



Published in final edited form as:

Nature. 2020 July ; 583(7815): 259–264. doi:10.1038/s41586-020-2347-0.

Insights about variation in meiosis from 31,228 human sperm genomes

Avery Davis Bell^{1,2,*}, Curtis J. Mello^{1,2}, James Nemesh^{1,2}, Sara A. Brumbaugh^{1,2}, Alec Wysoker^{1,2}, Steven A. McCarroll^{1,2,*}

¹Department of Genetics, Harvard Medical School

²Program in Medical and Population Genetics, Broad Institute of MIT and Harvard

Meiosis, while critical for reproduction, is also variable and error-prone: crossover rates vary among gametes, between the sexes, and among humans of the same sex, and chromosome mis-segregation leads to aneuploidy^{1–8}. To study diverse meiotic outcomes and how they co-vary across chromosomes, gametes, and humans, we developed Sperm-seq, a way to simultaneously sequence the genomes of thousands of individual sperm. We analyzed the genomes of 31,228 human gametes from 20 sperm donors, identifying 813,122 crossovers and 787 aneuploid chromosomes. Sperm donors had aneuploidy rates ranging from 0.01 to 0.05 aneuploidies per gamete; crossovers partially protected chromosomes from nondisjunction at meiosis I. Some chromosomes and donors underwent more-frequent nondisjunction during the meiosis I cell division, while other chromosomes and donors showed more segregation failures during meiosis II; many genomic anomalies that could not be explained by simple nondisjunction also occurred. Diverse recombination phenotypes – from crossover rates to crossover location and separation (a measure of crossover interference) – co-varied strongly across individuals and cells. Our results can be incorporated with earlier observations into a unified model in which a core mechanism, the variable physical

Users may view, print, copy, and download text and data-mine the content in such documents, for the purposes of academic research, subject always to the full Conditions of use:http://www.nature.com/authors/editorial_policies/license.html#terms

* Correspondence to: averydavisbell@gmail.com; mccarroll@genetics.med.harvard.edu. **Correspondence and requests for materials** should be addressed to S.A.M and A.D.B.

Author Contributions

A.D.B. and S.A.M. conceived and led the studies. A.D.B., S.A.M., and C.J.M. developed the experimental methods. A.D.B. and C.J.M. performed all experiments, generating all data. A.D.B. and S.A.M. designed the crossover and aneuploidy analysis strategies, and A.D.B. performed the crossover and aneuploidy analyses. A.D.B., J.N., and A.W. wrote the sequence and variant processing software, pipelines, and analytical methods. A.D.B. wrote the crossover calling and analysis software. A.D.B. and S.A.B. wrote the aneuploidy calling software. A.D.B. and S.A.M. wrote the manuscript with contributions from all authors.

Data Availability

Crossover and aneuploidy data (individual events and counts per donor and/or cell), including source data underlying Figs. 2, 3b-e and Extended Data Figs. 5–9, are available via Zenodo, <http://dx.doi.org/10.5281/zenodo.2581570>. Raw sequence data are available in the SRA via dbGaP for general research use upon application and approval (study accession number phs001887.v1.p1).

Code Availability

Analysis scripts and documentation are available via Zenodo, <http://dx.doi.org/10.5281/zenodo.2581595>.

Competing Interests

A.D.B. and S.A.M. are inventors on a United States Provisional Patent application (PCT/US2019/029427; applicant: President and Fellows of Harvard College) currently in PCT stage relating to droplet-based genomic DNA capture, amplification and sequencing that is capable of obtaining high-throughput single-cell sequence from individual mammalian cells, including sperm cells. A.D.B. is an occasional consultant for Ohana Biosciences since October 2019. The other authors declare no competing interests.

Supplementary Information (Supplementary Notes, Supplementary Discussion, and Supplementary methods) is included with this paper.

compaction of meiotic chromosomes, generates inter-individual and cell-to-cell variation in diverse meiotic phenotypes.

One way to learn about human meiosis has been to study how genomes are inherited across generations. Genotype data are available for millions of people and thousands of families; crossover locations are estimated from genomic segment sharing among relatives and linkage-disequilibrium patterns in populations^{2,4,7,9,10}. Although inheritance studies sample only the few gametes per individual that generate offspring, such analyses have revealed that average crossover number and crossover location associate with common variants at many genomic loci^{3-6,11,12}.

Another powerful approach to studying meiosis is to directly visualize meiotic processes in gametocytes, which has made it possible to see that homologous chromosomes usually begin synapsis (their physical connection) near their telomeres¹³⁻¹⁵; to observe double-strand breaks, a subset of which progress to crossovers, by monitoring proteins that bind to such breaks^{16,17}; and to detect adverse meiotic outcomes, such as chromosome mis-segregation^{18,19}. Studies based on such methods have revealed much cell-to-cell variation in features such as the physical compaction of meiotic chromosomes^{20,21}.

More recently, human meiotic phenotypes have been studied via genotyping or sequencing up to 100 gametes from one person, demonstrating that crossovers and aneuploidy can be ascertained from direct analysis of gamete genomes²²⁻²⁶. Despite these advances, it has not yet been possible to measure multiple meiotic phenotypes genome-wide in many individual gametes from many people.

Development of Sperm-seq

We developed a method (“Sperm-seq”) with which to sequence thousands of sperm genomes quickly and simultaneously (Fig. 1). A key challenge in developing Sperm-seq was to deliver thousands of molecularly accessible-but-intact sperm genomes to individual nanoliter-scale droplets in solution. Tightly compacted²⁷ sperm genomes are difficult to access enzymatically without loss of their DNA into solution; we accomplished this by decondensing sperm nuclei using reagents that mimic the molecules with which the egg gently unpacks the sperm pronucleus (Extended Data Fig. 1a-d). These sperm DNA “florets” were then encapsulated into droplets together with beads that delivered unique DNA barcodes for incorporation into each sperm’s genomic DNA; we modified three technologies so as to do this (Drop-seq²⁸, 10X Chromium Single Cell DNA, and 10X GemCode²⁹, which was used to generate the data in this study) (Extended Data Fig. 1e-f). We then developed, adapted, and integrated computational methods for determining the chromosomal phase of each donor’s sequence variants and for inferring the ploidy and crossovers of each chromosome in each cell.

We used this combination of molecular and computational approaches to analyze 31,228 sperm cells from 20 sperm donors (974–2,274 gametes per donor), sequencing a median of ~1% of the haploid genome of each cell (Extended Data Table 1). Deeper sequencing allows detection of ~10% of a gamete’s genome.

Sperm-seq enabled inference of donors' haplotypes along the full length of every chromosome: alleles from the same parental chromosome tend to appear in the same gametes, so the co-appearance patterns of alleles across many sperm enabled alleles to be assembled into chromosome-length haplotypes (Extended Data Fig. 2a, Methods). *In silico* simulations and comparisons to kilobase-scale haplotypes from population-based analyses indicated that Sperm-seq assigned alleles to haplotypes with 97.5–100% accuracy (Extended Data Fig. 2b,c, Supplementary Notes).

The phased haplotypes determined by Sperm-seq allowed us to identify cell “doublets” from the presence of both parental haplotypes at loci on multiple chromosomes (Extended Data Fig. 2d-f, Methods). We also identified surprising “bead doublets,” in which two beads' barcodes reported identical haplotypes genome-wide, through different SNPs, and thus appeared to have captured the same gamete genome (Extended Data Fig. 3a,b, Methods, Supplementary Methods). Bead doublets were useful for evaluating the replicability of Sperm-seq data and analyses (Extended Data Fig. 3c-e), which is usually impossible to do in inherently destructive single-cell sequencing.

Recombination rate in sperm donors, cells

We identified crossover (recombination) events in each cell as transitions between the parental haplotypes we had inferred analytically (Methods). We identified 813,122 crossovers in the 31,228 gamete genomes (Extended Data Table 1). Crossover locations were inferred with a median resolution of 240 kb, with 9,746 (1.2%) inferred within 10 kb (Extended Data Table 1, Supplementary Notes). Analysis of bead doublets indicated high accuracy of crossover inferences (Extended Data Fig. 3e). Estimates of crossover rate and location were robust to down-sampling to the same coverage in each cell (Extended Data Fig. 4, Supplementary Methods).

The 20 sperm donors' recombination rates ranged from 22.2–28.1 crossovers per cell, consistent with estimates from other methods^{3,5,6,10–12,24,26}, though with far more precision at the individual-donor level (95% confidence intervals of 22.0–22.4 to 27.9–28.4 crossovers per cell), due to the large number of gametes analyzed per donor (Extended Data Table 1, Extended Data Fig. 5a). Individuals with higher global crossover rates had more crossovers on average on each chromosome (Extended Data Fig. 5b). We generated genetic maps for each of the donors from their 25,839–62,110 observed crossovers; these maps were broadly concordant with a family-derived paternal genetic map⁶ (Extended Data Fig. 5c,d; Supplementary Notes and Supplementary Methods).

Much more variation was present at the single-cell level: cells routinely harbored 17 to 37 crossovers (1st and 99th percentiles, median across donors), with a standard deviation of 4.23 across cells (median across donors), vs. a standard deviation of 1.53 across donors' crossover rates. Among gametes from the same donor, gametes with fewer crossovers in half of their genome tended to have fewer crossovers in the other half of their genome (Pearson's $r = 0.09$, two-sided $p = 8 \times 10^{-54}$ with all gametes from all donors combined after within-donor normalization; Supplementary Notes). This relationship, predicted by earlier

observations in families⁵ and spermatocytes²¹, suggests that crossover number on each chromosome is partly shaped by factors that act nucleus-wide.

