

## SUPPLEMENTARY INFORMATION FOR

### **Somatic mutant clone selection in human skin varies with body site**

Joanna C Fowler<sup>1</sup>, Charlotte King<sup>1</sup>, Christopher Bryant<sup>1</sup>, Michael Hall<sup>1,3</sup>, Roshan Sood<sup>1</sup>, Swee Hoe Ong<sup>1</sup>, Eleanor Earp<sup>1</sup>, David Fernandez-Antoran<sup>1</sup>, Jonas Koeppel<sup>1</sup>, Stefan Dentre<sup>1,2</sup>, David Shorthouse<sup>3</sup>, Amer Durrani<sup>4</sup>, Kate Fife<sup>4</sup>, Edward Rytina<sup>4</sup>, Doreen Milne<sup>4</sup>, Amit Roshan<sup>4,5</sup>, Krishnaa Mahububani<sup>6</sup>, Kourosh Saeb-Parsy<sup>6</sup>, Benjamin A Hall<sup>3</sup>, Moritz Gerstung<sup>2,7</sup>, Philip H Jones<sup>1,3\*</sup>

<sup>1</sup> Wellcome Sanger Institute, Hinxton CB10 1SA, UK

<sup>2</sup> European Molecular Biology Laboratory, European Bioinformatics Institute, Cambridge CB10 1SD, UK

<sup>3</sup> MRC Cancer Unit, University of Cambridge, Hutchison-MRC Research Centre, Cambridge Biomedical Campus, Cambridge CB2 0XZ, UK

