

Figure S1. Inflammatory facial ulcer 11 months post gene therapy.
(A) Photograph of lesion; Histological H&E staining (B, C and D).

FIGURE LEGEND

Figure S1. Inflammatory facial ulcer 11 months post gene therapy.

(A) Photograph of lesion; (B and C) H&E histological appearance demonstrating pseudoepitheliomatous hyperplasia at the edge of the ulcer at x10 and x20 magnification, respectively; (D) H&E histological appearance demonstrating a dense inflammatory infiltrate of lymphocytes, plasma cells, eosinophils and neutrophils at x40 magnification.

Analysis of integration site distributions and relative clonal abundance for subject pWAS_UK11

Apr 28 2017

Introduction

The attached report describes results of analysis of integration site distributions and relative abundance for samples from gene therapy trials. For cases of gene correction in circulating blood cells, it is possible to harvest cells sequentially from blood to monitor cell populations. Frequency of isolation information can provide information on the clonal structure of the population. This report summarizes results for subject pWAS_UK11 over time points m6, m12, m20 in UCGS genome draft hg38.

The samples studied in this report, the numbers of sequence reads, recovered integration vectors, and unique integration sites available for this subject are shown below. We quantify population clone diversity using Gini coefficient, Shannon index, and UC50. The Gini coefficient provides a measure of inequality in clone abundance in each sample. The coefficient equals zero when all sites are equally abundant (poisson) and increases as fewer sites account for more of the total (oligopoisn). Shannon index is another widely used measure of diversity and accounts for both abundance and evenness of the integration events. Alternatively, the UC50 is the number of unique clones which make up the top 50% of the sample's abundance. For polyclonal samples, one may expect a low Gini coefficient, high Shannon Index, and high UC50 (proportional to the total number of unique sites identified in the sample).

Under most circumstances only a subset of sites will be sampled. We thus include an estimate of sample size based on frequency of unique integration from the SoniLength method (Barré, 2012). The 'Size' column denotes the estimated lower bound for population size derived using Chao estimate (Chao, 1987). If sample replicates were present then estimates were subjected to jackknife bias correction.

Trial	GTSF	Replicates	Patient	Timepoint	CellType	TotalReads	InferredCells	UniqueSites	Fragment	VEN	S.chao1	Gini	Shannon	UC50
WAS	GTSF277	4	pWAS_UK11	m6	Boleils	1015853	4543	3803	Shearing	NA	20148	0.1443	8.1681	1532
WAS	GTSF279	4	pWAS_UK11	m6	Boleils	212871	314	237	Shearing	NA	985	0.1991	5.3693	81
WAS	GTSF275	4	pWAS_UK11	m6	Neutrophils	784542	1963	1490	Shearing	NA	7251	0.2063	7.1815	509
WAS	GTSF278	4	pWAS_UK11	m6	NKcells	1088799	2949	2090	Shearing	NA	757	0.2371	7.5072	616
WAS	GTSF274	4	pWAS_UK11	m6	PBMC	885556	2061	1357	Shearing	NA	4628	0.2280	6.9992	327
WAS	GTSF276	4	pWAS_UK11	m6	Teils	1016780	5570	2996	Shearing	NA	10620	0.4152	7.2942	488
WAS	GTSF121	4	pWAS_UK11	m12	Boleils	898668	2159	1713	Shearing	NA	7071	0.1757	7.3530	643
WAS	GTSF123	4	pWAS_UK11	m12	Monocytes	355039	344	233	Shearing	NA	2088	0.2827	5.9202	62
WAS	GTSF119	4	pWAS_UK11	m12	Neutrophils	541561	707	476	Shearing	NA	8103	0.2963	5.2909	127
WAS	GTSF122	4	pWAS_UK11	m12	NKcells	498233	688	466	Shearing	NA	1850	0.2575	5.9939	127
WAS	GTSF118	4	pWAS_UK11	m12	PBMC	456621	2047	1721	Shearing	NA	7899	0.1468	7.3294	688
WAS	GTSF120	4	pWAS_UK11	m12	Teils	965593	3050	1999	Shearing	NA	8248	0.1797	7.3013	477
WAS	GTSF1491	4	pWAS_UK11	m20	Boleils	693499	975	562	Shearing	NA	2648	0.2827	6.1524	143
WAS	GTSF1493	4	pWAS_UK11	m20	Monocytes	791066	127	60	Shearing	NA	803	0.4867	3.3602	5
WAS	GTSF1489	4	pWAS_UK11	m20	Neutrophils	784120	1925	1226	Shearing	NA	4433	0.2796	6.9410	333
WAS	GTSF1492	4	pWAS_UK11	m20	NKcells	456570	523	361	Shearing	NA	6729	0.2813	5.6614	100
WAS	GTSF1488	4	pWAS_UK11	m20	PBMC	707262	6460	4335	Shearing	NA	14817	0.2963	7.5829	1106
WAS	GTSF1490	4	pWAS_UK11	m20	Teils	749603	1165	559	Shearing	NA	1980	0.4150	5.8776	105

Do any uniquely mapped clones account for greater than 20% of the total?

