

1 Supplemental Methods and Supplemental Figures

2 Discovery of an unusual high number of *de novo*  
3 mutations in sperm of older men using duplex  
4 sequencing

5 Renato Salazar, Barbara Arbeithuber, Maja Ivankovic, Monika Heinzl, Sofia Moura, Ingrid  
6 Hartl, Theresa Mair, Angelika Lahnsteiner, Thomas Ebner, Omar Shebl, Johannes Pröll, Irene  
7 Tiemann-Boege

8

9 Table of Contents

10 ***Supplemental Methods* ..... 1**

11 ***Discovery of an unusual high number of de novo mutations in sperm of older men***  
12 ***using duplex sequencing*..... 1**

13 ***Table of Contents* ..... 1**

14 ***Supplemental Methods* ..... 3**

15     **Collection of sperm and testes samples ..... 3**

16     **Sample preparation..... 3**

17     **Duplex library preparation ..... 4**

18     **Targeted enzymatic fragmentation..... 4**

19     **Adapter synthesis..... 5**

20     **End-repair/A-tailing and ligation ..... 7**

21     **DNA amplification and targeted capture ..... 7**

22     **Quality control ..... 9**

23     **Sequencing..... 10**

24     **Sequencing Data Availability ..... 10**

25     **Data processing and variant filtering ..... 10**

26     **Variant frequency comparison ..... 11**

|   |  |           |
|---|--|-----------|
| 1 | <b>Sensitivity evaluation .....</b>  | <b>12</b> |
| 2 | <b>Droplet Digital PCR.....</b>  | <b>12</b> |
| 3 | <b>DS vs ddPCR and BEA.....</b>  | <b>13</b> |
| 4 | <b><i>Supplemental References .....</i></b>  | <b>13</b> |
| 5 | <b><i>Supplemental Figures.....</i></b>  | <b>15</b> |
| 6 |  |           |
|   | Supplemental Figure S1. Types of substitutions found either in the SSCS or the DCS.  | 15        |
|   | Supplemental Figure S2. Age distribution of the donors used in the young and old sperm pool.   | 16        |
|   | Supplemental Figure S3. Mutation frequencies of the sequenced libraries.   | 17        |
|   | Supplemental Figure S4. Mutation frequencies of (non)-CpG transitions and transversions.   | 18        |
|   | Supplemental Figure S5. Mutation frequencies of each substitution type.  | 19        |
|   | Supplemental Figure S6. Mutational spectra (A) based on mutation frequencies, (B) transcription bias and (C) mutational signature.                                 | 20        |
|   | Supplemental Figure S7. Mutational spectra of variants compared to the mutational spectra extracted from public databases.   | 21        |
|   | Supplemental Figure S8. Mutational spectra (A) based on mutation frequencies, (B) transcription bias and (C) mutational signatures of the Younger and Older group. | 22        |
|   | Supplemental Figure S9. Mutational signature compared to a catalogue of signatures (version 3.2) in the COSMIC database.   | 23        |
|   | Supplemental Figure S10. DCS Coverage of each library.   | 24        |
|   | Supplemental Figure S11. Different adapter designs for DS.   | 25        |
|   | Supplemental Figure S12. Example of fragment length distribution of different library preparation steps.   | 26        |

7

8

# 1 Supplemental Methods

## 2 ***Collection of sperm and testes samples***

3 Sperm DNA from anonymous donors mainly from Central European ancestry aged between 26  
4 to 30 and 49 to 59 years were pooled using 5 donors per pool; the exact pool composition is  
5 shown in Table S4. Samples were collected from the Kinderwunsch Klinik, MedCampus IV,  
6 Kepler Universitätsklinikum, Linz following the protocol approved by the ethics commission of  
7 Upper Austria (Approval F1-11). Three human genomic DNA samples encoding the *FGFR3*  
8 mutations c.742C>T (NA00711), c.746C>G (NA08909), and c.749C>G (CD00002) were  
9 purchased from Coriell Cell Repositories (Camden, NJ). Similarly, genomic DNA with the *FGFR3*  
10 c.1620C>A transversion (p.540N>K) was extracted from the B-lymphocyte cell line (GM18666)  
11 purchased from the Coriell Cell depository.

## 12 ***Sample preparation***

13 Sperm genomic DNA (gDNA) was extracted using the Gentra Puregene Cell kit (QIAGEN) as  
14 follows: first, 25  $\mu\text{L}$  of semen ( $\sim 2 \times 10^6$  sperm cells) were centrifuged (8,000g, 20 seconds) and  
15 the resulting sperm cell pellet was incubated overnight with 150  $\mu\text{L}$  of cell lysis solution, 6  $\mu\text{L}$   
16 of 1M DTT, and 0.5  $\mu\text{L}$  of proteinase K (20mg/mL) at 37°C. RNase treatment was performed by  
17 adding 0.75  $\mu\text{L}$  of RNase A (4mg/mL) and incubating at 37°C for 15 minutes. For the protein  
18 precipitation step, the sample was put on ice for 5 minutes. Then, 50  $\mu\text{L}$  of protein  
19 precipitation solution were added followed by vigorous manual shaking for 1 minute. Next,  
20 the sample was centrifuged for 20 minutes at 13,000g to remove the protein precipitate. The  
21 supernatant was transferred into a fresh tube and a second centrifugation step was  
22 performed if the solution was still cloudy. The DNA was precipitated by adding 150  $\mu\text{L}$  of  
23 isopropanol (100%) and 0.25  $\mu\text{L}$  of glycogen solution to the supernatant, followed by gentle  
24 inversions. A centrifugation step was done (30 minutes, 13,000g at 4°C) and the supernatant  
25 was discarded. DNA was washed with 150  $\mu\text{L}$  of 70% ethanol, followed by a 3-minute  
26 centrifugation (13,000g, 4°C). Supernatant was removed and DNA was dried at room  
27 temperature for 3 minutes. Finally, DNA was resuspended in 25  $\mu\text{L}$  of DNA hydration solution.  
28 Genomic DNA from B-lymphocyte cell line was extracted using the same kit. Cells in cultured  
29 medium ( $\sim 2 \times 10^6$  cells) were added to a 1.5 mL microcentrifuge tube and centrifuged for 5  
30 seconds at 13,000g. Supernatant was discarded and pellet was resuspended in residual

1 supernatant by vortexing vigorously. Cell lysis solution was added (300  $\mu$ L) followed by a 10  
2 seconds vortexing step to lyse cells. Next, 100  $\mu$ L of protein precipitation solution were added  
3 followed by 20 seconds of vortexing. The mixture was then centrifuged for 1 minute at  
4 13,000g to pellet the proteins. Isopropanol (300  $\mu$ L) was added to a clean 1.5 mL  
5 microcentrifuge tube followed by the addition of the supernatant of the previous step. The  
6 tube was inverted 50 times, followed by a centrifugation step (1 minute, 13,000g). At this  
7 stage, DNA formed a pellet and the supernatant was carefully discarded. Then, 300  $\mu$ L of 70%  
8 ethanol were added and the tube was inverted several times to wash the extracted DNA. The  
9 tube was centrifuged for 1 minute at 13,000g and the supernatant was removed, followed by  
10 a drying step at room temperature for 5 minutes. Dried DNA was resuspended in 100  $\mu$ L of  
11 DNA hydration solution and was vortexed for 5 seconds at medium speed. DNA was incubated  
12 overnight at 22°C with gentle shaking. Finally, samples were transferred to a storage tube.

### 13 ***Duplex library preparation***

14 DS was based on previous protocols (Kennedy et al. 2014) with substantial modifications  
15 including an initial restriction digest to pre-select the target regions. In the different libraries  
16 we used some minor modifications to our main protocol, either starting with distinct input  
17 DNA amounts, different adapters, or amplification strategies; see Table S3. These  
18 modifications should not impact the variants called.

### 19 ***Targeted enzymatic fragmentation***

20 Genomic DNA (500-5,000 ng, Table S3) was subject to an overnight restriction enzyme digest  
21 that targeted 5 regions of the *FGFR3* gene (Table S8). Resulting fragments were expected to  
22 have a length between ~300 to ~600 bp. A double size selection was performed using  
23 SPRIselect beads (Beckman Coulter) in order to exclude fragments outside of this size range.  
24 To the 100  $\mu$ L of restriction digest solution, 93.5  $\mu$ L of PCR grade water were added. Then,  
25 106.5  $\mu$ L of beads were added (0.55 volumes of beads), mixed by pipetting thoroughly and  
26 incubated at room temperature for 5 minutes. Tubes were then placed on a magnet for 5  
27 minutes and 290  $\mu$ L of supernatant were transferred to a new tube. Next, 84.2  $\mu$ L of beads  
28 (1.0 volumes in total considering the initial bead solution) were added to the solution and  
29 mixed by pipetting. The mixture was incubated at room temperature for 5 minutes. Tubes  
30 were placed on a magnet and supernatant was discarded. Beads were washed twice with 80%  
31 ethanol while still on the magnet. Beads were dried at room temperature and 50  $\mu$ L of PCR  
32 grade water were added. Beads were resuspended by pipetting and incubated at room

1 temperature for 5 minutes. Finally, tubes were placed on a magnet and the clear supernatant  
2 containing the size-selected DNA was transferred to a new tube.