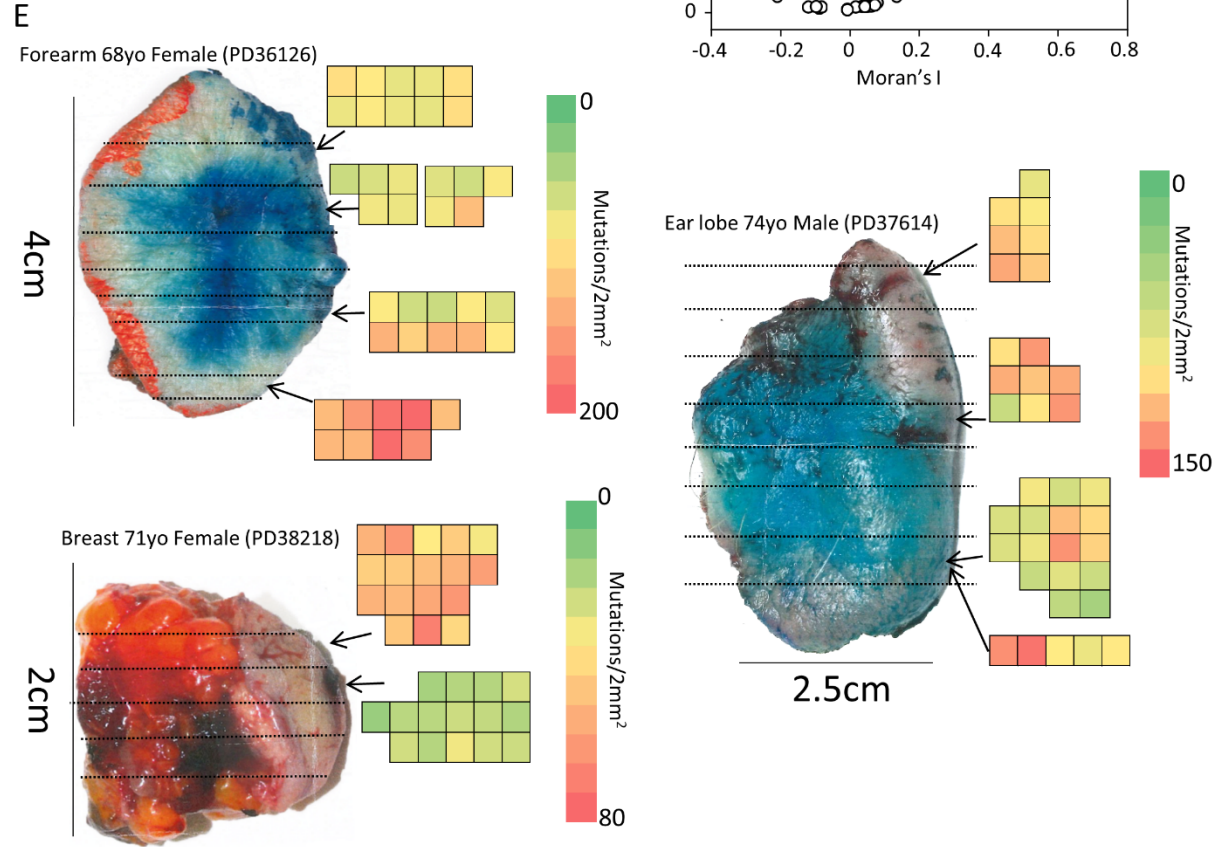
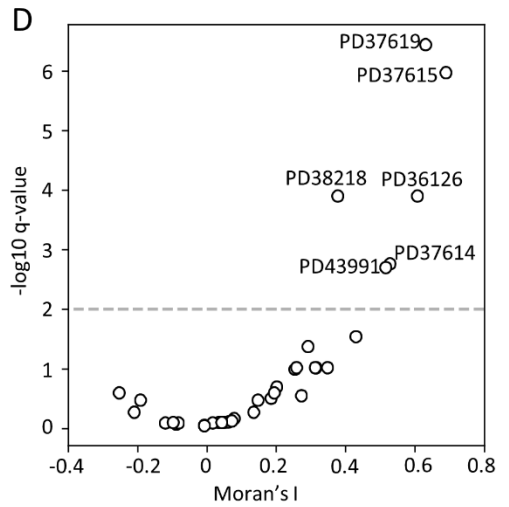
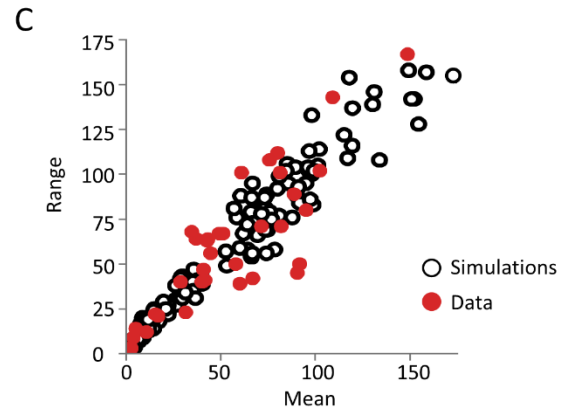
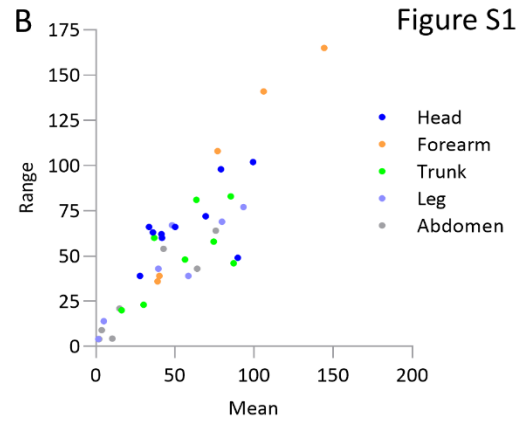
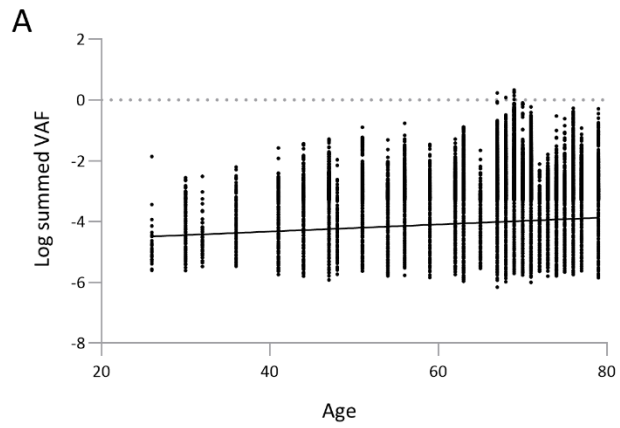
<sup>4</sup> Cambridge University Hospitals NHS Foundation Trust, Cambridge Biomedical Campus Cambridge, CB2 0QQ, UK.

