## Supplementary Methods

### Supplementary references

DSB (Denoised and Scaled by Background) normalisation method[1]

Harmony package for batch correction (v1.0)[2]

### Sample preparation

*Fetal developmental stage assignment and chromosomal assessment*

Developmental age was estimated from standardized measurements of foot length and heel-to-knee length[3]. Quantitative Fluorescence PCR of chromosomes X, Y, 13, 15, 16, 18, 21 and 22 was performed on fetal skin or chorionic villi to assign gender and exclude common chromosomal abnormalities. In Down syndrome FBM samples, *GATA1* mutation was excluded as previously described[4].

*Plate-based scRNA-seq*

Two FBM suspensions (both 17 PCW) were prepared for FACS sorting (see Methods; antibody details in **Supplementary Table 26;** metadata in **Supplementary Tables 1, 15**). Target populations were gated as shown in **Extended Data Fig. 1e**. From live, CD45$^+$ single cells, CD123$^+$HLA-DR$^-$ basophils and CD123$^+$HLA-DR$^+$ pDCs were gated. From the remaining cells, CD34$^+$CD117$^{mid-hi}$ progenitors and CD117$^{hi}$ mast cells were gated. Next, CD125$^+$FSC$^{hi}$ eosinophils were gated. Subsequently, HLA-DR$^+$CD79a$^+$ B cells were separated. As CD79 is weakly expressed on the cell surface, a significant number of B cells fell in the HLA-DR$^+$CD79$^-$ gate, forming the CD14$^-$CD204$^-$CLEC9A$^-$CD1c$^-$ population, thus did not enter subsequent sort gates. From HLA-DR$^+$ cells, CD14$^+$CD204$^-$ monocytes were gated. Within the CD14$^-$CD204$^-$ population, CLEC9A$^+$ DC1 and CD1c$^+$ DC2 were identified. From HLA-DR$^-$ cells, CD11b$^-$ CD52$^+$ T and NK cells were excluded. From the remaining cells CD11b$^-$CD66b$^-$ promyelocytes, CD11b$^-$CD66b$^+$ metamyelocytes/myelocytes and CD11b$^+$CD66b$^+$ neutrophils were selected.

Note, the CD11b gating was as per published reports, mature neutrophils being CD11b[+] and immature neutrophils being CD11b[-5]. The number of sorted cells per subset were: 68 neutrophils, 68 myelocytes/metamyelocytes, 56 promyelocytes, 60 pDCs, 72 HSC/MPPs, 60 pDCs, 72 B cells, 70 eosinophils, 70 mast cells, 72 monocytes, 70 DC1 and 70 DC2. Plates containing lysed single cells were processed using a modified Smart-seq2 protocol[6]. Libraries were generated using the Nextera XT kit (Illumina) with 384 cells per library. Cells were barcoded using Index v.2 sets A, B, C and D (Illumina). Libraries were sequenced using an Illumina NextSeq 550 on High-output mode to achieve a minimum of 1 million reads per cell.

*CITE-seq experiment optimisation*

The 'TotalSeqA' 198 oligo-conjugated panel was generated in collaboration with Biolegend and includes 4 isotype controls (**Supplementary Table 27**). Biolegend undertook extensive optimisation of this panel, including pooling, titration experiments, extensive testing on in-house cell lines and optimisation of the stability of the final product. We performed pilot experiments isolating MNCs and CD34[+] cells from cord blood, to be stained either before or after FACS sorting with the TotalSeqA panel, i) to optimise the staining protocol with primary cells and ii) to identify antibodies that still needed titrating. These samples were taken all the way through to Illumina sequencing. Bioinformatic analysis of these CITE-seq data was performed independently by Biolegend in La Jolla, as well as our in-house team in Cambridge. This analysis suggested that any antibody which had more than 3.5% of the total reads would need to be 'competed' in subsequent experiments. Moreover, staining after sorting resulted in high background and non-specific staining, and was therefore abandoned as an approach for subsequent samples. The following antibodies contributed to more than 3.5% of the total reads: CD7, CD5, CD47, CD24 and CD325. Biolegend supplied 'cold antibodies' (antibodies without the oligo conjugation) to be included as a competition cocktail in the staining protocol with the TotalSeqA panel (**Supplementary Table 26**). A second pilot experiment was then performed

where primary cells were stained with the TotalSeqA panel and the cold competition antibodies before FACS sorting. Analysis of this experiment by Biolegend showed that the antibodies were balanced and the sequencing space well allocated amongst all antibodies. All subsequent experiments were performed as described in the CITE-seq section (see Methods), with staining of the TotalseqA and competition cocktail performed in parallel before sorting.

**Data analysis**

*Differentially expressed gene statistics*

Differential gene expression analysis referenced in text and shown in violin plots were run using the two-sided Wilcoxon rank-sum statistical test with Benjamini-Hochberg procedure for multiple testing correction. *p*-values are shown in the relevant Supplementary Tables.