### 3 *Adapter synthesis*

4 We used 3 different adapter synthesis protocols to produce Adapter 1, Adapter 2 and Adapter  
5 3. Supplemental Figure S11 shows the structural differences between adapters: Adapter 1 is  
6 identical to the original open-hairpin structure, whereas adapters 2 and 3 have a closed loop  
7 hairpin structure that is opened up after the ligation step by the cleavage of the uracil with  
8 the USER enzyme mix (New England Biolabs). It is suggested that closed loops adapters reduce  
9 adapter dimers during ligation (New England Biolabs). Structurally, adapter 3 is similar to  
10 adapter 2, but is 9 bp longer and contains a phosphorothioate bond before the T-overhang on  
11 the 3' end such that removal of this overhang is reduced. T-tailed DS adapters were  
12 synthesized with modifications to the original protocol (Kennedy et al. 2014). Oligonucleotides  
13 used for each adapter are listed in Table S12.

14 Adapter 1 was synthesized as described by (Kennedy et al. 2014), except that the amounts  
15 used for each component were downscaled. Oligonucleotides mws51 and mws55 (for each,  
16 10  $\mu\text{L}$  of 100 pmole/ $\mu\text{L}$ ) were added to a single PCR tube and heated in a PCR machine for 5  
17 minutes at 95°C. The PCR machine was immediately turned off and the mixture was left inside  
18 for one hour, resulting in the annealing of both oligonucleotides. The extension of annealed  
19 oligonucleotides was performed by adding 23.65  $\mu\text{L}$  of PCR grade water, 5.6  $\mu\text{L}$  of 10 $\times$  NEB  
20 buffer 2 (New England Biolabs), 5.6  $\mu\text{L}$  of 2.5mM each dNTPs (Biozym), and 1.15  $\mu\text{L}$  of NEB  
21 Klenow Exo- (5U/ $\mu\text{L}$ ) (New England Biolabs). The mixture was incubated for 1 hour at 37°C.  
22 Purification was done by adding 28  $\mu\text{L}$  (0.5 volumes) of ammonium acetate (5 M) (Invitrogen)  
23 followed by the addition of 168  $\mu\text{L}$  (2 volumes) of pre-cooled 100% ethanol. The solution was  
24 inverted several times until it got cloudy. Precipitation occurred while the solution was  
25 incubated for 30 minutes at -20°C. The solution was then centrifuged for 30 minutes at 4°C in  
26 a microcentrifuge at 14,000 rpm. Supernatant was then discarded and 1 mL of 80% ethanol  
27 was added followed by a 5-minute centrifugation at 4°C, 14,000 rpm. The supernatant was  
28 discarded and the tube was dried at room temperature until the pellet became transparent.  
29 The pellet was resuspended in 40  $\mu\text{L}$  of PCR grade water. The T-overhang was generated by  
30 digestion with HpyCH4III enzyme: 47  $\mu\text{L}$  of PCR grade water, 10  $\mu\text{L}$  of 10 $\times$  CutSmart buffer  
31 (New England Biolabs), and 3  $\mu\text{L}$  of HpyCH4III (5U/ $\mu\text{L}$ ) (New England Biolabs) were added to  
32 the extended adapters and incubated for 16 hours. To purify the adapters, 400  $\mu\text{L}$  of PCR  
33 grade water were added, followed by 250  $\mu\text{L}$  of ammonium acetate (5 M). Then, 1.5 mL of

1 pre-cooled 100% ethanol were added and the solution was mixed by inverting until it got  
2 cloudy. Adapters were precipitated while incubated at -20°C for 30 minutes. The solution was  
3 then centrifuged for 30 minutes at 4°C, 14,000 rpm. Supernatant was discarded and 1 mL of  
4 80% ethanol was added followed by a 5-minute centrifugation at 4°C, 14,000 rpm. The  
5 supernatant was again discarded and the tube was dried at room temperature until the pellet  
6 became transparent. Adapters were resuspended by adding 22 µL of 10mM mM TE<sub>low</sub> (Tris-  
7 HCl pH 8.0, 0.1 mM EDTA).

8 Adapter 2 was synthesized using the adapter 1 protocol. The only difference is that instead of  
9 using oligonucleotides mws51 and mws55 in the first step, only oligonucleotide DS\_Hairpin\_U  
10 was used, and consequently, an extra 10 µL of PCR grade water were added to the extension  
11 reaction.

12 The synthesis of adapter 3 started by using oligonucleotides DS\_Hairpin\_U and TA-Overhang.  
13 First, 10 µL of each oligonucleotide (each at 100 pmole/µL) were added to a PCR tube. The  
14 phosphorylation of 5' ends was done by adding 23.5 µL of PCR grade water, 5.2 µL of T4 DNA  
15 Ligase Buffer (10×) with ATP (10 mM) (New England Biolabs), and 3.3 µL of T4 Polynucleotide  
16 Kinase (10U/µL), followed by an incubation at 37°C for 30 minutes. Then, the solution was  
17 heated in a PCR machine at 95°C for 5 minutes. The machine was then turned off and the  
18 mixture stayed inside for 1 hour to anneal the oligonucleotides. The extension of annealed  
19 oligonucleotides was done by mixing 50 µL from the previous step with 26.7 µL of PCR grade  
20 water, 5µL of T4 DNA Ligase Buffer (10×) with ATP (10 mM), 10 µL of dNTPs (2.5 mM), 1 µL of  
21 BSA (10 mg/mL) (New England Biolabs), 0.6 µL of T4 DNA Ligase (2,000 U/µL) (New England  
22 Biolabs), and 1.7 µL of T4 DNA Polymerase (New England Biolabs). The solution was incubated  
23 at 12°C for 30 minutes followed by 16°C for another 30 minutes. The purification of adapters  
24 was accomplished by adding 400 µL of PCR grade water followed by 250 µL of ammonium  
25 acetate (5 M). Then, 1.5 mL of pre-cooled 100% ethanol were added and the solution was  
26 mixed by inverting until it got cloudy. Adapters were precipitated while incubated at -20°C for  
27 30 minutes. The solution was then centrifuged for 30 minutes at 4°C, 14,000 rpm. The  
28 supernatant was discarded and 1 mL of 80% ethanol was added followed by a 5-minute  
29 centrifugation at 4°C, 14,000 rpm. The supernatant was again discarded and the tube was  
30 dried at room temperature until the pellet became transparent. Adapters were resuspended  
31 by adding 22 µL of 10mM mM TE<sub>low</sub> (Tris-HCl pH 8.0, 0.1 mM EDTA).

1 *End-repair/A-tailing and ligation*

2 End-repair/A-tailing was done by mixing 50  $\mu\text{L}$  of size-selected genomic DNA with 7  $\mu\text{L}$  of  
3 NEBNext Ultra II End Prep Reaction Buffer (New England Biolabs), and 3  $\mu\text{L}$  of NEBNext Ultra II  
4 End Prep Enzyme Mix (New England Biolabs). The mixture was incubated at 20°C for 30  
5 minutes followed by 65°C for 30 minutes.

6 The ligation of adapters occurred by adding 30  $\mu\text{L}$  of NEBNext Ultra II Ligation Master Mix  
7 (New England Biolabs), 1  $\mu\text{L}$  of NEBNext Ultra II Ligation Enhancer (New England Biolabs), and  
8 2.5  $\mu\text{L}$  of adapter dilution to the 60  $\mu\text{L}$  of End-repair/A-tailing mix. The solution was then  
9 incubated at 20°C for 15 minutes, followed by 4°C overnight. Next, 3  $\mu\text{L}$  of USER (1 U/ $\mu\text{l}$ ) (New  
10 England Biolabs) were added and the solution was incubated at 37°C for 15 minutes,  
11 immediately followed by AMPure XP beads (Beckman Coulter) purification. To prepare the 2.5  
12  $\mu\text{L}$  of adapter dilution, we first estimated the number of genomic DNA molecules was  
13 estimated based on the ng present after size selection, and then the amount of adapter  
14 needed to have a 20-fold excess. PCR grade water was added to fill the volume up to 2.5  $\mu\text{L}$ .

