# The intersectional genetics landscape for human

Andre Macedo* and Alisson M. Gontijo*

Chronic Diseases Research Center (CEDOC), NOVA Medical School | Faculdade de Ciências Médicas, Universidade Nova de Lisboa, Lisbon, Portugal.
*   Correspondence: andre.macedo@nms.unl.pt (A.M.) and alisson.gontijo@nms.unl.pt (A.M.G.)

**Supplementary Materials**

This file contains Supplementary Figures S1-4, Supplementary Tables S1-S4, and Supplementary Data S1. The latter is available at https://github.com/AndreMacedo88/VEnCode.
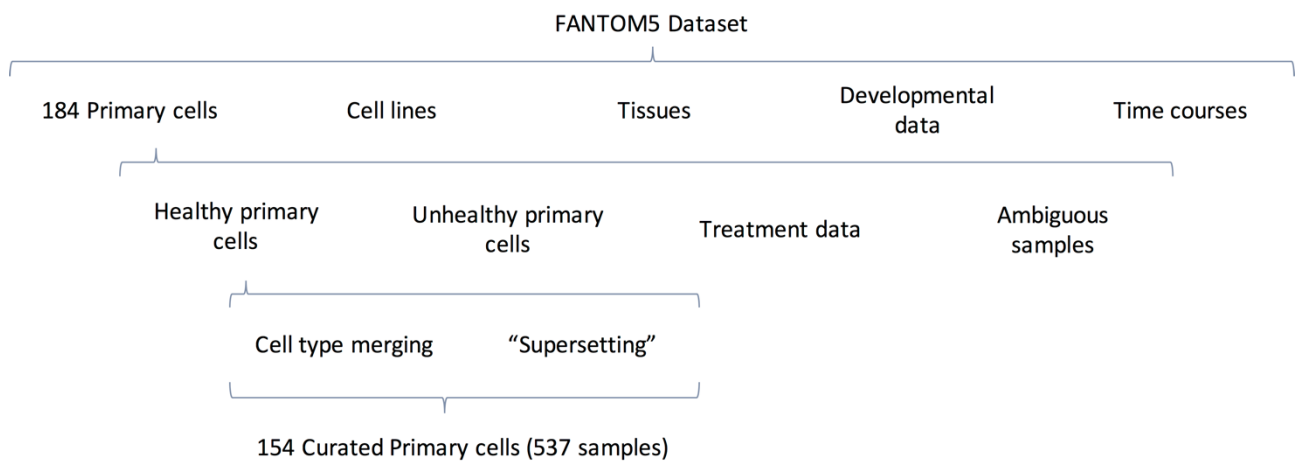
## Supplementary Figures



**Figure S1. Pipeline for FANTOM5 data preparation and curation.** Further details on cell type merging, "supersetting", and excluded primary cells are provide Supplementary Table S1.

**1. Calculate $E_{raw}$ from various samples of inactive promoters**

Vary number cell types

RE activity

$RE1$ 0 0 0 0 0 0 0 0 0 0 0 ...
$RE2$ 0 0 0 0 0 0 0 0 0 0 0 ...
... 0 0 0 0 0 0 0 0 0 0 0 ...
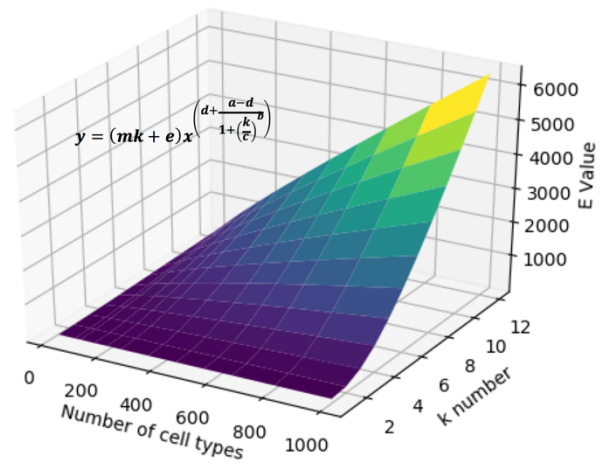0 0 0 0 0 0 0 0 0 0 0 ...
0 0 0 0 0 0 0 0 0 0 0 ...