*Statistics from barplots*

For cell type proportion analysis across gestational ages and different tissues, proportions were modelled as a quasibinomial distribution (see further information below in '*Calculating differences in cell type proportions across gestational stages and organ*'). *P*-values for the significance of change in proportion between conditions were assessed and values subsequently detailed in **Supplementary Tables 18, 22, 23, 38, 45**. For all barplot statistics, asterisks were used to indicate significant changes in proportion, with *, **, *** and **** representing *p*-values of <0.05, <0.01, <0.001 and 0.0001 respectively. Directionality of trend between proportions and stage were accessed by the Spearman's rho coefficient computed between proportions and stage. Increasing or decreasing trends were denoted with 'up' and 'down' arrows respectively.

*Statistics from colony experiments*

Statistical analysis of culture wells producing colonies between FBM and FL HSC/MPPs by Mann Whitney test yielded *p*=0.011 (* with 7 96-well plate replicates from 3 biologically independent samples). Comparison between FBM and FL committed progenitors by the same method yielded *p*=0.0006 (***). Comparison of number of colony types per well for paired progenitor types was performed by binomial test, comparing 1 colony type with >1 colony type. 2-sided *p*-values were 0.0008 for FBM and FL HSC/MPPs (**, n=164) and 0.27 for FBM and

FL committed progenitors (replicates as above). Comparison of the number of myeloid-only colonies for paired FBM and FL HSC/MPPs was performed by a binomial test, comparing 'myeloid only' with 'myeloid+other ($p$=0.0001 from k=77 myeloid-colony wells).

*Logistic regression for label transfer of annotations*

Label transfer and probability of label correspondence between GEX matrices in single cell datasets were carried out using a Logistic Regression (LR) model trained against the FBM 10x data. The LR model was built utilising the 'sklearn.linear_model.LogisticRegression' module in the sklearn package (v0.22). The LR model was trained on normalized gene expression data (x_var =33,712) and annotations (x_sample=64) of the whole discovery FBM scRNA-seq dataset (103,228 cells), and used to predict the probability of correspondence between labels in the target dataset. The model used was L2 Ridge Regression regularised with penalty strength of 0.2 using the 'lbfgs' solver. Predicted label probability distribution within pre-computed clusters in the target data were used to assign cluster identity by majority vote. We further assessed the predicted labels by computing Adjusted Rand index and Mutual information scores from the modules 'sklearn.metrics.adjusted_rand_score' and 'sklearn.metrics.mutual_info_score' between original cell labels and predicted comparative cluster labels in overlapping cell states in each dataset.

*Kernel density estimation of cellular abundance*

Gaussian kernel density functions (scipy v1.4.1) were estimated in the FDG space (bandwidth of 0.1) separately using 7,500 sampled cells from each tissue-specific dataset. To visualise cell densities on the common landscape, values were computed for all cells with respective kernel density functions.

*Comparisons between scRNA-seq datasets*

To compare gene expression programs across tissue compartments (and scRNA-seq datasets), we combined raw GEX matrices and performed joint matrix transformation and dimension reduction as previously described in Methods. Datasets were integrated using Harmony with tissue source or sample as a covariate. DEG analysis was performed as described in '*Annotation of clusters*', unless otherwise stated in tool-specific description in Methods (using two-sided Wilcoxon rank sum testing with Benjamini-Hochberg for multiple testing adjustment).

*Calculating differences in cell type proportions across gestational stages and organ* GraphPad Prism (v8.1.0) was used for plotting and statistical comparison. Statistically significant differences in cell type proportions across tissue were conducted using a one-way ANOVA with Tukey's multiple comparison tests. Significance was noted on corresponding scatter plots using asterisks, where scatter plots display proportion per biological replicate. Cell-type proportions per sample were obtained by adjusting observed proportions by $CD45^+$/$CD45^-$ sort gate (**Supplementary Table 19**).

For cell type proportion statistical analysis across gestational stage and tissue, proportions were modelled as a quasibinomial distribution. For both analyses the condition (gestational stage or tissue) was provided as a covariate for the proportion of the cell type being assessed. The quasibinomial model was fitted using glm from the MASS R package (v7.3-54). The *p*-value for the significance of the change in proportion at 95% CI between conditions was assessed using a one-sided likelihood-ratio test, computed using the *anova.glm* function and Bonferroni adjusted for multiple comparisons using the *p.adjust* function (both functions from the stats (v3.6.2) R package). Significant changes in cell type proportion were highlighted on bar plots using asterisks. The sign of Spearman's rho and Pearson's R coefficients computed between proportions and stage were used to assess directionality of monotonic trends of flux between

analogous cell states of different developmental stages. Increasing or decreasing trends were denoted with 'up' and 'down' arrows respectively.

