

1 **Whole genome and exome sequencing reference datasets from a multi-**
2 **center and cross-platform benchmark study**

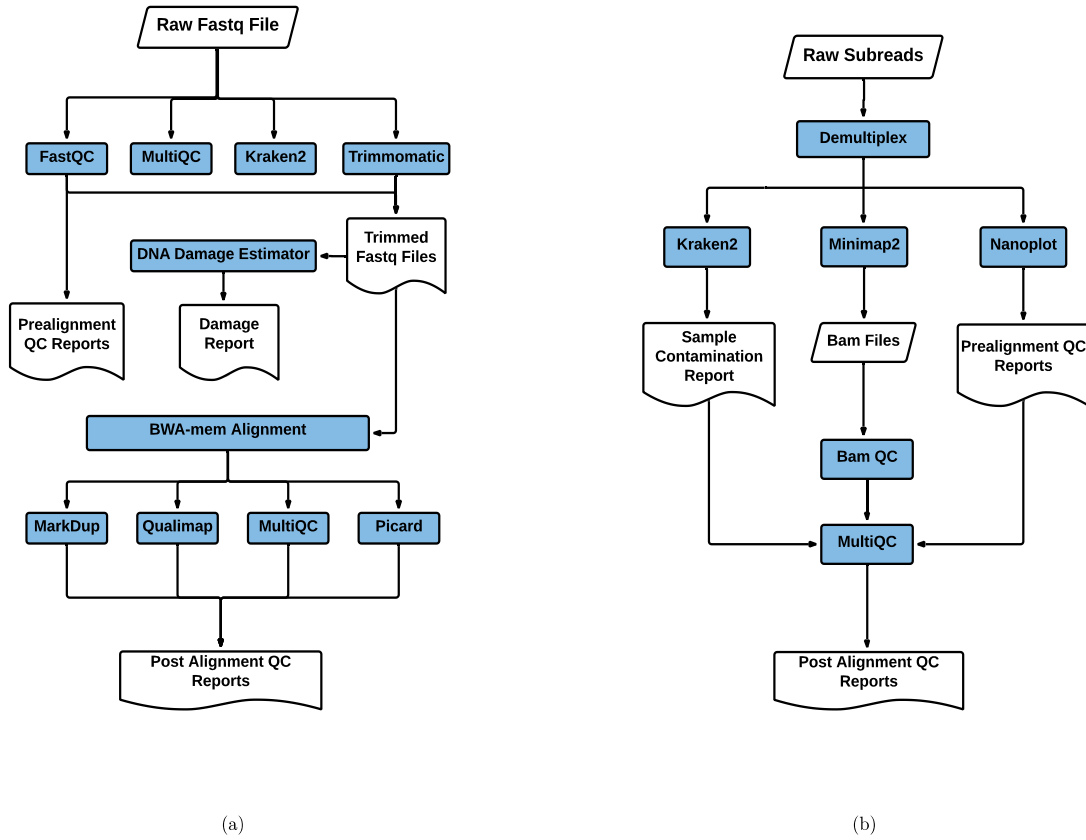
3
4 Yongmei Zhao^{1#}, Li Tai Fang², Tsai-wei Shen¹, Sulbha Choudhari¹, Keyur Talsania¹, Xiongfong
5 Chen¹, Jyoti Shetty³, Yuliya Kriga³, Bao Tran³, Bin Zhu⁴, Zhong Chen⁵, Wanqiu Chen⁵, Charles
6 Wang⁵, Erich Jaeger⁶, Daoud Meerzaman⁷, Charles Lu⁸, Kenneth Idler⁸, Luyao Ren⁹, Yuanting
7 Zheng⁹, Leming Shi⁹, Virginie Petitjean¹⁰, Marc Sultan¹⁰, Tiffany Hung¹¹, Eric Peters¹¹ Jiri
8 Drabek^{12,13}, Petr Vojta^{12,13}, Roberta Maestro^{13,14}, Daniela Gasparotto^{13,14}, Sulev Köks^{13,15,16}, Ene
9 Reimann^{13,17}, Andreas Scherer^{13,18}, Jessica Nordlund^{13,19}, Ulrika Liljedahl^{13,19}, Jonathan Foox²⁰,
10 Christopher E. Mason²⁰, Chunlin Xiao²¹, Huixiao Hong²², Wenming Xiao^{23#}