Vary $k$

**2. Calculate correlation best-fitting curve**

Vary number cell types: $y = sk^h$

Coefficient variation with k:

$s = mk + e$  $R^2 = 0.9993$

$h = d + \dfrac{a-d}{1 + \left(\frac{k}{c}\right)^b}$  $R^2 = 1$

| | | Coefficients | | | |
|---|---|---|---|---|---|
| a | b | c | d | m | e |
| -164054,1 | 0,9998811 | 6,08895E-06 | 1,00051 | 0,9527 | -0,1131 |

$y = (mk + e)x^{\left(d + \frac{a-d}{1 + \left(\frac{k}{c}\right)^b}\right)}$

**3. Generate function that generates expected best $E$**

$y = (mk + e)x^{\left(d + \frac{a-d}{1 + \left(\frac{k}{c}\right)^b}\right)}$

**Figure S2. Generating the function that returns the best possible $E$ ($E_{best}$).** We generated $E_{raw}$ (as described in Figure 6) for simulated data reflecting a best-case-scenario for a VEnCode – where all $k$ REs are inactive in the non-target cell – varying the number of cell types and $k$ in the dataset (1.). By changing cell type number ranging from 20 to 1000, we calculated the best fitting curves that predicted $E_{raw}$ and generated the general equation $y = sk^h$ (2.). However, $s$ and $h$ depend on the number of REs ($k$). So, varying $k$ from 1 up to 10, we obtained the best fitting equations that explain this variation (2.). With this data we then generated a function that best accounts for the variation of both $k$ and cell type number in the dataset (3.). The effectiveness of this equation can be seen in Table S2.
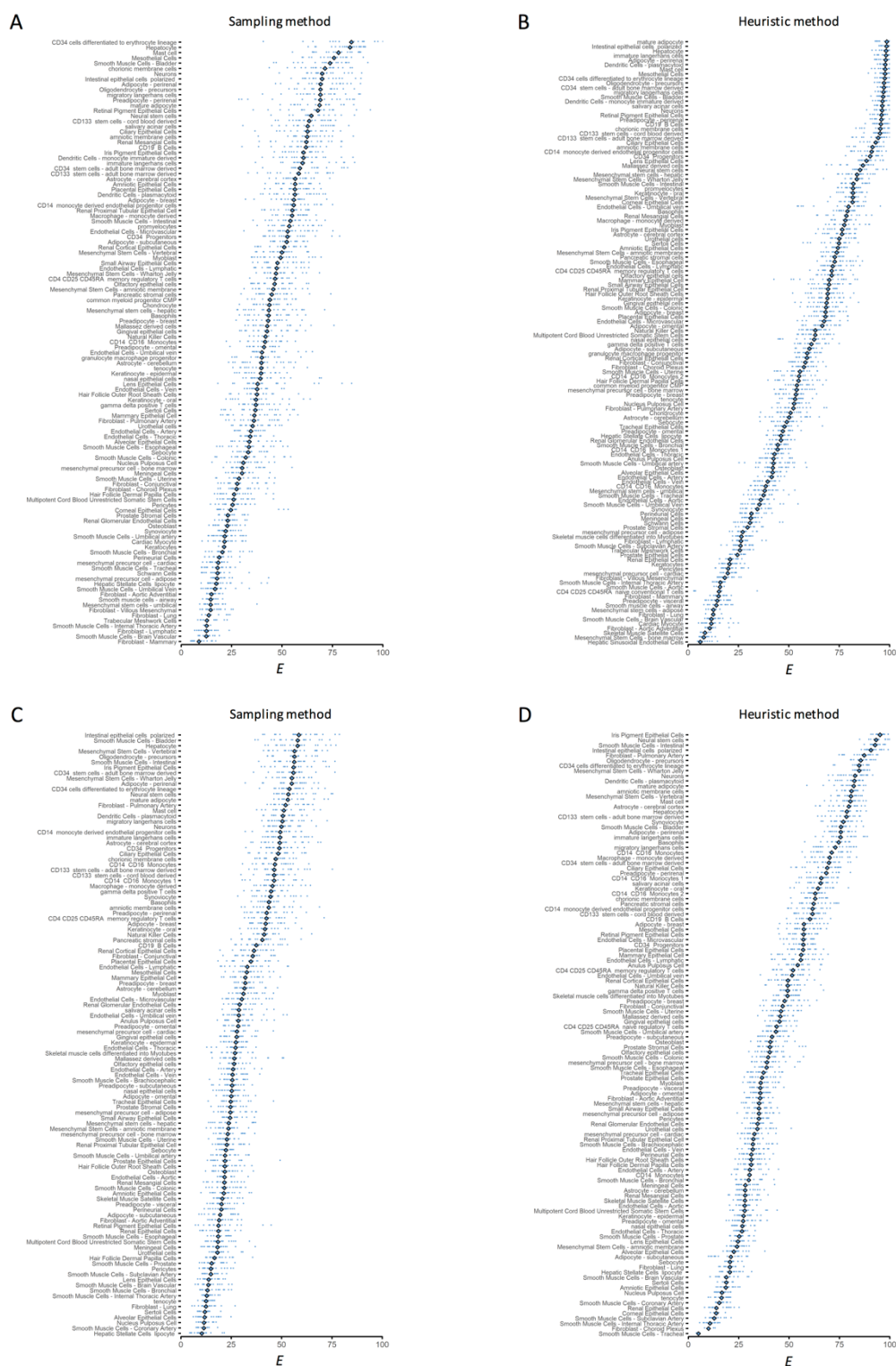