15 Bead purification was performed as follows: 77.2  $\mu\text{L}$  (0.8 volumes) of beads were added to the  
16 96.5  $\mu\text{L}$  of ligation reaction and mixed by pipetting; the mixture was incubated at room  
17 temperature for 15 minutes and then placed on a magnet for 5 minutes; the supernatant was  
18 discarded and beads were washed twice with 400  $\mu\text{L}$  of 80% ethanol; the beads were dried at  
19 room temperature and resuspended by adding 20  $\mu\text{L}$  of PCR grade water; the mixture was  
20 incubated at room temperature for 5 minutes and then placed on a magnet until the solution  
21 was clear; the supernatant containing ligated DNA was transferred to a new tube.

22 *DNA amplification and targeted capture*

23 The amplification of ligated DNA was done using KAPA HiFi HotStart ReadyMix (KAPA  
24 Biosystems). Components and volumes used are shown in Table S9 (PCR 1) and primer  
25 sequences are available in Table S10. Each library was split into multiple reactions, each with  
26 an input material of approximately 240 ng. With the exceptions of 4 libraries (see Table S3),  
27 the first step of amplification was 12 cycles of single primer extensions (Table S9) followed by  
28 the addition of the primer NEBNext Universal and a standard PCR amplification.

29 For purification, multiple PCR reactions originating from the same ligation reaction were  
30 pooled. Depending on the number of reactions and the corresponding final volume, 1.2  
31 volumes of AMPure XP beads were mixed with each sample. The purification procedure  
32 continued as described above. DNA was eluted in 20  $\mu\text{L}$  of PCR grade water.

1 The target capture was modified from Schmitt *et al.* (Schmitt et al. 2015). First, the purified  
2 PCR 1 output was mixed with 5 µg of Cot-I DNA (Invitrogen) and 1 nmol of blocking  
3 oligonucleotides mws60 and mws61 (Table S11). The mixture was then lyophilized and  
4 resuspended in 3.1 µL water, 8.5 µL NimbleGen 2× Hybridization Buffer (Roche) and 3.4 µL  
5 NimbleGen Hybridization Component A (Roche). Next, the mixture was heated for 10 minutes  
6 at 95°C, and 3 pmol of pooled biotinylated oligos specifically designed for the targeted regions  
7 (Table S11) were added at 65°C, followed by an incubation of 4 hours and 30 minutes.  
8 Meanwhile, 22.5 µL of NimbleGen 10× Wash Buffer I, 15 µL of 10× Wash Buffer II, 15µL of 10×  
9 Wash Buffer III, 30 µL of 10× Stringent Wash Buffer, and 150 µL of 2.5× Bead Wash Buffer  
10 (Roche) were diluted in PCR grade water in order to obtain a 1× working solution. Wash Buffer  
11 I and Stringent Wash Buffer were incubated at 65°C. Ten minutes before the end of the  
12 hybridization reaction, streptavidin beads were washed as following: 75 µL of M-270  
13 streptavidin beads (Life Technologies) were placed on a magnet and the supernatant was  
14 discarded; 150 µL of Bead Wash Buffer were added and the mixture was vortexed for 10  
15 seconds and placed on a magnet; the supernatant was discarded and the wash was repeated;  
16 75 µL of Bead Wash Buffer were added and the mixture was vortexed until homogeneous;  
17 shortly before the hybridization reaction was completed, beads were placed on a magnet, and  
18 supernatant was removed. With the beads still on the magnet, the hybridization reaction was  
19 added to the beads, quickly mixed by pipetting (it is important to be fast to maintain the  
20 temperature), transferred back to the PCR tube, and placed back in the PCR machine also at  
21 65°C for 45 minutes. Sample were vortexed for 3 seconds every 10 minutes. After incubation,  
22 captured DNA was bound to the beads and the next step is the washing which was performed  
23 as following: with the tube still on the PCR machine, 75 µL of Wash Buffer I was added to the  
24 mixture and then vortexed for 10 seconds; the mixture was placed on a 1.5 mL tube and then  
25 placed on a magnet; the supernatant was removed and 150 µL of Stringent Wash Buffer were  
26 added and mixed by pipetting; the mixture was incubated for 5 minutes at 65°C; Stringent  
27 Wash Buffer wash was repeated; after the second incubation, the mixture was placed on a  
28 magnet and the supernatant was removed followed by the addition of 150 µL of room  
29 temperature Wash Buffer I; the mixture was vortexed for 2 minutes and placed on a magnet;  
30 150 µL of Wash Buffer II were added and the mixture was vortexed for 1 minute and then  
31 placed on a magnet; the supernatant was removed and 150 µL of Wash Buffer III was added  
32 followed by 30 seconds of vortexing; the mixture was then placed on a magnet, the  
33 supernatant was removed and 44 µL of PCR grade water were added to resuspend the beads  
34 with the bound DNA.

1 The amplification of hybridized targets was done as shown in Table S9 (column PCR 2-4). PCR  
2 Product was purified using 1.2 volumes of AMPure XP beads (100  $\mu$ L of PCR product and 120  
3  $\mu$ L of beads). Apart from the different initial volumes, purification was done as described  
4 before. For the final elution, 20  $\mu$ L of PCR grade water was used.

5 For the second and third target captures, the same procedure was applied but only half of the  
6 amounts of Cot-I DNA, blocking oligos and biotinylated oligos were used (the remaining  
7 volume of the latter was replaced by PCR grade water). Also, the PCR reaction was slightly  
8 different because NEBNext<sup>®</sup> Multiplex Oligos for Illumina<sup>®</sup> (New England Biolabs) were used  
9 instead of NEB\_mws20 primer in order to add indexes for sequencing multiple libraries in one  
10 sequencing run.

#### 11 *Quality control*

12 Dilutions of 2  $\mu$ L aliquots of every step of the library preparation were subject to dsDNA  
13 concentration measurements and fragment size distribution analysis to monitor the quality of  
14 the library preparation. To measure the concentration of dsDNA, the DeNovix High Sensitivity  
15 Assay (DeNovix Inc.) and a DeNovix DS-11 FX fluorometer (DeNovix Inc.) were used following  
16 the exact instructions of the manufacturer. The fragment size distribution was obtained using  
17 the High Sensitivity DNA Kit (Agilent) and a 2100 Bioanalyzer instrument (Agilent) according to  
18 the exact instructions of the manufacturer. Analysis of results was carried on with the 2100  
19 Expert Software (Agilent). Figure S12 shows an example of fragment length distribution for a  
20 four of the different library preparation steps. The final distribution of fragment lengths is  
21 distinctive for our approach when compared with random fragmentation methods because all  
22 fragments of each targeted region have identical length.

23 The last step before pooling of different libraries for sequencing was a qPCR using the Library  
24 Quantification kit (KAPA Biosystems), following the procedure detailed by the manufacturer.  
25 The results were used to adjust the concentrations previously obtained with the fluorescence  
26 method. The adjusted concentrations were then used to calculate the volumes of each library  
27 that should be pooled together and sent for sequencing, following the guidelines of the  
28 sequencing provider used.

## 1 **Sequencing**

2 Sequencing was performed on the Illumina MiSeq platform using the MiSeq Reagent v3 600  
3 cycles kit (Illumina) at the Center for Medical Research of the Johannes Kepler University, Linz,  
4 Austria, and at the VBCF NGS Unit, Vienna, Austria.

## 5 **Sequencing Data Availability**

6 The raw sequencing data generated in this study have been submitted to the NCBI BioProject  
7 database (<https://www.ncbi.nlm.nih.gov/bioproject/>) under accession number PRJNA684907.