11
12 **TABLE OF CONTENTS**

13
14 **Supplementary Figures** Page 2
15 **Suppl. Figure 1:** Preprocessing and QC analysis pipelines Page 2
16 **Suppl. Figure 2:** WGS and WES cross-site data quality metrics Page 3
17 **Suppl. Figure 3:** GC coverage bias Page 3
18 **Suppl. Figure 4:** Cumulative genome coverage for WGS and WES Page 4
19 **Suppl. Figure 5:** HCC1395 tumor genome ploidy and heterogeneity Page 5
20 **Suppl. Figure 6:** Source of DNA damage artifacts Page 6
21 **Suppl. Figure 7:** Mutation calling repeatability and O_Score distribution Page 7
22

23
24
25

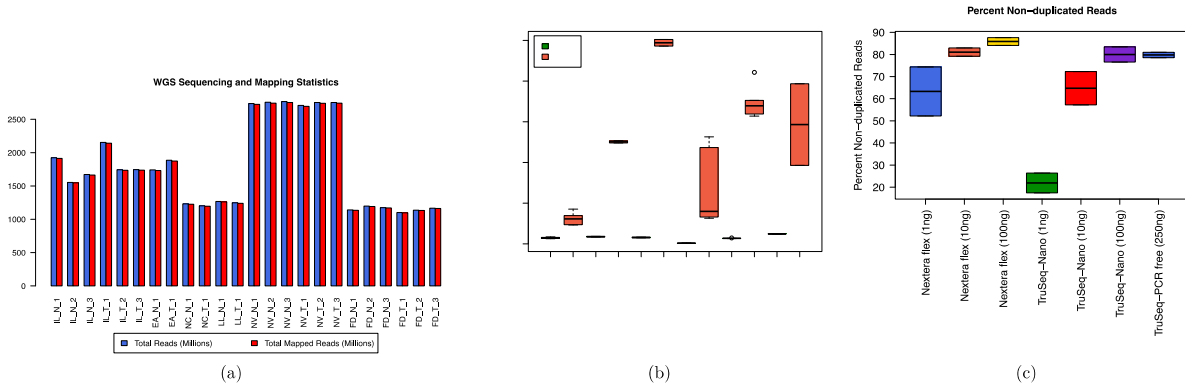
Supplementary Figures



26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43