**Figure S3. *E* value variation by cell type.** 5 to 20 VEnCodes at *k* = 4 were obtained for every possible primary cell type and their *E* values were determined as described in Figure 6. **A, B.** Results for VEnCodes generated using the promoter dataset. In (**A**) the sampling method (see Figure 3) was used to obtain VEnCodes and determine *E* for 114 cell types. In (**B**), the heuristic method (see Figure 4) was used for the same purpose, allowing us to analyze *E* for up to 20 VEnCodes for 131 cell types. **C, D.** Results for VEnCodes generated using the enhancer dataset. **C.** *E* value distribution of VEnCodes obtained with the sampling method for 112 primary cell types. **D.** Similarly, 117 cell types had 5-20 enhancer VEnCodes retrieved using the heuristic method and their *E* was determined.
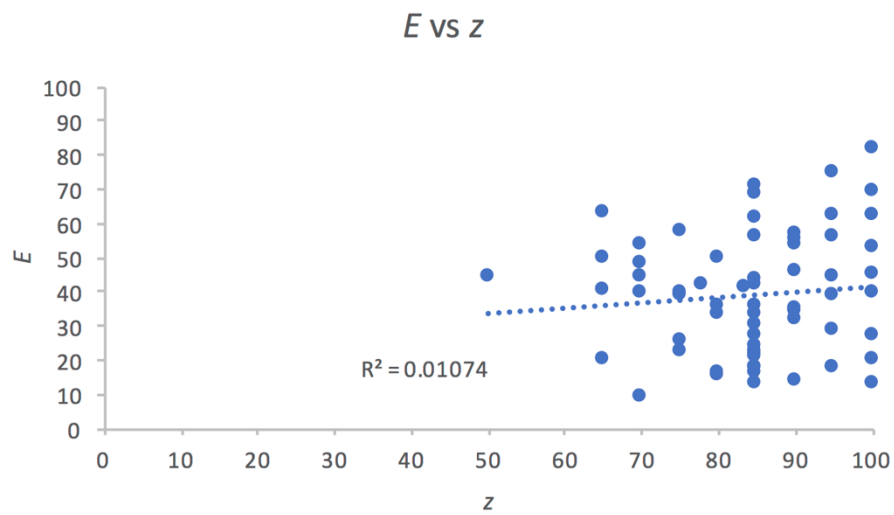
**Figure S4. *E* versus *z* scores.** Plotted is the *E* and *z* scores for each cell type of a list of 64 cell types with three donors and which we managed to retrieve both *E* and *z* values. Also plotted is the linear regression attempt to model the relationship between these two values.