## 8 **Data processing and variant filtering**

9 FASTQ files were analyzed in Galaxy (on both public ([usegalaxy.org](https://usegalaxy.org)) and private ([zusie.jku.at](https://zusie.jku.at))  
10 servers) to process the data according to a DS specific workflow that is available at  
11 <https://usegalaxy.org/u/jku-itb-lab/w/ds-fgfr3-2021>. The workflow link contains information  
12 on versions used for each tool. This workflow combines the toolset *Du Novo* (Afgan et al.  
13 2016; Stoler et al. 2016; Stoler et al. 2020) with a recently published set of tools. These are  
14 designed not only to monitor the quality and efficiency of each library, but also to re-examine  
15 variant calls with a tier-based strategy, giving the user a better overview and control over the  
16 mutations called (Povysil et al. 2021). The first step of the workflow is grouping together all  
17 reads that share the same barcodes, which is done by the tool *Du Novo: Make families*. Next,  
18 the *Du Novo: Correct barcodes* tool will group together barcodes with differences up to 3  
19 nucleotides (this setting can be changed) since those are probably originated by PCR or  
20 sequencing errors and would otherwise be lost (Stoler et al. 2020). The tool *Du Novo: Align*  
21 families will then perform a multiple-sequencing alignment of every read within each family of  
22 reads. The consensus sequences are then built by *Du Novo: Make consensus reads*. After  
23 trimming, consensus reads are mapped to the human genome (hg38) using Map with *BWA-*  
24 *MEM* (Li 2013) followed by a left realignment with *BamLeftAlign* (Garrison and Marth 2012).  
25 Variants are called using *Naive Variant Caller* (Blankenberg et al. 2014) and annotated by the  
26 Variant Annotator (Blankenberg et al. 2014).  
27 Then, the *Variant Analyzer* considers the called mutations, and re-examines the raw reads in  
28 order to assign them to different tiers (Table S5). This tier classification is based mainly on the  
29 number of reads per family and the proportion of the alternative allele in the consensus  
30 sequence (Povysil et al. 2021).

1 High-quality tiers (tier 1.1) require at least 3 reads per family ( $FS \geq 3$ ); whereas, second order  
2 tiers (tier 1.2-2.4) allow also reads without a family ( $FS \geq 1$ ). Both assume that at least 75% of  
3 the reads in the family carry the mutant allele. Lower quality tiers (tier 3 and 4) include also  
4 other parameters such as the read position (e.g., mutants in the first or last 10bp are assigned  
5 to tier 4) or lower the percentage of reads carrying the mutant allele to  $>50\%$  (tier 3).

6 The *Variant Analyzer* output was inspected and variants with DCS coverage below 1000,  
7 intronic variants, and SNPs were discarded. Only exonic variants of tier 1 were kept, together  
8 with tier 2 variants that were detected more than once in different libraries.

9 Table S14 shows an example of a partial *Variant Analyzer* output with variants found in  
10 different libraries. For this hypothetical case, we would immediately exclude variants chr4-  
11 1803735-C-A and chr4-1804278-A-T because all mutants have low quality tiers (tier 3 or 4);  
12 the high frequency and a further inspection on population databases indicates variant chr4-  
13 1804314-G-A is a SNP, so it would also be excluded; variant chr4-1803672-G-T is tier 2.1, but  
14 since it is only found once and it is not found in any other library, it would be also excluded  
15 from our data; for the remaining variants we inspect if they are exonic or intronic, and we  
16 would exclude variant chr4-1803861-G-T because it is an intronic variant for the *FGFR3*  
17 transcript we are focusing on; the remaining exonic variants would then be used for  
18 downstream analysis.

19 Selected variants were annotated using Variant Effect Predictor (VEP) (McLaren et al. 2016)  
20 and wANNOVAR (Yang and Wang 2015).

### 21 ***Variant frequency comparison***

22 The variant frequency was calculated by dividing the number of DCS calling the variant by the  
23 DCS coverage at the position of the variant within the library it was detected. Poisson  
24 confidence intervals (CI) were calculated for each variant according to Garwood, 1936  
25 (Garwood 1936; Patil and Kulkarni 2012), where the lower limit is  $\chi^2_{2x, 0.025} / 2$  and the upper  
26 limit is  $\chi^2_{2x+1, 0.975} / 2$  and  $x$  is the number mutants detected for a particular position. CIs for  
27 each mutation frequency were obtained by dividing each of the above-mentioned limits by  
28 the number of molecules analyzed for a particular position (DCS coverage).

29 Available data for the *FGFR3* variants were extracted from gnomAD v2.1.1 (transcript  
30 ENST00000440486.2, <https://gnomad.broadinstitute.org>) (Karczewski et al. 2020) and COSMIC  
31 v90 (transcript ENST00000440486.7, <https://cancer.sanger.ac.uk>) (Tate et al. 2018) databases  
32 on January 23<sup>rd</sup> 2020 and January 21<sup>st</sup> 2020, respectively. Variants from gnomAD were

1 remapped from GRCh37/hg19 to GRCh38/hg38 using NCBI Genome Remapping Service  
2 (<https://www.ncbi.nlm.nih.gov/genome/tools/remap>). Exonic single-nucleotide substitutions  
3 retrieved from gnomAD and COSMIC were annotated with VEP and wANNOVAR.  
4 Deleteriousness analysis was performed using CADD raw scores (Rentzsch et al. 2018)  
5 extracted from VEP annotation.  
6 Pairwise comparison between every group of variants was done with the Mann-Whitney *U*  
7 test using GraphPad Prism 8.2.1 (GraphPad Software). A multiple comparison correction using  
8 a False Discovery Rate (FDR) approach (two-stage set-up method of Benjamini, Krieger and  
9 Yekutieli) was done using the same software.  
10 Germline association was investigated using ClinVar data extracted from wANNOVAR output.  
11 Tumor association was investigated by consulting the presence or absence of variants in the  
12 extracted COSMIC data.

### 13 ***Sensitivity evaluation***

14 In order to assess the ability of our adapted DS approach to detect low- and ultralow-  
15 frequency variants, we sequenced four libraries with a mixture of genomic DNA containing the  
16 *FGFR3* variants c.742C>T, c.746C>G, c.749C>G, c.1620C>A, and wild-type DNA at different  
17 ratios. The ratio of each variant to WT was 1:10, 1:100, 1:1000, and 1:10,000 (one ratio per  
18 library) and it was calculated considering the heterozygosity of the genomic samples  
19 containing the variants analyzed. Variant calling and frequency calculations were done as  
20 described above.

### 21 ***Droplet Digital PCR***

22 Site-specific assays were designed using the Droplet digital PCR Assay design tool from BioRad  
23 (available at: <https://www.bio-rad.com/digital-assays>). For site-specific assay information see  
24 Table S13.

25 Individual 20  $\mu$ L ddPCR reactions were prepared by adding 10  $\mu$ L of 10 $\times$  SuperMix for probes  
26 (no dUTP), 2  $\mu$ L of genomic DNA ( $\sim$ 125ng/ $\mu$ L; total of  $\sim$ 72500 genomes), 6.7  $\mu$ L of Nucleic Acid  
27 free water, 1  $\mu$ L of primers and probes mix (900nM/probe and 100nM/primer), and 0.3  $\mu$ L of  
28 restriction enzyme (CviQI or MseI; 10U/ $\mu$ L). Note that for each sample 4 replicates were  
29 carried out (total of  $\sim$ 300,000 genomes). The digest occurred at room temperature for 15  
30 minutes prior to transferring the samples into the cartridges. 70  $\mu$ L of droplet generation oil

1 for probes and a sealing gasket were added and droplets were formed in the droplet  
2 generator. Next, the newly formed droplets were transferred to a 96-well plate and sealed for  
3 5 seconds at 180°C. PCR was ran as follows: 10 minutes at 95°C, 40 cycles of: 30 seconds at  
4 94°C and 1 minute at 53/55°C (depending on the target site, see Table S13), and a final step of  
5 10 minutes at 98°C. Each step had a ramp rate of 2°C and the lid was heated to 105°C. End-  
6 point analysis was done in the droplet reader. Results were analyzed using QuantaSoft  
7 Analysis Pro Software version 1.7.4 (Bio-Rad Laboratories Inc.).

## 8 **DS vs ddPCR and BEA**

9 Six mutations detected with DS were measured with ddPCR or BEA in the same donors to  
10 confirm the authenticity of DS mutations. Frequencies were calculated by simply dividing the  
11 number of mutants by the number of total molecules analyzed. CIs were calculated for each  
12 variant as described above. To investigate whether there are significant differences between  
13 the results obtained with DS and ddPCR/BEA, a Fischer exact test was performed for all  
14 comparisons using the GraphPad Prism 8.2.1 software.

