Libraries were prepared using Chromium Single Cell 3' Reagent Kits v2 according to the manufacturer's protocol. HTO libraries were prepared as described in¹⁸. Libraries were sequenced together on Illumina HiSeq 4000.

Single-Cell RNA-seq Data analysis

Gene counts were obtained by aligning reads to the mm10 genome using CellRanger software (v1.3) (10X Genomics). HTO counts were obtained by using Cite-Seq-Count package³⁴. Single cell data were initially analyzed using the Seurat package⁵⁵. HTOs were assigned to cells using the HTODemux function and doublets were eliminated from analysis. Cells with less than 1000 UMIs or 400 genes and more than 4000 UMIs or 0.05% reads mapped to mitochondrial genes were also left out of the analysis. Tregs and naive CD4⁺ cells were removed from analysis by using the SubsetData function. Data were normalized using the NormalizeData function and scaled using the ScaleData function, regressing out number of UMIs and percentage of expressed mitochondrial genes. Variable genes were found by FindVariableGenes function, using genes with mean expression over 0.0125 and 4 UMIs per cell. Dispersion cutoff was calculated based on the fano factor distribution per gene. This way, 550-950 variable genes were selected in different Teff datasets. PCs were calculated using the RunPCA function and significant PCs were selected using the JackStraw function. tSNE and KNN clusters were computed on significant PCs using the RunTSNE and FindClusters functions, respectively. UMAP dimensionality reduction was calculated on significant PCs using the RunUMAP function. T_H signature scores were computed by the mean expression of signature genes per cell.

Diffusion maps are useful for identifying differentiation trajectories, as they allow for pseudotemporal ordering of single cells in a high-dimensional gene expression space⁵⁶. Diffusion maps were generated using the Seurat package RunDiffusion function with default settings.

Imputation can denoise the cell count matrix and fill in missing transcripts by data diffusion⁵⁷ Imputation was performed using the built-in Seurat AddImputedScore function with default parameters on all variable genes. PCs and tSNE were then r-recomputed based on the imputed values.

Principal Components were identified and plotted using the Seurat PCHeatmap function using default parameters.

Correlation Coefficient Analysis (CCA)⁵⁵ was performed by running RunMultiCCA function on 500 variable genes between the 4 samples. 20 significant Correlation Coefficients (CC) were selected for alignment using the AlignSubSpace function. tSNE and KNN clustering were run as previously, based on 20 CCs.

To compute Euclidean distances within groups of cytokine expressing cells, cytokine positive cells were identified as expressing one or more normalized UMIs. Distances between each selected cell to other cells were calculated for 1000 top variable genes using the dist function in R. P values were computed using the Mann-Whitney test. For Dendrogram analysis, cytokine expressing cells were identified as above, and distances between different samples expressing different cytokines were computed by the dist function in R with default settings on top 1000 variable genes. Hierarchical clustering (hclust function in R) was then employed to generate the dendrogram.

Highly Th-specific genesets (Supplementary Table 2) were generated by manual curation, starting mainly from published signatures as well as other scRNAseq datasets^{19,32,58–61} and selecting genes that were reproducibly present in these signatures. We removed transcripts that overlapped between resulting T_H genesets, often simple markers of cell activation frequent in such signatures, as well as some non-T transcripts which are not infrequently found to contaminate published signatures (e.g. *Cd19*, *Cd79a*, *Cd8a*). We also added several transcripts known to correlate with *Ifng*, *Ill17a* or *Ill4/il13* (*Cxcr3*, *Tmem176a*, *Areg*). The gene signature average for these genes was then calculated via the AddModuleScore function in Seurat V3. Expression of cell cycle genes was calculated based on CellCycleScoring function in Seurat V3 (cc.genes based on⁶²)

Gene Module Generation

After filtering transcripts for robust expression (those that appear in more than 10 cells in any one of the infected or SPF samples), gene-gene correlation (Pearson, cor function in R) was calculated within each dataset. The ten matrices (one for each replicate and condition) of pairwise gene-gene correlations were then averaged for Extended Data Fig. 5c.

To select the genes with the highest correlations, a threshold correlation score in the 98% percentile was calculated for each gene, and 588 genes with correlation scores higher than 0.05 were selected for further analysis. Gene modules were then identified by Affinity Propagation⁶³ using the APcluster R package with negative distance similarity function and number of input similarities (q) set to 0. Gene modules were overlaid on the tSNE by computing the mean expression of module genes per each cell.

Clustering approaches

BackSPIN: The data was normalized with Seurat parameters and then subset to the top 588 most variable genes according to the Seurat pipeline. In order to know if significant clusters would emerge from more elaborate clustering methods, we used BackSPIN, an unsupervised biclustering method that sorts both genes and cells into clusters⁶⁴. The motivation behind backSPIN was that by iterative partitioning, the algorithm would be able to cluster true cell subsets and gene subsets together. One important parameter for BackSPIN involves defining

the partitioning ‘rate’ (i.e., how much to subset the groups at each iterative process). This was set at the default of 0.1. Other parameters specified were the number of levels (numLevels) to partition by, set at 2, the number of top variable genes to cluster, set at 596, the initial number of iterations (first_run_iters), set to 10, and subsequent number of iterations (runs_iters) set to 8. The default initial decrease rate of 0.1 (first_run_iters) and the default subsequent decrease rate (runs_step) of 0.3 were used. The decrease rates helped determine the precision of clusters. Finally, threshold values were set at the default value of 2 for both minimum number of cells (split_limit_c) and genes (split_limit_g). A threshold score of 1.15 was used to determine when to stop partitioning the data (stop_const), and the default threshold for determining which group a gene would fall into was kept at 0.015.

BISCUIT: BISCUIT iteratively learns to identify features in each cluster and create clusters with these specific features by imputing and normalizing the data⁴⁵. The motivation behind BISCUIT is that by imputing the data, it catches variation provided by genes that may have dropped out. The major parameter for BISCUIT is the dispersion parameter (alpha) which allows the algorithm to sort cells into more clusters or less clusters, which was set to 1. The following parameters were used for running BISCUIT: The default setting of 20 genes/batch, the default number of 20 iterations, and 100 as the number of cells in each batch. Once complete, the final clusters were projected onto the tSNE plot of Fig. 2c computed by Seurat. Cell cluster outputs from BISCUIT were projected onto the tSNE computed by Seurat.

Dip test

Data were normalized with standard Seurat parameters as described before⁵⁵. The same number of variable genes, defined by Seurat, was used in the continuity analysis. In order to test for “discontinuity” in transcriptomic-based representation of a set of cells, we used the Hartigan’s Dip Test of multi-modality²⁰. The Dip Test asks whether the pairwise distances between all pairs of cells can be best supported by a unimodal or a multi-modal distribution. The intuition behind this test comes from the fact that if there are two or more clear sub-populations of cells that cluster together with clear boundaries, then given a high-dimensional representation of these cells (i.e., vectors of length g consisting of gene expression levels for g genes), then there would one more region of low density in between highly dense regions in this space. These low-density regions would thus create a “dip” in the distribution of pairwise distances between all cells in this space. One important parameter here is the representation of gene expression data used in computing the pairwise distances between the cells. To support the ability of the Dip Test to identify regions of low density, we first applied a projection defined by minimum separation hyperplane²¹ to gene expression data from variable genes (defined by Seurat) and then applied the Dip Test on the distances computed on the projected data.

Binary classification IL17 or IFN- γ producing cells

We trained a deep neural network run on the Keras platform (www.keras.io). The input geneset was the 500 most variable genes across the entire scRNAseq dataset of Fig. 2 (naturally leaving out *Il17* and *Ifng* transcripts), and the network was trained to classify *Ifng*- or *Il17a*-producing cells (randomly assigned to 80% training set, 20% test set). The data

matrix was normalized by the mean of the expression of each gene across the 2,885 cells (otherwise the transcripts with highest expression levels dominate the output). The DNN was composed of 3 hidden layers with the following features: size of the hidden layers: 512-128-64 with random weights initialization; activation function: Sigmoid; optimizer for backward propagation: ADAM optimizer; number of epochs: 50; training and testing on CPU; Batch size: 100. We added downstream a decision function with the possibility of NoCall (for non-producing cells), the classification as Il17-producing was accepted if the output softmax score of the cell was above 0.95 (and below 0.05 for *Ifng*-producing cells), otherwise NoCall. We voluntarily over-fitted the model to fit the distribution of the output softmax score with the decision function constraints. We used a Keras-based (version 2.2.4) neural network (www.keras.io) on Python2. The Integrated Gradients library was used to compute the overall contribution score of each gene as the mean of its contribution scores across the whole data set. To test the reproducibility of the Integrated Gradients we split randomly the data set in two sub datasets on which we independently trained models, repeating the operation 100 times on each dataset and taking the mean of these 100 scores. As a positive control, the same architecture was used to distinguish Teff from Treg cells, which it could do with 98.8% (Teff) and 89.7% (Treg) accuracy on average, per below).

Accuracy: 97.58 - 98.05%	Actual Tregs	Actual Teff	
Prediction Treg (in 10 independent runs)	1120-1170	42-64	
Prediction Teff (in 10 independent runs)	78-128	5700-5722	
	1248	5764	nCells=7012

Clonotype analysis and CITE-Seq

Mice were infected with *Salmonella*, as above and colon single-cell suspensions prepared as above Antibody staining (cell hashing and CITE-Seq) was performed simultaneously by adding TotalSeq-C™ hashtags 1-7 (Day0: hashtag7, Day3: hashtag6, Day5: hashtag 5, Day7: hashtag 4, Day10: hashtag 3, Day17: hashtag 1), anti-CD14 (C0424) and anti-I-A/I-E (C0117) (Biolegend) to the cells at a ratio of 1:100 in RPMI 2% FCS and incubating on ice for 15 minutes. Cells were then washed twice with RPMI 2% FCS and sorted as described above before encapsulation (10X Genomics). Gene expression, Feature and TCR V(D)J libraries were prepared using the 5' V(D)J v1 kit (10x Genomics). Rearranged TCRs were identified by running Cellranger v3.0, and TCR chains and N, P nucleotides per clonotypes were determined with the help of the IMGT database (http://imgt.org/imgt_vquest). Repeated clonotypes were defined by shared TCR alpha and beta with identical CDR3s at the nucleotide level. For UMAP and clonotype analysis, cells in cycle were taken out. *Ifng/Il17a* expressing cells were defined as cells that had reads for either transcript. Euclidean distances between cells expressing the same repeated TCR clonotype were measured using the `dist()` function on either the T_H geneset (Supplementary Table 2) or the 1000 most variable genes.