For some trials, a reporting criteria is whether any cell clones expand to account for greater than 20% of all clones. This is summarized below for subject pWAS_UK11. Abundance is estimated using the SoniLength method. Data such as this must, of course, be interpreted in the context of results from other assays. Distances reported refer to transcription start sites (5').

No sites found in this patient which are greater than 20% of the total data.

Do any multithit event account for greater than 20% of the total?

Up until now, the analysis has been looking at integration sites that can be uniquely mapped. But it is also helpful to look at reads finding multiple equally good alignments in the genome which can be referred to as 'Multithit' if an integration site occurred within a repeat element (i.e. Alu, LINE, SINE, etc), then it might be helpful to access those sites for potential detrimental effects. These collections of sequences are analyzed separately due to their ambiguity. To make some sense of these multithits, we bin any sequence(s) which share 1 or more genomic locations hence forming pseudo collections which can be referred to as OTUs (operation taxonomic units). Once the OTUs are formed, we compare breakdowns of unique sites and multithits. There is also a case for any multithits which higher abundance than a unique site in a given sample. Below is a table similar to the one shown previously except we show any events instead of genomic locations which are OTUs as opposed to clones in the data.

No multithits sites found in this patient which are greater than 20% of the total data.

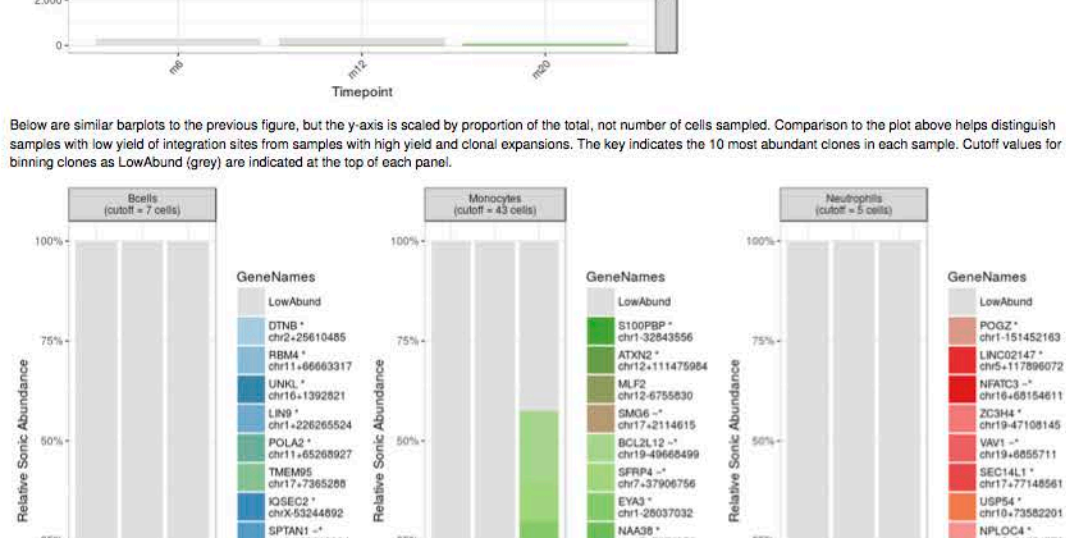
Relative abundance of cell clones

The relative abundance of cell clones is summarized in the attached bar graphs. The cell fraction studied is named at the top, the time points are marked at the bottom. The different bars in each panel show the major cell clones, as marked by integration sites. A key to the sites is shown at the right. Throughout the whole report, each integration site is assigned a GeneName given by:

