<u>**Supplementary information:**</u>

**Computational Methods:**

**Gene expression analysis**

For in vitro differentiated $T_H$ cells, samples were normalized using DESeq[1] within the Strand NGS software suite. A gene was retained for further analysis if it had > 20 reads in all three replicates in at least one condition, resulting in 12,742 genes (Naive, $T_H$0+block, $T_H$0, $T_H$1, $T_H$1+IL-27, $T_H$2, $T_H$17, VitD3/Dex at times 0, 0.5, 2 and 4 hours post re-stimulation in vitro, and Foxp3RFP[+]IL-10GFP[+] or IL-10GFP[-] ex vivo) have values within cut-off). To identify genes of interest, we assessed the Pearson correlation coefficients of their expression across 34 samples with that of IL10 (Strand NGS); we identified transcription factors among these genes using a manually curated list using GO annotations and GeneSpring GX, Agilent Technologies).

For ex vivo CD4[+] T cells, all analyses were performed with the R statistical package version 3.3.1 (2016) and Bioconductor libraries version 3.3[2].  For each sample, expressed genes were identified by fitting a two-component Gaussian mixture to the $\log_2$ (raw count+1) value with mclust[3], using a probability threshold of belonging to the expressed class of 0.1. A gene was considered reliably expressed if it belonged to the expressed class in $\geqq$ 3 samples. The $\log_2$ intensity values of the identified 11,769 reliably expressed genes were normalized across all samples using limma[4]. Unsupervised hierarchical clustering analysis of the samples was performed using the Spearman correlation as a distance measure and the complete-linkage clustering using the R package gplots[5] (Figure 3a).

**Singular Value Decomposition (SVD) analysis**

SVD was performed on the filtered and quantile normalized mRNA expression data set (11,769 genes) to characterize the overall structure of the data and identify major sources of gene expression variation. Three linear models were fitted to each right singular vector: the full linear model, in which the design formula contains both the disease and the strain factors, and two reduced models, in which either the strain or the disease factor was used. To test the association of each principal component with the disease and/or the strain factor we performed an analysis of variance between the full linear model (strain and disease) and each of the two individual reduced models; both the *P*-value of the Chi-squared-test of ANOVA and the

Akaike Information Criterion (AIC) were considered. Finally, the most informative components were identified using the following criteria: (1) the fraction of explained variance in gene expression for a given component is above 4% (visual identification of the threshold which corresponds to the lower part of the elbow), (2) *P*-value of the Chi-test of ANOVA < 0.01 between full and the reduced model, and (3) AIC of the reduced model is lower than AIC of the full model (Figure 3b-c). To visualize the right singular vectors, we plotted the average expression of right singular vectors per sample group coloring all samples corresponding to *Maf*<sup>fl/fl</sup> black, and those corresponding to *Maf*<sup>fl/fl</sup> Cd4-cre white (Figure 3d). The biological interpretation of the principal components was facilitated by the identification of the genes whose expression profiles correlate and contribute most strongly (either positively or negatively) with the expression profile of the singular vector. The highest (most positive scores in both projection and correlation) and lowest (most negative scores in both correlation and projection) genes were selected for each singular vector using the K-mean clustering method allowing 10 clusters per component, and selecting those genes belonging to the most positive and negative cluster.

GO analysis of the genes contributing most to each component of the SVD analysis and the differentially up- and down-regulated genes was performed using a Fisher test with topGO Bioconductor package[8]. Only GO terms containing at least 10 annotated genes were considered. A *P*-value of 0.05 was used as the level of significance. The top significant GO terms were manually selected by removing redundant GO terms and terms which contain fewer than 5 significant genes (Figure 3d).

**Differential gene expression**

For differential gene expression analysis, samples were normalized using DESeq[1] within Strand NGS. For each infection condition, genes were then filtered (>20 reads in all three replicates), leaving 12,037 (malaria), 13,554 (HDM allergy) and 12,053 genes (EAE). with more than 20 reads where at least 100 percent of samples in any 1 out of 2 conditions (*Maf*<sup>fl/fl</sup> and *Maf*<sup>fl/fl</sup> Cd4-cre) have values within cut-off. Differentially expressed genes were determined by two-sided moderated *t*-test (Avadis NGS; cut off P<0.05 and absolute fold change>1.5); 2,635 (malaria), 1,073 in (HDM allergy) and 265 (EAE). Proportional Venn diagrams were generated using euler*APE*<sup>6</sup> (Figure 4a). Ingenuity Pathway Analysis (IPA) (QIAGEN Redwood City, www.qiagen.com/ingenuity) was used to retrieve the following annotations: transcription regulator, ligand dependent nuclear receptor, transmembrane

receptor + G-protein coupled receptor and cytokine + growth factor (249, 53 and 138 genes respectively = 440 in total of 3,967 differentially expressed genes). We used the IPA annotated gene-gene interactions to generate networks visualized with Cytoscape[7] (Figure 4b-d).