<sup>5</sup> Cancer Research UK Cambridge Institute, University of Cambridge, Li Ka Shing Centre, Robinson Way, Cambridge CB2 0RE, UK

<sup>6</sup> Department of Surgery and Cambridge NIHR Biomedical Research Centre, Biomedical Campus, Cambridge CB2 0QQ

<sup>7</sup> European Molecular Biology Laboratory, Genome Biology Unit, Heidelberg, DE

\*correspondence to [pj3@sanger.ac.uk](mailto:pj3@sanger.ac.uk)



### Supplemental figure 1 - Assessing mutational variability in human skin

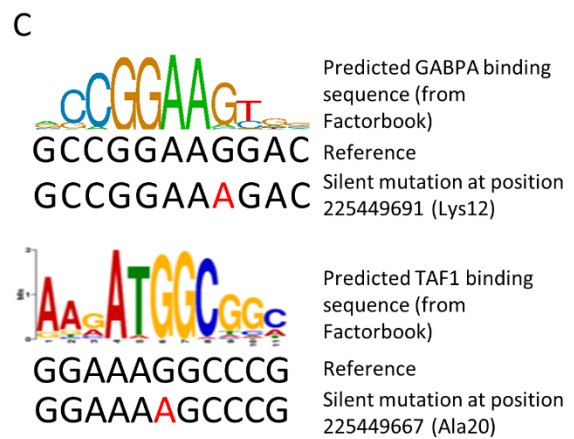
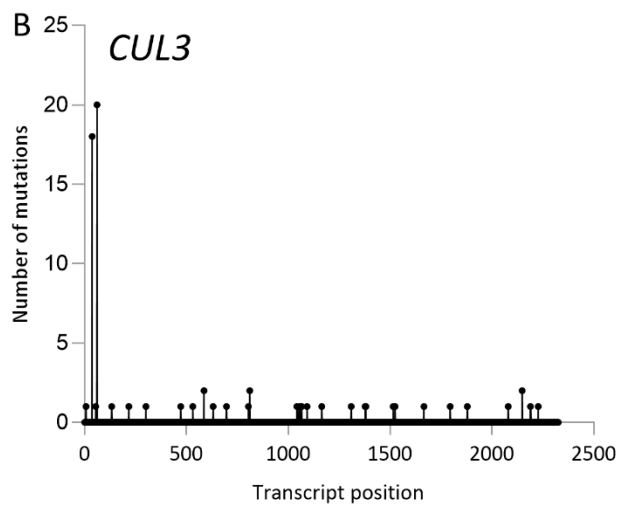
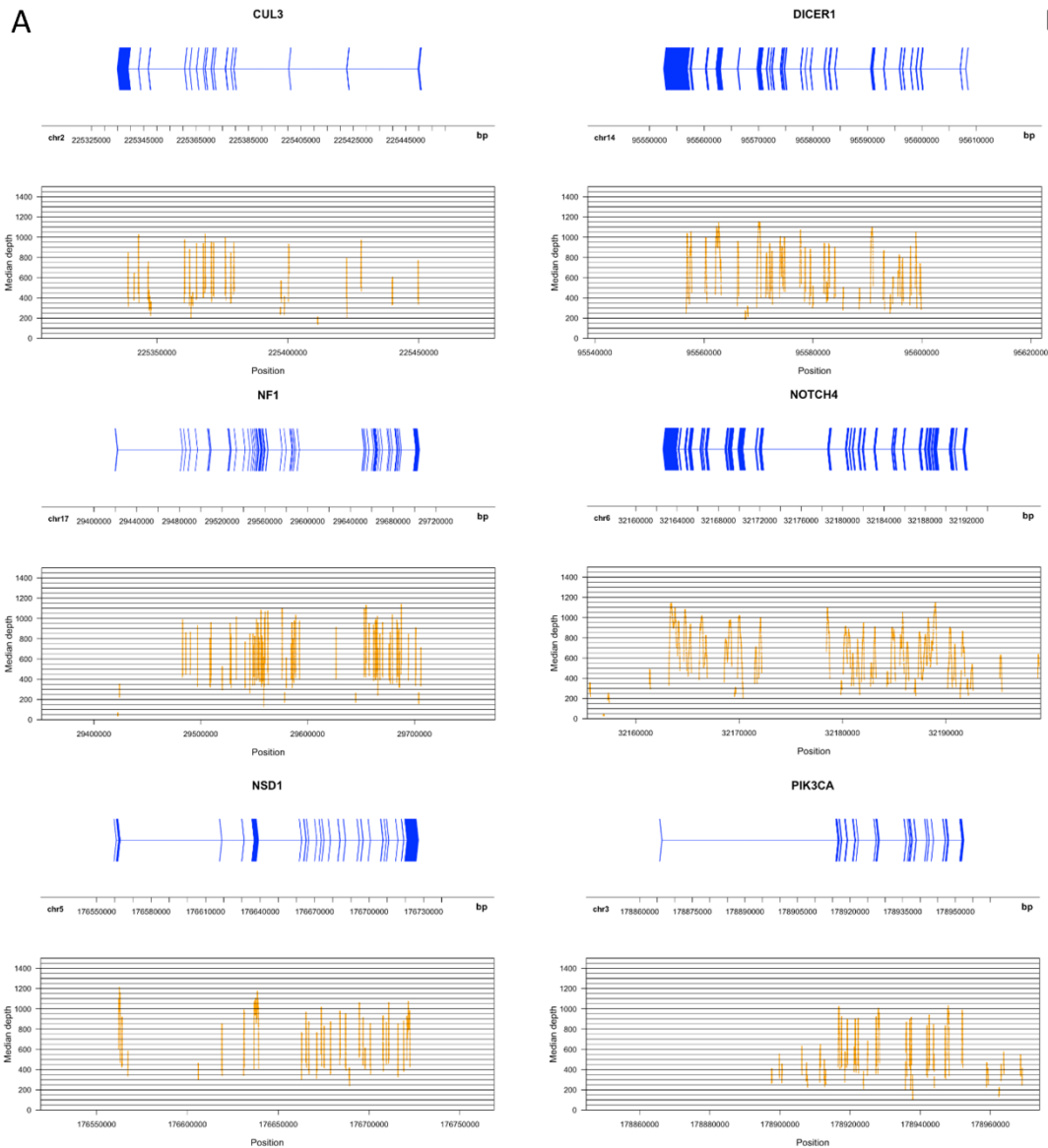
- a. Log of summed VAF. A simple linear regression finds that there is an increase in clone size (VAF) with age ( $p < 2.2 \times 10^{-16}$ , slope = 0.01161, intercept = -4.79356), however, this should be interpreted with caution because, even after log transformation, the data does not meet the normal distribution assumption. Note for this analysis PD38219 was removed as an outlier because of their predominantly non-UV signature
- b. Mean mutation number per sample vs the range (difference between the highest and lowest mutation count in single biopsies) for each patient. Different colours represent different body sites.
- c. Mean mutation number per sample vs the range for each patient (red) and simulations (open black markers).
- d. Patients whose epithelia shows a greater than expected difference in the number of mutations/  $2\text{mm}^2$  gridded compared to their adjacent neighbours. Moran's I and  $-\log_{10}$  q-value for each donor are shown. Values of Moran's I closer to 1 show increased positive autocorrelation (more clustering of high values/clustering of low values).
- e. Further examples of epidermis which show a gradient in the number of mutations across the tissue. Heat maps show the number of mutations in each  $2\text{mm}^2$  grid.