## 15 Supplemental References

- 16 Afgan E, Baker D, van den Beek M, Blankenberg D, Bouvier D, Cech M, Chilton J, Clements D,  
17 Coraor N, Eberhard C et al. 2016. The Galaxy platform for accessible, reproducible and  
18 collaborative biomedical analyses: 2016 update. *Nucleic Acids Res* **44**: W3-w10.
- 19 Bergstrom EN, Huang MN, Mahto U, Barnes M, Stratton MR, Rozen SG, Alexandrov LB. 2019.  
20 SigProfilerMatrixGenerator: a tool for visualizing and exploring patterns of small  
21 mutational events. *BMC Genomics* **20**: 685.
- 22 Blankenberg D, Von Kuster G, Bouvier E, Baker D, Afgan E, Stoler N, Taylor J, Nekrutenko A,  
23 Galaxy T. 2014. Dissemination of scientific software with Galaxy ToolShed. *Genome*  
24 *Biology* **15**: 403.
- 25 Garrison E, Marth G. 2012. Haplotype-based variant detection from short-read sequencing.  
26 arXiv:1207.3907.
- 27 Garwood F. 1936. Fiducial Limits for the Poisson Distribution. *Biometrika* **28**: 437-442.
- 28 Karczewski KJ Francioli LC Tiao G Cummings BB Alföldi J Wang Q Collins RL Laricchia KM Ganna  
29 A Birnbaum DP et al. 2020. The mutational constraint spectrum quantified from  
30 variation in 141,456 humans. *Nature* **581**: 434-443.
- 31 Kennedy SR, Schmitt MW, Fox EJ, Kohn BF, Salk JJ, Ahn EH, Prindle MJ, Kuong KJ, Shen JC,  
32 Risques RA et al. 2014. Detecting ultralow-frequency mutations by Duplex Sequencing.  
33 *Nat Protoc* **9**: 2586-2606.
- 34 Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.  
35 arXiv:1303.3997.
- 36 McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GRS, Thormann A, Flicek P, Cunningham F. 2016.  
37 The Ensembl Variant Effect Predictor. *Genome Biology* **17**.
- 38 Patil VV, Kulkarni HV. 2012. Comparison of confidence intervals for the poisson mean: some  
39 new aspects. *REVSTAT - Statistical Journal* **10**: 211+.

1 Povysil G, Heinzl M, Salazar R, Stoler N, Nekrutenko A, Tiemann-Boege I. 2021. Increased yields  
2 of duplex sequencing data by a series of quality control tools. *NAR Genomics and*  
3 *Bioinformatics* **3**.

4 Rentzsch P, Witten D, Cooper GM, Shendure J, Kircher M. 2018. CADD: predicting the  
5 deleteriousness of variants throughout the human genome. *Nucleic Acids Research* **47**:  
6 D886-D894.

7 Schmitt MW, Fox EJ, Prindle MJ, Reid-Bayliss KS, True LD, Radich JP, Loeb LA. 2015. Sequencing  
8 small genomic targets with high efficiency and extreme accuracy. *Nat Methods* **12**:  
9 423-425.

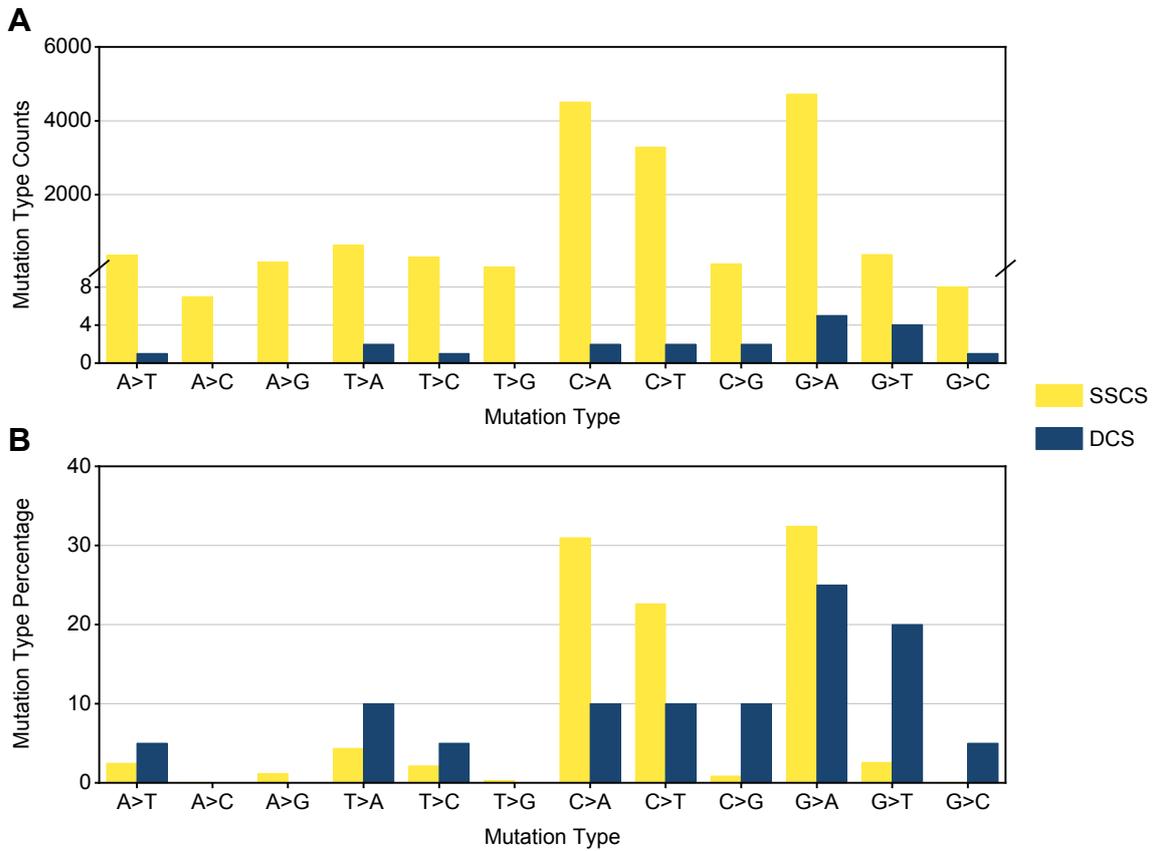
10 Stoler N, Arbeithuber B, Guiblet W, Makova KD, Nekrutenko A. 2016. Streamlined analysis of  
11 duplex sequencing data with Du Novo. *Genome Biol* **17**: 180.

12 Stoler N, Arbeithuber B, Povysil G, Heinzl M, Salazar R, Makova KD, Tiemann-Boege I,  
13 Nekrutenko A. 2020. Family reunion via error correction: an efficient analysis of duplex  
14 sequencing data. *Bmc Bioinformatics* **21**.

15 Tate JG, Bamford S, Jubb HC, Sondka Z, Beare DM, Bindal N, Boutselakis H, Cole CG, Creatore C,  
16 Dawson E et al. 2018. COSMIC: the Catalogue Of Somatic Mutations In Cancer. *Nucleic*  
17 *Acids Research* **47**: D941-D947.

18 Yang H, Wang K. 2015. Genomic variant annotation and prioritization with ANNOVAR and  
19 wANNOVAR. *Nature Protocols* **10**: 1556-1566.  
20  
21

# 1 Supplemental Figures

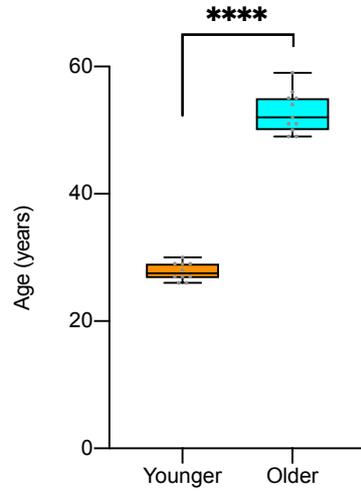


2

3 **Supplemental Figure S1.** Types of substitutions found in the SSCS and in the DCS. (A) Absolute number of counts of each mutation  
 4 type found in libraries "FGFR3 Up O Oct19 Re-seq" and "FGFR3 Down O BAT" using SSCS (in yellow) and DCS (in blue) data. (B)  
 5 Percentage of each mutation type found in the same libraries using SSCS data (n=14552) and DCS data (n=15).

6

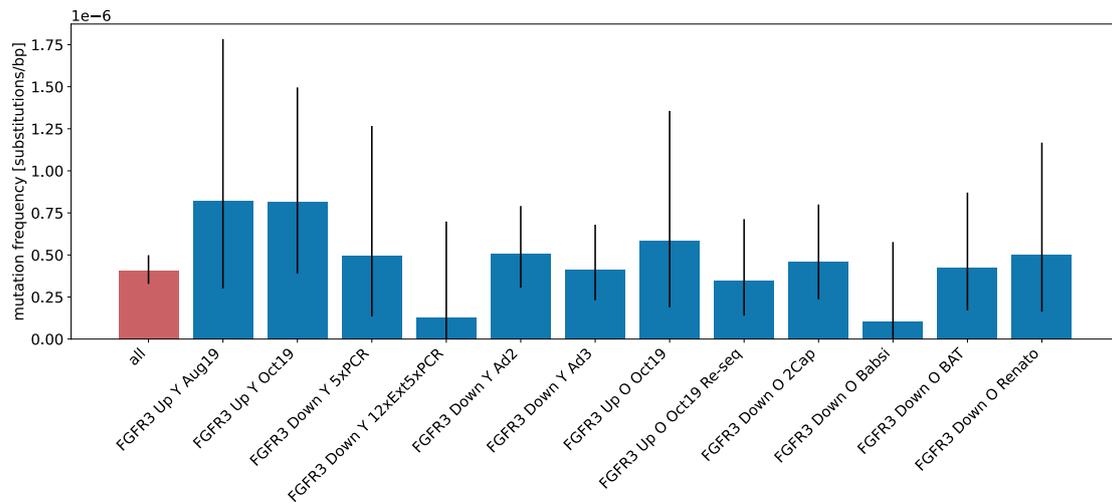
1  
2



3  
4  
5  
6

**Supplemental Figure S2.** Age distribution of the donors used in the young and old sperm pool. Each pool was composed by 5 different donors. Pools are significantly different as estimated by the Mann-Whitney *U* test ( $P < 0.0001$ ).