Single-cell ATAC-seq

Total CD4⁺ T cells were isolated from the colon of Salmonella-infected mice as described for scRNAseq, with the exception of using Collagenase II and Dispase instead of Collagenase VIII. 25K Cells were sorted directly into 2% FCS RPMI and subsequently spun down and reconstituted in 0.04% PBS. Nuclei isolation, GEM generation and library preparation were performed as described in the Chromium Single Cell ATAC (10X Genomics) manual (<https://support.10xgenomics.com/single-cell-atac>). Libraries were sequenced on Illumina NextSeq. OCR counts were obtained by aligning reads to the mm10 genome using Cell Ranger ATAC software (v1.1) (10X Genomics). scATAC-seq was analyzed using the Seurat/Signac pipeline (<https://satijalab.org/signac/index.html>). For QC, cells with less than 5000 peak calls and less than 20 percent reads mapped to peaks were filtered out. For normalization of peak counts used to drive the UMAP representation, we used the RunTFIDF function to calculate the Term Frequency Inverse Document Frequency (TFIDF). For dimensionality reduction data structure was learned via latent semantic indexing (LSI, RunLSI function) and Single Value Decomposition (SVD, RunSVD function). Contaminating non-T cells were taken out, and UMAP and cell clusters were then recalculated. Naïve T, Treg and Teff clusters were identified and attributed based on the Gene Activity matrix, constructed using the FeatureMatrix function and the Gencode v18 annotation; peaks that were found within the gene body and up to 2 kb upstream of TSS were assigned to the corresponding genes.

To calculate the *Tbx21* and *Rorc* scores shown in Fig. 5b, we counted the raw reads falling into 300 bp intervals centered on OCRs that were highly correlated to expression of corresponding genes in the Immgen compendium (per Supplementary Tables 3F and 2 of²⁷). For the *Rorc* locus, signals at three OCRs with TSS Gene-OCR correlation score >10 were used; for the *Tbx21* locus 11 OCRs with TSS Gene-OCR scores >15 (Fig. 5C). Read counts were then summed and averaged into a score per cell using the AddModuleScore function. Cells were assigned as *Rorc*⁺ or *Tbx21*⁺ if the average OCR score for these loci was greater than 0. Coverage maps were then generated using the CoveragePlot function, only applied to Teff cells.

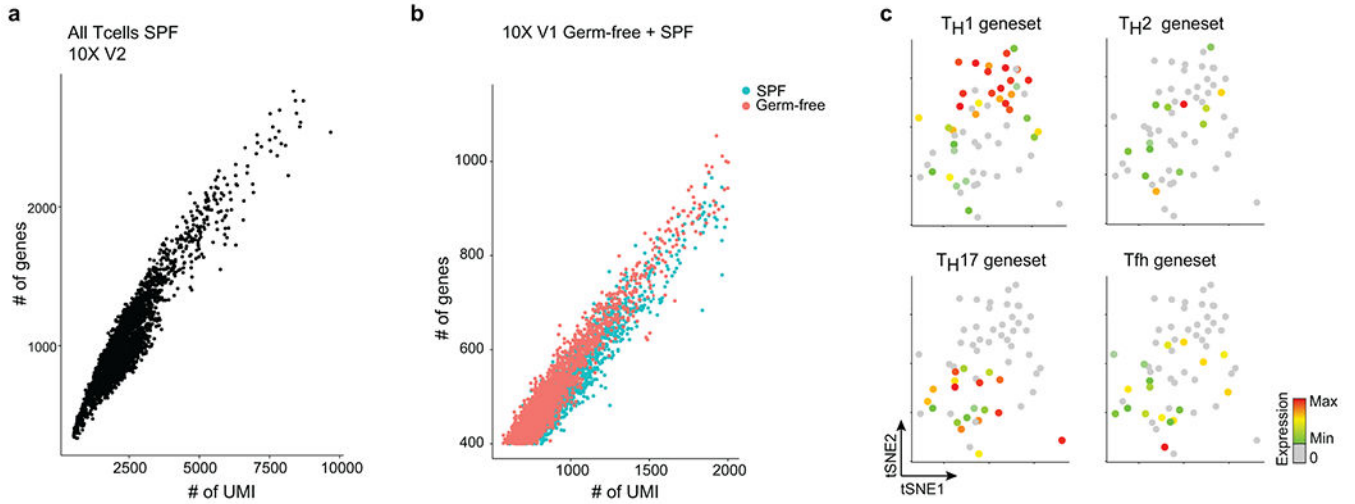
Raw bulk ATACseq data from *in vitro* differentiated Th0, T_H1 and T_H17 cells were generously provided by P. Thakore and A. Schnell⁶⁵. The *Tbx21* and *Rorc* chromatin scores were computed as above from read counts (normalized to the total read number for each biological replicate).

TF deviation and variability scores were calculated using the chromVAR package (v1.8)²⁸ with motifs from the JASPAR2018 database. The filtered Teff-only scATAC-seq counts matrix was used as input, with peaks overlapping motifs determined using the motifmatchr matchMotifs function. The chromVAR computeDeviations function was used to calculate the bias-corrected deviation scores for each TF motif. Briefly, this method computes the difference between observed fragments within peaks containing a given motif and the total expected number of fragments using the average of all cells. These 'raw deviation' scores are then normalized for technical biases using a set of background peaks matched for GC content and accessibility to yield the 'bias-corrected deviation scores.' Variability of TF motifs across the Teff data was calculated using the chromVAR computeVariability function.

DATA AVAILABILITY

The data reported in this paper have been deposited in the Gene Expression Omnibus (GEO) database under accession no. GSE 160055).

Extended Data

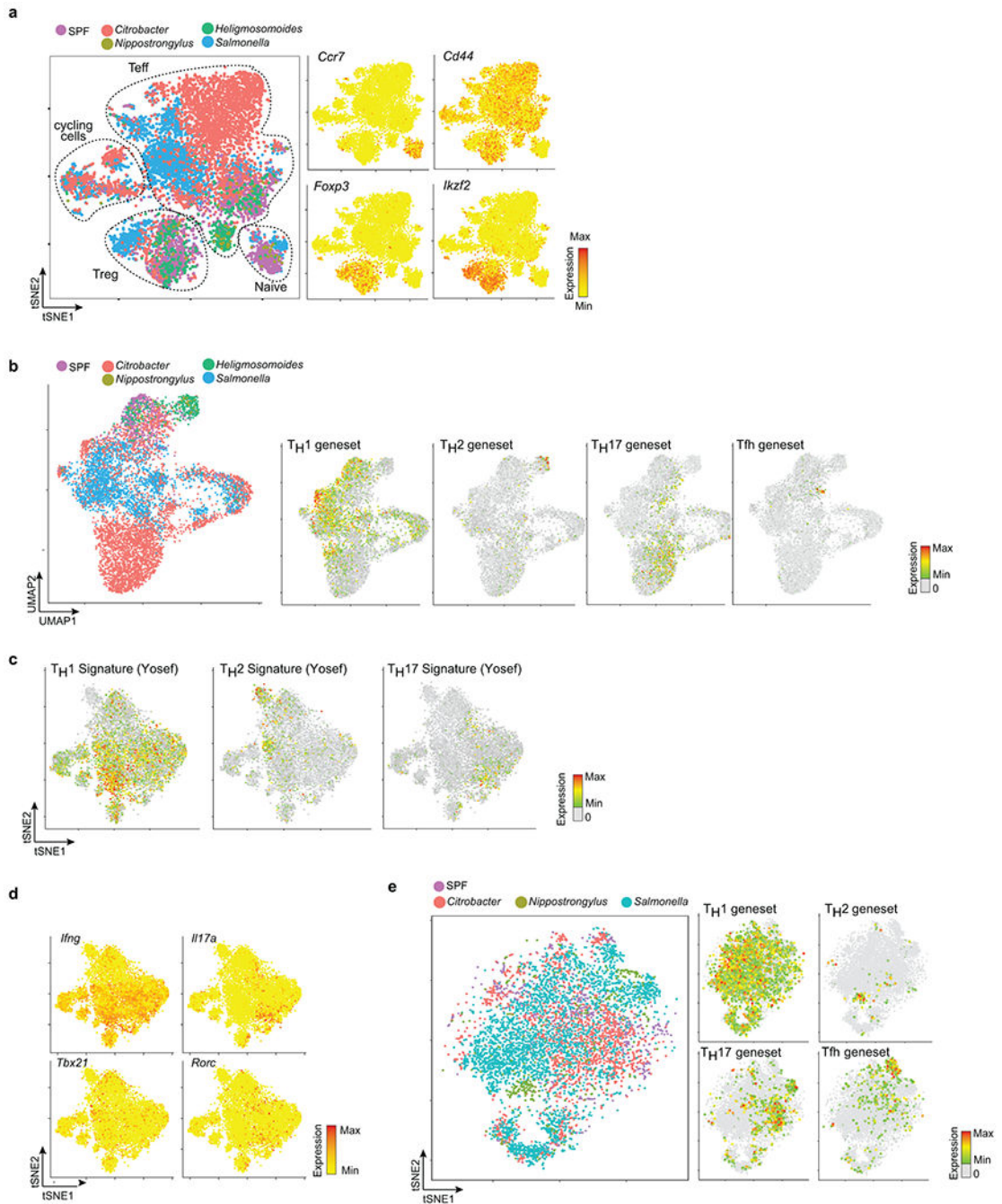


Extended Data Fig. 1. scRNAseq of Teff under normal conditions

a) Quality control plots (per-cell number of unique reads vs number of transcripts detected) for the scRNAseq data from total colonic CD4⁺ T cells (data from Fig. 1a).

b) Same plots as (a), for CD4⁺ QC of scRNAseq data from total colonic CD4⁺ T cells of germ-free and SPF mice.

c) SMART-SEQ2 single-cell data from colon T memory cells (from Miragaia et al. {11059}). Aggregate expression of Th-specific genesets (defined as for Fig. 1) are overlaid on the tSNE.



Extended Data Fig. 2. scRNAseq of Teff under infectious conditions.

a) tSNE representation of all CD4⁺ T cells in the scRNAseq data from the parallel infection experiment of Fig. 2. Left panel: each color represents cells from a different infection condition. Tregs, naive Tconvs, cycling cells and Teffs are circled; right panel: expression of key genes.

b) UMAP representation of Teff cells from the same experiment, colored by condition; Right panels: Overlay of T_H genesets (per Fig. 2).