Crossover location and interference

All 20 donors shared a tendency to concentrate their crossovers in the same regions of the genome, with large concentrations of crossovers in distal regions, as expected from earlier analyses of families^{4,6,9,11,30}, and more modest shared enrichments in many centromere-proximal regions (Fig. 2a, Extended Data Fig. 6). Guided by these empirical patterns, we divided the genome into “crossover zones,” each bounded by local minima in crossover density (Extended Data Fig. 6b, Supplementary Methods). These zones are much larger-scale than fine-scale-sequence-driven crossover hotspots^{7,31–33}, which the spatial resolution of most crossover inferences was not well suited for analyzing.

Intriguingly, the crossover zones with the most variable usage across people were all adjacent to centromeres; individuals with high recombination rates used these zones much more frequently (Fig. 2a, Extended Data Fig. 6a; with simulated equal SNP coverage, Extended Data Fig. 4c,e). The relative usage of distal and proximal zones varied greatly among donors and was correlated with donors’ recombination rates (Extended Data Fig. 7). These results were robust to alternative definitions of “distal” vs. “proximal” (Extended Data Fig. 7c, Supplementary Notes).

Positive crossover interference causes crossovers in the same meiosis to be further apart than they would be if crossovers were independent events^{26,30,34,35}. The effect of crossover interference was visible in each of the 20 sperm donors (Extended Data Fig. 8, Supplementary Methods). Crossover separation varied greatly among sperm donors and correlated inversely with recombination rate (Extended Data Fig. 7b), results that were robust to chromosome composition and that applied similarly to same-arm and opposite-arm crossover pairs (Extended Data Fig. 7e,f, Supplementary Notes).

The extremely strong correlations of donors’ crossover rates with crossover locations and interference could arise from an underlying biological factor that coordinates these phenotypes, or could arise trivially from the fact that chromosomes with more crossovers would also tend to have crossovers more closely spaced and in more regions. To distinguish between these possibilities, we focused on data from the 180,738 chromosomes with exactly two crossovers (here called “two-crossover chromosomes”; Supplementary Notes). Even in this two-crossover chromosome analysis, distal-zone usage (Fig. 2b) and crossover separation (Fig. 2c) correlated strongly and negatively with genome-wide recombination rate (additional control analyses described in Supplementary Notes and Extended Data Fig. 7d,g,h). These relationships indicate that a donor’s crossover-location and crossover-spacing phenotypes reflect underlying biological factors that vary from person to person, as opposed to resulting indirectly from the number of crossovers on a chromosome.

To test whether this co-variation of diverse meiotic phenotypes also governs variation at the single-gamete level, we asked whether cells with more crossovers than the average for their donor also exhibit the same kinds of crossover-spacing and crossover-location phenotypes

that donors with high crossover rates do (Supplementary Methods). Indeed, two-crossover chromosomes from cells with more crossovers tended to have closer crossover spacing and increased relative use of non-distal zones (Fig. 2d,e, Extended Data Figure 7i,j; unnormalized results in Supplementary Notes). This result indicates that the correlated meiotic-outcome biases that distinguish people from one another also distinguish the gametes within each individual (Discussion).

Chromosome and sperm donor aneuploidy

Aneuploidy generally arises from a chromosome mis-segregation that yields two aneuploid cells: one in which that chromosome is absent (a loss), and one in which it is present in two copies (a gain). Among the 31,228 gametes, we found 787 whole-chromosome aneuploidies and 133 chromosome arm-scale gains and losses (2.5% and 0.4% of cells, respectively, Fig. 3a, Methods). All chromosomes and sperm donors were affected. The sex chromosomes and acrocentric chromosomes had the highest rates of aneuploidy, consistent with fluorescence *in situ* hybridization analysis-based estimates^{18,19} (Fig. 3b).

The 20 young (18–38-year-old) sperm donors, considered by clinical criteria to have normal-range sperm parameters, exhibited aneuploidy frequencies ranging from 0.010 to 0.046 aneuploidy events per cell (Fig. 3c, Extended Data Table 1). Permutation tests indicated that this 4.5-fold variation in observed aneuploidy rates reflected genuine inter-individual variation (one-sided $p < 0.0001$, Supplementary Notes).

Under the prevailing model for the origins of aneuploidy, sperm with chromosome losses and gains should be equally common. However, we observed 2.4-fold more chromosome losses than chromosome gains (554 losses vs. 233 gains, proportion test two-sided $p = 2 \times 10^{-30}$). This asymmetry did not appear to reflect technical ascertainment bias (Extended Data Fig. 9a, Supplementary Notes). This surprising result is further considered in the Supplementary Discussion.

Errors in chromosome segregation can occur at meiosis I (MI), when homologs generally separate, or at meiosis II (MII), when sister chromatids separate. Because recombination occurs in MI prior to disjunction but does not occur at centromeres, errors during MI result in chromosomes with different (homologous) haplotypes at their centromeres, whereas sister chromatids nondisjoined in MII have the same (sister) haplotype at their centromeres (Fig. 3a). (Sex chromosomes X and Y disjoin in MI, and the sister chromatids of X and Y disjoin at MII.) Encouragingly, for chromosome 21 – the principal chromosome for which earlier estimates were possible – our finding of 33% MI events and 67% MII events matched previous estimates from trisomy 21 patients with paternal-origin gains³⁶.

Across all chromosomes, MI gains and MII gains had very different relative frequencies in different individuals and on different chromosomes (Fig. 3d,e). For example, sex chromosomes were 2.2 times more likely to be affected in MI than MII, whereas autosomes were 2.0 times more likely to be affected in MII than MI (proportion test two-sided $p = 1.3 \times 10^{-6}$). The lack of correlation between MI and MII vulnerabilities (Fig. 3d,e) indicated that MI and MII are differentially challenging to different chromosomes and to different people.

Although crossovers are required for proper chromosomal segregation³⁷ and seem protective against nondisjunction in maternal meiosis, in which chromosomes are maintained in diplotene of meiosis I for decades⁸, the relationship of crossovers to aneuploidy is less clear in paternal meiosis^{24,36,38–41}. We found that chromosome gains originating in MI – when recombination occurs – had 36% fewer total crossovers than matched, well-segregated chromosomes did (Supplementary Methods), suggesting that crossovers protected against MI nondisjunction of the chromosomes on which they occurred (Extended Data Fig. 9b, Supplementary Notes). No similar relationship was observed for MII gains (though the simulated control distribution for MII is inherently less accurate, Supplementary Notes) or at other levels of aggregation (Extended Data Fig. 9b-d, Supplementary Notes).

Other chromosome-scale genomic anomalies

Many sperm had complex patterns of aneuploidy that could not be explained by the canonical single-chromosome mis-segregation event. We detected 19 gametes that had three, instead of one, copies of entire or nearly entire chromosomes (2, 15, 20, and 21) (Fig 3f, Extended Data Fig. 10a,b). Chromosome 15 was particularly likely to be present in two extra copies; in fact, sperm with three copies of all or most of chromosome 15 ($n = 10$) outnumbered sperm with two copies of chromosome 15 ($n = 2$) (Fisher's exact test vs. Poisson two-sided $p = 2 \times 10^{-7}$, Supplementary Notes).

Other gametes carried anomalies encompassing incomplete chromosomes. These included: one cell that gained the *p* arm of chromosome 4 while losing the *q* arm; cells with gains of two copies of a chromosome arm; and cells with losses of chromosome arms (Fig. 3f; Extended Data Fig. 10c,d). One cell carried at least eight copies of most of the *q* arm of chromosome 4 (Fig. 3f). This gamete – which we estimate contained almost a billion base pairs of extra DNA – carried both parental haplotypes of chromosome 4, though almost all of the ~8 copies came from just one of the parental haplotypes (93% of observed alleles in the amplified region were haplotype 2). Diverse mutational processes likely generate these genomic anomalies (Supplementary Discussion).

Discussion

Inter-individual variation in crossover rates has previously been inferred from SNP data from families^{2–7,9–12}. Here, highly parallel single-gamete sequencing revealed that donors with high crossover rates also exhibit closer crossover spacing, even when controlling for the number of crossovers actually made on a chromosome. Based on these analyses, we consider it most likely that inter-individual variation in crossover interference is the true driver of variation in crossover rate and placement.

These same constellations of correlated meiotic crossover phenotypes – low interference, high rates, use of centromere-proximal zones – tended to characterize the same gametes from any donor. Cells with more crossovers in half of their genome tended to have more crossovers in the other half, tended to have made consecutive pairs of crossovers closer together in genomic distance – even when making just two crossovers on a chromosome –

and tended to have placed proportionally more of their crossovers in non-distal chromosomal regions.

What could cause these meiotic phenotypes to covary across chromosomes, in individual cells, and among people? The physical length of chromosomes during meiosis, which reflects their compaction, has been observed to vary up to two-fold among individual spermatocytes while being strongly correlated across chromosomes in the same spermatocyte; spermatocytes with more-compacted chromosomes also generally have fewer incipient crossovers^{20,21,42}. A unifying model (Extended Data Fig. 11) explains the covariance of these meiotic phenotypes while providing a candidate mechanism for inter-individual variation: cell-to-cell variation in the compaction of meiotic chromosomes – and person-to-person variation in the average degree of this compaction – would cause these phenotypes to co-vary in the manner observed in Fig. 2b-e.

Our enthusiasm about this model relies on multiple additional earlier observations (Extended Data Fig. 11). Firstly, at a cellular level, crossover interference occurs as a function of physical (micron) distance along the meiotic chromosome axis or synaptonemal complex rather than as a function of genomic (base pair) distance^{43–45}. Secondly, the first crossover on a chromosome is more likely to occur distally^{13–15}. Such a model also predicts a shared mechanism for sex differences in recombination rates and inter-individual variation among individuals of the same sex: oocytes have a longer synaptonemal complex, more crossovers, and decreased crossover interference as measured in genomic distances than spermatocytes, but have the same synaptonemal complex length extent of crossover interference^{22,42,46,47}.