1

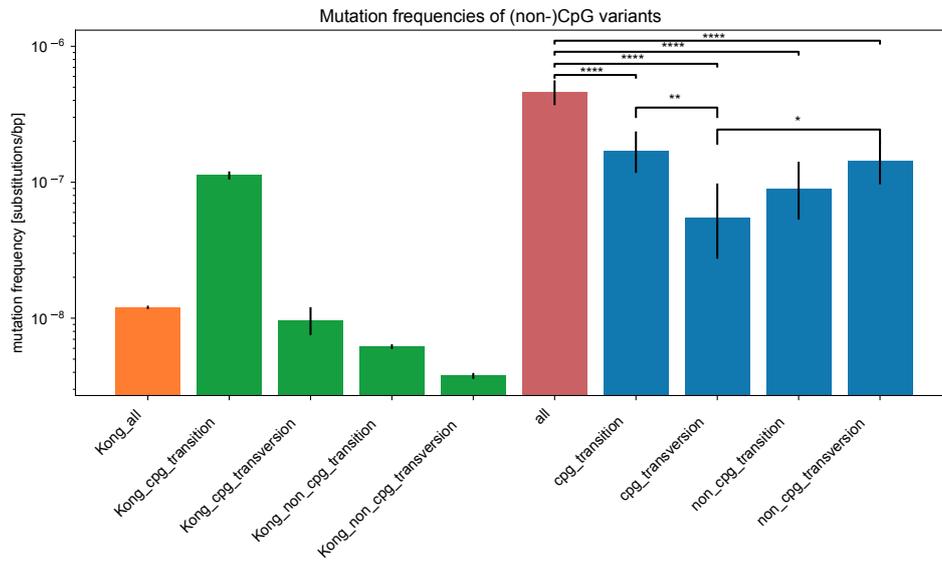


2

3 **Supplemental Figure S3.** Mutation frequencies of the different libraries estimated as the number of different variants divided by  
4 the number of total sequenced nucleotides (estimated by the mean coverage of the region multiplied with the size of the region)  
5 of that particular library. The mutation frequency of all variants (n=92) is compared to the individual libraries (FGFR3 Up Y Aug19  
6 n=6, FGFR3 Up Y Oct19 n=10, FGFR3 Down Y 5xPCR n=4, FGFR3 Down Y 12xExt5xPCR n=1, FGFR3 Down Y Ad2 n=19, FGFR3 Down Y  
7 Ad3 n=15, FGFR3 Up O Oct19 n=5, FGFR3 Up O Oct19 Re-seq n=7, FGFR3 Down O 2Cap n=12, FGFR3 Down O Babsi n=1, FGFR3  
8 Down O BAT n=7, FGFR3 Down O Renato n=5). Error bars are confidence intervals of a Poisson distribution. Pairwise testing  
9 between the libraries is performed with Chi-square testing with Bonferroni-Holm correction and no significance difference is  
10 found.

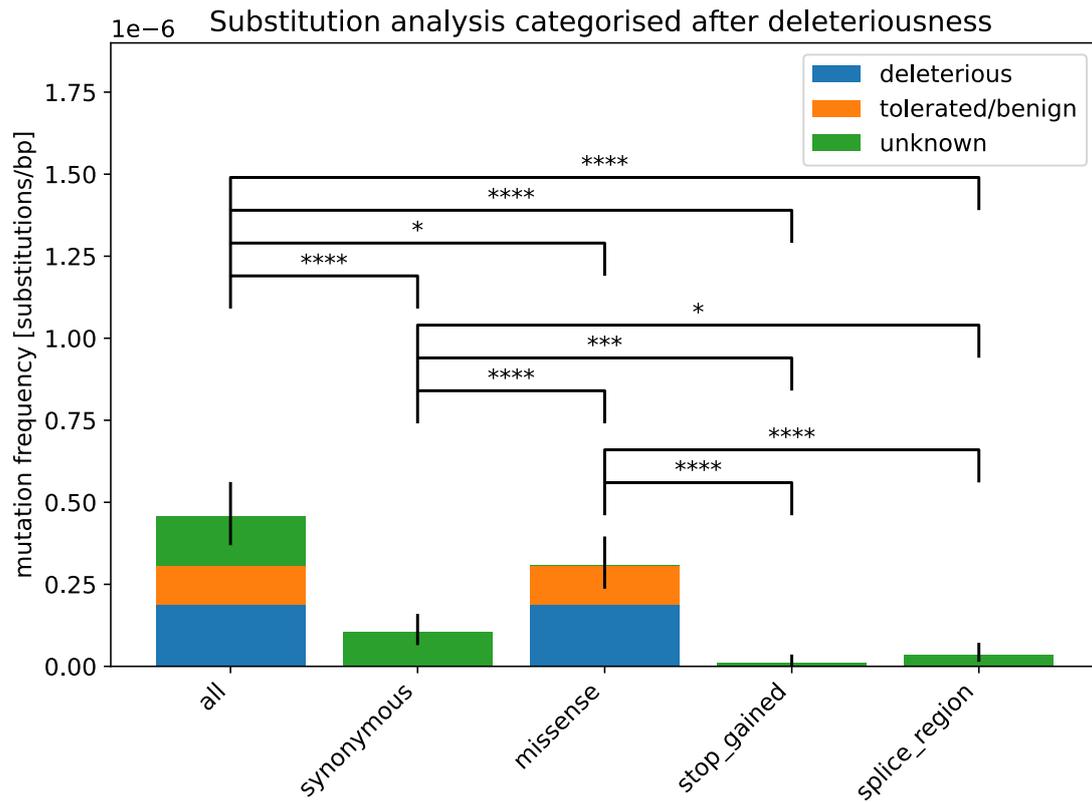
11

1  
2



3

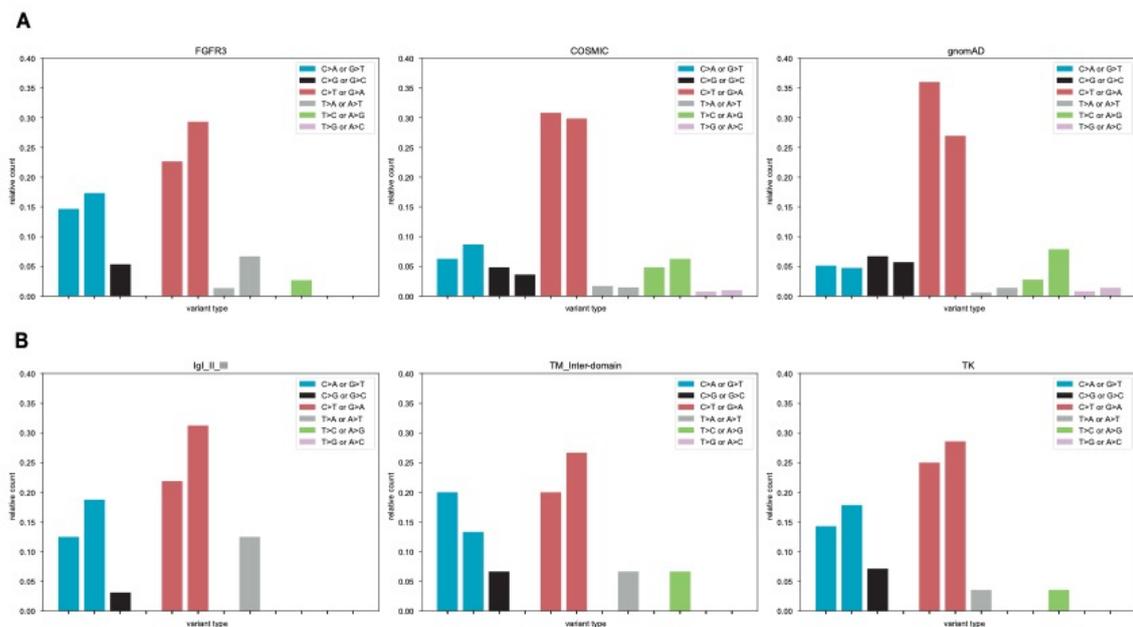
4 **Supplemental Figure S4.** Mutation frequencies of CpG and non-CpG variants (all n=92, CpG transitions n=34, CpG transversions  
5 n=11, non-CpG transitions n=18, non-CpG transversions n=29) compared to the pedigree sequencing of trios (Kong et al., 2012).  
6 Error bars are confidence intervals of a Poisson distribution. Pairwise testing is performed with Chi-square testing with  
7 Bonferroni-Holm correction and only significant differences are shown ((\* p-value < 0.05; (\*\*) p-value < 0.01; (\*\*\*) p-value <  
8 0.0001).