**ChIP-seq data analysis**

Raw sequencing reads for c-Maf ChIP-seq were obtained from GEO GSE40918 (single end, read length 36 nt)[9] and given the nature of the library was analyzed as follows. Reads were trimmed using Trimmomatic 0.36 (parameters HEADCROP:2 TRAILING:25 MINLEN:26)[10] and then mapped to the mouse genome mm10 using Bowtie 1.1.2 (parameters y -m2 --best --strata -S)[11]. Peaks were called for each replicate using MACS2 2.1.1 (default parameters; q-value < 0.01)[12] and a consensus peak set was generated from the union of both replicates; for overlapping peaks, the one with the best confidence score was kept. This resulted in 45,727 c-MAF ChIP-seq peaks (Supplementary figure 5). A consensus c-Maf binding motif was inferred from the ChIP-seq dataset using the CRUNCH suite[13] and validated using the ~2000 most confident ChIPseq peaks, as determined by q-value, using the MEME-ChIP[14] software (data not shown). The motif is shown in Supplementary figure 5. All queries for motif matches on both DNA strands within the ATAC-seq peak sequences were performed using FIMO[15].

**ATAC-seq data analysis**

ATAC-seq libraries were sequenced using Illumina HiSeq 2500 (paired end, lengths ranging from 50 to 100 nt) and given the nature of the library was analyzed as follows (method is distinct from ChIP-seq analysis owing to differences in data content). Adapters and low-quality bases were removed from reads using Skewer 0.2.2[16] (parameters -m pe -q 26 -Q 30 -e -l 30 -L 50). Reads were mapped to the mouse genome mm10 using BWA-MEM[17] with default parameters. Duplicates were removed using Picard 2.1.1[18]; discordant alignments, and/or with a mapQ<30 were discarded using SAMtools 1.3.1[19]. Mapped reads were shifted by +4 and -5 bp on the forward or reverse strands respectively to account for the transposase insertion. Fragments spanning nucleosomes (>99bp length) were removed as performed by Buenrostro *et al*[20]. Peaks representing open chromatin regions were identified for each sample using MACS2 2.1.1 using parameters designed for finding enrichment in cutting sites[12] (parameters --keep-dup all --nomodel --shift -100 --extsize 200; q-value < 0.01).

We used DiffBind 2.0.2[21] (parameters dba.count:minOverlap=0, score= DBA_SCORE_RPKM, bRemoveDuplicates=FALSE, bUseSummarizeOverlaps= TRUE;

dba.analyze: method=DBA_DESEQ2, bFullLibrarySize=T) to normalize for library sizes across all samples, and we calculated the Spearman correlation coefficients of normalized read counts between each pair of ATAC-seq sample. Samples were hierarchically clustered using the pairwise correlation coefficients and visualized using the BioConductor ComplexHeatmap library[22] (Figure 6a).

For each disease model, we defined a consensus set of ATAC-seq peaks as the union of peaks found in the *Maf* [fl/fl] Cd4-cre and *Maf* [fl/fl] samples (Malaria: 87,533; HDM: 54,745; EAE: 42,286 peaks). Diffbind 2.0.2 was also used to identify changes in ATAC-seq peaks between *Maf* [fl/fl] Cd4-cre and *Maf* [fl/fl], interpreted as chromatin remodeling events (Figure 6b; absolute fold-change in read coverage>1.5 and FDR<0.05). The sequences underlying the 1,273 remodeled peaks belonging to the malaria dataset were subjected to *de novo* motif discovery using MEME-ChIP[14].

**Assigning direct and indirect targets of c-Maf regulation**

ATAC-seq peaks were defined as c-Maf-associated if they overlapped with a c-Maf ChIP-seq peak or contained a c-Maf-motif match (Supplementary figure 5). The distance distributions between ATAC-seq peaks and annotated transcription start sites (TSS) show that c-Maf-associated peaks tend to occur much closer to genes, with most within 3kb of the TSS. A gene was assigned to an ATAC-seq peak if the peak overlapped or fell within +/- 3kb of the gene body boundaries; assignments were performed using the ChIPseeker BioConductor library[23]. In each treatment condition, a differentially expressed gene was defined as a direct c-Maf target if it was assigned to a c-Maf-associated ATAC-seq peak (1,828 genes in Malaria, 631 in HDM, 149 in EAE; Supplementary figure 5). All others were defined as indirect targets.