Human genetics research has revealed that recombination phenotypes are heritable and associate with common variants at many genomic loci^{3–6,11,12}. A recent genome-wide association study found that variation in crossover rate and placement is associated with variants near genes that encode components of the synaptonemal complex, which connects and compacts meiotic chromosomes, and with genes involved in the looping of homologs along the chromosome axis³. Our model predicts that inherited genetic variation at these loci may bias the average degree of compaction of meiotic chromosomes; the fact that this same property varies among cells from the same donor^{20,21} shows that variance is well-tolerated and compatible with diverse-but-successful meiotic outcomes.

The sharing of co-varying phenotypes between the single-cell and person-to-person levels suggests that a core biological mechanism shapes both inter- and intra-individual (single-cell) variation in meiotic outcomes. Such parallelisms between cell-biological and human-biological variation could in principle exist in a wide variety of biological contexts.

Methods

A companion protocol for generating single-sperm libraries using the methods presented here is available via Protocol Exchange⁴⁸. Custom scripts (available via Zenodo⁴⁹) are referenced by name in the sections describing analyses they perform. Recombination and aneuploidy data generated via the described methods are also publicly available⁵⁰. All

statistical analyses were performed in R unless otherwise noted. Details on further analysis methods are provided in the Supplementary Methods.

Sample information

Sperm samples from 20 anonymous, karyotypically normal sperm donors were obtained from New England Cryogenic Center under a Not Human Subjects determination from the Harvard Faculty of Medicine Office of Human Research Administration (protocols M23743–101 and IRB16–0834). Donors consented at the time of initial donation for samples to be used for research purposes. The Not Human Subjects determination was based on the use of discarded biospecimens that had been consented for research and the fact that researchers had no interactions with the biospecimen donors and no access to identifiable information about the biospecimens. The reviewing committee also reviewed and approved our deposition of the data into an NIH repository. All experiments were performed in accordance with all relevant guidelines and regulations. (Specimens can be obtained from New England Cryogenic Center upon IRB approval.)

Samples arrived in liquid nitrogen in “egg yolk buffer” or “standard buffer with glycerol” (no further buffer information provided) and were aliquoted and stored in liquid nitrogen in the same buffers.

Per sperm bank policy, donors are 18–38 years old at the time of donation and precise age of donors is not released. Donor identifiers used in the paper were created specifically for this study and are not linked to any external identifiers.

ddPCR to evaluate genome accessibility

To evaluate how often regions from two different chromosomes co-occurred (as would be expected from cells), we performed droplet digital PCR with naked DNA, untreated sperm cells, or sperm cells decondensed as described subsequently but with variable heat incubation times. For each assay targeting each chromosome, a 20× assay mix was created by combining 25.2 μL of 100 μM forward primer (IDT), 25.2 μL of 100 μM reverse primer (IDT), and 7 μL of 100 μM probe (IDT for FAM-labelled probes, Life Technologies for VIC-labelled probes) with 82.6 μL ultrapure water. ddPCR was performed as described previously⁵¹, following section 3.2 steps 4–12, but with untreated sperm or sperm DNA florets as input instead of DNA.

For this analysis, chromosome 7 was targeted with an assay to intergenic region chr7: 106552149–106552176 (hg38); forward primer sequence: CGTAATGGGGCACAGGGATATA; reverse primer sequence: CTGTGAGAGGTAGAGAATCGCC; probe sequence: CACAGAGTCCATTTGCAGCACCTCAGT; probe fluorophore: FAM. Chromosome 10 was targeted with an assay to *RPP30* at chr10:92631759–92631820; forward primer sequence: GATTGGACCTGCGAGCG; reverse primer sequence: GCGGCTGTCTCCACAAGT; probe sequence: CTGACCTGAAGGCTCT; probe fluorophore: VIC.

We calculated the percentage of molecules expected to be linked from each reaction following Regan *et al.*⁵².

Sperm cell library generation

Accessible sperm nuclei “florets” were generated using a combination of published decondensation protocols^{53,54} with some modifications. Sperm aliquots containing >200,000 cells were thawed on ice and then washed by spinning for 10 minutes at 400 g at 4°C. The pellet was resuspended in 10 µL phosphate-buffered saline (PBS, Gibco/LifeTechnologies) and re-centrifuged under the same conditions. The sperm pellet was resuspended in 2.5 µL of a sucrose buffer containing 250 mM sucrose (Sigma), 5 mM MgCl₂ (Sigma), and 10 mM Tris HCl (pH 7.5, Thermo Scientific). Sperm aliquots were submerged in liquid nitrogen and immediately quick-thawed by holding them in a warm fist; three such freeze-thaw cycles were performed.

Freeze-thawed sperm solution was combined with 22.5 µL decondensation buffer (113 mM KCl [Sigma], 12.5 mM KH₂PO₄ [Sigma], 2.5 mM Na₂HPO₄ [Sigma], 2.5 mM MgCl₂ [Sigma], and 20 mM Tris [Thermo Scientific] freshly supplemented with 150 µM heparin [sodium salt from porcine, Sigma H3393] and 2 mM beta-mercaptoethanol [Sigma]). The reaction was incubated at 37°C for 45 minutes. To allow enzymatic DNA amplification, heparin was inactivated by mixing the sperm solution with 0.5 U heparinase I (Sigma H2519) by gently pipetting and incubating at room temperature for 2 hours⁵⁵.

The sperm solution was moved to ice, and sperm floret concentration was determined by diluting 1:100 with PBS and staining with 1X SYBR I (Thermo Scientific) and counting using the green fluorescence channel at 10x magnification.

Droplets were prepared using the following modifications to 10X Genomics’ GemCode (version 1²⁹) User Guide Revision C (in place of steps 5.1–5.3.9); all reagents come from the 10X Genomics GemCode kit. Ultrapure water was combined with 10,833 sperm to a final volume of 5 µL; 10,000 sperm were used for library generation. To each sperm sample was added 60 µL of a master mix containing 32.5 µL GemCode reagent mix, 1.5 µL primer release agent, 9.2 µL GemCode polymerase, and 16.8 µL ultrapure water.

GemCode beads were vortexed at full speed for 25 seconds, and then diluted 1:11 with ultrapure water to a total volume of at least 90 µL per sample. Per 10X’s protocol, 60 µL of sample-master mix combination was added to the droplet generation chip, followed by 85 µL of freshly pipette-mixed 1:11-diluted bead mixture and 150 µL of droplet generation oil.

Droplets were generated and processed through library generation following 10X Genomics’ GemCode (version 1) User Guide Revision C (step 5.3.10 through the end of section 6).

Sequencing and sequence data processing

Two libraries were generated per sperm donor and additional libraries were generated for four initial samples with low cell counts. Four or five libraries were sequenced at a time on S2 200 cycle flow cells on an Illumina NovaSeq. The read structure was 178 cycles read 1, 8 cycles read 2 (index read one), 14 cycles read 3 (index read two containing the cell barcode;

later treated as the reverse read), and 5 cycles read 4 (unused; included to fulfill the NovaSeq's paired-end requirement).

To convert the data to mapped BAM files with cell and molecular barcodes encoded as read tags, we used Picard Tools v2.2 (<http://broadinstitute.github.io/picard>) and Drop-seq Tools v2.2 (<https://github.com/broadinstitute/Drop-seq/releases>; see https://github.com/broadinstitute/Drop-seq/blob/master/doc/Drop-seq_Alignment_Cookbook.pdf for details on running many of the tools)²⁸:

Illumina BCL files were converted to unmapped BAM files using Picard's ExtractIlluminaBarcodes and IlluminaBasecallsToSam with read structure 178T8B14T (cell barcodes, present in the i5 index, were incorporated as read 2 for ease of downstream processing). BAMs were processed to include unique molecular identifiers (UMIs) and cell barcodes as read tags, and to exclude reads with poor-quality cell barcodes or UMIs; consequently, each read was retained as single-end with 14-bp cell barcode stored in tag XC and 10-bp molecular barcode/unique molecular identifier (UMI) stored in tag XM. The first 10 bp of read 1 were used as the UMI. First, DropSeq Tools' TagBamWithReadSequenceExtended was called with `BASE_RANGE=1-14`, `BASE_QUALITY=10`, `BARCODED_READ=2`, `DISCARD_READ=true`, `TAG_NAME=XC`, `NUM_BASES_BELOW_QUALITY=1`. Subsequently, TagBamWithReadSequenceExtended was called again with `BASE_RANGE=1-10`, `BASE_QUALITY=10`, `HARD_CLIP_BASES=true`, `BARCODED_READ=1`, `DISCARD_READ=false`, `TAG_NAME=XM`, `NUM_BASES_BELOW_QUALITY=1`. Finally, DropSeq Tools' FilterBAM was called with parameter `TAG_REJECT=XQ`.

Reads were aligned to hg38 using `bwa mem`⁵⁶ v0.7.7-r441. BAMs were converted to FastQ using Picard's SamToFastQ, FastQ reads were aligned using `bwa mem -M`, and then unmapped BAMs were merged with mapped BAMs using Picard's MergeBamAlignment, with non-default options `INCLUDE_SECONDARY_ALIGNMENTS=false` and `PAIRED_RUN=false`. Reads were marked PCR duplicates using Drop-seq Tools' SpermSeqMarkDuplicates (part of Drop-seq tools v2.2 and above) with options `STRATEGY=READ_POSITION`, `CELL_BARCODE_TAG=XC`, `MOLECULAR_BARCODE_TAG=XM`, `NUM_BARCODES=20000`, `CREATE_INDEX=true`. BAM files for all lanes and index sequences from the same sample were merged using Picard's MergeSamFiles prior to alignment and/or during duplicate marking with all BAMs given as input to SpermSeqMarkDuplicates.