1

2 **Supplemental Figure S5.** Mutation frequencies of the different substitution types categorized after the deleteriousness based on  
 3 the SIFT score (all n=92, synonymous n=21, missense n=62, stop\_gained n=2, splice\_region n=7). Error bars are confidence  
 4 intervals of a Poisson distribution. Pairwise testing is performed with Chi-square testing with Bonferroni-Holm correction and only  
 5 significant differences are shown (\* p-value < 0.05; \*\* p-value < 0.01; (\*\*\*) p-value < 0.001; (\*\*\*\*) p-value < 0.0001).

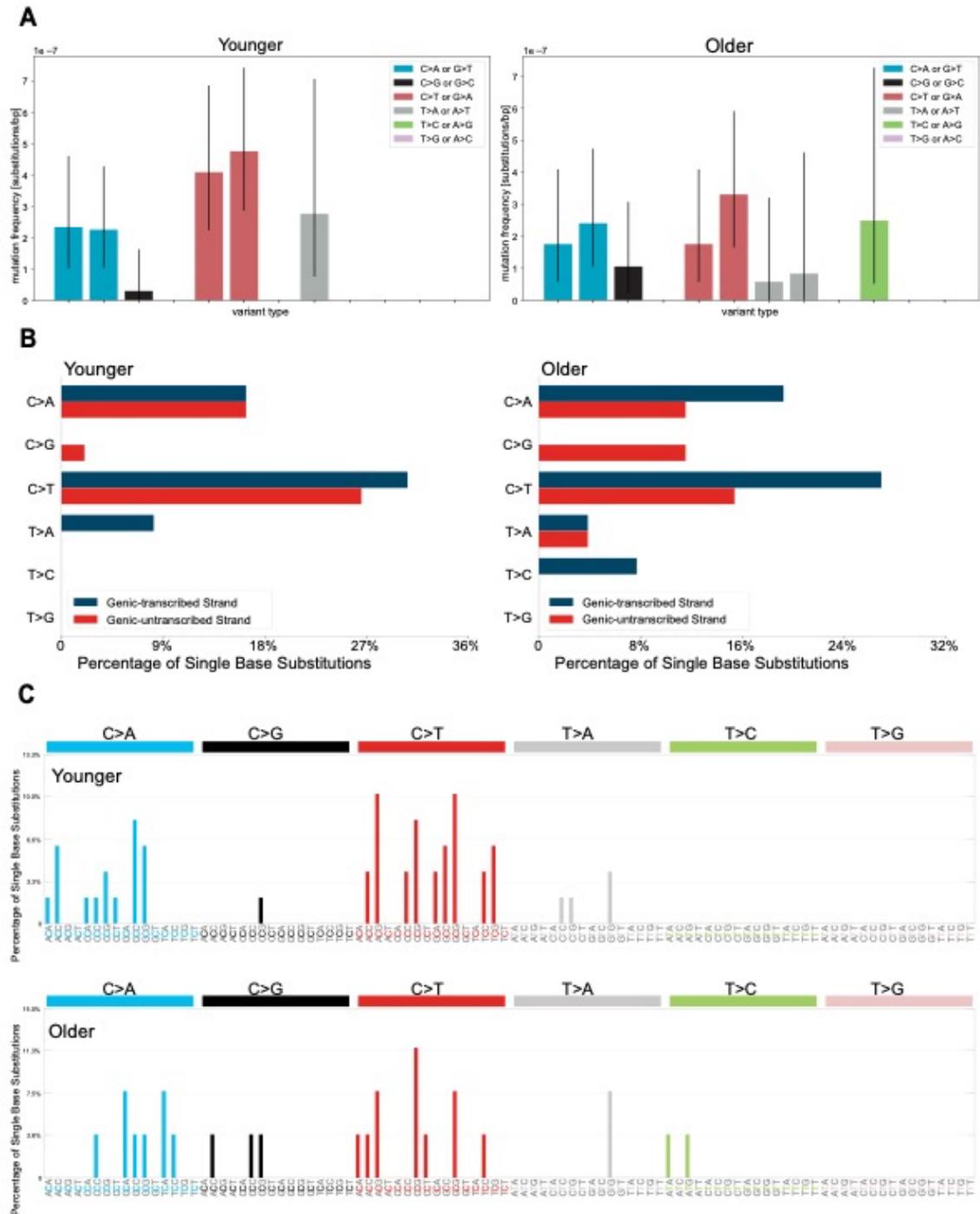




1  
2 **Supplemental Figure S7.** Mutational spectra of variants compared to the mutational spectra extracted from public databases. (A)  
3 Mutational spectra of all variants (n=75) compared to COSMIC v94 (n=415) and gnomAD v3.1.1 (n=508). (B) Mutational spectra for  
4 domains (Igl-III n=32, TM\_Inter-domain n=15, TK n=28). Note that the mutational spectra are based on the relative mutation  
5 count, therefore every variant is considered only once although they can be present in multiple libraires. The cosine similarities  
6 between the spectra can be found in Supplemental Table S7. For the mutational spectra of the databases COSMIC v94 and  
7 gnomAD v3.1.1, all variants except for indels of the targeted regions were downloaded. Next, we compared them to the spectra of  
8 the duplex sequencing data using the measure of relative counts.