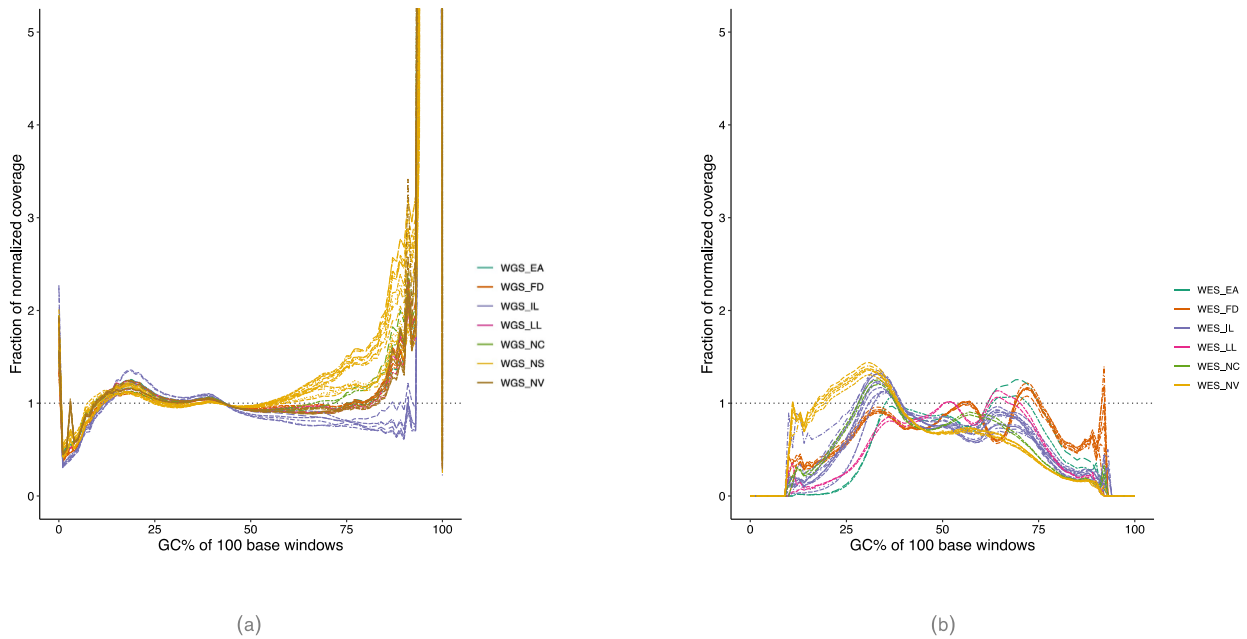
Suppl. Figure 1: Preprocessing and QC analysis pipelines used for whole genome and whole exome-seq data analysis **(a)** Short-read Illumina sequencing QC pipeline. **(b)** PacBio long-read sequencing QC pipeline.

44
45



46
47
48
49
50
51
52
53

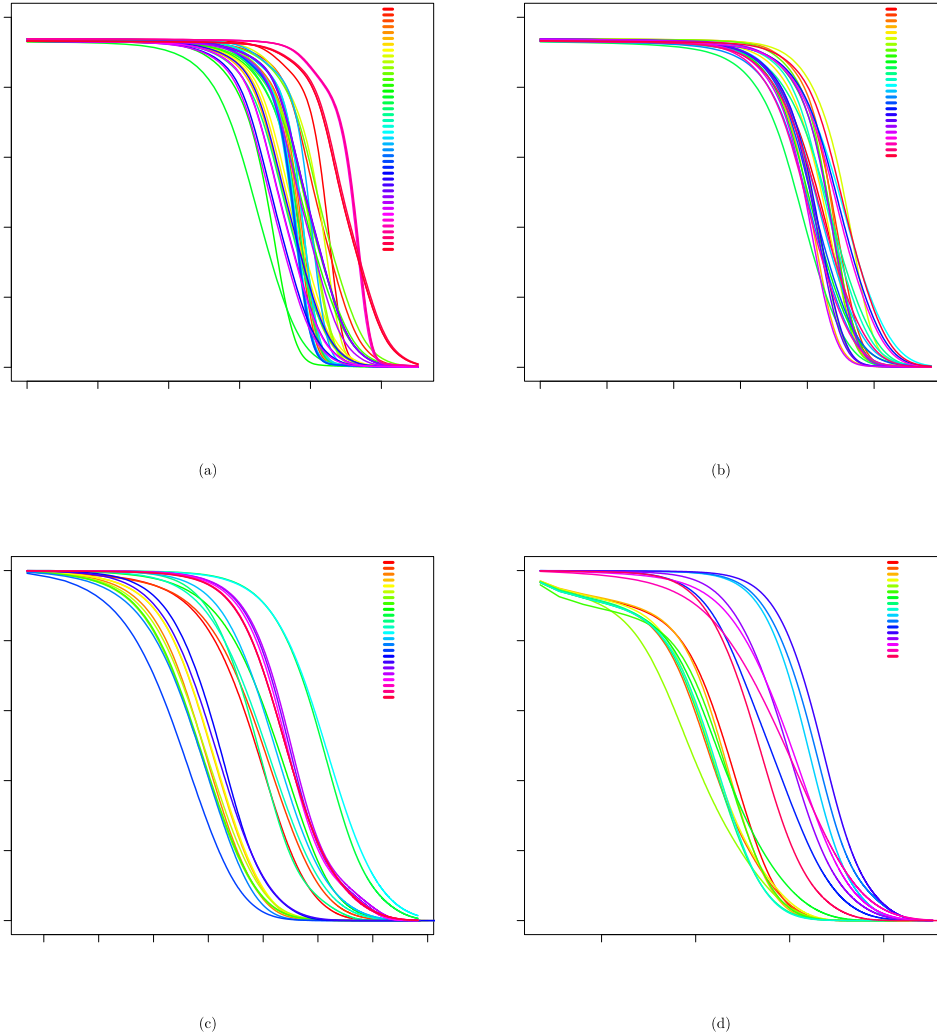
Suppl. Figure 2: WGS and WES cross-site data quality metrics **(a)** WGS cross-site sequencing yields (Millions) mapped reads (Millions) statistics. **(b)** Adapter contents in WGS and WES Illumina short-read data set across 6 data centers. **(c)** None duplicate mapped reads in Nextera Flex, TruSeq Nano and TruSeq PCR free protocols with different input amount.