Variant calling, sperm cell genotyping

For each donor, we pooled all reads from all libraries, including reads that did not derive from a barcode associated with a complete sperm cell. Using GATK v3.7^{57,58} in hg38, we followed GATK's best practices documentation for base quality score recalibration, gVCF generation using HaplotypeCaller (in DISCOVERY mode with `-stand_call_conf 20`), and joint genotyping with GenotypeGVCFs. We filtered variants with `SelectVariants -selectType SNP` and `VariantFiltration (--filterExpression "QD<3.0")`. We then performed VQSR following GATK's best practices, except that we excluded annotations MQ and DP (VariantRecalibrator with GATK provided resources; `-an QD, MQRankSum`,

ReadPosRankSum, FS, and SOR; -mode SNP; --trustAllPolymorphic; and tranches 90, 99.0, 99.5, 99.9, and 100.0). We applied tranche 99.9 recalibration using ApplyRecalibration -mode SNP and obtained the names of SNPs from dbSNP 146⁵⁹ using VariantAnnotator --dbsnp. We filtered our sites to contain only biallelic SNPs present in Hardy–Weinberg equilibrium in 1000 Genomes Phase 3⁶⁰ using SelectVariants --concordance with a VCF containing only these sites (from GATK’s resource bundle). We excluded SNPs in centromeric regions or acrocentric arms as defined by the UCSC Genome Browser’s cytoband track^{61,62} (<http://genome.ucsc.edu>; the same centromere boundaries were used in all analyses) and those in known paralogous regions as lifted over from Genovese et al 2014⁶³. We selected only heterozygous SNPs using SelectVariants -selectType SNP --selectTypeToExclude INDEL --restrictAllelesTo BIALLELIC --excludeFiltered --setFilteredGtToNocall --selectexpressions ‘vc.getGenotype(““““<sample name>““““).isHet()’.

We identified SNPs present in each sperm cell and which allele was present using GenotypeSperm (part of Drop-seq Tools v2.2 and above). For downstream analyses, we generated a file with columns cell, pos, and gt, with gt having the value 0 for the reference allele and 1 for the alternate allele for SNPs that had one or more UMIs covering only one base matching the reference or alternate allele. (See our script *gtypesperm2cellsbyrow.R*.)

Chromosome-scale phasing

We identified barcodes potentially associated with cells by plotting the cumulative fraction of reads associated with each ranked barcode and identifying the inflection point of this curve (see Extended Data Fig. 1f). We then included only barcodes with substantial read depth on either the X or the Y chromosome but not both, as the vast majority of sperm cells should contain only one sex chromosome. (We later added these barcodes back in before formally identifying and excluding cell doublets).

To phase sperm donors’ genomes, we used all quality-controlled heterozygous sites in these cell barcodes expected to correspond to sperm cells, excluding observations of SNPs where the observed allele was not the reference or alternate allele in the parental genome or where more than one allele was observed. For each chromosome, we converted per-cell SNP calls into “fragments” for input into the HapCUT phasing software^{64,65} by considering each consecutive pair of SNPs observed in a cell to be a fragment (see our script *gtypesperm2fmf.R*). We then used HapCUT with parameter --maxiter 100 to generate chromosomal phase. After identifying and removing cell doublets (see below), we repeated phasing with only non-doublet cell barcodes.

To validate our phasing method, we simulated single-cell SNP observations from known haplotypes, including 2% genotype errors and a variable percentage of cell doublets. Briefly, sites were randomly sampled from one known haplotype of chromosome 17 until a crossover location probabilistically assigned based on the deCODE recombination map⁶, then sampled from the other haplotype (one crossover was simulated per cell). To simulate PCR or sequencing errors, 2% of the sites were randomly assigned to an allele. Doublets were simulated by combining two cells and retaining 70% of the observed sites at random. We performed five random simulations for each doublet proportion, mean proportion of sites

“observed” in each cell, and number of cells simulated, and then followed our phasing protocol using each simulation. (See our script *simulatespermseqfromhaps.py*.)

To further validate phasing, we used Sperm-seq data to phase one donor’s genome and compared these phased haplotypes to this donor’s Eagle^{66,67}-generated haplotypes. We compared the phase relationship between each consecutive pair of SNPs (identifying the proportion of switch errors between the two phased sets). We also compared the Sperm-seq allele-allele phase of all pairs of alleles in perfect linkage disequilibrium in 1000 Genomes Phase 3⁶⁰ in the populations matching the donor’s ancestry.

Cell doublets

To identify cell barcodes associated with more than one sperm cell (cell doublets), we detected consecutively observed SNP alleles that appeared on different parental haplotypes, which could occur because of crossover, error, or the presence of two haplotypes in the same droplet (doublet). We ranked barcodes by the proportion of consecutive SNPs that spanned haplotypes using all SNPs from all autosomes except the autosome with the most haplotype-spanning consecutive SNPs (so as to avoid mistakenly identifying cells with chromosome gains as doublets); this resulted in a clear inflection point wherein cell doublets had a quickly accelerating proportion of haplotype-spanning consecutive SNPs (Extended Data Fig. 2d-f). All cell barcodes below this inflection point (identified with the function *ede* from the R package *inflection* <https://CRAN.R-project.org/package=inflection>) were considered non-doublet (Extended Data Fig. 2f). (See our script *computeSwitchesandInflThresh.R*.) Even though we exclude the autosome with the most haplotype-spanning consecutive SNPs from doublet identification, any cells with multiple chromosome gains (especially more than two) or whole-genome diploidy would be excluded by this method.

Crossover events

We identified crossover events on all autosomes (but excluded the *p* arms of acrocentric chromosomes where SNPs were excluded from analysis) by finding transitions between tracts of SNPs with alleles matching different parental haplotypes using a Hidden Markov Model written in R with package *HMM* (<https://CRAN.R-project.org/package=HMM>). To ensure that we detected crossovers located near the ends of SNP coverage (sub-telomeric regions are frequently used for crossovers in spermatogenesis), we ran the HMM both in the forward chromosomal and reverse-chromosomal directions, with start probability for one haplotype equal to 1 if the first two SNPs observed were of that haplotype. In addition to two states for parental haplotypes, we included a third “error” state to capture cases in which a haplotype 1 allele is observed in a haplotype 2 region (and vice versa), *e.g.*, due to PCR or sequencing error, gene conversion, or cases in which a small piece of off-haplotype ambient DNA was captured in a droplet. Crossovers were where one haplotype transitioned to another, or where one haplotype transitioned to the error state and then to the other haplotype. Crossover boundaries were the last SNP in the first haplotype and the first in the next. The key parameters for this algorithm are the transition probability between haplotypes (set to 0.001, from the per-cell median 26 crossovers divided by the per-cell median 24,710 heterozygous SNPs) and transition probability into and out of the “error” state (we set transition probability into this state to 0.03 from either haplotype, as only a few percent of

SNPs are off-haplotype; we set the probability of staying in error to 0.9 to allow for the occasional tract of SNPs from an ambient piece of off-haplotype DNA). Emission probabilities were 100% haplotype 1 alleles from haplotype 1, 100% haplotype 2 alleles from haplotype 2, and equal probability haplotype 1 or 2 alleles from the third “error” state. Crossover calling was robust to a range of low transition probabilities. (See our script *spseqHMMCO Caller_3state.R*, which calls crossovers on one chromosome.)

After aneuploidy identification, we marked aneuploid chromosomes as having no crossovers for all crossover analyses (absent chromosomes have no crossovers and crossovers are called differently on gained chromosomes, described subsequently).

Identifying even-coverage cell barcodes

We used Genome STRiP v2.0 (GS) (<http://software.broadinstitute.org/software/genomestrip/>)^{68,69} to determine sequence read depth (observed number of reads divided by expected number of reads) in bins of 100 kb of uniquely mappable sequence across the genome in each sperm cell, using GS’s default GC bias correction and repetitive region masking for gr38. We divided read depth by 2 to obtain read depth per haploid rather than diploid genome. Input to GS was a BAM file containing only cells of interest with read groups set to <sample name>:<cell barcode> (created using Drop-seq Tools’ ConvertTagToReadGroup with options CELL_BARCODE_TAG=XC, SAMPLE_NAME=<name of sample/donor>, CREATE_INDEX=true, and CELL_BC_FILE=list of barcodes potentially associated with cells, described above).

A minority of cell barcodes were associated with eccentric read depth across many chromosomes, with wave-like read depth vacillating between 0 and 2. (We hypothesize that these cell barcodes were associated with sperm nuclei that did not properly decondense, such that some regions of the genome were more accessible than others, leading to undulating read depth across more and less accessible chromatin.) To identify and exclude such barcodes, we treated read depths across each chromosome as a time series and used Box-Jenkins Autoregressive Integrated Moving Average (ARIMA) modelling to model how read depth observations relied on their previous values and their overall averages (implemented via the R package *forecast*^{70,71}, excluding differencing). By visual inspection, we determined that chromosomes with certain ARIMA criteria were likely to have undulating read depth, and that cell barcodes with five or more such identified chromosomes were likely to have eccentric read depth globally. We flagged individual chromosomes if 1) the sum of AR1 and AR2 coefficients was greater than 0.7, the AR1 coefficient was greater than 0.9, or the net sum of all AR and MA coefficients was greater than 1.25 and 2) either the net sum of AR and MA coefficients was greater than 0.4 or the intercept was less than 0.8 or greater than 1.2. If both criteria in (2) were met, this signified an exceedingly odd chromosome, which we counted twice. Cell barcodes with five or more chromosomes flagged in this way were excluded from downstream analyses. (Because gains of large amounts of the genome cause artificially depressed read depths on non-gained chromosomes, we manually examined any cells with a large range of ARIMA intercepts and over five chromosomes denoted as unstable. Any such cells that had simply gained a large proportion of the genome, *e.g.*, 3 copies of chromosome 2, were included rather than

excluded.) We cross-referenced all cell exclusions with called aneuploidies, confirming that cells were not excluded simply on the basis of having lost or gained a chromosome.