*Differential abundance testing using Milo*

To identify cell subpopulations enriched or depleted in DS vs disomic FBM scRNA-seq samples, we used Milo (development version) for differential abundance testing on KNN graph neighbourhoods. This was implemented in the R package miloR (https://github.com/MarioniLab/miloR). Briefly, we performed KNN graph embedding on DS and disomic FBM samples at matched age (12-13 PCW) using the reduced dimension space derived from Harmony integration (using the same number of dimensions and value of K used for UMAP embedding). We used a refined sampling algorithm to select a subset of cells spanning the KNN graph (defined as index cells) and we counted the number of cells from each sample in the neighbourhoods of index cells, where a neighbourhood is defined as the group of cells connected by an edge to an index cell. We tested for differences in abundance between the cells from DS and non-DS samples in each neighbourhood using the Quasi-likelihood test implemented by edgeR, controlling the FDR across the graph neighbourhoods. To assign cell type annotations to neighbourhoods, we took the most frequent annotation between cells in each neighbourhood. Neighbourhoods are generally homogeneous, retaining neighbourhoods with at least 60% of cells belonging to the most abundant cell type label.

*Visualising for single cell protein and RNA expression*

Gene and protein expression dot plots were visualised using the *sc.pl.dotplot* function in the Scanpy package in Python; dot colour indicated the mean expression values and dot size indicated the proportion of cells in each category that expressed the given marker. Violin plots were also visualised using the Scanpy package in Python, using the *sc.pl.violin* function. Heatmaps were visualised using the *sc.pl.heatmap* function, and any accompanying hierarchical

clustering, or, dendrograms were produced within this function using *sc.tl.dendrogram* with default parameters. All expression values were DSB-normalised for protein and log-transformed, normalised, and scaled for RNA. When upper limits were placed on either expression or population in GEX plots, this is indicated in the given figure legend.

*Gene enrichment scores*

To conduct gene enrichment scores against a reference published blood dataset[6], the top 100 DEGs (log2 fold change) of blood DC and monocyte cell types were input into *sc.tl.score_genes* function in Scanpy. Gene enrichment value for the blood reference cell type was then calculated as the average expression of the top DEG from the reference dataset, minus the average expression of another reference set of genes (randomly sampled from each binned expression value). Gene enrichment scores were visualised using a heatmap in the seaborn (v.0.9.0) package.

Cell cycle gene enrichment scores were calculated through use of a publicly available curated list of genes implicated in the human cell cycle[7] as the input for *sc.tl.score_genes_cell* function in Scanpy. The G2/M and S phase score for each cell thus represented low to high enrichment for a particular phase's genes. In order to serve as a proxy for a 'proliferative phase' score, the mean of the G2/M and S phase scores were calculated and plotted in UMAP space. Cells with G2/M/S cell cycle score greater than the mean were assigned as 'cycling' cells, else assigned as 'not cycling'.

Gene enrichment scores calculated using *sc.tl.score_genes* function in Scanpy were used to ascertain: i) apoptotic gene enrichment through use of genes implicated in the KEGG apoptotic pathway (GSEA:M8492), ii) NK cytotoxic gene enrichment through use of genes implicated in the KEGG NK cytotoxicity pathway (GSEA:M5669), and iii) TNF response gene enrichments using genes from the GO biological process database (GO:0034612).

*Direction of Transition (DoT)-score analysis*

To define the suitable origin for the DoT-score method, we projected CD34+ CITE-seq HSC/MPP I cells from each tissue (FBM, FL or CB) onto the ABM scRNA-seq landscape. Both fetal and adult data were scaled together and fetal cells from each tissue were projected onto 50 PC vectors computed using the ABM data. For each fetal cell, 15 nearest neighbours in the ABM data were identified. The normalised sum of neighbours identified per cell in the ABM served as a similarity score. We used all cells with similarity scores >0.05 (around 60 cells for each tissue) to compute the average gene expression, subsequently used as the origin points for the DoT-score analysis. DoT-score was computed as previously described[8], using the dotscore package publicly available at https://github.com/Iwo-K/dotscore. Weights were derived from genes with significant differential expression between: HSC/MPP I cells from FL compared to FBM (FDR < 0.05, abs($\log_2$(Fold Change) > 0.5); FBM compared to combined CB and FL cells (FDR < 0.05, abs($\log_2$(Fold Change) > 1); CB compared to combined FBM and FL cells (FDR < 0.05, abs($\log_2$(Fold Change) > 1 (**Supplementary Table 40**). DoT-score values correspond to z-scores estimated against simulated data. For clarity, colour-scales in DoT-score plots are clipped at 0.1% and 99.9% percentiles to avoid plotting extreme outliers.

*TCR- and BCR-enriched VDJ repertoire analysis using pyVDJ*

Using the pyVDJ Python package (v0.1.2), lanes of BCR-enriched and TCR-enriched 10x data were integrated with their corresponding 10x GEX lane data in Scanpy. Filtered CellRanger output files were then imported into the Scanpy workflow to investigate

productivity of chains, presence of heavy and light chains and clonal assignment. VDJ metadata by cell type was then exported from Scanpy and plotted in GraphPad Prism.

*Trajectory inference using Monocle3*

GEX matrices for cell-types of interest (filtered by highly variable genes, as defined in Scanpy) were loaded into the Monocle workflow as CellDataSet objects using Monocle3 (v0.2.1). GEX values were then normalised by log and size factor to address depth differences using the *preprocess_cds* function. For known lineages, cells were clustered with a resolution parameter of 1e-07 in order that one partition was returned (to ensure pseudotime with incorporated all cells). Cells were then ordered along pseudotime and with root state provided using the *order_cells* function. DEGs across pseudotime were calculated using a one-sided Moran's I statistical test (*graph_test* function) and DEGs grouped into 'modules' by their Moran's-derived correlation across pseudotime using the *find_gene_modules* function. Dynamically expressed genes across a given pseudotime were then plotted as heatmaps, with normalised logged and scaled gene expression values. Paired heat maps across conditions were the product of combined processing (log-transforming, normalising, scaling) of GEX counts and plotting gene expression over independently derived pseudotime trajectories.