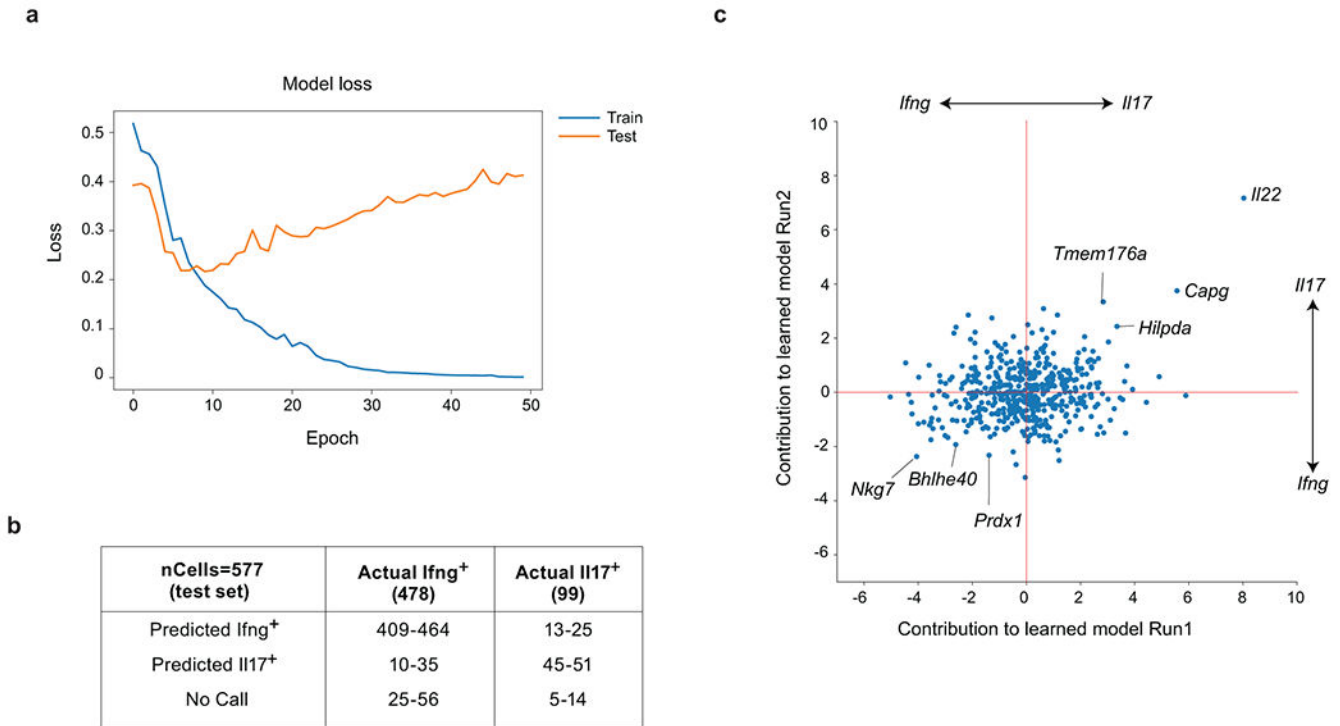
- c) Data from the same parallel-infection experiment as Fig. 2c and displayed using the same tSNE coordinates, highlighted with aggregate expression of T_H signature genes from ref{9203}.
- d) Expression of key cytokines and transcription factors in the same scRNAseq data as Fig. 2c.
- e) Independent parallel infection experiment. Samples were not hash-tagged, and processed in parallel encapsulations, and cell data were aligned by canonical correlation analysis (CCA) for tSNE representation, color-coded by sample. Right: expression of Th-specific genesets, defined as for Fig. 2c.



Extended Data Fig. 3. Different clustering approaches and signatures do not parse out the data into TH subsets.

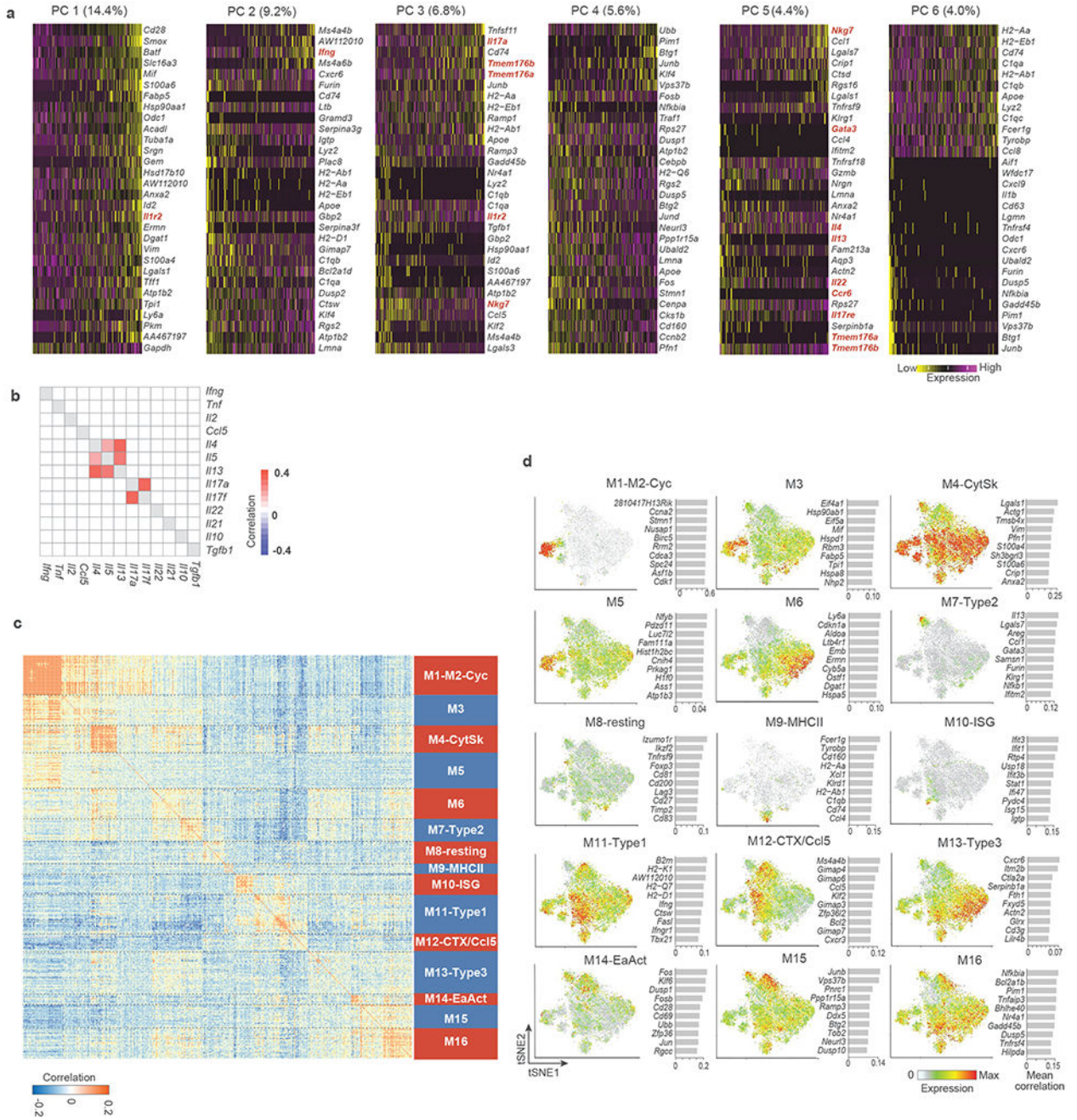
- a) KNN clusters shown on hash-tagged tSNE. Percentages of cells corresponding to each signature in each KNN cluster are shown in the table.
- b) Biscuit clusters shown on hash-tagged tSNE. Percentages of cells corresponding to each signature in each Biscuit cluster are shown in the table.
- c) Backspin clusters shown on hash-tagged tSNE. Percentages of cells corresponding to each signature in each Backspin cluster are shown in the table.

- d) Overlay of pathogenic T_H17 signatures from {9915, 8270}. Left panel: all Teff; right panel: only Il17a⁺ Teff.
 e) Overlay of *Citrobacter* T_H17 signature from {11132} on the tSNE plot.



Extended Data Fig. 4. Neural Network prediction of IFN- γ and Il17-producing phenotypes

- a) A Keras neural network was trained to use as input the expression of 500 most variable genes in Teff single-cell RNAseq data to predict *Ifng* or *Il17a* expression in each cell. Loss as a function of training epochs plotted here. Note the overfitting beyond 10 epochs (representative of >50 independent training runs with random 80/20 training/test).
- b) Accuracy of DNN-predicted cytokine expression by individual Teff cells, relative to their actual expression in the test scRNAseq data (non-expressing cells were not included as input, since there is uncertainty as to their real nature given drop-out frequencies in scRNAseq data). Numbers shown represent the range observed in 10 independent training runs (with different training/test sets).
- c) Contribution of each transcript to the prediction of *Il17a* or *Ifng* expression, as score in the Integrated Gradients, comparing the model learned in two independent runs. A positive score indicates influence on predicting *Il17a* expression, a negative score influence in predicting *Ifng* expression.

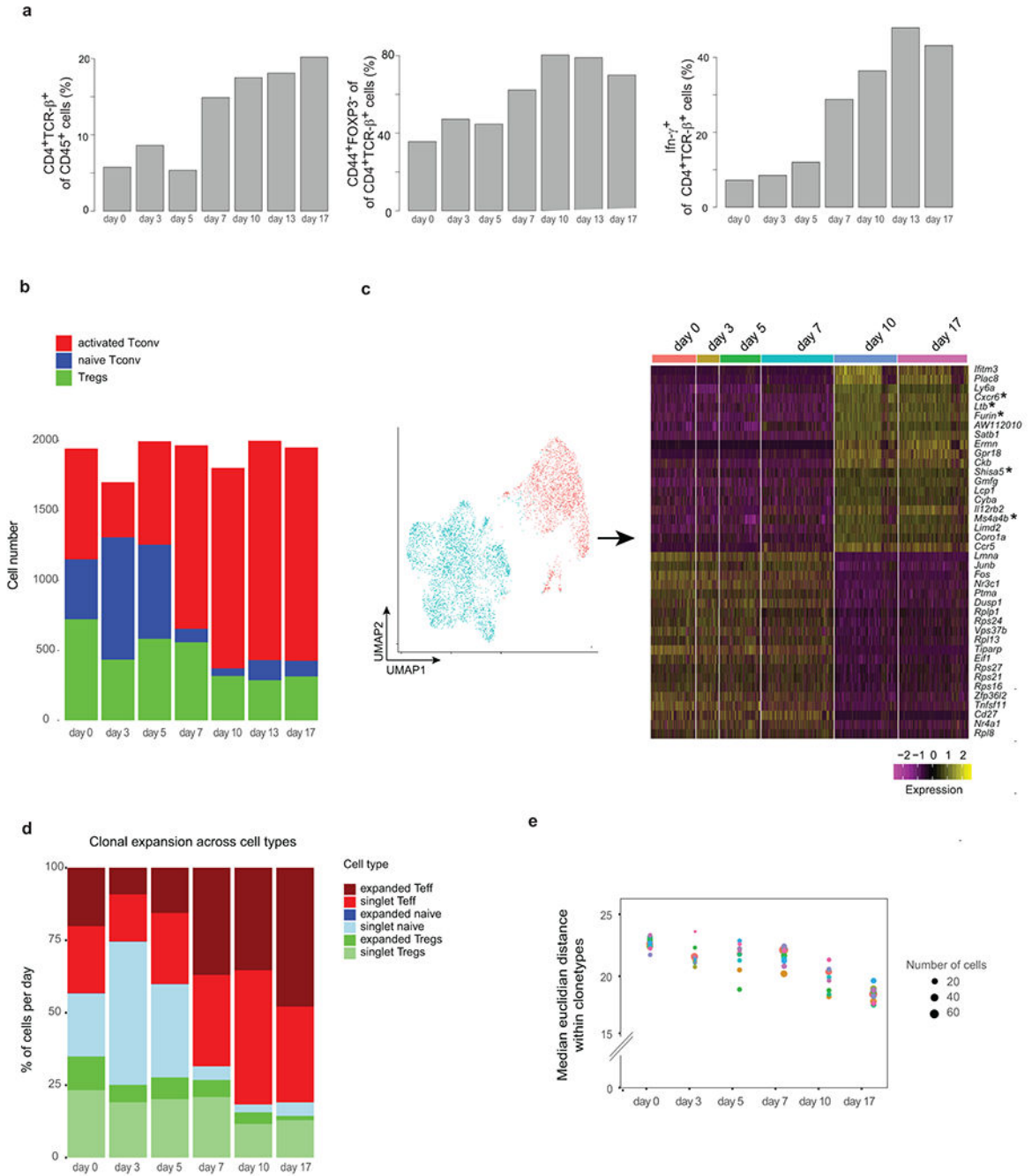


Extended Data Fig. 5. Th-associated genes are not the main drivers of Teff heterogeneity.