(See our scripts *setupsreaddepth.R*, *exclbadreaddepth_arima_1.R*, *exclbadreaddepth_initid_2.R*, and *exclbadreaddepth_finalize_3.R*)

Replicate barcodes (“bead doublets”)

One sperm cell can be encapsulated in a droplet with more than one barcoded bead. To identify such cases, where pairs of sperm genomes were identical, we determined the proportion of SNPs that were of the same haplotype for each pair of barcodes. We imputed the haplotype of all heterozygous SNPs based on the haplotype of surrounding observed SNPs and locations of recombination events and compared SNP haplotypes across sperm cell pairs. SNP observations between boundaries of crossovers were excluded from analysis. Sperm cells shared on average 50% of their genomes, but a few sets of barcodes shared nearly 100% of their SNP haplotypes (Extended Data Fig. 3a). We considered these pairs “bead doublets” or replicate barcodes. In all downstream analyses, only one barcode (chosen randomly) from a set corresponding to the same cell was used. (See our scripts *imputeHaplotypeAllSNPs.R*, *compareSpermHapsPropSNPs.R*, *combineChrsSpermHapsPropSNPs.R*, and *curateNonRepBCList.R*)

Crossover zones

To define regions of recombination use, we found local minima of the density (built-in function in R) of all crossovers’ median positions across all samples on each chromosome. Minima were identified using the *findPeaks* function (from <https://github.com/stas-g/findPeaks>) on the inverse density with $m=3$. Crossover zones run from the beginning of the chromosome (including the whole p arm for acrocentric chromosomes) to the location of the first local minimum, from the location of the first local minimum plus one basepair to the next local minimum, *etc.*, with the last zone on each chromosome ending at the chromosome end. (See our script *findcozones_peaks.R*.)

Aneuploidy and chromosome arm loss/gain

As described previously (see “Restricting to cell barcodes with coverage of the entire genome”), we used Genome STRiP (<http://software.broadinstitute.org/software/genomestrip/>)^{68,69} to determine read depth in each sperm cell in 100 kb bins. We located chromosomes or chromosome arms with aberrant read depth to identify aneuploidy.

We excluded genomic regions that had outlying read depth across all cells, defined as those with $p < 0.05$ in a one-sided one-sample t -test (looking for increased read depth) against the expected mean read depth of $2\#$ (defined below). To identify gains of autosomes, we performed a one-sided one-sample t -test (expecting increased read depth in a gain) for each cell against expected read depth for a gain of one copy, $2\#$. For each cell, this analysis compared the distribution all bins’ read depth across a region of interest to the gain expectation $2\#$, and flagged any cells whose read depth distributions were not significantly different ($p < 0.05$) We used the same approach to identify losses, comparing a cell’s read

depth distribution across bins to 0.1 and flagging any that were not significantly higher ($p < 0.05$).

The expected copy number for gains is 2, but the expected read depth for gains depends on the size of the chromosome: a library corresponding to a cell with a chromosome gain has more reads than would be in that same library without a gain. This phenomenon pulls read depth down globally by increasing the total number of expected reads, causing the denominator in each read depth bin (the expected number of reads in that bin) to increase. Therefore, we computed a chromosome-specific critical read depth value for identifying gains: $2\# = 2 * (\text{the proportion of the genome in base pairs coming from all chromosomes other than the tested one})$. For losses, we used 0.1 rather than 0 as the expected read depth because a small number of reads generally align to every chromosome in every library.

For non-acrocentric chromosomes, we performed aneuploidy calling for the arms separately and for the whole chromosome. Because amplification of more than two copies of a chromosome arm could result in the whole chromosome passing the p -value threshold, we required a whole-chromosome event to pass the p -value threshold at the whole-chromosome level and to have rounded read depth of both arms ≥ 2 for a gain (or 0 for a loss). For the acrocentric chromosomes, only the q arm was considered and any q arm gain or loss was considered to be a whole-chromosome event (unless investigated further).

For the sex chromosomes, we followed a similar statistical framework, but a loss was only considered an aneuploidy if both the X and the Y chromosomes were flagged as lost. A gain was called if both the X and Y chromosomes were present. (See our scripts *setupgsreaddepth.R*, *idaneus_initialttests.R*, *curateaneudata_clean.R*, *getautosomalaneumatrix.R*, and *getxykaryos_aneus.R* for aneuploidy calling and output formatting; see our scripts *curateAnFreqFromCodeMatrix.R*, *curateInitAnalyzeXYKaryos.R*, and *combineAnFreq_AutXY.R* for conversion of outputs of aneuploidy calling to cross-donor aneuploidy frequency tables.)

Chromosome gains' division of origin

To see when chromosome gains originated, we determined whether the centromeres of the multiple copies of the chromosomes were heterozygous and therefore from homologs, which typically disjoin in meiosis I (MI), or homozygous and therefore from sister chromatids, which typically disjoin in meiosis II (MII). We identified heterozygous regions for all cells using a Hidden Markov Model (HMM) in which the states are 1) heterozygous (emitting either haplotype's alleles) or 2) homozygous (emitting only one haplotype's alleles), with transition probability between the states equal to the recombination transition probability. For each gain, we determined whether heterozygous tracts overlapped the centromere. If a heterozygous tract 1) started before the start of the centromere and ended after the end of the centromere or 2) started at the first SNP observed on an acrocentric chromosome or within the first 10 SNPs and was more than 10 SNPs long, the chromosome was classified as an MI gain; if no heterozygous tract overlapped the centromere, it was classified as an MII gain. (See our scripts *getDiploidTracts_hmm.R*, *originOfGainID.R*, and *curateOriginMultSamps.R*.)

At the sex chromosomes, any XY sex chromosome gain derives from MI (X and Y are homologs), whereas an XX or YY gain derives from MII (sister chromatids duplicated).