*Decision tree construction using Rpart*

Decision tree construction was implemented using the Rpart package (v.4.1-15) to distinguish between the cell types (classes) in FBM (total) CITE-seq dataset using the 198 antibodies present in the ADT panel. The input for continuous decision tree construction was DSB-normalised protein matrices (closely following the computational approach used by Haas et al, 2021[9]). The aim of this decision tree was to identify discriminative markers for lineage-committed FBM immune cell populations, so the four rare progenitor cell types identified in this dataset were merged into CD38+/- progenitor populations: [CD38- pro.=HSC/MPP], [CD38+ pro.=CMP, GMP, ELP], leading to a reduction to 30 classes. To prevent class

imbalance, the data was subsampled to n=smallest_class_size per cell type using *scanpy.pp.subsample* function in scanpy with random_state set to 1 and 2 for production of training and test datasets, respectively. A decision tree was then built with 10-fold cross-validation using the training data as input into the *Rpart.rpart* function with complexity parameter - cp (used to define the cost-complexity measure of a given tree; based on misclassification rate of the terminal nodes)=-1 and min_split (minimum observations in node to attempt a split)=2 to ensure tree growth was not prematurely terminated in favour of reduced complexity. The resultant decision tree was evaluated for cross-validation error - xerror, cross-validation standard deviation - xstd and relative error - rerror across the full range of branching/complexity levels, and the tree was pruned by sequential filtering for complexity level with: i) identification of the smallest cp value that has an xerror smaller than the sum of the smallest xerror found across cp values and its xstd, ii) cp in the top 3 levels for smallest(complexity + cross-validation error) iii) cp with the lowest number of branch splits for which all classes are present as terminal leaves in the decision tree. To finally evaluate the sensitivity and specificity of the decision tree, class predictions were run for the test dataset and accuracy of the decision tree model was visualised and plotted using the *Rpart.confusionMatrix* function. Overall confusion matrix accuracy was computed with a 95% CI using a binomial test and checked with a one-sided test (see caret package documentation for *confusionMatrix* function).

*Prediction of cell-cell communication using CellPhoneDB*

To assign putative cell-cell interactions within our FBM scRNA-seq dataset, we used CellPhoneDB (v2.1.2). Log-transformed, normalised and scaled gene expression values for stromal and progenitor cell types of interest were exported from Scanpy along with their respective cell type metadata. Using the receptor-ligand database (v2.0.0), CellPhoneDB was

run using the statistical method, with *p*-value cut-off of 0.05 for significant receptor ligand pairs and a result precision of 3dp. To visualise spatiality of non-DS FBM niche interactions via Venn diagram, stromal cells were grouped into the following neighborhoods: 'endothelial' (EC-sinusoidal, EC-proliferating, EC-tip, Fb-arteriolar, Adipo-CAR), 'endosteal' (Fb-endosteal, Osteochondral precursor, Early osteoblast) and 'stromal' (Mac-stromal, Fb-fibroblast). Total cell numbers of interest were included in analysis unless stated otherwise in tool-specific Methods.

*Inference of transcription factors and their gene regulatory networks using PySCENIC*

The PySCENIC package (v0.9.19) and pipeline was used to identify transcription factors and their target genes in the combined non-DS and DS FBM scRNA-seq datasets. The ranking database (hg38 refseq-r80 500bp_up_and_100bp_down_tss.mc9nr.feather), motif annotation database (motifs-v9-nr.hgnc-m0.001-o0.0.tbl) and list of transcription factors (lambert2018.txt) were downloaded from the Aert's laboratory github page. An adjacency matrix of transcription factors and their targets was generated and pruned using the Aert's group suggested parameters. PySCENIC was used to calculate median TF activity in each cluster (from AUCell output), and DS vs non-DS were compared by t-test for each TF to calculate *p*-values. The regulons generated were used to predict which genes controlled by each transcription factor in downstream analysis.

*TNF response gene annotation, TNF superfamily interactions and TNFα-signalling pathway enrichment*

To produce the heatmap of differentially enriched inflammatory and cytokine production pathways in DS vs. non-DS FBM (**Supplementary Fig. 8i**), we derived DEGs between analogous cellstate compartments in non-DS and DS FBM scRNA-seq datasets. Differentially enriched inflammatory and cytokine production pathways in DS vs. non-DS FBM scRNA-seq stroma were defined by a two-sided Wilcoxon rank-sum test with Benjamini-Hochberg procedure for multiple testing correction. Genes were submitted to the fgsea package with

returned enriched pathway gene sets ranked by log fold change. Statistically significant DEGs (*p*-value < 0.05) were compared to intersect against TNF response associated genes acquired from the GO biological process database (GO:0034612). Intersecting TNF response genes in DEGs were ranked by log fold change between cell states in either condition. The full list of intersecting TNF response genes are shown in (**Supplementary Table 43**).