- a) Distribution of Top 6 PCs of Teffs from all hash-tagged samples, with cell cycle genes regressed out. Genes that are Th-associated are highlighted.
- b) Co-expression of key cytokines across all samples. Mean Pearson gene:gene correlation of cytokine genes across all samples. Only significantly correlated cytokines are colored ($p < 0.05$, χ^2 test). Significant P values: $IL4/IL13$ 6.3×10^{-3} , $IL4/IL5$ 1.8×10^{-98} , $IL5/IL13$ 5.5×10^{-129} , $IL17a/IL17f$ 1.3×10^{-4} .

c) Coregulated gene modules in Teff single-cells. Gene:gene correlation between 588 most variable genes was calculated independently within each condition/infection of the single-cell datasets, then averaged between conditions. 16 gene modules were determined by Affinity Propagation within this matrix, annotated at right.

d) Overlay of average expression of these gene modules on Teff tSNE (per 2c) with barplots showing genes with highest mean correlation (full list in Supplementary Table 3).



Extended Data Fig. 6. Unique clonotypes are not restricted to a T_H type and do not diversify over time.

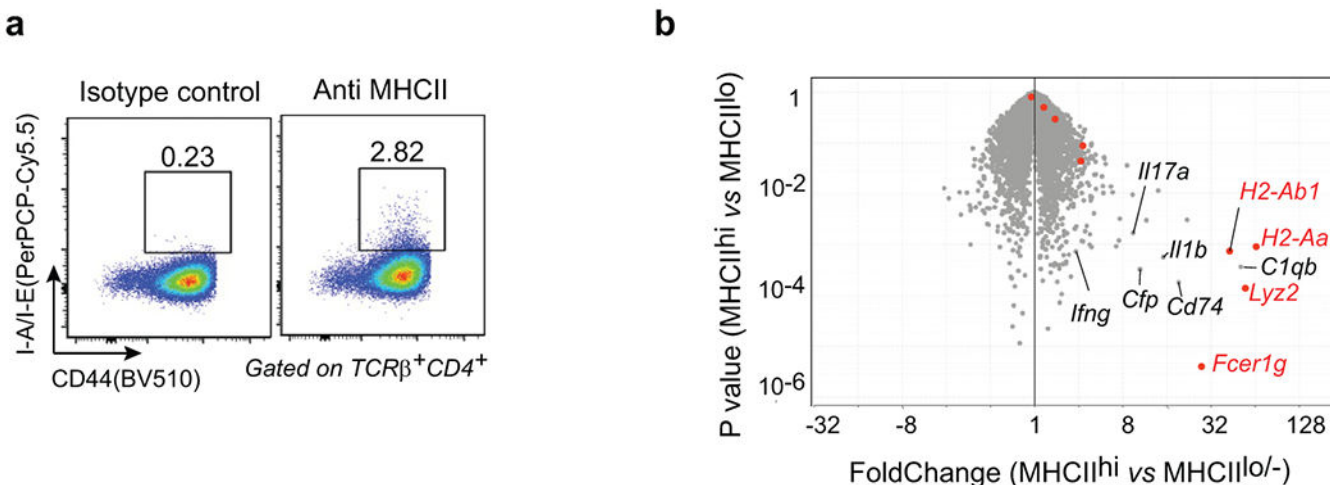
a) Quantification of flow cytometry data on cells from mouse LP at different timepoints of infection; Left: Proportion of CD4⁺ T cells within total CD45⁺; Middle: Proportion of Teff (CD44^{hi} Foxp3⁻) within total CD4⁺ T; Right: Proportion of IFN- γ ⁺ cells within total CD4⁺ T.

b) Cell numbers per scRNAseq clustering by day post infection. Treg clusters were identified as Foxp3⁺, naive cluster as Foxp3⁻ Ccr7⁺ and Teff clusters as Foxp3⁻ Cd44⁺.

c) Left: UMAP as in 5a, showing two groups of cell clusters: cells taken from mice after day 10 are colored in red, and cells taken prior to day 7 are colored in blue. Right: DEG analysis on top 20 differentially expressed genes between the two cluster groups. Asterisks represent genes that overlap with genes that are higher in Teff after *Salmonella* infection in Fig. 3a.

d) Bar graph representing proportions of cells belonging to singlet clones (clones that appear only once) or expanded clones (clones that appear more than once) in each of the clusters defined in S6b, grouped by day post infection.

e) Median Euclidean distances between cells within the same clonotype across the top 10 clonotypes for each timepoint. Euclidean distance was calculated based on the top 1000 variable genes. Each color dot represents a unique clonotype, and the size of the dot signifies the number of cells within each clonotype.



Extended Data Fig. 7. The unexpected MyT subset.

- a) Flow cytometric analysis (gated CD4⁺TCR β ⁺FOXP3⁻ Teff) cells from colonic LP of *Salmonella* infected mice.
- b) Volcano plot of bulk RNAseq from colonic Teff sorted as in C (LP of *Salmonella* infected mice). Genes highlighted in red belong to the myeloid genes listed in B.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

ACKNOWLEDGMENTS

We thank Drs. K. Murphy, N. Yosef, V. Kuchroo, R. Ramirez, D. Ramanan, M. Sassone-Corsi, E. Pamer, A. Anderson and J. Huh for insightful discussions, data and mouse lines; K. Hattori, C. Araneo, K. Seddu, D. Dionne and the Klarman Cell Observatory team, N. El-Ali and the Bauer Core Facility for help with mice, cell sorting and single-cell profiling. This work was supported by grants from the NIH to CB&DM (AI125603) and to the ImmGen

consortium (AI073072). EK was supported by a PhD Fellowship from Boehringer Ingelheim Fonds. KC was supported by NIGMS-T32GM007753 and a Harvard Stem Cell Institute MD/PhD Training Fellowship.

REFERENCES

1. Araneo BA, Marrack PC, & Kappler JW Functional heterogeneity among the T-derived lymphocytes of the mouse. II. Sensitivity of subpopulations to anti-thymocyte serum. *J Immunol* 114, 747–751 (1975). [PubMed: 1089717]
2. Bottomly K A functional dichotomy in CD4⁺ T lymphocytes. *Immunol Today* 9, 268–274 (1988). [PubMed: 2908229]
3. Mosmann TR et al. Two types of murine helper T cell clone. I. Definition according to profiles of lymphokine activities and secreted proteins. *J Immunol* 136, 2348–2357 (1986). [PubMed: 2419430]
4. Killar L et al. Cloned, Ia-restricted T cells that do not produce interleukin 4(IL 4)/B cell stimulatory factor 1(BSF-1) fail to help antigen-specific B cells. *J Immunol* 138, 1674–1679 (1987). [PubMed: 2950166]
5. Korn T, Bettelli E, Oukka M, & Kuchroo VK IL-17 and Th17 cells. *Annu. Rev. Immunol* 27, 485–517 (2009). [PubMed: 19132915]
6. Jabeen R & Kaplan MH The symphony of the ninth: the development and function of Th9 cells. *Curr. Opin. Immunol* 24, 303–307 (2012). [PubMed: 22365614]
7. Crotty S Follicular helper CD4 T cells (TFH). *Annu. Rev. Immunol* 29, 621–663 (2011). [PubMed: 21314428]
8. Murphy KM et al. Signaling and transcription in T helper development. *Annu. Rev. Immunol* 18, 451–494 (2000). [PubMed: 10837066]
9. O’Shea JJ & Paul WE Mechanisms underlying lineage commitment and plasticity of helper CD4⁺ T cells. *Science* 327, 1098–1102 (2010). [PubMed: 20185720]
10. Zhu J & Paul WE CD4 T cells: fates, functions, and faults. *Blood* 112, 1557–1569 (2008). [PubMed: 18725574]
11. Spits H et al. Innate lymphoid cells—a proposal for uniform nomenclature. *Nat. Rev. Immunol* 13, 145–149 (2013). [PubMed: 23348417]
12. Kelso A Th1 and Th2 subsets: paradigms lost? *Immunol. Today* 16, 374–379 (1995). [PubMed: 7546192]
13. Murphy KM & Stockinger B Effector T cell plasticity: flexibility in the face of changing circumstances. *Nat. Immunol* 11, 674–680 (2010). [PubMed: 20644573]
14. Geginat J et al. Plasticity of human CD4 T cell subsets. *Front Immunol* 5, 630 (2014). [PubMed: 25566245]
15. Stubbington MJT, Rozenblatt-Rosen O, Regev A, & Teichmann SA Single-cell transcriptomics to explore the immune system in health and disease. *Science* 358, 58–63 (2017). [PubMed: 28983043]
16. Miragaia RJ et al. Single-cell transcriptomics of regulatory T cells reveals trajectories of tissue adaptation. *Immunity* 50, 493–504 (2019). [PubMed: 30737144]
17. Lee JY et al. The transcription factor KLF2 restrains CD4(+) T follicular helper cell differentiation. *Immunity* 42, 252–264 (2015). [PubMed: 25692701]
18. Stoeckius M et al. Cell hashing with barcoded antibodies enables multiplexing and doublet detection for single cell genomics. *Genome Biol.* 19, 224 (2018). [PubMed: 30567574]
19. Yosef N et al. Dynamic regulatory network controlling TH17 cell differentiation. *Nature* 496, 461–468 (2013). [PubMed: 23467089]
20. Hartigan JA & Hartigan PM The dip test of unimodality. *The Annals of Statistics* 13, 70–84 (1985).
21. Pavlidis NG, Hofmeyr DP, & Tasoulis SK Minimum density hyperplanes. *The Journal of Machine Learning Research* 17, 5414–5446 (2016).
22. Lee Y et al. Induction and molecular signature of pathogenic TH17 cells. *Nat. Immunol* 13, 991–999 (2012). [PubMed: 22961052]