# Supplementary Tables

**Table S1. List of curated primary cell types and merged, supersets and excluded categories used in this study.** The list of curated cell types contains 154 primary cell types, encompassing a total of 537 cell line samples. Merging was done with the rationale of turning a vast data on diverse cell conditions into biologically relevant cell types. Merged cell types are then used as curated cell types and the original "cell type" data is not accessed independently. On the other hand, a superset cell type means that the superset data includes the subset data, but each subset is still included in the curated cell type list used in this study. The excluded cell type category lists the data not used at any point throughout the study.

| Curated cell types (154 Primary cell types) | |
|---|---|
| Adipocyte - breast | Mast cell |
| Adipocyte - omental | mature adipocyte |
| Adipocyte - perirenal | Melanocyte |
| Adipocyte - subcutaneous | Meningeal Cells |
| Alveolar Epithelial Cells | mesenchymal precursor cell - adipose |
| Amniotic Epithelial Cells | mesenchymal precursor cell - bone marrow |
| amniotic membrane cells | mesenchymal precursor cell - cardiac |
| Anulus Pulposus Cell | Mesenchymal stem cells - adipose |
| Astrocyte - cerebellum | Mesenchymal Stem Cells - amniotic membrane |
| Astrocyte - cerebral cortex | Mesenchymal Stem Cells - bone marrow |
| Basophils | Mesenchymal stem cells - hepatic |
| Bronchial Epithelial Cell | Mesenchymal stem cells - umbilical |
| Cardiac Myocyte | Mesenchymal Stem Cells - Vertebral |
| CD133+ stem cells - adult bone marrow derived | Mesenchymal Stem Cells - Wharton Jelly |
| CD133+ stem cells - cord blood derived | Mesothelial Cells |
| CD14+ monocyte derived endothelial progenitor cells | migratory langerhans cells |
| CD14+ Monocytes | Multipotent Cord Blood Unrestricted Somatic Stem Cells |
| CD14+CD16- Monocytes | Myoblast |
| CD14+CD16+ Monocytes | nasal epithelial cells |
| CD14-CD16+ Monocytes | Natural Killer Cells |
| CD19+ B Cells | Neural stem cells |
| CD34 cells differentiated to erythrocyte lineage | Neurons |
| CD34+ Progenitors | Neutrophil |
| CD34+ stem cells - adult bone marrow derived | Nucleus Pulposus Cell |
| CD4+ T Cells | Olfactory epithelial cells |
| CD4+CD25+CD45RA- memory regulatory T cells | Oligodendrocyte - precursors |
| CD4+CD25+CD45RA+ naive regulatory T cells | Osteoblast |
| CD4+CD25-CD45RA- memory conventional T cells | Pancreatic stromal cells |
| CD4+CD25-CD45RA+ naive conventional T cells | Pericytes |
| CD8+ T Cells | Perineurial Cells |
| Chondrocyte | Placental Epithelial Cells |
| chorionic membrane cells | Preadipocyte - breast |
| Ciliary Epithelial Cells | Preadipocyte - omental |
| common myeloid progenitor CMP | Preadipocyte - perirenal |
| Corneal Epithelial Cells | Preadipocyte - subcutaneous |
| Dendritic Cells - monocyte immature derived | Preadipocyte - visceral |
| Dendritic Cells - plasmacytoid | promyelocytes |
| Endothelial Cells - Aortic | Prostate Epithelial Cells |
| Endothelial Cells - Artery | Prostate Stromal Cells |
| Endothelial Cells - Lymphatic | Renal Cortical Epithelial Cells |
| Endothelial Cells - Microvascular | Renal Epithelial Cells |
| Endothelial Cells - Thoracic | Renal Glomerular Endothelial Cells |
| Endothelial Cells - Umbilical vein | Renal Mesangial Cells |
| Endothelial Cells - Vein | Renal Proximal Tubular Epithelial Cell |
| Eosinophils | Retinal Pigment Epithelial Cells |
| Esophageal Epithelial Cells | salivary acinar cells |
| Fibroblast - Aortic Adventitial | Schwann Cells |
| Fibroblast - Cardiac | Sebocyte |
| Fibroblast - Choroid Plexus | Sertoli Cells |
| Fibroblast - Conjunctival | Skeletal Muscle Cells |
| Fibroblast - Dermal | Skeletal muscle cells differentiated into Myotubes - multinucleated |
| Fibroblast - Gingival | Skeletal Muscle Satellite Cells |
| Fibroblast - Lung | Small Airway Epithelial Cells |
| Fibroblast - Lymphatic | Smooth muscle cells - airway |
| Fibroblast - Mammary | Smooth Muscle Cells - Aortic |
| Fibroblast - Periodontal Ligament | Smooth Muscle Cells - Bladder |
| Fibroblast - Pulmonary Artery | Smooth Muscle Cells - Brachiocephalic |
| Fibroblast - skin | Smooth Muscle Cells - Brain Vascular |
| Fibroblast - Villous Mesenchymal | Smooth Muscle Cells - Bronchial |
| gamma delta positive T cells | Smooth Muscle Cells - Carotid |
| Gingival epithelial cells | Smooth Muscle Cells - Colonic |
| granulocyte macrophage progenitor | Smooth Muscle Cells - Coronary Artery |
| Hair Follicle Dermal Papilla Cells | Smooth Muscle Cells - Esophageal |
| Hair Follicle Outer Root Sheath Cells | Smooth Muscle Cells - Internal Thoracic Artery |
| Hepatic Sinusoidal Endothelial Cells | Smooth Muscle Cells - Intestinal |
| Hepatic Stellate Cells (lipocyte) | Smooth Muscle Cells - Prostate |
| Hepatocyte | Smooth Muscle Cells - Pulmonary Artery |
| immature langerhans cells | Smooth Muscle Cells - Subclavian Artery |
| Intestinal epithelial cells (polarized) | Smooth Muscle Cells - Tracheal |
| Iris Pigment Epithelial Cells | Smooth Muscle Cells - Umbilical artery |
| Keratinocyte - epidermal | Smooth Muscle Cells - Umbilical Vein |
| Keratinocyte - oral | Smooth Muscle Cells - Uterine |
| Keratocytes | Synoviocyte |
| Lens Epithelial Cells | tenocyte |
| Macrophage - monocyte derived | Trabecular Meshwork Cells |
| Mallassez-derived cells | Tracheal Epithelial Cells |
| Mammary Epithelial Cell | Urothelial cells |