To produce the Sankey-plot of putative TNF superfamily interactions in DS FBM (**Extended Data Fig. 7j**), DEGs for equivalent cell states in DS vs non-DS (**Supplementary Table 21**) were filtered using CellPhoneDB (with cell type groupings as described in **Supplementary Table 7, 20**). To produce TNFα-signalling pathway enrichment dotplot in **Extended Data Fig. 7i**, the fgsea (v1.18.0) package and MSigDB database were used with non-DS/DS DEGs as input (**Supplementary Table 7, 20, 21**).

*Disease lists for interactive web portal*

The interactive web portal offers the capability to search expression profiles by genes associated with haematological disease. These gene lists were intended to compass inherited disorders with haematological phenotype and are derived from Genomics England clinical testing panels (https://panelapp.genomicsengland.co.uk/) (**Supplementary Table 50**).

**Statistics and reproducibility**

For all analysis of single cell datasets in this study, the complete set of biological replicates and total population of annotated cell types were used (unless otherwise noted below) in order to increase statistical power.

*Fetal liver, fetal yolk sac and thymus droplet-based scRNA-seq data*
For all re-analysis of FL (n=14, k=113,063, 7-17 PCW) and YS (n=3, 10,071, 4-7 PCW) scRNA-seq data (E-MTAB-7407), the complete set of biological replicates and total population of annotated cell types were shown in figures unless otherwise noted below. Original cell numbers can be found in the original publication[10].

For **Fig. 2a**, proportions of myeloid cell states arising from FBM haematopoiesis and their counterparts in other tissues were compared. The YS cell state originally assigned as DC progenitor in Popescu et al[4] was renamed as macrophage in line with further GEX exploration and re-annotation, and therefore not included in **Fig. 2a** or **Extended Data Fig. 3f**. Further YS progenitor (HSC/MPP, ELP, CMP, MEMP, MEP) and myeloid (GMP, promonocyte, MOP, macrophage) cell states were identified upon re-clustering and re-annotating the YS Lymphoid progenitor, MEMP, Myeloid progenitor and YS progenitor/MPP (metadata for reannotation available in **Supplementary Table 51**).

For **Fig. 2a**, the original FL Monocyte precursor and Neutrophil-myeloid progenitor were subclustered to resolve heterogeneity and revealed further myeloid states including MOP, monocyte, promonocyte and promyelocyte cells (metadata for reannotation available in **Supplementary Table 52**). For **Fig. 2a** and **Extended Data Fig. 4b**, the original FL HSC_MPP, MEMP and Pre pro B cell were sub clustered to further identify progenitor, myeloid and lymphoid cells, including (for future use in these plots): HSC/MPP, MEMP, GMP, ELP, MPP,

MEP, eo/baso/mast precursor, myeloid DC progenitor and pDC progenitor (metadata for reannotation available in **Supplementary Table 52)**.

For **Fig. 2a** and **Extended Data Fig. 3f,** FL myeloid nomenclature was updated such that monocyte precursor became promonocyte, and monocyte became CD14$^+$ monocyte. Due to the presence of a distinctive pDC precursor in FL, pDC and pDC precursor were merged into one 'pDC' grouping for purposes of the cross-tissue bar plot. For **Extended Data Fig. 8g,** FL endothelial populations were re-clustered to annotate sinusoidal endothelium, with metadata for reannotation available in **Supplementary Table 52**.

For all analysis of thymus scRNA-seq data (k=259,265), n=24 biologically independent samples were used and the total population of annotated cell types were shown in figures. Cell numbers can be found in the original publication[10].

*Fetal bone marrow droplet-based scRNA-seq data*

For analysis of combined FBM scRNA-seq 5' and 3' data (k=103,228), n=9 biologically independent samples were used, which spanned 12-19 PCW. Total population of annotated cell types and biological replicates were shown in figures - otherwise, downsampling methods (for example, for age-matching for proportional comparisons with Down syndrome FBM scRNA-seq data) are noted in relevant figure legends and/or methods.

The following refined cell number annotations were displayed in each of the figures (annotations available in **Supplementary Table 8**): CD4 T cell - 327, CD8 T cell - 171, CD14 monocyte - 8763, CD56 bright NK - 449, CMP- 425, DC1 - 50, DC2 - 598, DC3 - 705, DC precursor - 201, erythroid macrophage - 92, ELP - 1357, GMP - 1281, HSC/MPP - 92, ILC precursor - 67, LMPP - 34, MEMP - 16, MEP - 269, MK - 1000, MOP - 3838, MPP myeloid - 92, NK T cell - 111, NK progenitor - 26, Treg - 62, adipo-CAR - 353, arteriolar fibroblast - 83, basophil - 139,