23. Ghoreschi K et al. Generation of pathogenic T(H)17 cells in the absence of TGF- β signalling. *Nature* 467, 967–971 (2010). [PubMed: 20962846]
24. Omenetti S et al. The Intestine Harbors Functionally Distinct Homeostatic Tissue-Resident and Inflammatory Th17 Cells. *Immunity* 51, 77–89 (2019). [PubMed: 31229354]
25. Meredith M, Zemmour D, Mathis D, & Benoist C Aire controls gene expression in the thymic epithelium with ordered stochasticity. *Nat Immunol* 16, 942–949 (2015). [PubMed: 26237550]
26. Buenrostro JD et al. Single-cell chromatin accessibility reveals principles of regulatory variation. *Nature* 523, 486–490 (2015). [PubMed: 26083756]
27. Yoshida H et al. The cis-Regulatory Atlas of the Mouse Immune System. *Cell* 176, 897–912 (2019). [PubMed: 30686579]
28. Schep AN, Wu B, Buenrostro JD, & Greenleaf WJ chromVAR: inferring transcription-factor-associated accessibility from single-cell epigenomic data. *Nat. Methods* 14, 975–978 (2017). [PubMed: 28825706]
29. Aghaeepour N et al. GateFinder: projection-based gating strategy optimization for flow and mass cytometry. *Bioinformatics*. 34, 4131–4133 (2018). [PubMed: 29850785]
30. Becht E et al. Reverse-engineering flow-cytometry gating strategies for phenotypic labelling and high-performance cell sorting. *Bioinformatics*. 35, 301–308 (2019). [PubMed: 29931307]
31. Baron CS et al. Cell type purification by single-cell transcriptome-trained sorting. *Cell* 179, 527–542 (2019). [PubMed: 31585086]
32. Tibbitt CA et al. Single-Cell RNA Sequencing of the T Helper Cell Response to House Dust Mites Defines a Distinct Gene Expression Signature in Airway Th2 Cells. *Immunity* 51, 169–184 (2019). [PubMed: 31231035]
33. Arazi A et al. The immune cell landscape in kidneys of patients with lupus nephritis. *Nat Immunol* 20, 902–914 (2019). [PubMed: 31209404]
34. Stoeckius M et al. Simultaneous epitope and transcriptome measurement in single cells. *Nat Methods* 14, 865–868 (2017). [PubMed: 28759029]
35. Burel JG et al. The challenge of distinguishing cell-cell complexes from singlet cells in non-imaging flow cytometry and single-cell sorting. *Cytometry A*(2020).
36. Cosgrove D et al. Mice lacking MHC class II molecules. *Cell* 66, 1051–1066 (1991). [PubMed: 1909605]
37. Alexander TB et al. The genetic basis and cell of origin of mixed phenotype acute leukaemia. *Nature* 562, 373–379 (2018). [PubMed: 30209392]
38. Priceman SJ et al. Regulation of adipose tissue T cell subsets by Stat3 is crucial for diet-induced obesity and insulin resistance. *Proc Natl Acad Sci U S A* 110, 13079–13084 (2013). [PubMed: 23878227]
39. McGeachy MJ et al. TGF-beta and IL-6 drive the production of IL-17 and IL-10 by T cells and restrain T(H)-17 cell-mediated pathology. *Nat. Immunol* 8, 1390–1397 (2007). [PubMed: 17994024]
40. Krausgruber T et al. T-bet is a key modulator of IL-23-driven pathogenic CD4⁺ T cell responses in the intestine. *Nat Commun.* 7, 11627 (2016). [PubMed: 27193261]
41. La Manno G et al. RNA velocity of single cells. *Nature* 560, 494–498 (2018). [PubMed: 30089906]
42. Becattini S et al. Functional heterogeneity of human memory CD4⁺ T cell clones primed by pathogens or vaccines. *Science* 347, 400–406 (2015). [PubMed: 25477212]
43. Wei SC et al. Negative co-stimulation constrains T cell differentiation by imposing boundaries on possible cell states. *Immunity* 50, 1084–1098 (2019). [PubMed: 30926234]
44. Li H et al. Dysfunctional CD8 T cells form a proliferative, dynamically regulated compartment within human melanoma. *Cell* 176, 775–789 (2019). [PubMed: 30595452]
45. Azizi E et al. Single-cell map of diverse immune phenotypes in the breast tumor microenvironment. *Cell* 174, 1293–1308 (2018). [PubMed: 29961579]
46. Wu TD et al. Peripheral T cell expansion predicts tumour infiltration and clinical response. *Nature* 579, 274–278 (2020). [PubMed: 32103181]

47. Zilionis R et al. Single-cell transcriptomics of human and mouse lung cancers reveals conserved myeloid populations across individuals and species. *Immunity*. 50, 1317–1334 (2019). [PubMed: 30979687]
48. Bielecki P et al. Skin inflammation driven by differentiation of quiescent tissue-resident ILCs into a spectrum of pathogenic effectors. *bioRxiv* 10.1101/461228 (2018).
49. Cano-Gamez E et al. Single-cell transcriptomics identifies an effectorness gradient shaping the response of CD4⁺ T cells to cytokines. *Nat. Commun* 11, 1801 (2020). [PubMed: 32286271]
50. Mostafavi S et al. Parsing the interferon transcriptional network and its disease associations. *Cell* 164, 564–578 (2016). [PubMed: 26824662]

METHODS-ONLY REFERENCES

51. Sefik E et al. Individual intestinal symbionts induce a distinct population of ROR γ ⁺ regulatory T cells. *Science* 349, 993–997 (2015). [PubMed: 26272906]
52. Hess J, Ladel C, Miko D, & Kaufmann SH Salmonella typhimurium aroA- infection in gene-targeted immunodeficient mice: major role of CD4⁺ TCR-alpha beta cells and IFN-gamma in bacterial clearance independent of intracellular location. *J Immunol* 156, 3321–3326 (1996). [PubMed: 8617956]
53. Collins JW et al. *Citrobacter rodentium*: infection, inflammation and the microbiota. *Nat Rev Microbiol.* 12, 612–623 (2014). [PubMed: 25088150]
54. Camberis M, Le GG, & Urban J Jr. Animal model of *Nippostrongylus brasiliensis* and *Heligmosomoides polygyrus*. *Curr. Protoc. Immunol* Chapter 19, Unit (2003).
55. Butler A et al. Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat Biotechnol.* 36, 411–420 (2018). [PubMed: 29608179]
56. Haghverdi L, Buettner F, & Theis FJ Diffusion maps for high-dimensional single-cell analysis of differentiation data. *Bioinformatics.* 31, 2989–2998 (2015). [PubMed: 26002886]
57. van Dijk D et al. Recovering Gene Interactions from Single-Cell Data Using Data Diffusion. *Cell* 174, 716–729 (2018). [PubMed: 29961576]
58. Ciofani M et al. A validated regulatory network for Th17 cell specification. *Cell* 151, 289–303 (2012). [PubMed: 23021777]
59. Muranski P et al. Th17 cells are long lived and retain a stem cell-like molecular signature. *Immunity.* 35, 972–985 (2011). [PubMed: 22177921]
60. Nurieva RI et al. Generation of T follicular helper cells is mediated by interleukin-21 but independent of T helper 1, 2, or 17 cell lineages. *Immunity.* 29, 138–149 (2008). [PubMed: 18599325]
61. Yusuf I et al. Germinal center T follicular helper cell IL-4 production is dependent on signaling lymphocytic activation molecule receptor (CD150). *J Immunol.* 185, 190–202 (2010). [PubMed: 20525889]
62. Kowalczyk MS et al. Single-cell RNA-seq reveals changes in cell cycle and differentiation programs upon aging of hematopoietic stem cells. *Genome Res.* 25, 1860–1872 (2015). [PubMed: 26430063]
63. Bodenhofer U, Kothmeier A, & Hochreiter S APCluster: an R package for affinity propagation clustering. *Bioinformatics.* 27, 2463–2464 (2011). [PubMed: 21737437]
64. Zeisel A et al. Brain structure. Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq. *Science* 347, 1138–1142 (2015). [PubMed: 25700174]
65. Wang C et al. Metabolic and epigenomic regulation of Th17/Treg balance by the polyamine pathway. *BioRxiv* www.biorxiv.org/content/10.1101/2020.01.23.911966v1.full (2020).

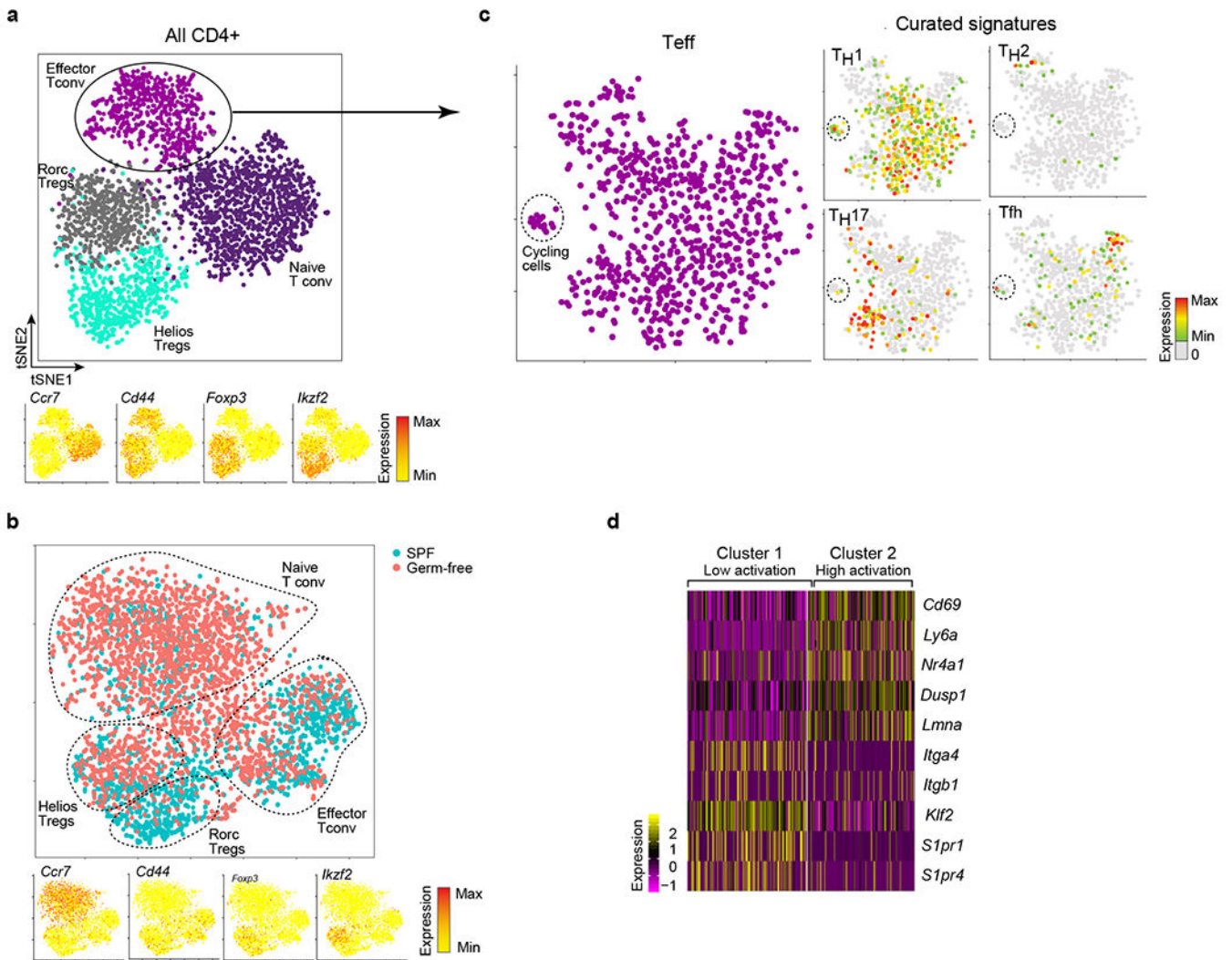


Figure 1: The transcriptional landscape of CD4⁺ T cells in the colon

a) scRNAseq analysis of total colonic LP CD4⁺ T cells of SPF mouse (computed from 658 most variable genes). Top: tSNE representation, color-coded by KNN cell clusters, identified based on expression of prototypic transcripts (bottom).