### Supplemental figure 2 – Positive selection across different body sites

- a. Number of different types of mutations for positively selected genes (global  $q < 0.01$ ,  $dN/dS > 2$ )
- b. Mutational spectrum of single nucleotide variants from all 2mm<sup>2</sup> grid samples from targeted sequencing. The bar plot shows the number of mutations in each of the 96 possible trinucleotide contexts.
- c. Percentage of mutant epithelia for all positively selected genes (global  $q < 0.01$ ,  $dN/dS > 2$ )
- d. Sliding window plot of missense mutations per codon in *NOTCH1*. Observed counts shown by the black line. Expected counts assuming that missense mutations were distributed across the gene according to the mutational spectrum shown in grey. Domains of NOTCH1 protein shown underneath the x-axis. EGF-repeats shown in blue, with ligand-binding site EGF repeats 11-12 shown in dark blue.

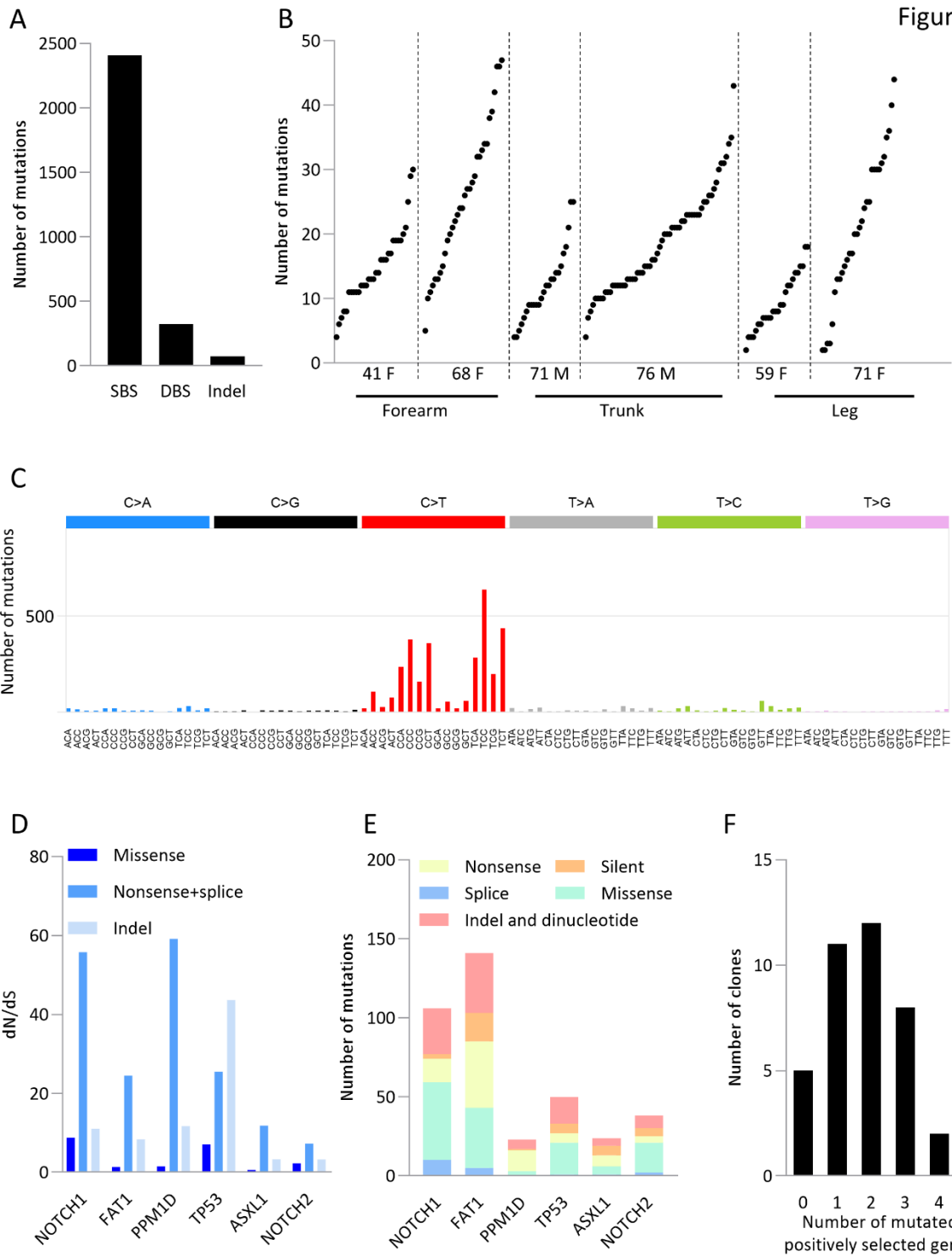
Figure S3



### **Supplemental figure 3 – Negative selection is not due to lower sequencing coverage**

- a. Sequence coverage per gene of all genes showing negative dN/dS selection. In each case gene structure is shown in the top panel and the median sequence depth per base is shown below. Sequence baits are designed to cover only the exonic regions explaining the lack of sequence depth in the introns.
- b. Distribution of number of *CUL3* silent mutations along the transcript length of the gene.
- c. Consequence of recurrent *CUL3* silent mutations on predicted transcription factor binding sites.

Figure S4



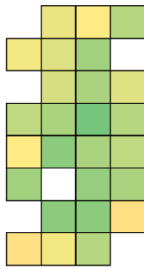
**Supplemental figure 4 – Sequence analysis of data from 0.25mm diameter punches is highly similar to that from 2mm<sup>2</sup> grid samples**

- a. Total number of mutations from targeted sequencing of 232 0.25mm diameter punches from 6 donor from different body sites. Number of single base (SBS), double base (DBS) and insertion/ deletions (indel) are shown. Proportions are comparable to that seen in 2mm<sup>2</sup> gridded sample sequencing.