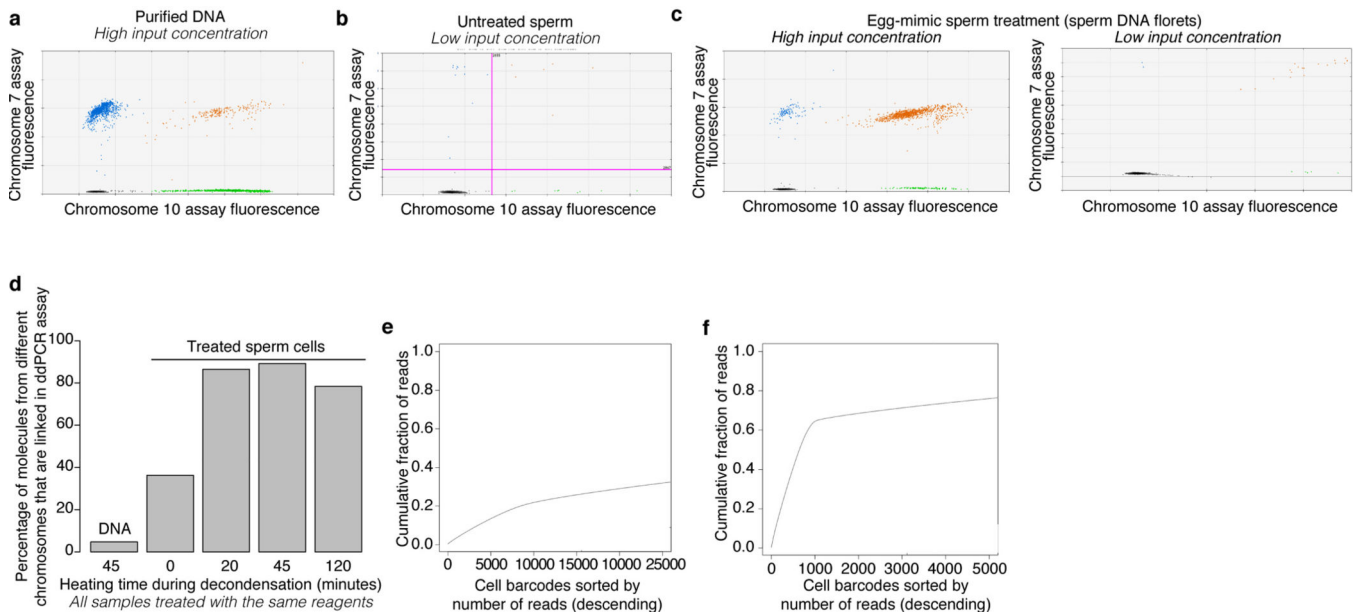
Extended Data

Author Manuscript

Author Manuscript

Author Manuscript

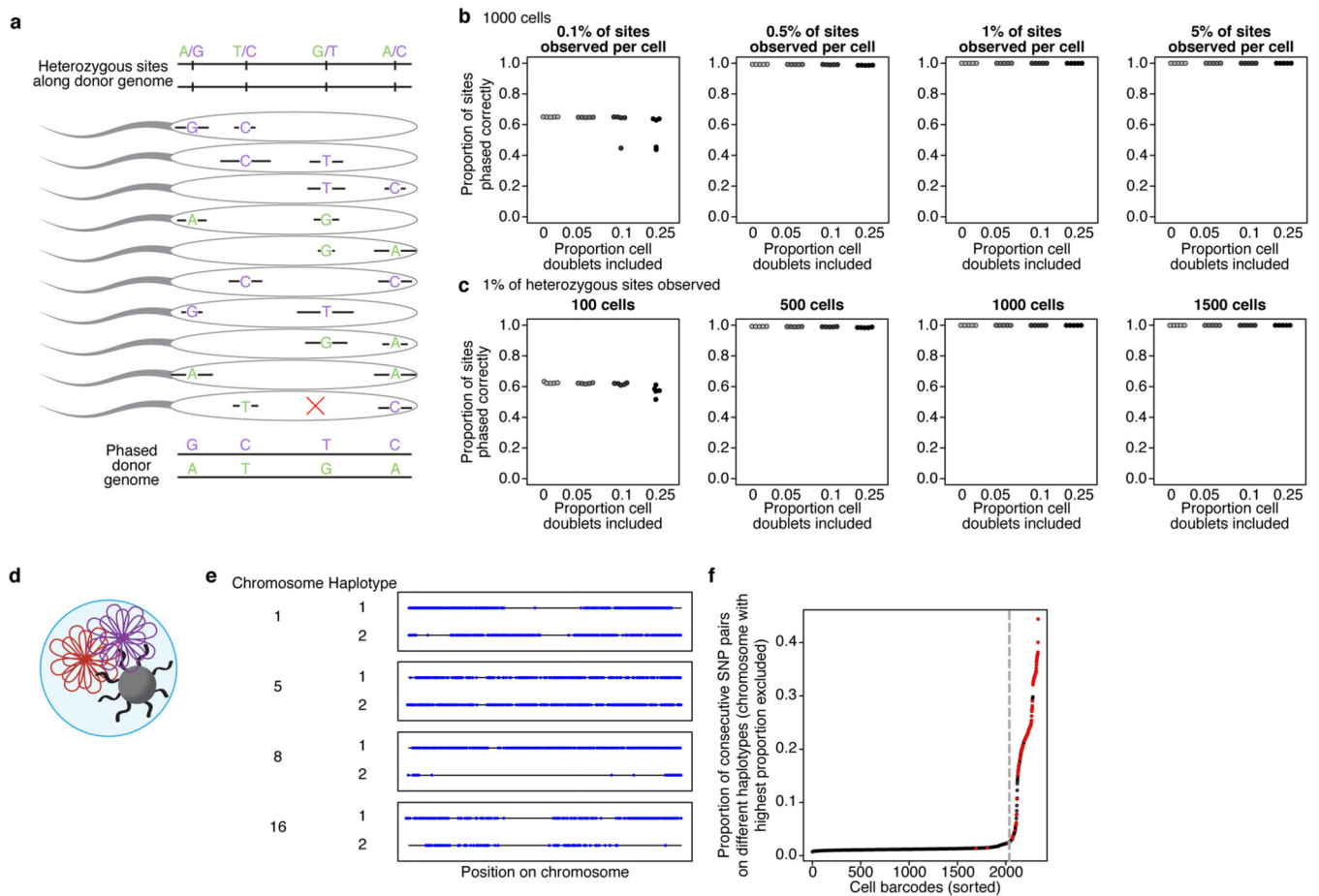
Author Manuscript



Extended Data Fig. 1. Characterization of egg-mimic sperm preparation and optimization of bead-based single-sperm sequencing.

a-c, Two-channel fluorescence plots showing the results of droplet digital PCR (ddPCR) with input template noted in each title, demonstrating that two loci (from different chromosomes) are detectable in the same droplet far more often when sperm DNA florets (rather than purified DNA) are used as input. Each point represents one droplet. Gray points in the bottom left quadrant represent droplets in which neither template molecule was detected; blue points in the top left quadrant represent droplets in which the assay detected a template molecule for the locus on chromosome 7; green droplets in the bottom right quadrant represent droplets in which the assay detected a template molecule for the locus on chromosome 10; and brown point in the top right quadrant represent droplets in which both loci were detected. With a high concentration of purified DNA as input (**a**), comparatively fewer droplets contain both loci than when untreated (**b**) or treated (**c**) sperm were used as input. Sperm “florets” treated with the egg-mimicking decondensation protocol had a much higher fraction of droplets containing both loci than purified DNA (compare **a** and **c**, right, high-input treated sperm) and had more-sensitive ascertainment and cleaner results (quadrant separation) than untreated sperm (compare **b** and **c**, left, low-input sperm and treated sperm). The pink lines in (**b**) delineate the boundaries between droplets categorized as negative or positive for each assay. **d**, Optimization of sperm preparation: Characterization of the effect of different lengths of 37°C incubation of sperm cells treated with egg-mimicking decondensation reagents on how often the loci on chromosomes 7 and 10 were detected in the same ddPCR droplet. Y axis, the percentage of molecules calculated to be linked to each other (*i.e.* physically linked in input) for assays targeting chromosomes 7 and 10. Extracted DNA (a negative control) gives the expected result of random assortment of the two template molecules into droplets (first bar). The 45-minute heat treatment was used for all subsequent experiments in this study. **e** and **f**, Distribution of sequence reads across cell barcodes from droplet-based single-sperm sequencing. Each panel shows the cumulative fraction (y-axis) of all reads from a sequencing run coming from

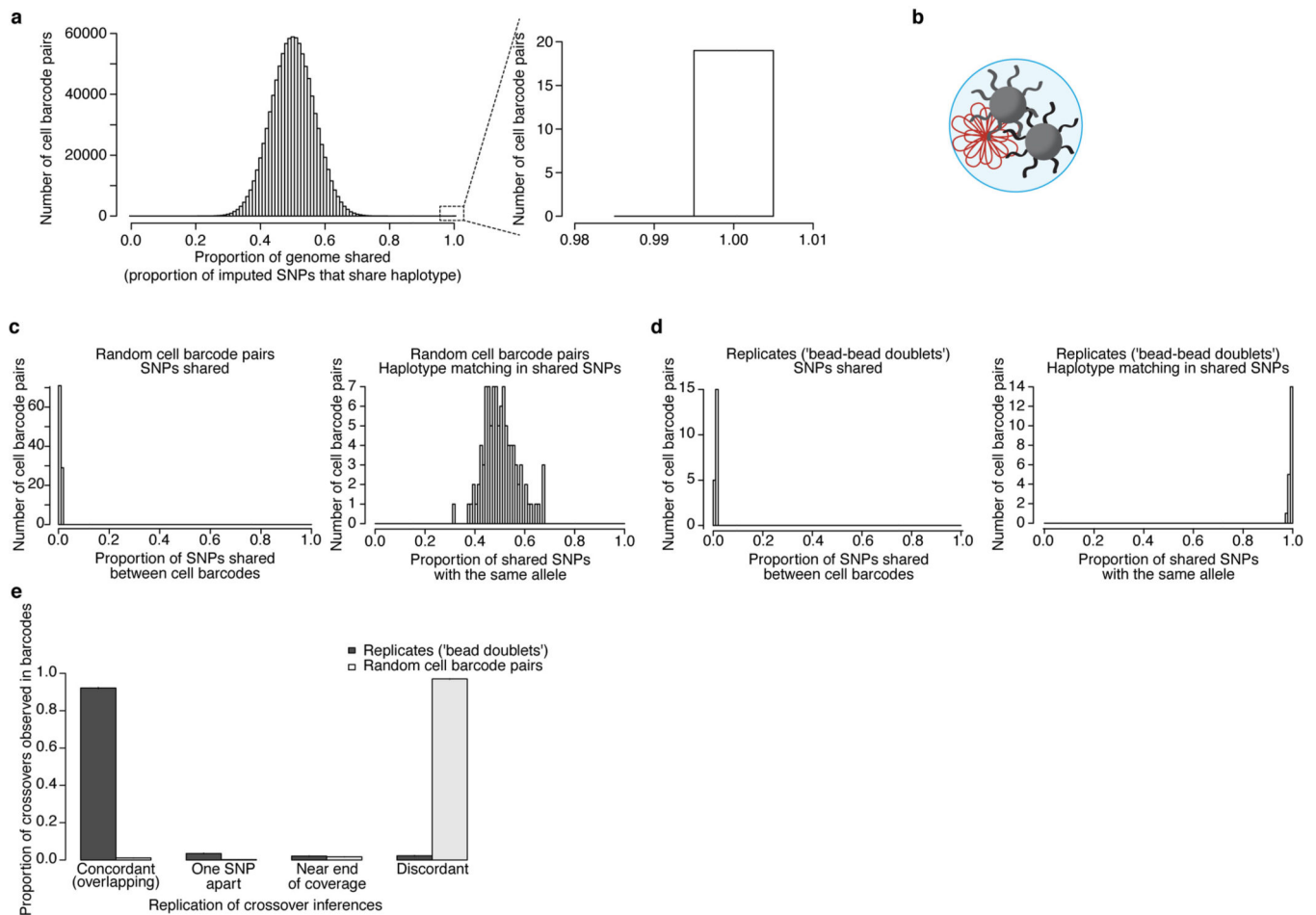
each read-number-ranked cell barcode; a sharp inflection point delineates the barcodes with many reads from those with few reads. Points to the left of the inflection point are the cell barcodes that associated with many reads (i.e., beads that co-encapsulated with cells); the height of the inflection point reflects the proportion of the sequence reads that come from these barcodes. Only reads that mapped to the human genome (hg38) and were not PCR duplicates are included. **e**, Data from an initial adaptation of 10X Genomics' GemCode linked reads system²⁹ where a small proportion of the reads come from cell barcodes associated with putative cells. **f**, Data from the final, implemented adaptation of 10X Genomics' GemCode linked reads system²⁹ for the same number of input sperm nuclei as in **e**. Note that this x-axis includes five times fewer barcodes than in (**e**).