**Table S1. Continued.**

| Merged data | |
|---|---|
| **Merged cell types** | **Original cell types** |
| **CD14+ Monocytes** | CD14+ monocytes - mock treated |
| | CD14+ monocytes - treated with BCG |
| | CD14+ monocytes - treated with B-glucan |
| | CD14+ monocytes - treated with Candida |
| | CD14+ monocytes - treated with Cryptococcus |
| | CD14+ monocytes - treated with Group A streptococci |
| | CD14+ monocytes - treated with IFN + N-hexane |
| | CD14+ monocytes - treated with lipopolysaccharide |
| | CD14+ monocytes - treated with Salmonella |
| | CD14+ monocytes - treated with Trehalose dimycolate (TDM) |
| | CD14+ Monocytes |
| **CD19+ B Cells** | CD19+ B Cells (pluriselect) |
| | CD19+ B Cells |
| **CD4+CD25+CD45RA- memory regulatory T cells** | CD4+CD25+CD45RA- memory regulatory T cells expanded |
| | CD4+CD25+CD45RA- memory regulatory T cells |
| **CD4+CD25+CD45RA+ naive regulatory T cells** | CD4+CD25+CD45RA+ naive regulatory T cells expanded |
| | CD4+CD25+CD45RA+ naive regulatory T cells |
| **CD4+CD25-CD45RA- memory conventional T cells** | CD4+CD25-CD45RA- memory conventional T cells expanded |
| | CD4+CD25-CD45RA- memory conventional T cells |
| **CD4+CD25-CD45RA+ naive conventional T cells** | CD4+CD25-CD45RA+ naive conventional T cells expanded |
| | CD4+CD25-CD45RA+ naive conventional T cells |
| **CD8+ T Cells** | CD8+ T Cells (pluriselect) |
| | CD8+ T Cells |
| **Chondrocyte** | Chondrocyte - de diff |
| | Chondrocyte - re diff |
| **Fibroblast - skin** | Fibroblast - skin dystrophia myotonica |
| | Fibroblast - skin normal |
| | Fibroblast - skin spinal muscular atrophy |
| | Fibroblast - skin walker warburg |
| **Mast cell** | Mast cell - stimulated |
| | Mast cell |
| **Melanocyte** | Melanocyte - dark |
| | Melanocyte - light |
| **Neutrophil** | neutrophil PMN |
| | Neutrophils |
| **Prostate Epithelial Cells** | Prostate Epithelial Cells (polarized) |
| | Prostate Epithelial Cells |