b) scRNAseq analysis of total colonic LP CD4⁺ T cells of Germ-free and SPF mice. Top: tSNE representation, color-coded by cell of origin. Marked clusters are identified based on expression of prototypic transcripts (bottom).

c) tSNE representation, restricted to CD4⁺ Teff selected in Fig. 1A (tSNE computed from 584 most variable genes). Right: overlay of combined expression of prototypic T_H genesets (Supplementary Table 2).

d) Heatmap of Teff divided into two clusters by KNN clustering. Representative genes overexpressed in each cluster are shown.

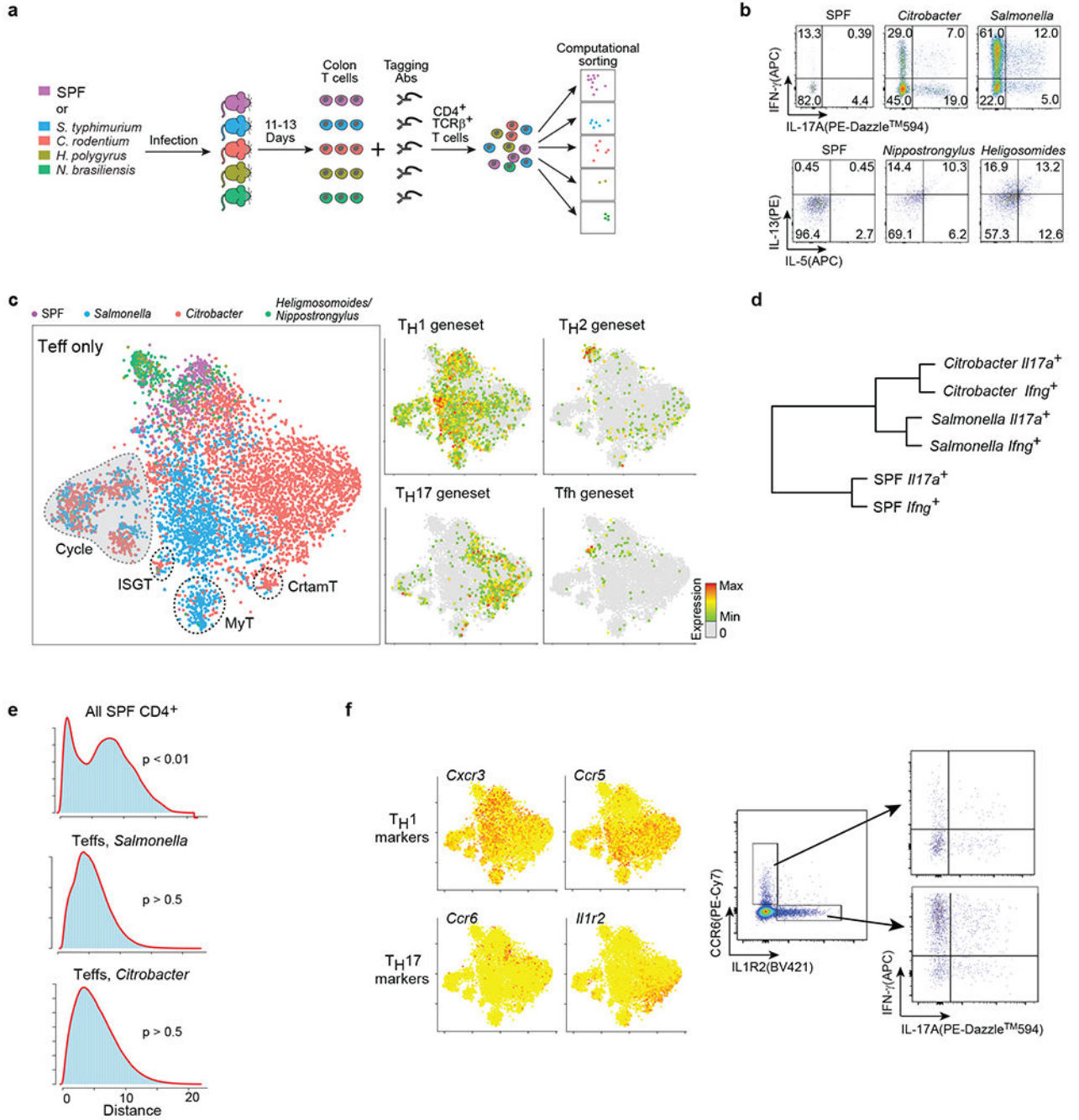


Figure 2: Variation in Teff transcriptomes shows continuous distribution that is not dictated by “T_H subsets”

a) Schematic of hash-tagging experiment. Mice were infected with different pathogens, their colonic LP cells extracted, labeled with hash-tagging antibodies, sorted as CD4⁺ T cells and processed as a single batch on the 10X chromium controller. Sample demultiplexing was done computationally.

b) Flow cytometric confirmation of the intestinal infections after intracellular staining for the cytokines shown (gated on CD4⁺TCRβ⁺FOXP3⁻CD44^{hi} cells).

- c) tSNE representation of Teff scRNAseq data from mice under different infection conditions (computed from 930 variable genes). Left: color-coded by condition/infection; right: overlay of combined expression of prototypic TH genesets.
- d) Dendrogram of Euclidean distances between cells in the scRNAseq dataset of Fig. 2c, splitting cells that express *Ifng* or *Il17a* in each of the infection conditions.
- e) Hartigan's Dip test applied to whole colonic CD4⁺ T cells from SPF mouse (top), or only to Teff from *Salmonella*-infected (middle) or *Citrobacter*-infected mice (bottom). MyT and cycling cells were taken out for this analysis.
- f) Expression of commonly-used markers of T_H subsets. Top: RNA expression in the scRNAseq data (overlaid on the tSNE from 2c). Bottom: Protein expression by flow cytometry in CD4⁺ Teff (gated CD4⁺TCRβ⁺FOXP3⁻CD44^{hi} cells) from *Salmonella*-infected mice.

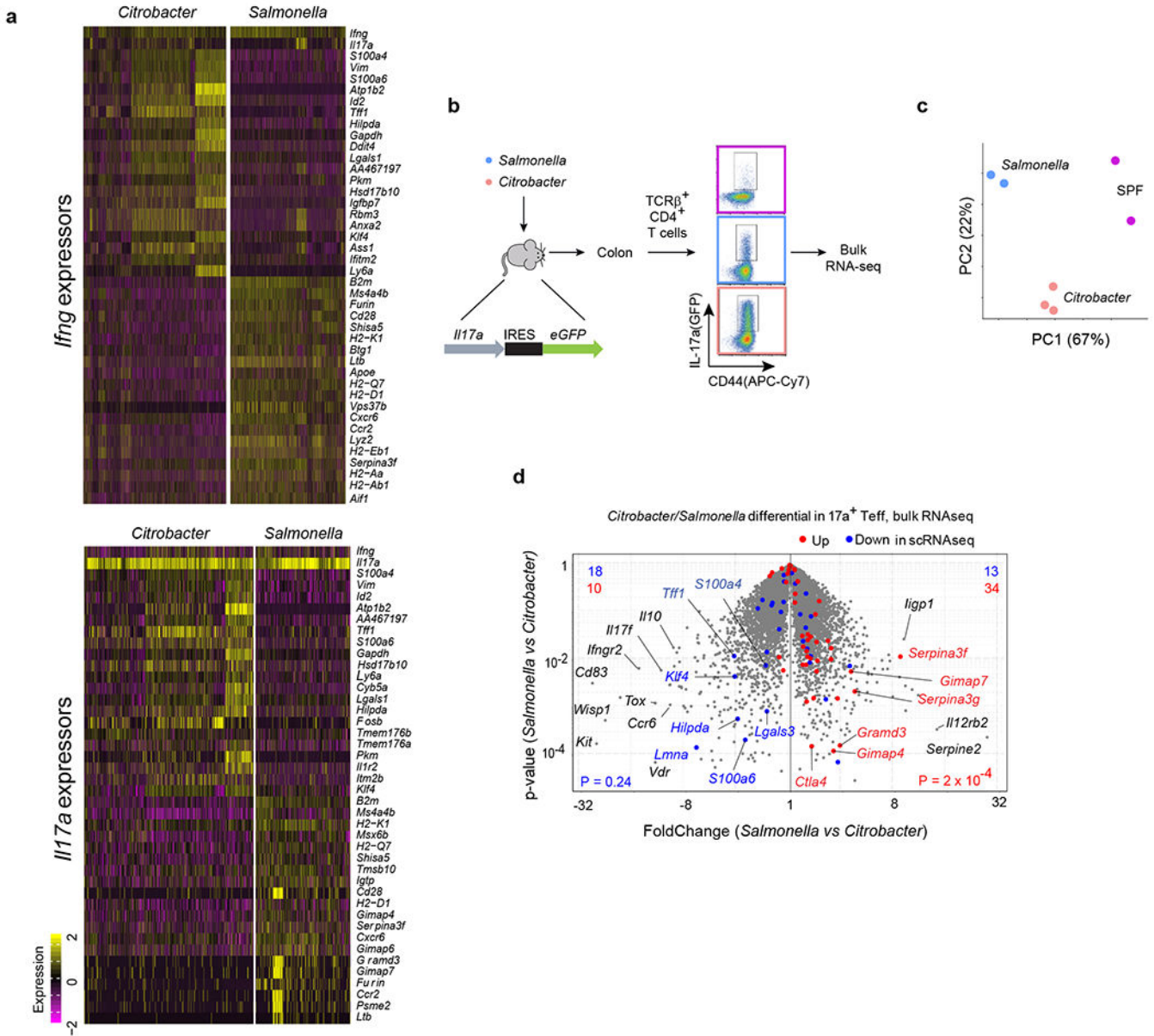


Figure 3: Teff phenotypes are distinguishable by infection rather than by TH type
 a) Heatmap of transcripts most differentially expressed within *Ifng*⁺ (top) or *Il17a*⁺ (bottom) Teff from *Citrobacter* or *Salmonella* infected mice.
 b) Sort of IL17a-expressing CD4⁺ CD44⁺ Teff cells from *Il17a-IRES-gfp* reporter mice (infected or not with *Salmonella* or *Citrobacter*) for expression profiling by ultra-low input RNAseq.
 c) PCA analysis of datasets from b.
 d) Volcano plot of bulk RNAseq data from C, comparing IL-17A⁺ Teff from *Salmonella* or *Citrobacter* infected mice. Red and blue highlights: transcripts found differentially expressed by scRNAseq.