chondrocyte - 80, early MK - 1624, early erythroid - 7474, early osteoblast - 280, endosteal fibroblast - 54, eo/baso/mast precursor - 175, eosinophil - 321, erythroid macrophage - 92, immature B cell - 1988, immature EC - 42, late erythroid - 4636, mast cell - 648, mature NK - 136, mid erythroid - 14297, monocytoid macrophage - 290, muscle - 131, muscle stem cell - 254, myelocyte - 3794, myeloid DC progenitor - 31, myofibroblast - 78, naive B cell - 1411, neutrophil - 4501, osteoblast - 363, osteoblast precursor - 456, osteochondral precursor - 191, osteoclast - 1221, pDC - 712, pDC progenitor - 23, pre B progenitor - 14229, pre pro B progenitor - 5427, proliferating EC - 26, promonocyte - 7437, promyelocyte - 2191, schwann cells - 9, sinusoidal EC - 550, stromal macrophage - 1464, tDC - 193, tip EC - 362, pro B progenitor - 5528.

When cell types were grouped into broad lineages (e.g, **Fig. 1a**), cell numbers were as follows: HSC/MPP and pro. - 3795, erythroid - 26407, MK - 2624, B_lineage - 28583, DC - 2459, eo/baso/mast - 1108, neutrophil - 10486, monocyte - 20038, T_NK - 1349, stroma - 6379. Other broad groupings used are detailed in **Supplementary Table 7** or otherwise noted in relevant figure legends and/or methods.

*Fetal bone marrow plate-based scRNA-seq data*
For analysis of FBM Smart-seq2 scRNA-seq validation data (k=486), n=2 biologically independent samples both 17 PCW were used and the following cell numbers were shown: HSC/MPP = 32, B cell = 52, DC1 = 34, DC2 = 15, monocyte = 32, PMN = 65, basophil = 20, eosinophil = 54, mast cell = 47, myelocyte = 61, pDC = 30, promyelocyte = 44. Metadata (including annotations) are available in **Supplementary Tables 1 and 15**.

*Down syndrome FBM scRNA-seq droplet-based scRNA-seq data*
For all analysis of fetal DS FBM scRNA-seq 5' data (k=16,743), n=4 biologically independent

samples were used, spanning 12-13 PCW. We performed two independent experiments on the same 4 DS FBM biological replicates in order to capture sufficient cells (see Supplementary Table 1; the latter independent experiment relates to 10x lanes with IDs: 'DSOXPool_GEX' and 'DSOX19_GEX'). All subsequent statistical analyses were run with lane ID as covariate (rather than biological replicate). All biological replicates and total population of annotated cell types were shown in figures unless otherwise noted in relevant figure legends and/or methods.

The following refined cell number annotations were displayed in each of the figures (annotations are available in **Supplementary Table 53**): CAR - 4, CD14 monocyte - 320, CD56 bright NK - 79, CD8 T cell - 181, CMP - 50, DC1 - 45, DC2 - 228, DC3 - 108, HSC/MPP - 105, ILC precursor -13, MEMP - 130, MK - 83, MOP - 422, MSC -53, Treg - 8, chondrocyte - 4, early B cell - 42, early MK - 34, early erythroid - 1,348, endothelium - 111, eo/baso/mast precursor - 53, eosinophil - 63, late erythroid - 6,336, macrophage - 113, mast cell - 66, mature B cell - 31, mature NK - 147, mid erythroid - 5,230, myelocyte - 243, neutrophil - 273, osteoblast -11, osteoclast - 57, pDC - 14, pre B cell - 115, promonocyte - 395, pre pDC - 110, promyelocyte - 107, transitional NK cell -11.

When cell types were grouped into broad lineages (e.g, **Extended Data Fig. 7a**), cell numbers were as follows: HSC/MPP and pro. = 338, Erythroid = 12,914, MK = 117, B lineage = 188, DC = 505, Neutrophil = 623, Eo/baso/mast = 129, Monocyte = 1,137, TNK = 439, Stroma = 353. Other broad groupings are detailed in **Supplementary Table 20**.

*Adult bone marrow scRNA-seq data*

For all analyses of adult BM 10x (k=142,026) data, n=4 biologically independent samples and total population of annotated cell types were used, unless otherwise noted in relevant figure legends and/or methods. Lanes from four donors (BM1, BM2, BM5, BM6, ranging from 26-52 years old) were downloaded from the Human Cell Atlas Data Coordination Portal 'Census of

Immune Cells' project;

https://data.humancellatlas.org/explore/projects/cc95ff89-2e68-4a08-a234-480eca21ce79/).

The following refined cell number annotations were displayed in each of the figures (annotations are available in **Supplementary Table 5)**: CD14 monocyte - 3670, CD16 monocyte - 1938, CD56 bright NK - 1228, CLP - 882, CMP - 288, DC1 - 135, DC2 - 481, DC3 - 550, DC precursor - 462, HSC/MPP - 497, LMPP - 80, MEMP - 785, MK - 577, MOP - 1440, MPP - 365, Treg - 6327, early MK - 136, early erythroid - 5441, erythroid macrophage - 77, immature B cell - 2728, late erythroid -1150, mature CD8 T cell - 15725, mature NK - 6074, memory B cell - 4106, memory CD4 T cell - 22197, mid erythroid - 2192, monocyte-DC - 515, myelocyte - 6675, myeloid DC progenitor - 110, naive B cell - 19265, naive CD4 T cell - 5873, naive CD8 T cell - 8965, neutrophil - 2482, pDC - 1134, pDC progenitor - 63, plasma cell - 2074, pre B cell - 971, pro B progenitor - 1390, promonocyte - 7448, promyelocyte - 2197, stroma - 161, tDC - 75, transitional B cell - 2151, transitional NK - 946.