**Table S1. Continued.**

| | Curated supersets | | Excluded cell types |
|---|---|---|---|
| **superset** | **subset** | | |
| **CD14+ Monocytes** | CD14+CD16- Monocytes | | mesenchymal precursor cell - ovarian cancer left ovary |
| | CD14+CD16+ Monocytes | | mesenchymal precursor cell - ovarian cancer metastasis |
| | | | mesenchymal precursor cell - ovarian cancer right ovary |
| **CD4+ T Cells** | CD4+CD25+CD45RA- memory regulatory T cells | | Osteoblast - differentiated |
| | CD4+CD25+CD45RA+ naive regulatory T cells | | Peripheral Blood Mononuclear Cells |
| | CD4+CD25-CD45RA- memory conventional T cells | | Whole blood (ribopure) |
| | CD4+CD25-CD45RA+ naive conventional T cells | | |

**Table S2. Normalized $E$ ($E = E_{raw}/E_{best}$) values for a range of number of cell types in data and $k$ REs used to generate a VEnCode.** Values were generated, as described in Figure 6, for a intraindividual robust best-case possible VEnCode and normalized using the function described in Figure S2. Thus, the values reflect the average $E$ expected for the most intraindividual robust VEnCodes.

| | | Number of Cell types | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | **20** | **80** | **100** | **154** | **200** | **250** | **350** | **450** | **550** | **650** | **800** | **1000** |
| | **1** | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 |
| | **2** | 99.5 | 98.8 | 96.5 | 96.1 | 97.9 | 100.0 | 97.2 | 97.3 | 97.4 | 100.0 | 100.0 | 100.0 |
| | **3** | 96.3 | 98.9 | 97.7 | 96.7 | 98.6 | 97.9 | 95.8 | 98.6 | 95.2 | 97.1 | 99.5 | 95.9 |
| | **4** | 99.2 | 98.4 | 97.3 | 98.4 | 97.6 | 98.4 | 99.1 | 98.3 | 98.5 | 98.6 | 98.7 | 97.0 |
| $k$ | **5** | 99.5 | 98.9 | 99.1 | 99.1 | 97.7 | 99.0 | 99.0 | 97.2 | 98.8 | 98.5 | 99.1 | 99.0 |
| | **6** | 98.5 | 99.9 | 99.4 | 99.3 | 98.5 | 99.7 | 99.3 | 99.4 | 99.1 | 99.0 | 99.1 | 99.8 |
| | **7** | 100.0 | 99.9 | 100.0 | 98.9 | 99.5 | 99.9 | 100.0 | 99.3 | 99.7 | 99.2 | 98.0 | 100.0 |
| | **8** | 99.9 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 99.6 |
| | **9** | 100.0 | 100.0 | 100.0 | 99.9 | 100.0 | 100.0 | 100.0 | 99.9 | 100.0 | 100.0 | 100.0 | 100.0 |
| | **10** | 100.0 | 100.0 | 100.0 | 99.7 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 |

**Table S3. List of curated cancer cell types used in this study.** The list contains 158 biologically relevant cell types, encompassing a total of 274 cancer cell line samples.