We tested the enrichment of c-Maf-associated ATAC-seq peaks among differentially expressed genes, compared with non-differentially expressed genes (Extended Data Table 1). The enrichment is statistically significant for HDM and EAE (p<2.2e-16 and p<1.704e-03 respectively; Chi-squared test). There is also an enrichment for malaria though it does not meet the threshold for statistical significance (p=0.06; Chi-squared test); this is in line with observations that malaria samples display a much broader set of differentially expressed genes (Figure 4a).

In order to highlight genes with high c-Maf abundance within accessible regions from those with few c-Maf sites we calculated a score for every gene $g$:

$$Score(g) = \frac{1}{k}\sum_{i=1}^{j} -log_{10}(C_i)$$

where $k$ is the number of ATAC-seq peaks assigned to gene $g$; $j$ is the number of ChIP-seq peaks that intersect any of the $k$ ATAC-seq peaks; $C$ is the q-value confidence score for a ChIP-seq peak. These scores were converted to rank-based quantiles. Same methodology was applied for motif data, using the $P$-value of the match as $C$. These scores are used to display the heatmaps in (Figure 6c).

We cross-checked the direct and indirect target assignments using, the Binding and Expression Target Analysis (BETA) software[24] (parameters -g mm10 --da 1 --df 0.05 -c 1). BETA takes as input TF-binding and gene expression data, modelling the regulatory potential of a binding site according to its distance to the TSS. BETA does not accept fold-change cutoffs to denominate differentially expressed genes, therefore to ensure the same set of differentially expressed genes we set the fold-changes of non-differentially expressed genes to 0, upregulated genes to 1 and down-regulated to -1, and left the $P$-values unchanged (used by BETA to rank the expression changes). The ChIP-seq data was intersected with the ATAC-seq data, thus, only ChIP-seq peaks within accessible regions in each context would affect the outcome of the software. Heatmap visualization of these scores was done using the ComplexHeatmap BioConductor library (Figure 6c).

**Genome-wide differential footprints**

To identify regulators with potential differences in TF-binding in $Maf^{fl/fl}$ Cd4-cre and $Maf^{fl/fl}$ samples, we applied the BaGFoot software using all ATAC-seq peaks identified in each treatment condition[25]. BaGFoot predicts these changes by searching for TF-binding motif matches in regions with altered ATAC-seq insertion patterns between two conditions. We used all 129 motifs of class A and B quality in the HOCOMOCO database v10[26]. Since BaGFoot currently does not consider replicates we performed three $Maf^{fl/fl}$ Cd4-cre and $Maf^{fl/fl}$ pair-wise comparisons for each disease model and calculated the average changes in accessibility and footprint-depth. Results are displayed as bagplots, using a fence of factor 2 (Figure 7a). We identified TFs with potentially altered binding by identifying the outliers of the multivariate

distribution, as assessed by the Mahalanobis distance of each TF to the multivariate distribution. The statistical significance of these distances was tested using a Chi-square distribution followed by a Benjamini-Hochberg correction for multiple-testing, the recommended approach by BaGFoot.

We also assessed if any of the TFs identified by BaGFoot could explain the expression changes of the indirect c-Maf targets. For this, we tested whether the corresponding motif is enriched within the accessible neighbourhood of differentially expressed genes compared with non-differentially expressed genes using a Fisher's exact test, with Benjamini-Hochberg correction for multiple-testing (q-value<0.05) (Table 2 and Figure 7a)

The displayed metaprofile of Tn5 insertions, the footprint, was corrected for Tn5 insertion bias obtained from BaGFoot software. The footprint shown depicts the average of the three biological replicates, the dashed lines correspond to the average Tn5 insertions in such metaprofile.

**Visualization of sequencing data**

All sequencing data presented in Genome Browser views were normalized to RPKMs using the bamCoverage software in DeepTools 2.4.2[27]. Tracks were visualized using IGV 2.3.89 [28], with replicates overlaid on top of each other. The fold-change values of ATAC-seq peaks were retrieved by DiffBind 2.0.2, these represent changes in chromatin accessibility (negative and positive values being a reduction or gain in accessibility, respectively, and 0 means no change). c-Maf ChIP-seq peak q-values were retrieved with MACS2 (q-values were -log10 transformed, thus, the greater the number the higher the confidence of existence of a peak). A bedgraph file was generated for each data type and treatment. The resulting bedgraph files were imported to IGV and visualized using "heatmap" option.

**References**

1.  Anders, S. & Huber, W. Differential expression analysis for sequence count data. *Genome Biol.* **11,** R106, doi:10.1186/gb-2010-11-10-r106 (2010).

2.  R: A language and environment for statistical computing. R Foundation for Statistical Computing. Vienna, Austria (2014).

3.  Fraley, C., Raftery, A. E., Murphy, T. B. & Scrucca, L. mclust Version 4 for R: Normal Mixture Modeling for Model-Based Clustering, Classification, and Density Estimation. 2012.