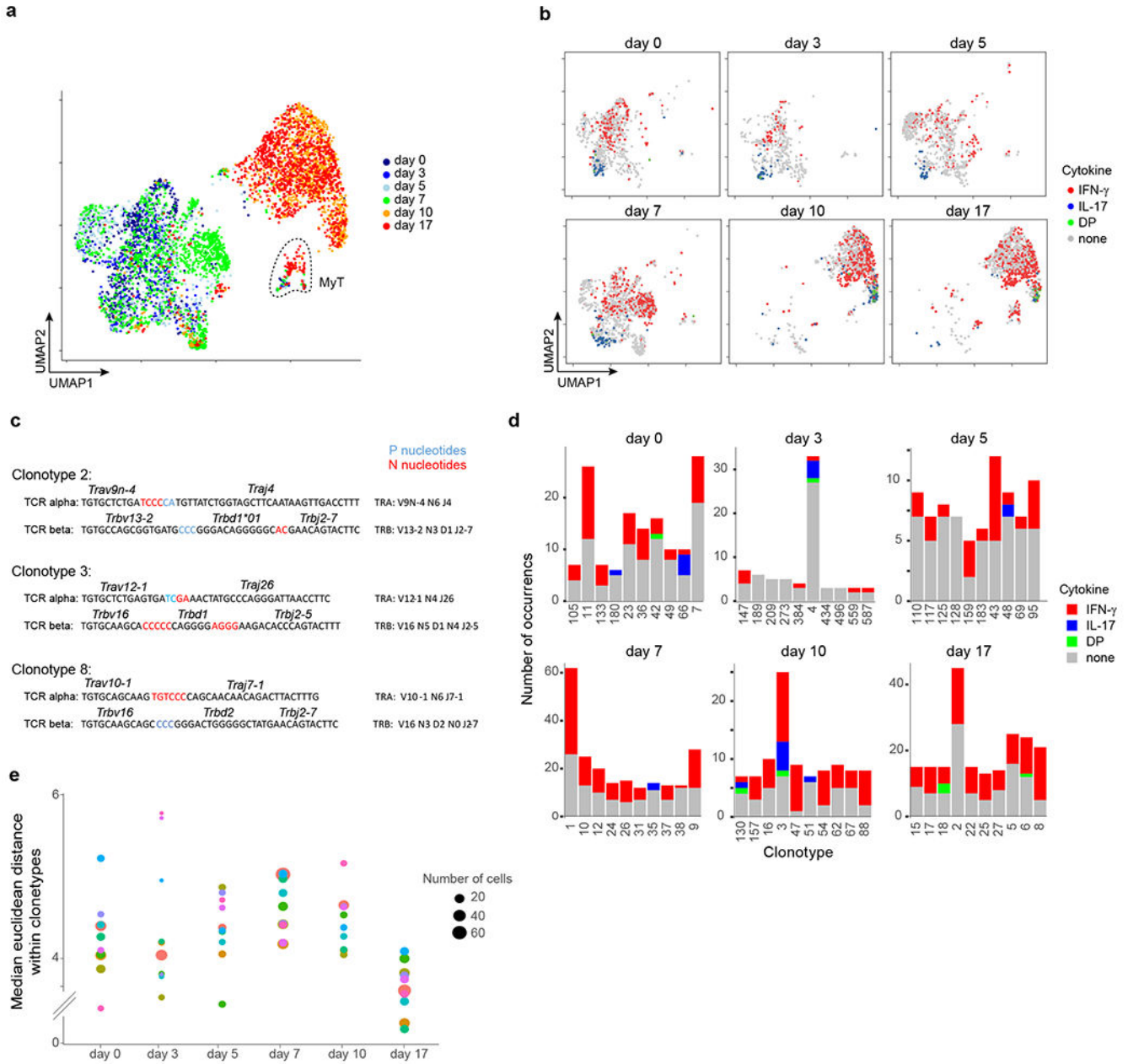


Figure 4: Repeated clonotypes can adopt different phenotypes, and do not diverge over time
 a) UMAP representation of Teff from mouse LP at different timepoints post-infection with *Salmonella*. MyT cells are circled.
 b) Cells from different time points, cytokine-producing cells highlighted as shown
 c) Representative examples of clonotypes with unique CDR3 identified by scTCR sequencing (non-germline N and P nucleotides shown).
 d) Numbers of *Ill17a*-, *Ifng*-, or both (DP) -expressing cells within the 10 most frequent clonotypes identified in each individual timepoint.
 e) Median Euclidean distances between cells within the same clonotype across the top 10 clonotypes for each timepoint. Euclidean distance was calculated based on T_H genes from

Supplementary Table 2. Clonotypes are color-coded, and the size denotes the number of cells that express each clonotype.

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

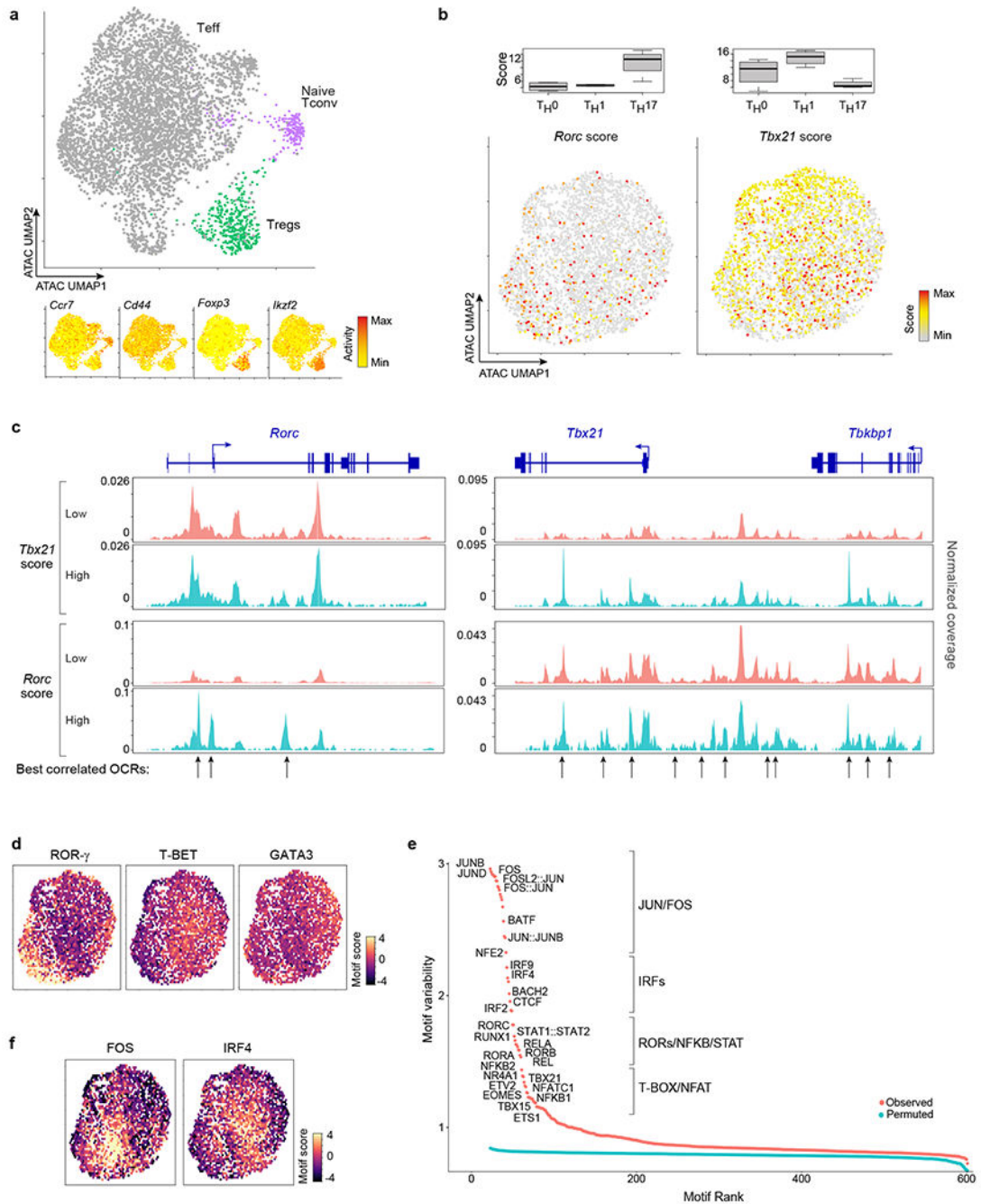


Figure 5: The chromatin states of TefFs are found on a continuum

a) scATACseq of total LP CD4⁺ T cells from Salmonella-infected mice. Top: UMAP representation, with Tregs and naive Tconv identified based on gene activity at prototypic loci (bottom).

b) Cell chromatin scores for *Rorc* and *Tbx21* loci, computed from accessibility of expression-correlated OCRs. Top: score in *in vitro* differentiated Th0, Th1 and Th17 cells. N=4 biological replicates for each condition. Centre, median; box limits, first and third

- percentiles; whiskers, 1.5× interquartile range (IQR). Bottom: scores for each cell in the scATACseq data from (a) (Teff only) in a UMAP plot.
- c) Aggregated coverage maps around *Rorc* and *Tbx21* loci in Teffs split based according to their chromatin score at each locus (shown at left); arrows: location of the best expression-correlated OCRs used to compute the scores.
 - d) Over-representation in each Teff cell (data from a) of TF-binding motifs in accessible chromatin (chromVAR bias-corrected TF motif deviation scores²⁸) for classic master regulators (UMAP framework from b)
 - e) Combined variability across the Teff scATACseq data for OCRs that contain motifs for different TFs (blue: null distribution for permuted dataset). TF families across the ranking shown at right.
 - f) TF motif deviation scores per Teff cell (as in d) for FOS and IRF4 motifs.