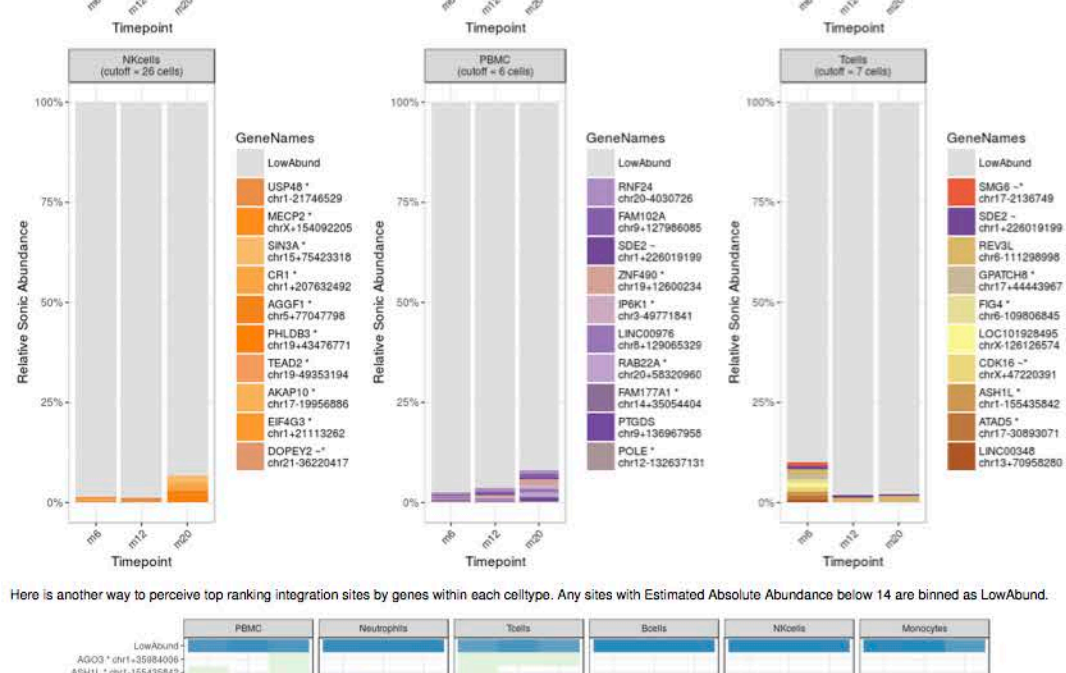
- GeneName refers to the closest gene to either end and strand,
- L indicates the site is within a L1 repeat element,
- S indicates the site is within a SINE repeat element,
- I indicates the gene was associated with lymphoma in humans.

Integration sites were recovered using Ligation mediated PCR after random fragmentation of genomic DNA, which reduces recovery bias compared with restriction enzyme cleavage. Relative abundance was not measured from read counts, which are known to be inaccurate, but from marks introduced into DNA specimens prior to PCR amplification using the SoniLength method (Barré, 2012).

In the barplots below, the x-axis indicates each sample type and time point, the y-axis is scaled by proportion of the total, not number of cells sampled. Comparison to the plot above helps distinguish samples with the yield of integration sites from samples with high yield and clonal expansions. The key indicates the 10 most abundant clones in each sample. Cutoff values for binning clones as LowAbund (grey) are indicated at the top of each panel.



Below are similar barplots to the previous figure, but the y-axis is scaled by proportion of the total, not number of cells sampled. Comparison to the plot above helps distinguish samples with the yield of integration sites from samples with high yield and clonal expansions. The key indicates the 10 most abundant clones in each sample. Cutoff values for binning clones as LowAbund (grey) are indicated at the top of each panel.

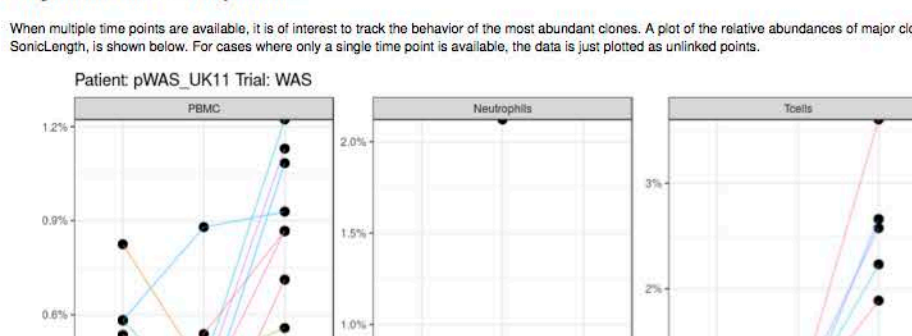


Here is another way to perceive top ranking integration sites by genes within each celltype. Any sites with Estimated Absolute Abundance below 14 are binned as LowAbund.