When cell types were grouped into broad lineages (e.g, **Extended Data Fig. 6a**), cell numbers were as follows: HSC/MPP and pro. = 3,007, Erythroid = 8,783, MK = 713, B lineage = 32,685, DC = 3,415, Neutrophil = 11,354, Monocyte = 14,496, TNK = 67,335, Stroma = 238. Other broad annotations are available in **Supplementary Table 2.**

*Cord blood scRNA-seq data*

For all analysis of CB 10x (k=148,442) data, n=4 biologically independent samples and total population of annotated cell types were used, unless otherwise noted in relevant figure legends and/or methods. Lanes from four donors (CB1, CB2, CB5, CB6 at 40-42 PCW) were downloaded from the Human Cell Atlas Data Coordination Portal 'Census of Immune Cells' project; https://data.humancellatlas.org/explore/projects/cc95ff89-2e68-4a08-a234-480eca21ce79/).

The following refined cell number annotations were displayed in each of the figures (annotations are available in **Supplementary Table 4**): CD8 T cell - 16345, CD14 monocyte - 13324, CD16 monocyte - 888, CD56 bright NK - 4066, CMP - 272, DC1 - 67, DC2 -155, DC precursor - 169, GMP - 203, HSC/MPP - 194, ILC precursor -1519, MEMP - 338, MK - 1262, early MK - 496, early erythroid - 532, late erythroid - 878, mature NK - 7860, mid erythroid - 2627, myelocyte - 3726, naive B cell - 19516, naive CD4 T cell - 69338, neutrophil - 3458, pDC - 242, preDC - 269, promonocyte - 607, tDC - 91.

When cell types were grouped into broad lineages, cell numbers were as follows: HSC/MPP and pro. - 1007, erythroid - 4037, MK - 1758, B cells - 19516, DC - 993, neutrophil - 7184, monocyte - 14819, T/NK - 99128. Broad annotations are available in **Supplementary Table 3.**

*CD34+ (FBM, fetal liver, cord blood) CITE-seq data*
For analysis of non-lineage committed progenitors in the CD34+ CITE-seq data, lanes from FBM (n=3, k=8,829, 14-17 PCW), fetal liver (n=4, k=18,904, 14-17 PCW) and cord blood (n=4, k=7,540, 40-42 PCW) were run using both 3GEX and ADT technology. The total population of annotated cell types and biological replicates were shown in figures unless otherwise noted in relevant figure legends and/or methods. The following refined cell number annotations were derived from RNA-based annotations of the data and displayed in each of the RNA-based figures (with protein-based analysis containing fewer cells by virtue of further filtering post-protein-QC and post intersect). RNA-based annotations available in **Supplementary Table 16** (see table description and **Supplementary Table 1** for further information): DC progenitor I_CB - 54, DC progenitor I_FBM - 247, DC progenitor I_FL - 50, DC progenitor II_CB - 76, DC progenitor II_FBM - 298, DC progenitor II_FL - 344, Early LyP_CB - 183, Early LyP_FBM - 301, Early LyP_FL - 440, EoBasoMC_CB - 153, EoBasoMC_FBM - 153, EoBasoMC_FL -

568, EryP I_CB - 224, EryP I_FBM - 172, EryP I_FL - 1,927, EryP II_CB - 10, EryP II_FBM - 70, EryP II_FL - 849, EryP III_CB - 86, EryP III_FBM - 133, EryP III_FL - 1,020, EryP IV_CB - 223, EryP IV_FBM - 319, EryP IV_FL - 1,864, HSC/MPP I_CB - 1,455, HSC/MPP I_FBM - 284, HSC/MPP I_FL - 1,699, HSC/MPP II_CB - 1,086, HSC/MPP II_FBM - 378, HSC/MPP II_FL - 1,298, HSC/MPP III_CB - 83, HSC/MPP III_FBM - 159, HSC/MPP III_FL - 544, HSC/MPP IV_CB - 1,020, HSC/MPP IV_FBM - 307, HSC/MPP IV_FL - 269, Late EryP I (Pro-erythroblast)_CB - 174, Late EryP I (Pro-erythroblast)_FBM - 229, Late EryP I (Pro-erythroblast)_FL - 1,598, Late EryP II (Erythroblast)_CB - 26, Late EryP II (Erythroblast)_FBM - 31, Late EryP II (Erythroblast)_FL - 868, LyP I (CLP)_CB - 787, LyP I (CLP)_FBM - 520, LyP I (CLP)_FL -