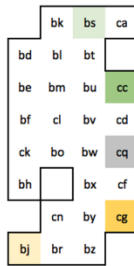


- b. Number of mutations per 0.25mm diameter punch sample. Each point represents a sample. Patients are ordered by age and divided by body site. Age and gender of the patient is shown below.
- c. Mutational spectrum of single nucleotide variants from all 0.25mm diameter punch samples from targeted sequencing. The bar plot shows the number of mutations in each of the 96 possible trinucleotide contexts.
- d. dN/dS ratios for missense, nonsense/ splice substitutions and insertions/deletions (indel) for all 0.25mm diameter punch samples. Only genes with global  $q < 0.01$  are shown.
- e. Number of different types of mutations for positively selected genes. Proportions are similar to that seen from targeting sequencing from 2mm<sup>2</sup> gridded samples.
- f. Number of positively selected driver genes carrying non-synonymous mutations from whole genome sequencing of clonal 0.25mm diameter punch samples.

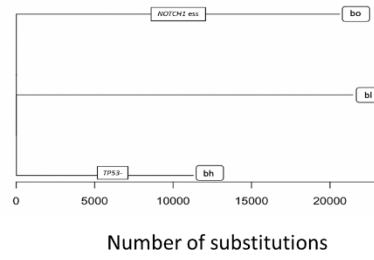
**A**  
59yo Female leg (PD38331)



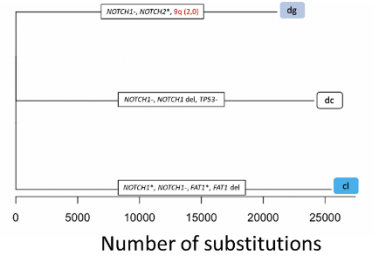
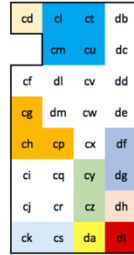
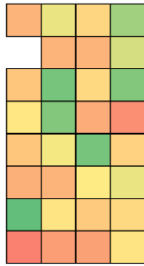
**B**



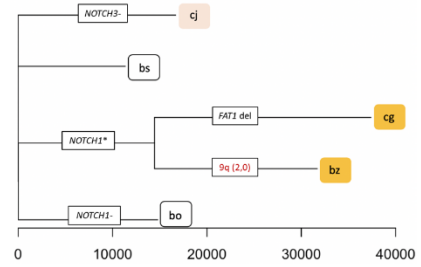
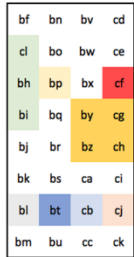
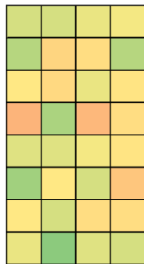
**C**



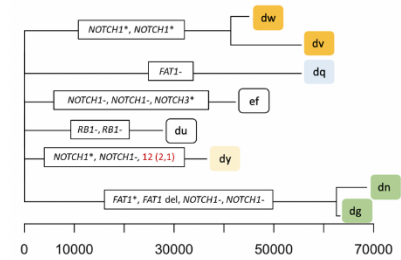
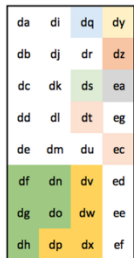
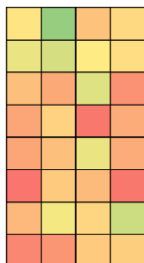
71yo Female leg (PD38334)



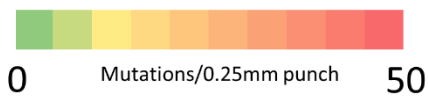
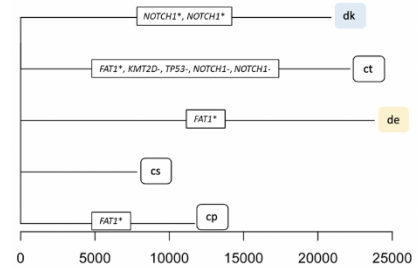
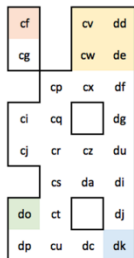
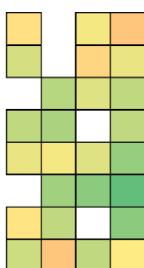
41yo Female forearm (PD37577)