Longitudinal behavior of major clones

When multiple time points are available, it is of interest to track the behavior of the most abundant clones. A plot of the relative abundances of major clones, based on output from SoniLength, is shown below. For cases where only a single time point is available, the data is just plotted as unsorted points.

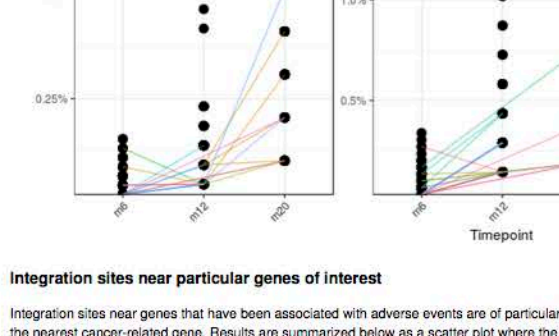


Integration sites near particular genes of interest

Integration sites near genes that have been associated with adverse events are of particular interest. Thus, we have cataloged all integration sites for which a gene of interest is the nearest cancer-related gene. Results are summarized below as a scatter plot where the y-axis shows relative abundance of the integration site and x-axis is distance to the nearest oncogene's 5' end.

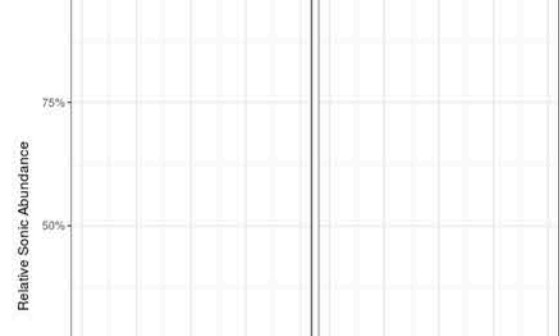
Negative distances indicate that the integration site is downstream from 5' end after the TSS. Positive distances indicate that the integration site is upstream from (i.e. before) the TSS. Note that all RefSeq splicing isoforms are used for this analysis, so the reference TSS may not be the same for each listed integration site.

LMO2



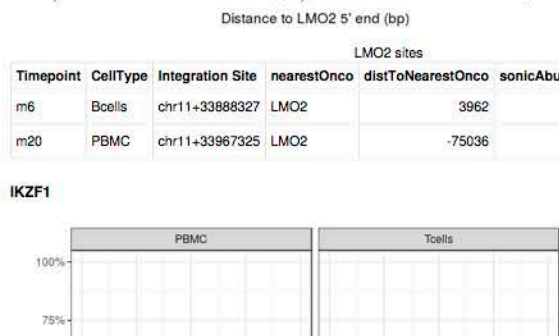
Timepoint	CellType	Integration Site	nearestOnco	distToNearestOnco	sonicAbundance	sonicAbundanceRank
m6	Boleils	chr11-3388872	LMO2	3962	1	3803
m20	PBMC	chr11-33987325	LMO2	-75036	1	4335

IKZF1



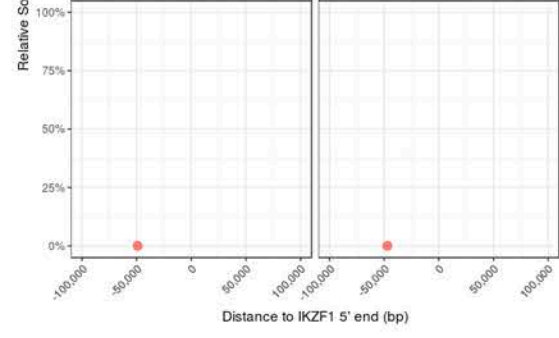
Timepoint	CellType	Integration Site	nearestOnco	distToNearestOnco	sonicAbundance	sonicAbundanceRank
m12	Teils	chr7-5020818	IKZF1	-43454	1	1999
m6	Boleils	chr7-42525087	IKZF1	-48995	2	549
m6	NKcells	chr7-50258668	IKZF1	-47116	1	2090
m20	PBMC	chr7-50319108	IKZF1	62	1	4335