9  
10

1



2

3

4

5

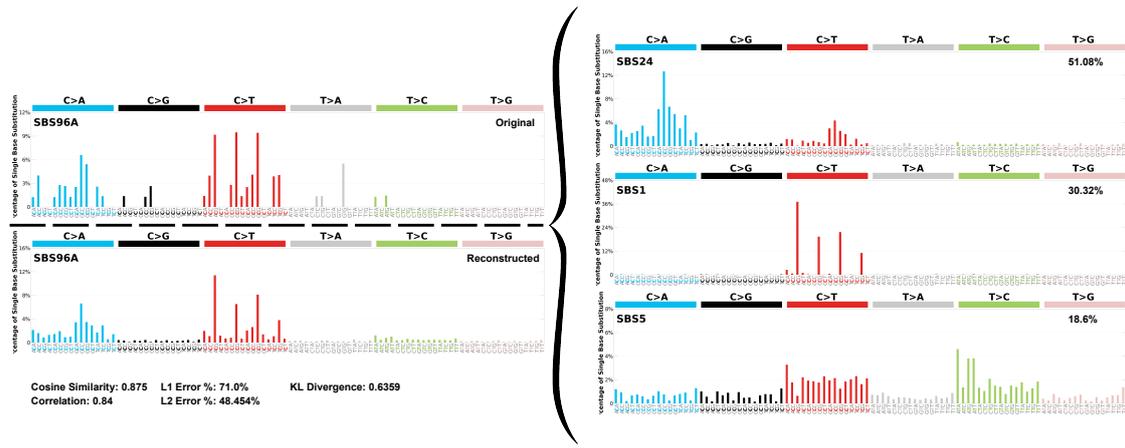
6

7

8

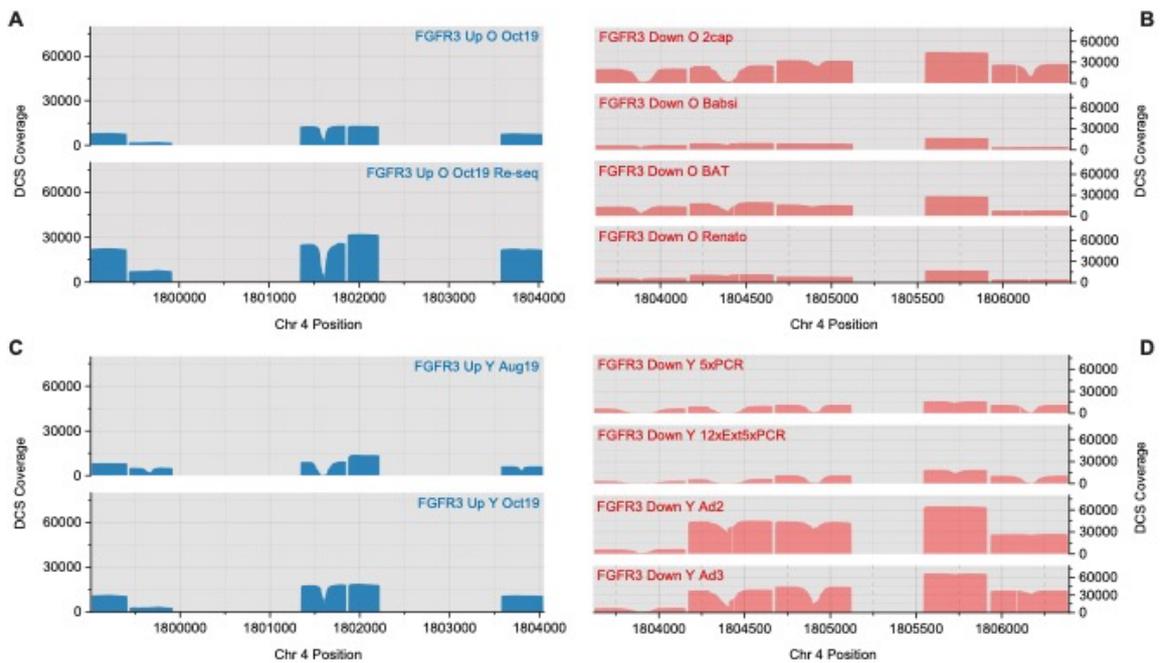
9

**Supplemental Figure S8.** Mutational spectra of the Younger and Older group (A) based on mutation frequencies (Younger n=55, Older n=37) estimated as the mutation count divided by all sequenced nucleotides. The cosine similarities between the groups can be found in Table S13. (B) Mutational spectra categorized by the transcribed and un-transcribed strand. (C) Mutational signatures (includes the 5' and 3' adjacent nucleotide to the variant). Note that for (B, C) every variant is considered only once (Younger n=49, Older n=26) although they can be present in multiple libraires. This is required for creating the mutational spectra based on relative counts by the tools of Bergstrom et al.,2019 ((Bergstrom et al. 2019)).



1

2 **Supplemental Figure S9.** Mutational signature compared to a catalogue of signatures (version 3.2) in the COSMIC database. The  
 3 original signature can be explained by three signatures (SBS24, SBS1 and SBS5) where SBS24 (Aflatoxin exposure) contributes  
 4 51.08%, SBS1 (spontaneous or enzymatic deamination of 5-methylcytosine to thymine) 30.32% and SBS5 (unknown) 18.6% of the  
 5 variants.



1

2 **Supplemental Figure S10.** DCS Coverage of each library. (A) Libraries from younger donors that targeted regions up 1-5. (B)  
 3 Libraries from younger donors that targeted regions down 1-5. (C) Libraries from older donors that targeted regions up 1-5. (D)  
 4 Libraries from older donors that targeted regions down 1-5. The variation in DCS coverage among libraries can be explained by  
 5 differences in read number per library or modification in the library preparation (e.g., input DNA amount, targeted region,  
 6 adaptor, hybridization capture or PCR strategy; for more details see Supplemental Methods).

7

8

9

10

11

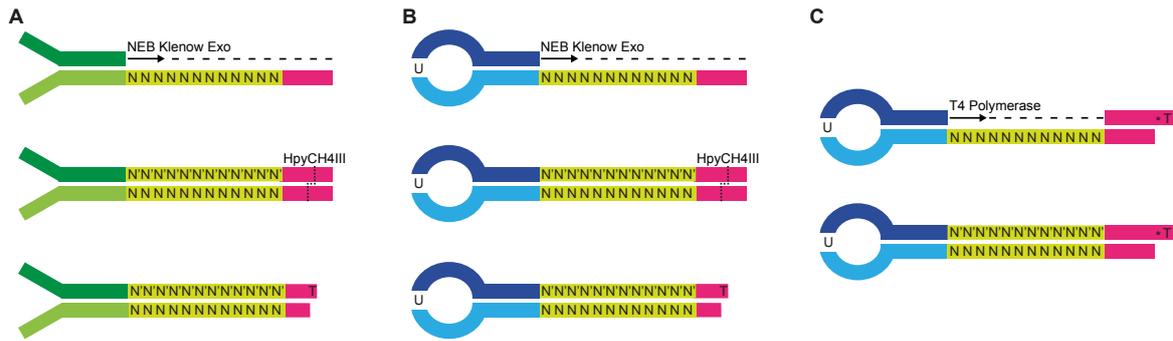
12

13

14

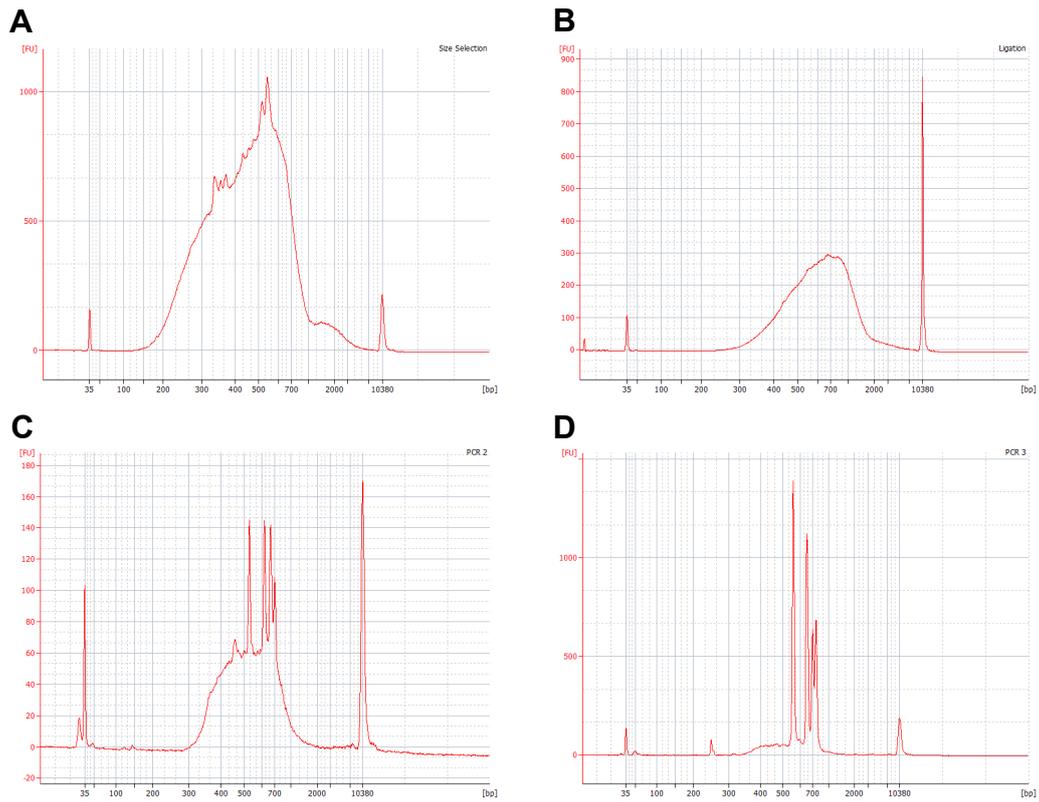
15

16



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16

**Supplemental Figure S11.** Different adapter designs for DS. (A) Duplex adapter synthesis and final structure as described by Kennedy et al. (Kennedy et al. 2014). (B) Adapter 2 synthesis and final structure. The use of a closed loop hairpin oligo instead of the annealing of two different oligos is the only difference to the Kennedy et al. design. (C) Adapter 3 synthesis and final structure. This adapter uses the same looped-hairpin oligo as adapter 2 but it adds a second oligo with a phosphorothioate bond at the 3' end just before the T-overhang. After the self-annealing of the hairpin oligo and the annealing of the oligo containing the T-overhang, T4 polymerase is used to extend the hairpin oligo and to create the complement of the random barcode.



1

2 **Supplemental Figure S12.** Example of fragment length distribution of different library preparation steps. (A) After double size  
 3 selection (~300 - 700 bp). (B) After ligation of adapters. (C) After the second PCR (performed after the first hybridization capture).  
 4 (D) Final fragment length distribution. These examples are from library 'FGFR3 Down' O Renato and were originated by the 2100  
 5 Expert software (Agilent)

6