Extended Data Fig. 2. Evaluation of chromosomal phasing and identification of cell doublets.

a, Phasing strategy. Green and purple denote the chromosomal phase of each allele (unknown before analysis). Each sperm cell carries one parental haplotype (green or purple) except where a recombination event separates consecutively observed SNPs (red “X” in bottom sperm). Because alleles from the same haplotype will tend to be observed in the same sperm cells, the haplotype arrangement of the alleles can be assembled at whole-chromosome scale. **b**, Evaluation of our phasing method using 1,000 simulated single-sperm genomes (generated from two *a priori* known parental haplotypes and sampled at various levels of coverage). Since cell doublets (which combine two haploid genomes and potentially two haplotypes at any region) can in principle undermine phasing inference, we included cell doublets in the simulation (in proportions shown on the X axis, which bracket the observed doublet rates). Each point shows the proportion of SNPs phased concordantly with the correct (*a priori* known) haplotypes (Y axis) for one simulation (five simulations were performed per proportion of cell doublets-percentage of observed sites condition pair). **c**, Relationship of phasing capability to number of cells analyzed. Data are as in (**b**), but for different numbers of simulated cells. All simulations had an among-cell mean of 1% of heterozygous sites observed. **d**, A cell doublet: when two cells (here, sperm DNA florets) are co-encapsulated in the same droplet, their genomic sequences will be tagged with the same barcode; such events must be recognized computationally and excluded from downstream

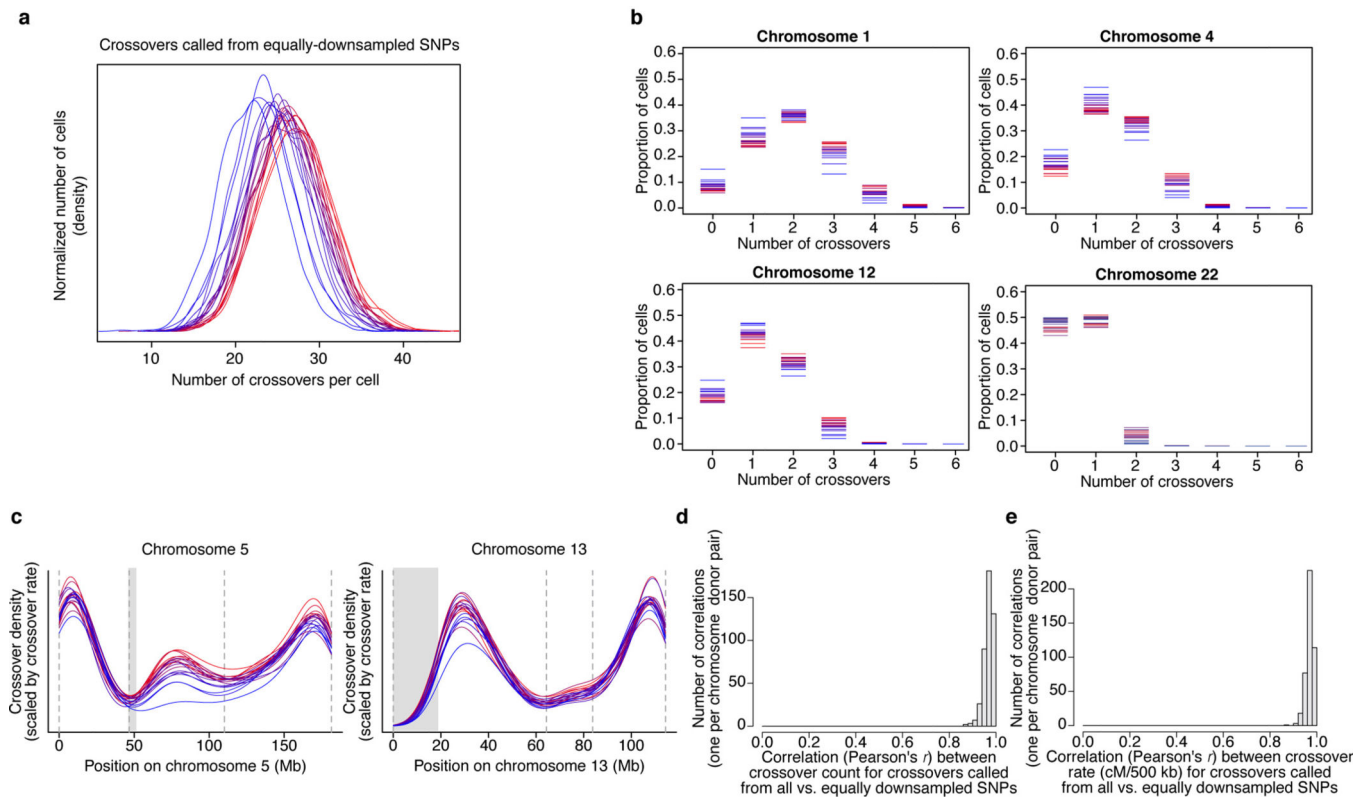
analyses. **e**, Four example chromosomes from a cell barcode associated with two sperm cells (a cell doublet). Black lines: haplotypes; blue circles: observations of alleles, shown on the haplotype from which they derive. Both parental haplotypes are present across regions of chromosomes where the cells inherited different haplotypes. **f**, Computational recognition of cell doublets in Sperm-seq data (from an individual sperm donor, NC11). The proportion of consecutively observed SNP alleles derived from different parental haplotypes is used to identify cell doublets; this proportion is generally small (arising from sparse crossovers, PCR/sequencing errors, and/or ambient DNA) but is much higher when the analyzed sequence comes from a mixture of two distinct haploid genomes. We use 21 of the 22 autosomes to calculate this proportion, excluding the autosome with the highest such proportion given the possibility that a chromosome is aneuploid. The dashed gray line marks the inflection point beyond which sperm genomes are flagged as potential doublets and excluded from downstream analysis. Red points indicate barcodes with coverage of both the X and Y chromosome (potentially X+Y cell doublets or XY aneuploid cells); black points indicate barcodes with one sex chromosome detected (X or Y). The red (XY) cells below the doublet threshold are XY aneuploid but appear to have just one copy of each autosome.



Extended Data Fig. 3. Identification and use of “bead doublets.”

a, SNP alleles were inferred genome-wide (for each sperm genome) by imputation from (i) the subset of alleles detected in each cell and (ii) Sperm-seq-inferred parental haplotypes. For each pair of sperm genomes (cell barcodes), the proportion of all SNPs at which they shared the same imputed allele was estimated. A small but surprising number of such pairwise comparisons (19 of 984,906 from the donor shown, NC14) indicate essentially identical genomes (ascertained through different SNPs). **b**, We hypothesize that this arises from a heretofore undescribed scenario we call “bead doublets”, in which two barcoded beads have co-encapsulated with the same gamete and whose barcodes therefore tagged the same haploid genome. **c**, Random pairs of cell barcodes (here 100 pairs selected from donor NC10) tend to interrogate few of the same SNPs (left), and tend to detect the same parental haplotype on average at the expected 50% of the genome (right). **d**, “Bead doublet” barcode pairs (here 20 pairs from donor NC10, who had the median number of bead doublets, left) also interrogate few of the same SNPs, yet detect identical haplotypes throughout the genome (right). Results were consistent across donors. **e**, Use of “bead doublets” to characterize the concordance of crossover inferences between distinct samplings of the same haploid genome by different barcodes. The bead doublets (barcode pairs) were compared to 100 random barcode pairs per donor. Crossover inferences were classified as “concordant” (overlapping, detected in both barcodes), as “one SNP apart” (separated by just one SNP,

detected in both barcodes), as “near end of coverage” (within 15 heterozygous SNPs of the end of SNP coverage at a telomere, where power to infer crossovers is partial), or as discordant. Error bars (with small magnitude) show binomial 95% confidence intervals for the number of crossovers per category divided by number of crossovers total in both barcodes (32,714 crossovers total in 1,201 bead doublet pairs; 67,862 crossovers total in 2,000 random barcode pairs; some barcodes are in multiple bead doublet or random barcode pairs).



Extended Data Fig. 4. Numbers and locations of crossovers called from down-sampled data (equal number of SNPs in each cell, randomly chosen).

To eliminate any potential effect of unequal sequence coverage across donors and cells, down-sampling was used to create data sets with equal coverage (numbers) of heterozygous SNP observations in each cell. Crossovers were called from these random equally sized sets of SNPs from all cells. **a** and **b**, Crossover number per cell globally (**a**) and per chromosome (**b**) (785,476 total autosomal crossovers called from down-sampled SNPs included, 30,778 cells included, aneuploid chromosomes excluded). **c**, Density plots of crossover location with crossover midpoints plotted and area scaled to be equal to per-chromosome crossover rate. Gray rectangles mark centromeric regions; coordinates are in hg38. **d**, Similar numbers of crossovers were called from full data and equally down-sampled SNP data: we performed correlation tests across cells for each donor and chromosome to compare the number of crossovers called from all data to the number of crossovers called from equal numbers of randomly down-sampled SNPs. The histogram shows Pearson's r values for all 460 (20 donors x 23 chromosomes [total number plus number for 22 autosomes]) tests (n per test = 974–2,274 cells per donor as in Extended Data Table 1, all chromosome comparisons Pearson's $r > 0.83$, all two-sided $p < 10^{-300}$). **e**, Crossovers called from equally down-sampled SNP data were in similar locations to those called from all data: we performed correlation tests comparing crossover rate in 500 kb bins (cM/500 kb) from all data vs. equally down-sampled SNP data for each donor and chromosome. The histogram shows Pearson's r values for all 460 (20 donors x 23 chromosomes [genome-wide rate plus rate for 22 autosomes]) tests (n per test = number of 500 kb bins per chromosome [genome-wide: 5,739, chromosomes 1 through 22: 497, 484, 396, 380, 363, 341, 318, 290, 276, 267, 270,