CND2



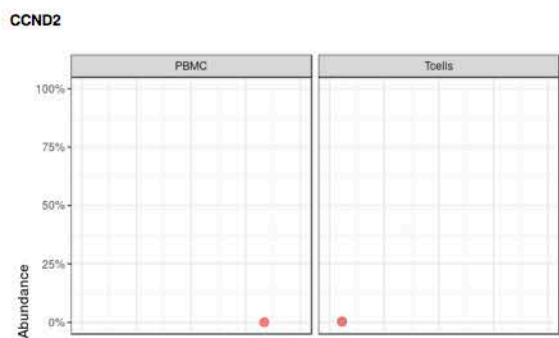
Timepoint	CellType	Integration Site	nearestOnco	distToNearestOnco	sonicAbundance	sonicAbundanceRank
m12	Teils	chr12-4185175	CND2	-8850	10	14
m12	Boleils	chr12-4251002	CND2	-22723	1	1713
m6	PBMC	chr12-4340473	CND2	6878	1	1353
m6	Teils	chr12-4185175	CND2	-8850	17	29
m6	NKcells	chr12-4207289	CND2	-56454	1	2090

HMG2



Timepoint	CellType	Integration Site	nearestOnco	distToNearestOnco	sonicAbundance	sonicAbundanceRank
m12	Boleils	chr12-6564758	HMG2	23068	1	1713
m6	Boleils	chr12-6540316	HMG2	5957	1	3803
m20	PBMC	chr12-6578992	HMG2	-4537	2	897

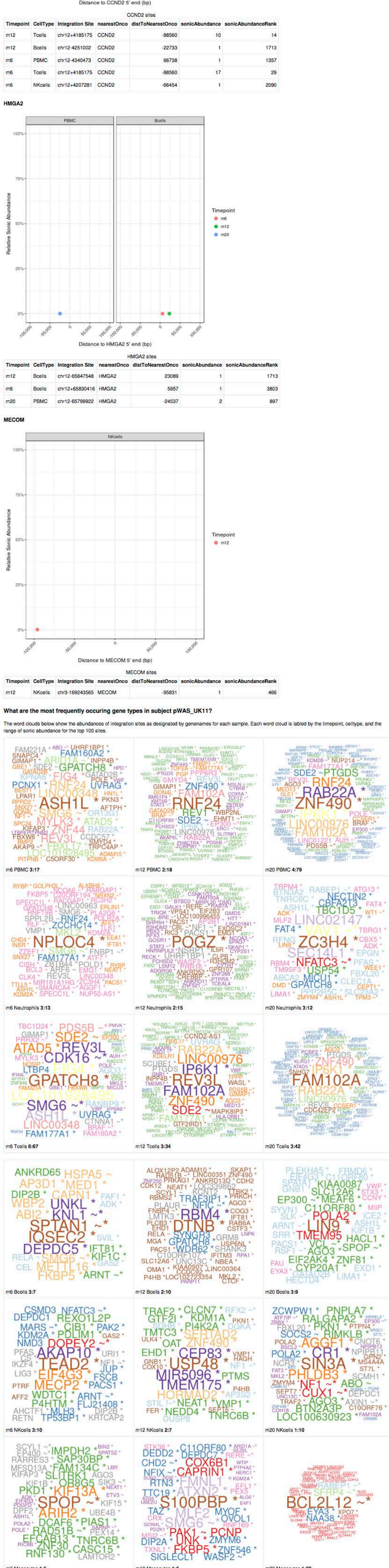
MECOM



Timepoint	CellType	Integration Site	nearestOnco	distToNearestOnco	sonicAbundance	sonicAbundanceRank
m12	NKcells	chr186243565	MECOM	-9581	1	466

What are the most frequently occurring gene types in subject pWAS_UK11?

The word clouds below show the abundances of integration sites as designated by genes for each sample. Each word cloud is labeled by the timepoint, celltype, and the range of clone abundance for the top 100 sites.



Integration sites shared between cell types.

Month 6

	Bcells	Monocytes	Neutrophils	NKcells	PBMC	Tcells
Bcells	3821	0	2	5	8	5
Monocytes	0	237	5	2	2	1
Neutrophils	2	5	1499	14	43	63
NKcells	5	2	14	2104	6	3
PBMC	8	2	43	6	1362	317
Tcells	5	1	63	3	317	3016

Month 12

	Bcells	Monocytes	Neutrophils	NKcells	PBMC	Tcells
Bcells	1725	4	6	5	17	3
Monocytes	4	238	7	3	6	0
Neutrophils	6	7	478	8	16	2
NKcells	5	3	8	468	14	0
PBMC	17	6	16	14	1732	42
Tcells	3	0	2	0	42	2014

Month 20

	Bcells	Monocytes	Neutrophils	NKcells	PBMC	Tcells
Bcells	583	0	8	2	10	0
Monocytes	0	60	0	0	1	0
Neutrophils	8	0	1237	9	112	13
NKcells	2	0	9	362	8	2
PBMC	10	1	112	8	4370	110
Tcells	0	0	13	2	110	564