68yo Female forearm (PD36126)



71yo Male trunk (PD38215)

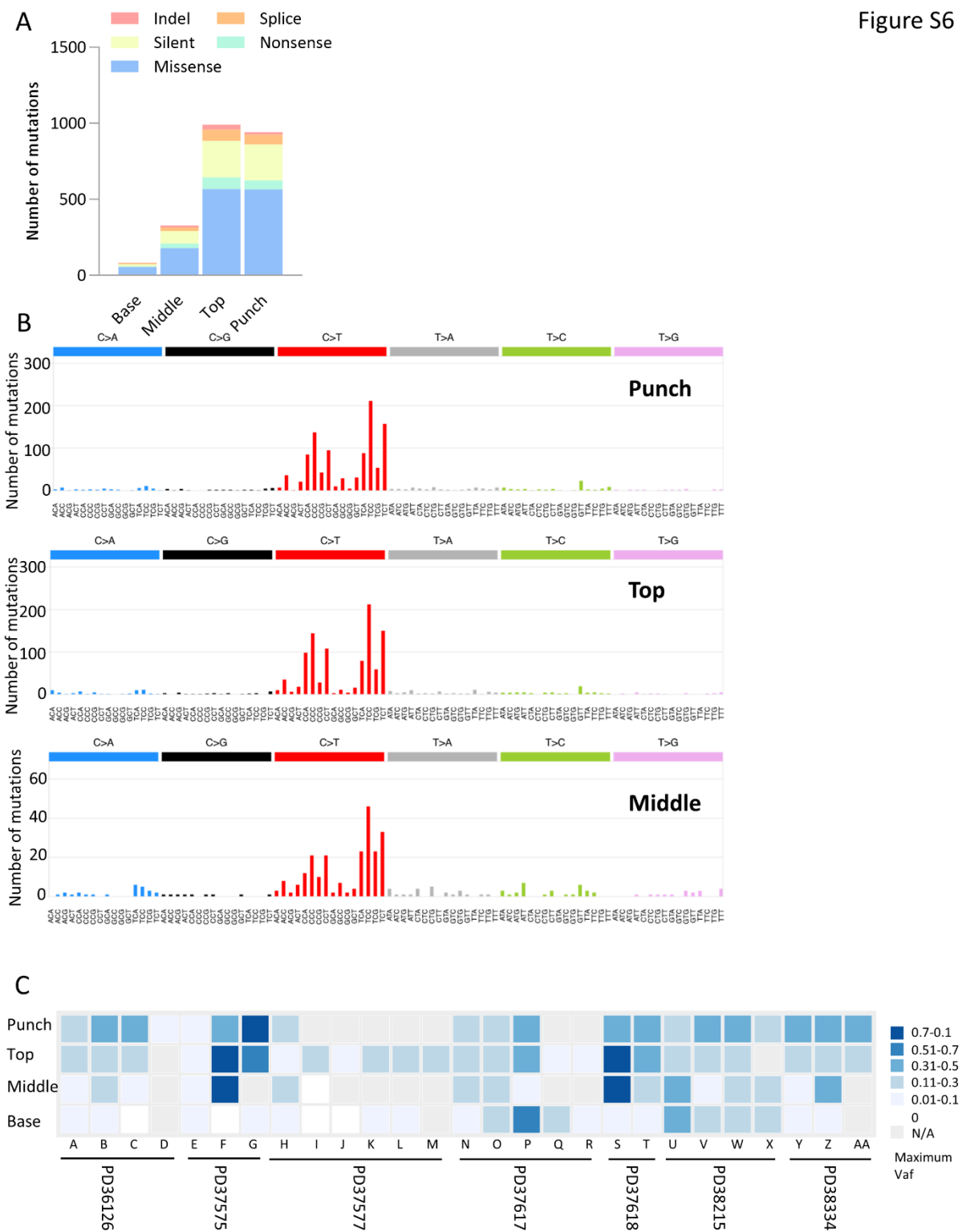


\* Nonsense    Ess Essential splice  
- Missense    CNA

**Supplemental figure 5 – Heatmaps and mutational variability of 0.25mm diameter punch samples**

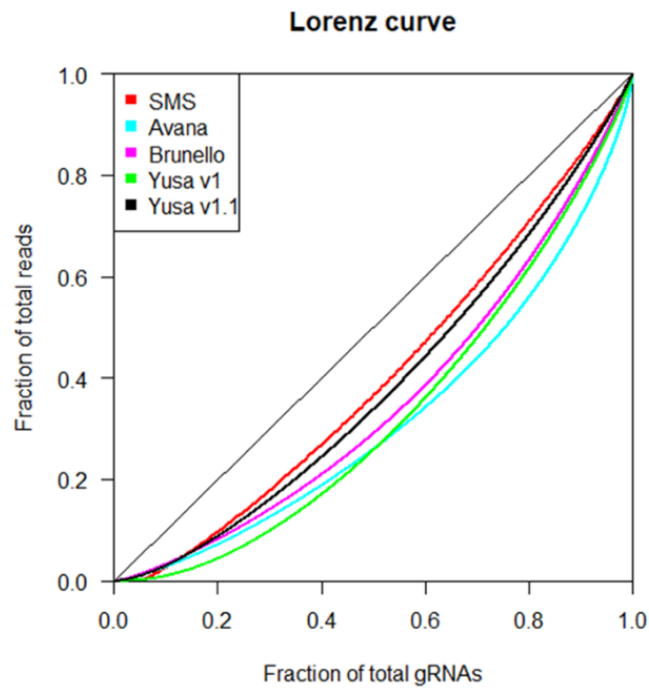
- a. Heat map of number of mutations per punch across other five donors. Each square represents an individual 0.25mm diameter sample.
- b. Epidermal clone map area for each donor. Mutations with a VAF of at least 0.2 are shown, with each clone represented in a different colour. Polyclonal samples are shown in white.
- c. Maximum parsimony tree of clonal substitutions detected in 24 whole-genome 0.25mm diameter punch samples from five additional donors. Branch lengths are equivalent to the number of clonal single and double base substitutions and are annotated with clonal non-synonymous mutations detected in the 13 genes found to be under positive selection. Within each branch, driver mutations are arbitrarily ordered. Copy number alterations are shown in red.

Figure S6



**Supplemental figure 6 – Mutational spectrum of epidermal 0.25mm diameter punch, top or middle regions of the hair follicle**

- a. Number of different types of mutations for each part of the follicle.
- b. Mutational spectrum for different parts of the follicle, excluding the base where there were too few mutations (<200) inhibiting mutational signature decomposition. The bar plot shows the number of mutations in each of the 96 possible trinucleotide contexts.