891, LyP II (pre pro-B)_CB - 460, LyP II (pre pro-B)_FBM - 1,499, LyP II (pre pro-B)_FL - 554, LyP III (pro-B)_CB - 3, LyP III (pro-B)_FBM - 450, LyP III (pro-B)_FL - 36, LyP IV (pre-B)_CB - 479, LyP IV (pre-B)_FBM - 777, LyP IV (pre-B)_FL - 619, MEP/MkP_CB - 348, MEP/MkP_FBM - 207, MEP/MkP_FL - 1,184, MEP_CB - 39, MEP_FBM  - 324, MEP_FL - 384, MyP_CB - 427, MyP_FBM - 951, MyP_FL - 845, Cycling LyP_CB - 144, Cycling LyP_FBM - 1,020, Cycling LyP_FL - 1,053. Additional sinusoidal EC were also captured in this dataset by virtue of being CD34+ (k=281 FBM), (k=1,296 FL). These cells are shown in **Extended Data Fig. 8g** protein analysis (total RNA-based annotations available in **Supplementary Table 47**, see table description for further information on post-QC post- protein intersect filtering).


*Fetal bone marrow (total) CITE-seq data*

For analysis of FBM CITE-seq data (n=3, k=8,978, 16-17 PCW), the total population of annotated cell types and biological replicates were shown in figures unless otherwise noted in relevant figure legends and/or methods. The following refined cell number annotations were displayed in each of the figures (annotations available in **Supplementary Table 11**): basophil -

15, CD14 monocyte - 1,384, CD4 T cell - 39, CD56 bright NK - 66, CMP - 78, DC1 - 13, DC2 - 87, DC3 - 20, early erythroid - 517, early MK - 91, ELP - 177, eosinophil - 22, GMP - 108, HSC/MPP - 36, immature B cell - 403, late erythroid - 670, mast cell - 57, mid erythroid - 466, MK - 31, MOP - 280, naive B cell - 249, neutrophil - 294, osteoclast - 58, pDC - 139, pre B progenitor - 2,241, pre pro B progenitor - 248, pro B progenitor - 366, promonocyte - 620, promyelocyte - 103, sinusoidal EC - 42, stromal macrophage - 47, tip EC - 11.

*Blood monocyte and DC 10x data*

Monocyte-DC blood SS2 scRNA-seq data (GSE94820) were downloaded from a published study[6]. The available RPKM counts for 1140 monocytes (k=768) and DCs (k=372) were logged and scaled (in line with 10x analysis), in preparation for DEG analysis conducted as described below. Refined celltype population frequency can be found in the original study.

*Murine fetal bone marrow 10x data*

Mouse fetal bone marrow scRNA-seq data (GSE122467) were downloaded from a published study[11]. The available scRNA-seq count matrix was subsetted to stromal cell types and log-transformed, normalised and scaled (in line with human fetal 10x analysis), in preparation for DEG analysis conducted as described below. Refined celltype population frequency can be found in the original study. Genes differentially expressed between mouse BM ECs were taken from Baryawno et al[12] and used as a reference for FBM ECs in **Extended Data Fig. 8h.**

**Supplementary references**

1.  Mulè, M. P., Martins, A. J. & Tsang, J. S. Normalizing and denoising protein expression data from droplet-based single cell profiling. doi:10.1101/2020.02.24.963603.

2.  Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with Harmony. *Nature Methods* vol. 16 1289–1296 (2019).

3.  Pandey, V. D., Singh, V., Nigam, G. L., Usmani, Y. & Yadav, Y. Fetal foot length for assessment of gestational age: A comprehensive study in north India. *Journal of the Anatomical Society of India* vol. 65 S19 (2016).

4.  Tunstall-Pedoe, O. *et al.* Abnormalities in the myeloid progenitor compartment in Down syndrome fetal liver precede acquisition of GATA1 mutations. *Blood* **112**, 4507–4511 (2008).

5.  Matsushima, H. *et al.* Neutrophil differentiation into a unique hybrid population exhibiting dual phenotype and functionality of neutrophils and dendritic cells. *Blood* **121**, 1677–1689 (2013).

6.  Villani, A.-C. *et al.* Single-cell RNA-seq reveals new types of human blood dendritic cells, monocytes, and progenitors. *Science* **356**, (2017).

7.  Macosko, E. Z. *et al.* Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell* **161**, 1202–1214 (2015).

8.  Kucinski, I. *et al.* Interactions between lineage-associated transcription factors govern haematopoietic progenitor states. *EMBO J.* e104983 (2020).

9.  Triana, S. H. *et al.* Single-cell proteo-genomic reference maps of the hematopoietic system enable the purification and massive profiling of precisely defined cell states. doi:10.1101/2021.03.18.435922.

10. Popescu, D.-M. *et al.* Decoding human fetal liver haematopoiesis. *Nature* **574**, 365–371 (2019).

11. Baccin, C. *et al.* Combined single-cell and spatial transcriptomics reveal the molecular,

cellular and spatial bone marrow niche organization. *Nat. Cell Biol.* **22**, 38–48 (2020).

12. Baryawno, N. *et al.* A Cellular Taxonomy of the Bone Marrow Stroma in Homeostasis and Leukemia. *Cell* **177**, 1915–1932.e16 (2019).