54
55
56
57
58
59
60
61

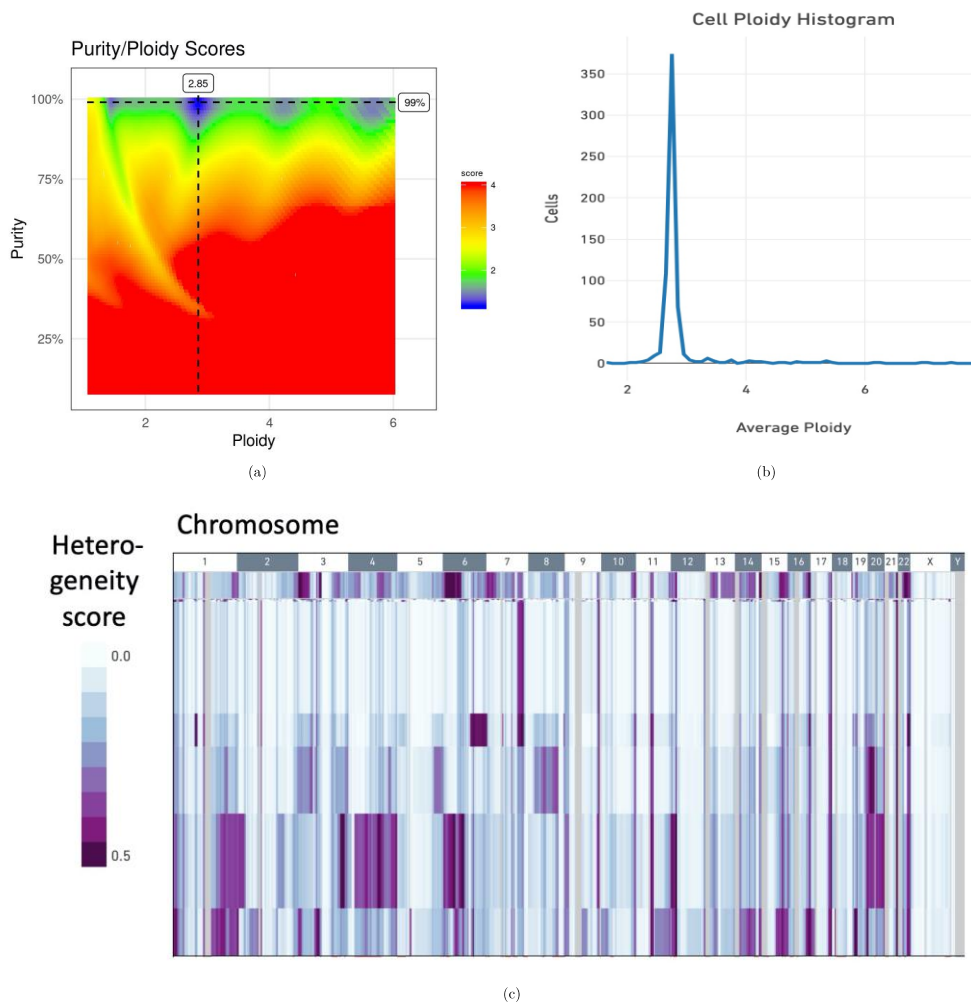
Suppl. Figure 3: GC Coverage Bias - using the Picard CollectGCBiasMetrics to measure the relative GC against the sequencing coverage to show bias in coverage across regions of the genome with varying GC content. A perfect library would be a flat line at $y = 1$. The plot X-axis denotes the GC content bin from the corresponding reference sequence ranging from 0 - 100%. Y-axis denotes the normalized coverage measurement of sequence depth for the particular GC bin for WGS data **(a)** and WES data **(b)**.