| Curated cancer cell types | |
|---|---|
| acantholytic squamous carcinoma cell line:HCC1806 | large cell non-keratinizing squamous carcinoma cell line:SKG-II-SF |
| acute lymphoblastic leukemia (B-ALL) cell line | leiomyoblastoma cell line:G-402 |
| acute lymphoblastic leukemia (T-ALL) cell line | leiomyoma cell line |
| acute myeloid leukemia (FAB M0) cell line | leiomyosarcoma cell line:Hs 5 |
| acute myeloid leukemia (FAB M1) cell line | lens epithelial cell line:SRA |
| acute myeloid leukemia (FAB M2) cell line | liposarcoma cell line |
| acute myeloid leukemia (FAB M3) cell line | lung adenocarcinoma cell line |
| acute myeloid leukemia (FAB M4) cell line | lung adenocarcinoma papillary cell line:NCI-H441 |
| acute myeloid leukemia (FAB M4eo) cell line | lymphangiectasia cell line:DS-1 |
| acute myeloid leukemia (FAB M5) cell line | lymphoma malignant hairy B-cell cell line:MLMA |
| acute myeloid leukemia (FAB M6) cell line | malignant trichilemmal cyst cell line:DJM-1 |
| acute myeloid leukemia (FAB M7) cell line | maxillary sinus tumor cell line:HSQ-89 |
| adenocarcinoma cell line:IM95m | medulloblastoma cell line |
| adrenal cortex adenocarcinoma cell line:SW-13 | melanoma cell line |
| adult T-cell leukemia cell line:ATN-1 | meningioma cell line:HKBMM |
| alveolar cell carcinoma cell line:SW 1573 | merkel cell carcinoma cell line |
| anaplastic carcinoma cell line:8305C | mesenchymal stem cell line:Hu5/E18 |
| anaplastic large cell lymphoma cell line:Ki-JK | mesodermal tumor cell line:HIRS-BM |
| anaplastic squamous cell carcinoma cell line:RPMI 2650 | Epithelioid mesothelioma cell line |
| argyrophil small cell carcinoma cell line:TC-YIK | Sarcomatoid mesothelioma cell line |
| astrocytoma cell line:TM-31 | Biphasic mesothelioma cell line |
| b cell line:RPMI1788 | mixed mullerian tumor cell line:HTMMT |
| B lymphoblastoid cell line: GM12878 ENCODE | mucinous adenocarcinoma cell line:JHOM-1 |
| basal cell carcinoma cell line:TE 354.T | mucinous cystadenocarcinoma cell line:MCAS |
| bile duct carcinoma cell line | myelodysplastic syndrome cell line:SKM-1 |
| biphenotypic B myelomonocytic leukemia cell line:MV-4-11 | myeloma cell line:PCM6 |
| bone marrow stromal cell line:StromaNKtert | myxofibrosarcoma cell line |
| breast carcinoma cell line | neuroblastoma cell line |
| bronchial squamous cell carcinoma cell line:KNS-62 | neuroectodermal tumor cell line |
| bronchioalveolar carcinoma cell line | neuroepithelioma cell line:SK-N-MC |
| bronchogenic carcinoma cell line:ChaGo-K-1 | neurofibroma cell line:Hs 53 |
| Burkitt lymphoma cell line | NK T cell leukemia cell line:KHYG-1 |
| carcinoid cell line | non T non B acute lymphoblastic leukemia cell line:P30/OHK |
| carcinosarcoma cell line:JHUCS-1 | non-small cell lung cancer cell line:NCI-H1385 |
| cervical cancer cell line | normal embryonic palatal mesenchymal cell line:HEPM |
| cholangiocellular carcinoma cell line:HuH-28 | normal intestinal epithelial cell line:FHs 74 Int |
| chondrosarcoma cell line:SW 1353 | oral squamous cell carcinoma cell line |
| choriocarcinoma cell line | osteoclastoma cell line:Hs 706 |
| chronic lymphocytic leukemia cell line:SKW-3 | osteosarcoma cell line |
| chronic megakaryoblastic cell line:MEG-01 | pagetoid sarcoma cell line:Hs 925 |