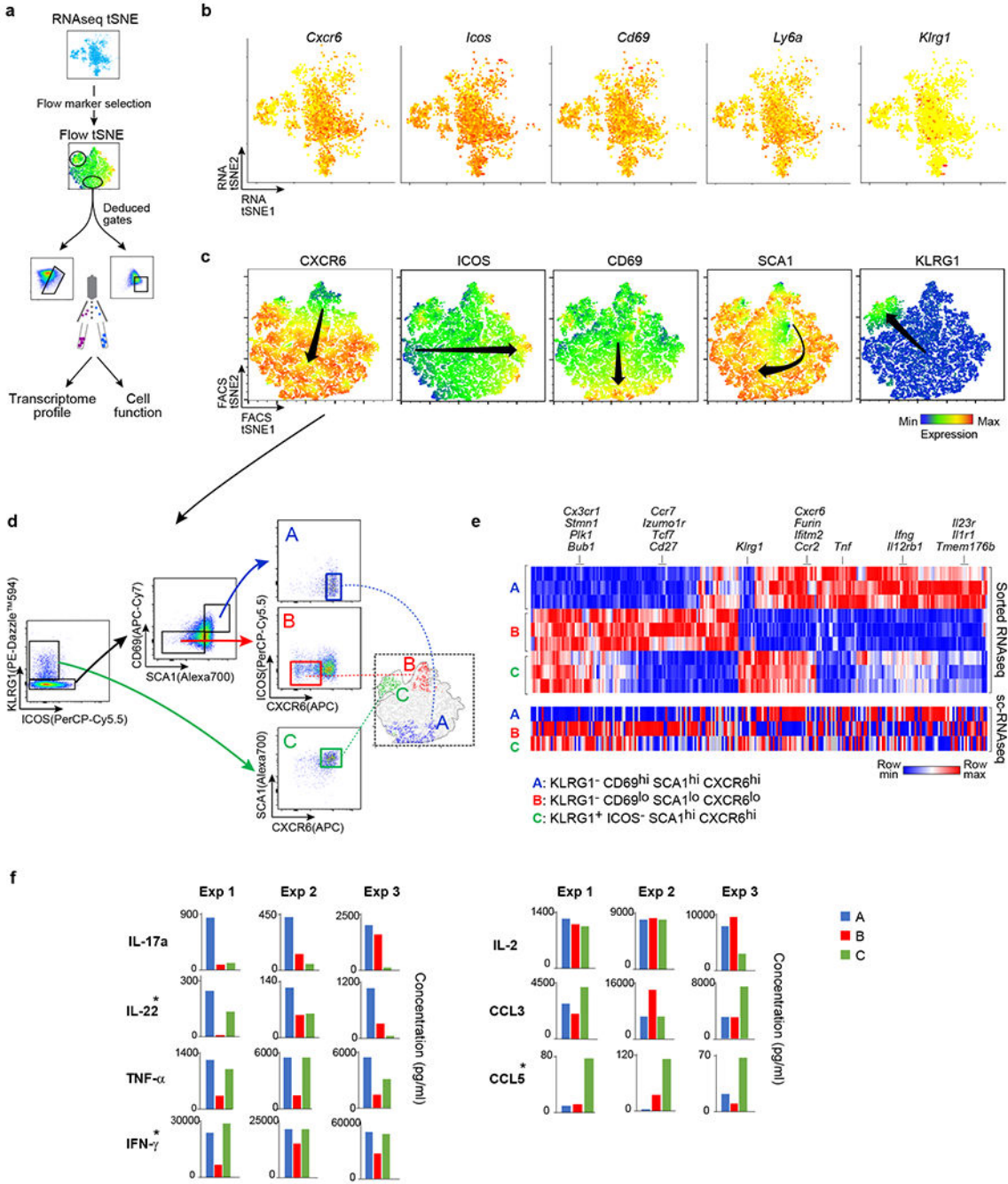


Figure 6: Transcriptional and functional validation of Teff continuity

a) Experiment schematic. Surface markers with continuous distribution in the scRNAseq are selected, and cells stained with corresponding antibodies for flow cytometry. tSNE plots are computed from the cytometry data, from which sorting gates are set to prepare cells for transcriptional and functional analysis.

b) Gene expression of selected surface markers in colonic Teffs from *Salmonella* infected mice (scRNAseq plot from 2c).

- c) Flow cytometry tSNE generated from fluorescence intensities of CD4⁺ Tregs stained for these markers.
- d) Sorting strategy, corresponding to the poles of the flow tSNE from C. The tSNE positions of the sorted cells are shown at right.
- e) Heatmap comparing differentially expressed genes in the bulk RNAseq profiling of populations A, B and C, sorted in D. Hierarchically clustered and row-mean normalized.
- f) Multiplex ELISA comparing secretion of cytokines and chemokines from populations A, B and C, where each bar is an independent biological replicate. *: cytokines with significant differences (at $p < 0.05$) between any two populations per paired student t-test (IL-17A: A vs B $p=4.0 \times 10^{-4}$, IL-22: A vs B $p=1.3 \times 10^{-2}$ and A vs C $p=3.6 \times 10^{-2}$, IFN- γ : A vs B $p=1.5 \times 10^{-2}$, CCL-5: A vs C $p=2.7 \times 10^{-2}$).

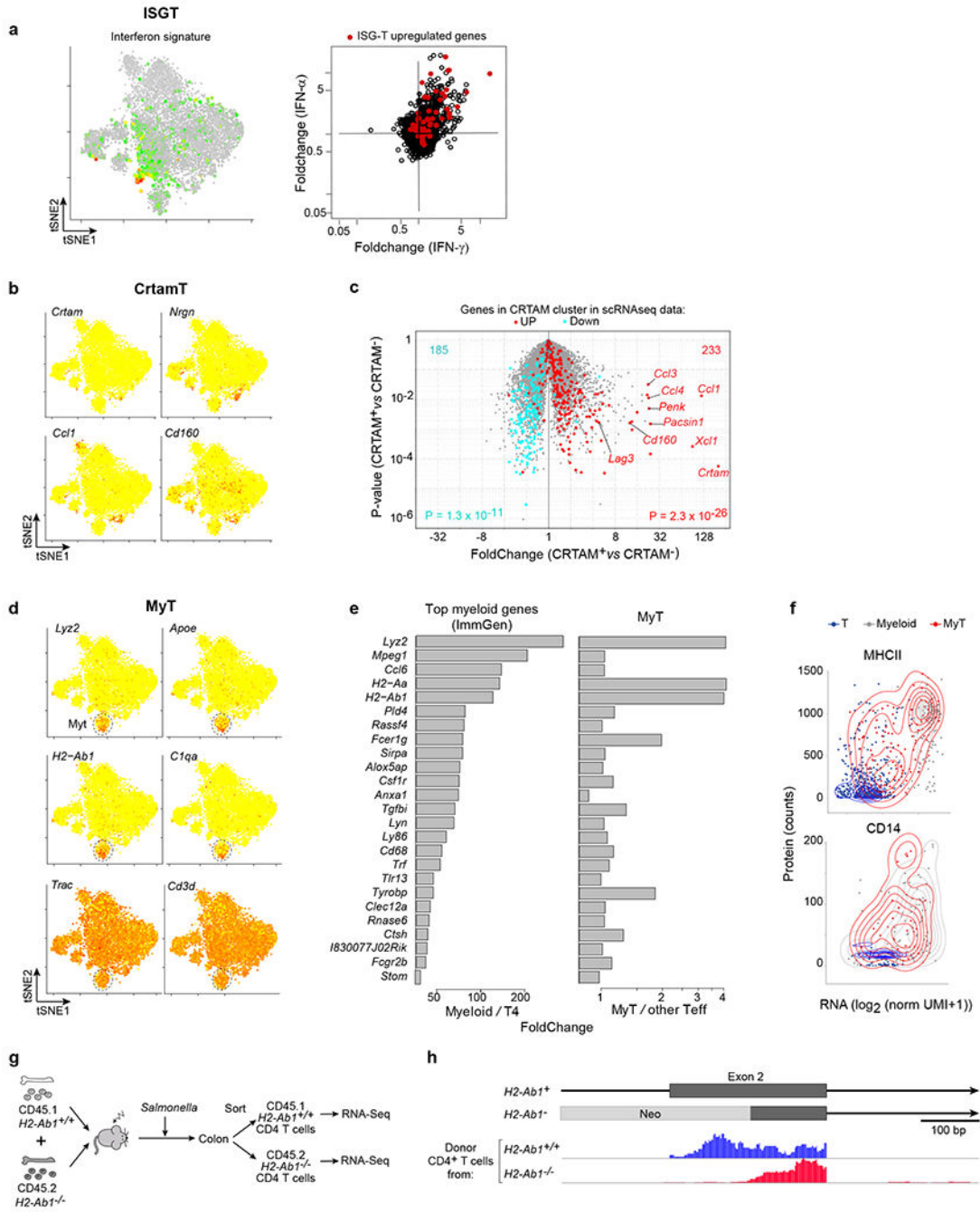


Figure 7: Novel Teff populations

a) The ISG-T subset. Left: Interferon-type-I signature⁵⁰ overlaid on the Teff tSNE. Right: Genes overexpressed in cluster ISG-T overlaid on top of genes upregulated in CD4⁺ T cells upon IFN- α or IFN- γ administration⁵⁰.

b) scRNAseq expression data of genes in *Crtam*⁺ cluster.

c) Volcano plot from RNA-seq of sorted CRTAM⁺ vs CRTAM⁻ colon Teff cells; over/under-expressed genes in the *Crtam*⁺ T cluster in scRNAseq data are shown in red and blue, with significance of overlap.

- d) Expression in MyT cells of genes overlaid on the general tSNE plot of Fig. 2c. top panels: typical myeloid cell transcripts; bottom: typical T cell transcripts.
- e) FoldChange histograms of myeloid-specific genes. Left, in myeloid vs CD4⁺ T cells (ImmGen RNAseq data); right: in MyT vs other colon Teff (*Salmonella*-infected, data from Fig. 2c). X axis on logarithmic scale.
- f) Contour plot representing RNA and protein expression in the single-cell data from Fig. 5 (x-axis: normalized scRNAseq; y-axis; raw CITE-seq counts) for MHCII (top) or CD14 (bottom). Individual cells are represented by dots, and are colored by their classification based on unsupervised clustering.
- g) Experiment schematic. Bone marrow from WT CD45.1 and CD45.2 *H2-Ab1*^{-/-} was mixed and transferred to irradiated CD45.1/CD45.2 hosts. After 8 weeks mice were infected with *Salmonella* and 13 days later the WT or KO LP CD4⁺ T cells were sorted for RNAseq.
- h) Top: schematic representation of the WT or KO *H2-Ab1* loci [neomycin resistance gene inserted into the second exon³⁶]. Bottom: position of RNAseq reads in colonic CD4⁺ Teff stemming from WT or *H2-Ab1* KO stem cells in mixed bone marrow chimeras infected with *Salmonella*.