62
63
64



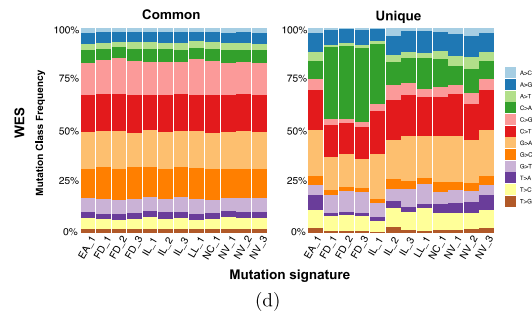
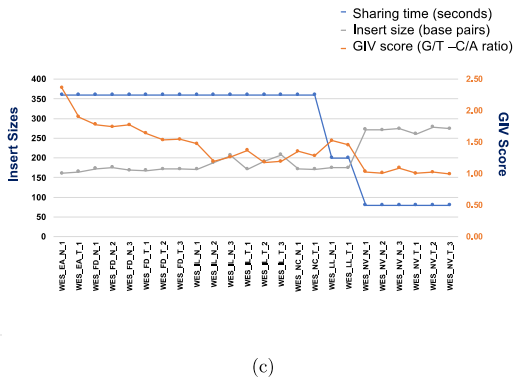
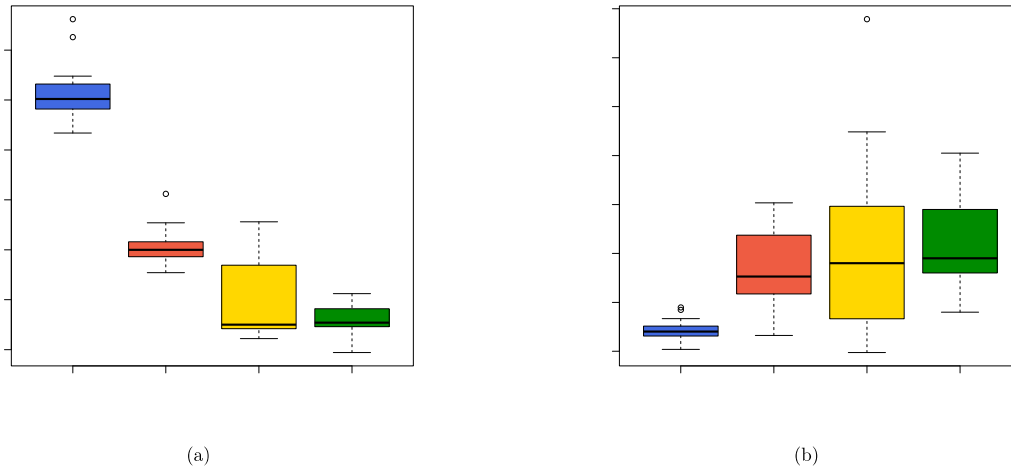
65
66
67
68
69
70
71
72
73
74
75