| chronic myeloblastic leukemia cell line:KCL-22 | pancreatic carcinoma cell line:NOR-P1 |
| chronic myelogenous leukemia cell line | papillary adenocarcinoma cell line:8505C |
| clear cell carcinoma cell line | papillotubular adenocarcinoma cell line:TGBC18TKB |
| colon carcinoma cell line | peripheral neuroectodermal tumor cell line:KU-SN |
| cord blood derived cell line:COBL-a untreated | pharyngeal carcinoma cell line:Detroit 562 |
| diffuse large B-cell lymphoma cell line:CTB-1 | plasma cell leukemia cell line:ARH-77 |
| ductal cell carcinoma cell line | pleomorphic hepatocellular carcinoma cell line:SNU-387 |
| embryonic kidney cell line: HEK293/SLAM untreated | prostate cancer cell line |
| embryonic pancreas cell line | rectal cancer cell line:TT1TKB |
| endometrial carcinoma cell line:OMC-2 | renal cell carcinoma cell line |
| endometrial stromal sarcoma cell line:OMC-9 | retinoblastoma cell line:Y79 |
| endometrioid adenocarcinoma cell line:JHUEM-1 | rhabdomyosarcoma cell line |
| epidermoid carcinoma cell line | sacrococcigeal teratoma cell line:HTST |
| epithelioid sarcoma cell line | schwannoma cell line:HS-PSS |
| epithelioid carcinoma cell line: HelaS3 ENCODE | serous adenocarcinoma cell line |
| Ewing sarcoma cell line:Hs 863 | serous cystadenocarcinoma cell line:HTOA |
| extraskeletal myxoid chondrosarcoma cell line:H-EMC-SS | signet ring carcinoma cell line |
| fibrosarcoma cell line:HT-1080 | small cell cervical cancer cell line:HCSC-1 |
| fibrous histiocytoma cell line:GCT TIB-223 | small cell gastrointestinal carcinoma cell line:ECC10 |
| gall bladder carcinoma cell line | small cell lung carcinoma cell line |
| gastric adenocarcinoma cell line | small-cell gastrointestinal carcinoma cell line:ECC4 |
| gastric cancer cell line | somatostatinoma cell line:QGP-1 |
| gastrointestinal carcinoma cell line:ECC12 | spindle cell sarcoma cell line:Hs 132 |
| giant cell carcinoma cell line | splenic lymphoma with villous lymphocytes cell line:SLVL |
| glassy cell carcinoma cell line:HOKUG | squamous cell carcinoma cell line:EC-GI-10.CNhs11252.10463-106H4 |
| glioblastoma cell line | squamous cell carcinoma cell line:JHUS-nk1.CNhs11749.10646-109A7 |
| glioma cell line:GI-1 | squamous cell carcinoma cell line:T3M-5.CNhs11739.10616-108G4 |
| granulosa cell tumor cell line:KGN | squamous cell lung carcinoma cell line |
| hairy cell leukemia cell line:Mo | synovial sarcoma cell line:HS-SY-II |
| Hep-2 cells mock treated | T cell lymphoma cell line:HuT 102 TIB-162 |
| hepatic mesenchymal tumor cell line:LI90 | teratocarcinoma cell line |
| hepatoblastoma cell line:HuH-6 | testicular germ cell embryonal carcinoma cell line |
| hepatocellular carcinoma cell line: HepG2 ENCODE | thymic carcinoma cell line:Ty-82 |
| hepatoma cell line:Li-7 | thyroid carcinoma cell line |
| hereditary spherocytic anemia cell line:WIL2-NS | transitional cell carcinoma cell line |
| Hodgkin lymphoma cell line:HD-Mar2 | tridermal teratoma cell line:HGRT |
| keratoacanthoma cell line:HKA-1 | tubular adenocarcinoma cell line:SUIT-2 |
| Krukenberg tumor cell line:HSKTC | Wilms tumor cell line |
| large cell lung carcinoma cell line | xeroderma pigentosum b cell line:XPL 17 |