- c. Maximum variant allele frequency (VAF) of any mutation in each section of the follicle. Each column is a follicle with the patient noted below. Each row is a different part of the follicle.



#### Supplemental figure 7: Gini index of Crispr/CAS9 screen libraries

The 'ineq' R package was used to generate Lorenz curves and to calculate a Gini index for the generated library and widely used genome-wide libraries. The Gini index is a measurement of statistical dispersion and reflects the evenness of read count distribution. A lower Gini index is indicative of a more even distribution (**Figure S7**).



	Mean	Median	Stdev
<b>Targeted</b>			
2mm2 grids			
All	668.7	667.0	111.0
Abdomen	682.7	660.7	146.1
Forearm	595.3	595.4	141.7
Head_melanoma	681.5	677.9	62.5
Head non-melanoma	613.5	603.4	78.6
Leg	727.4	720.4	77.1
Trunk	691.2	695.5	74.8
Punches			
All	44.9	42.3	20.3
Hair follicles			
All	372.3	410.7	158.7
<b>Whole genome sequencing</b>			
2mm2 grids			
All	49.5	45.7	8.1
Punches			
All	33.1	29.8	7.8

**Table S2 - Mean sequence coverage of all experiments**

Punches are 0.25mm diameter samples



**List of excel Tables S3-S8**

**Table S3 – Lists of genes on targeted sequencing bait sets**

**Table S4 - List of all mutations detected in all 2mm<sup>2</sup> gridded samples**

**Table S5 – dN/dS ratios of all genes with q value <0.01 from targeted sequencing of 2mm<sup>2</sup> gridded samples**

**Table S6 – Fold change in guide RNAs in targeted CRISPR-Cas9 knockout screen**

**Table S7 – List of all mutations detected in all 0.25mm diameter samples**

**Table S8 – List of all mutations detected in hair follicle samples**