Suppl. Figure 4: Cumulative genome coverage for WGS and WES cross-site data sets and FFPE WGS and FFPE WES data sets. The graph displays the percentages of the reference genome with at least the given depth of coverage in log scale for each sample **(a)** Genome coverage for each fresh DNA prepared Illumina WGS libraries for cross-site comparison. **(b)** Genome coverage for FFPE WGS libraries. **(c)** Genome coverage for each Fresh DNA prepared Agilent SureSelect V6+UTR exome capture libraries for cross-site comparison. **(d)** Genome coverage for FFPE whole exome capture libraries.



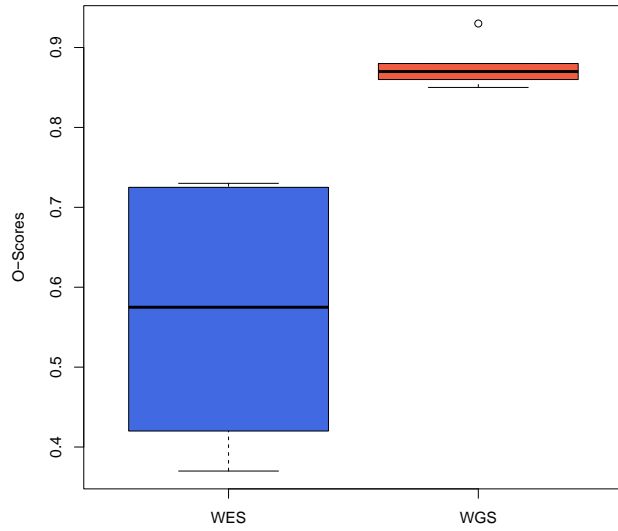
76
77
78
79
80
81
82
83
84
85
86
87
88

Suppl. Figure 5: HCC1395 tumor genome ploidy and heterogeneity measured from whole genome Illumina sequencing data and single cell CNV data. **(a)** Sample purity and ploidy estimated based on WGS for HCC1395 from Purple software, the tumor purity is above 99% with ploidy of 2.85. **(b)** Using 10X Genomics Single Cell CNV Solution, based on the analysis of 1270 cells for HCC1395 from 10x Single Cell CNV data set, Cellranger ploidy histogram displayed the vast majority of cells have ploidy of 2.8 as shown. **(c)** Heterogeneity analysis from 10x Single Cell CNV data. Each row represents a cell being sequenced. Integer-scaled CNA profiles across the genome of 1270 HCC1395 cells were obtained. Similar cells were clustered together based on CNAs. Subclonal populations are marked in tracks. The chromosome-scale gains showed in darker purple, and losses displayed as light color in heatmap.



90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107

Suppl. Figure 6: Source of DNA damage artifacts and effect on cancer genome mutation calls. One type of DNA damage artifacts is introduced during DNA fragment mechanical shearing. Longer library inserts have lower DNA damage scores (GIV scores) as shown in **(a)** library insert sizes for WGS cross-site libraries, WGS FFPE libraries, WES cross-site libraries, and WES FFPE libraries. **(b)** GIV score across different data sets; **(c)** correlation between library shearing time, insert sizes, and GIV scores. **(d)** percentage of mutation types for WES. The shared mutation across replicates is shown in left plot or sample specific unique mutations are shown in right plot. Note, high percentage of C/A mutation in the sites that also have high G/T_C/A GIV scores as displayed in **Suppl. Figure 6c**.



108
109
110
111

Suppl. Figure 7: Mutation calling repeatability and O_Score distribution for 12 WES and WGS runs. One-tailed t-tests of WES and WGS two groups has P-value of 3.32354E-07.