4.  Smyth, G. K., Ritchie, M., Thorne, N. & Wettenhall, J. LIMMA: linear models for microarray data. In Bioinformatics and Computational Biology Solutions Using R and Bioconductor. 397-420 (Springer, New York 2005).

5.  Warnes, G. R., Bolker, B., Bonebakker, L. & Gentleman, R. gplots: Various R programming tools for plotting data. *R package version* (2009).

6.  Micallef, L. & Rodgers, P. eulerAPE: drawing area-proportional 3-Venn diagrams using ellipses. *PLoS One* **9,** e101717, doi:10.1371/journal.pone.0101717 (2014).

7.  Shannon, P. *et al.* Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res.* **13,** 2498–2504, doi:10.1101/gr.1239303 (2003).

8.  Alexa, A. & Rahnenfuhrer, J. topGO: enrichment analysis for gene ontology. *R package version* (2010).

9.  Ciofani, M. *et al.* A validated regulatory network for Th17 cell specification. *Cell* **151,** 289–303, doi:10.1016/j.cell.2012.09.016 (2012).

10. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* **30,** 2114–2120, doi:10.1093/bioinformatics/btu170 (2014).

11. Langmead, B., Trapnell, C., Pop, M. & Salzberg, S. L. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol.* **10,** R25, doi:10.1186/gb-2009-10-3-r25 (2009).

12. Zhang, Y. *et al.* Model-based analysis of ChIP-Seq (MACS). *Genome Biol.* **9,** R137, doi:10.1186/gb-2008-9-9-r137 (2008).

13. Berger, S. *et al.* Crunch: Completely Automated Analysis of ChIP-seq Data. doi:10.1101/042903 (2016).

14. Ma, W., Noble, W. S. & Bailey, T. L. Motif-based analysis of large nucleotide data sets using MEME-ChIP. *Nat. Protoc.* **9,** 1428–1450, doi:10.1038/nprot.2014.083 (2014).

15. Grant, C. E., Bailey, T. L. & Noble, W. S. FIMO: scanning for occurrences of a given motif. *Bioinformatics* **27,** 1017–1018, doi:10.1093/bioinformatics/btr064 (2011).

16. Jiang, H., Lei, R., Ding, S.-W. & Zhu, S. Skewer: a fast and accurate adapter trimmer for next-generation sequencing paired-end reads. *BMC Bioinformatics* **15,** 182, doi:10.1186/1471-2105-15-182 (2014).

17. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv [q-bio.GN]* (2013).

18. Picard Tools - By Broad Institute. Available at: http://broadinstitute.github.io/picard/.

19. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25,** 2078–2079, doi:10.1093/bioinformatics/btp352 (2009).

20. Buenrostro, J. D., Giresi, P. G., Zaba, L. C., Chang, H. Y. & Greenleaf, W. J. Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome position. *Nat. Methods* **10,** 1213–1218, doi:10.1038/nmeth.2688 (2013).

22. Stark, R. & Brown, G. DiffBind: differential binding analysis of ChIP-seq peak data. *Bioconductor* Available at: http://bioconductor.org/packages/release/bioc/html/DiffBind.html.

22. Gu, Z., Eils, R. & Schlesner, M. Complex heatmaps reveal patterns and correlations in multidimensional genomic data. *Bioinformatics* **32,** 2847–2849, doi:10.1093/bioinformatics/btw313 (2016).

23. Yu, G., Wang, L.-G. & He, Q.-Y. ChIPseeker: an R/Bioconductor package for ChIP peak annotation, comparison and visualization. *Bioinformatics* **31,** 2382–2383 (2015).

24. Wang, S. *et al.* Target analysis by integration of transcriptome and ChIP-seq data with BETA. *Nat. Protoc.* **8,** 2502–2515, doi:10.1093/bioinformatics/btv145 (2013).

25. Baek, S., Goldstein, I. & Hager, G. L. Bivariate Genomic Footprinting Detects Changes in Transcription Factor Activity. *Cell Rep.* **19,** 1710–1722, doi:10.1016/j.celrep.2017.05.003 (2017).

26. Kulakovskiy, I. V. *et al.* HOCOMOCO: expansion and enhancement of the collection of transcription factor binding sites models. *Nucleic Acids Res.* **44,** D116–25, doi:10.1093/nar/gkv1249 (2016).

27. Ramírez, F. *et al.* deepTools2: a next generation web server for deep-sequencing data analysis. *Nucleic Acids Res.* **44,** W160–5, doi:10.1093/nar/gkw257 (2016).

28. Robinson, J. T. *et al.* Integrative genomics viewer. *Nat. Biotechnol.* **29,** 24–26, doi:10.1038/nbt.1754 (2011).