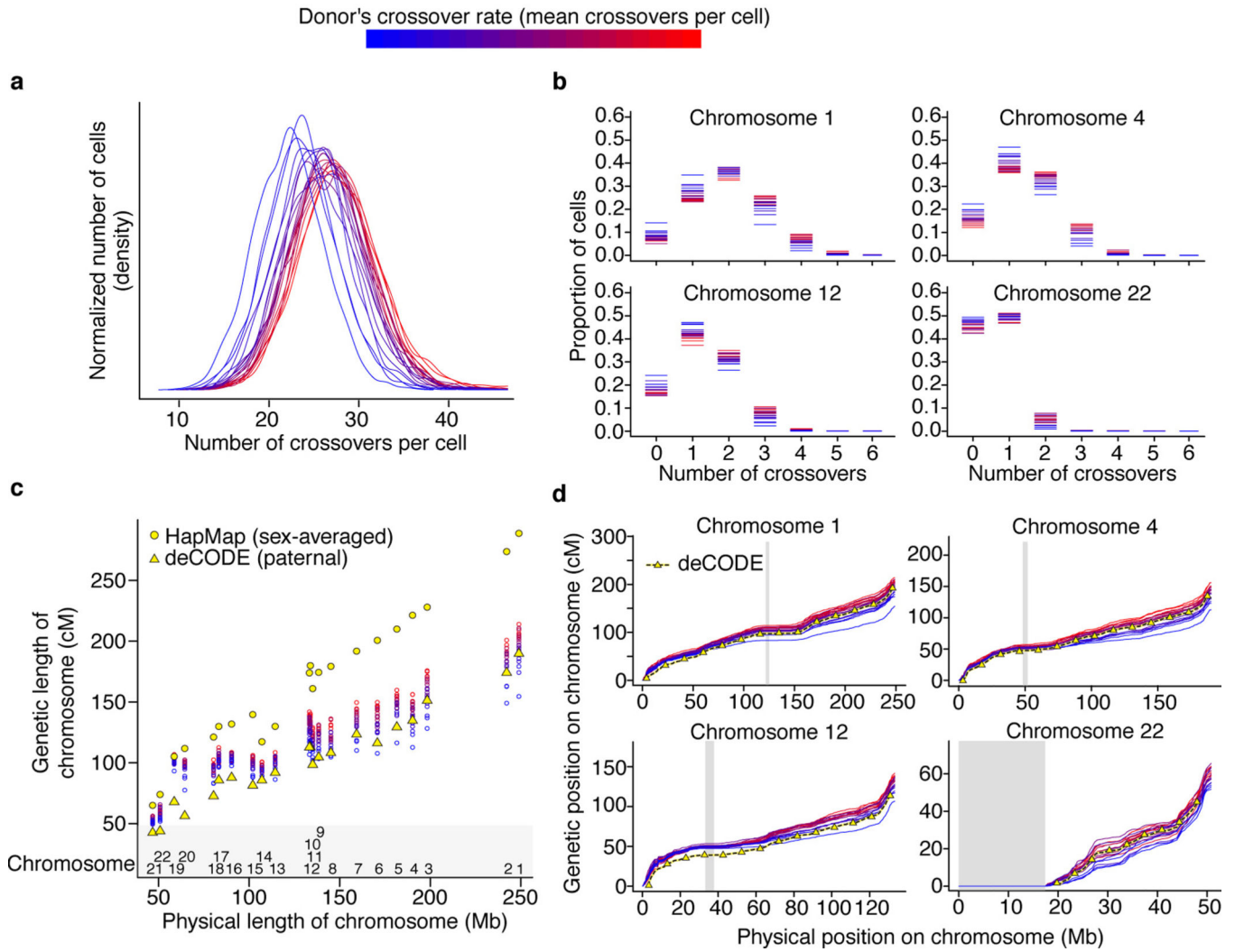
266, 228, 214, 203, 180, 166, 160, 117, 128, 93, 101], all chromosome comparisons
Pearson's $r > 0.87$, all two-sided $p < 10^{-300}$).

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript



Extended Data Fig. 5. Inter-individual and inter-cell recombination rate from single-sperm sequencing.

a, Density plot showing per-cell number of autosomal crossovers for all 31,228 cells (813,122 total autosomal crossovers) from 20 sperm donors (per-donor cell and crossover numbers as in Extended Data Table 1; aneuploid chromosomes were excluded from crossover analysis). Colors represent a donor's mean crossover rate (crossovers per cell) from low (blue) to high (red). This same mean recombination rate-derived color scheme is used for donors in all figures. Recombination rate differs among donors ($n = 20$, Kruskal-Wallis chi-squared = 3,665, $df=19$, $p < 10^{-300}$). **b**, Per-chromosome crossover number in each of the 20 sperm donors (data as in **a**) but shown for individual chromosomes). **c**, Per-chromosome genetic map lengths for: (i) each of the 20 sperm donors, as inferred from Sperm-seq data (colors from blue to red reflect donors' individual crossover rates as described above); (ii) a male average, as estimated from pedigrees by deCODE⁶ (yellow triangles); (iii) a population average (including female meioses, which have more crossovers), as estimated from HapMap data⁷ (yellow circles). The deCODE genetic maps stop 2.5 Mb from the ends of SNP coverage. **d**, Physical vs. genetic distances (for

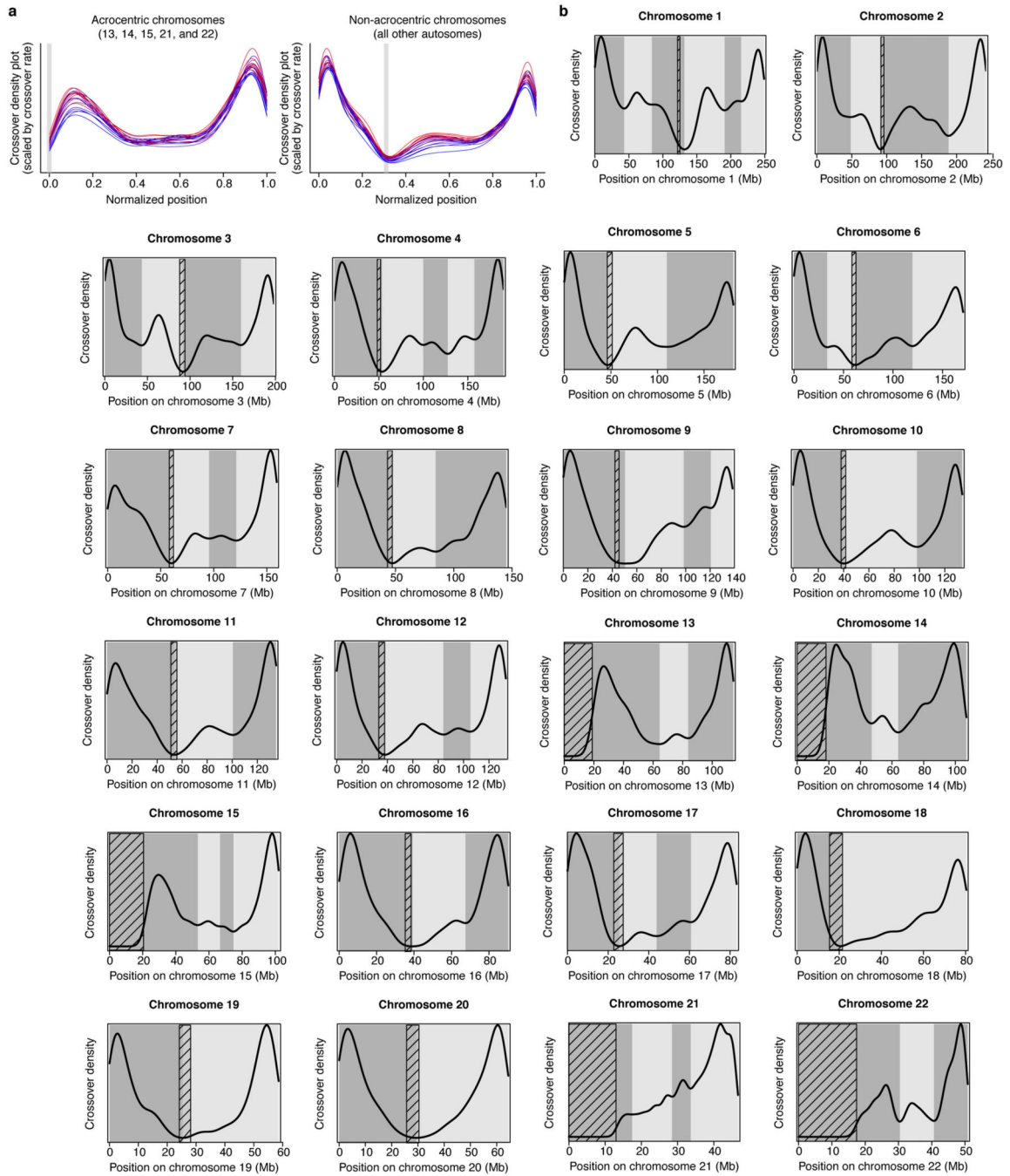
individualized sperm donor genetic maps and deCODE's paternal genetic map) plotted at 500 kb intervals (hg38). Gray boxes denote centromeric regions (or centromeres and acrocentric arms). Sperm-seq maps are broadly concordant with deCODE maps (correlation test results in Supplementary Notes) except at subtelomeric regions not included in deCODE's map.

Author Manuscript

Author Manuscript

Author Manuscript

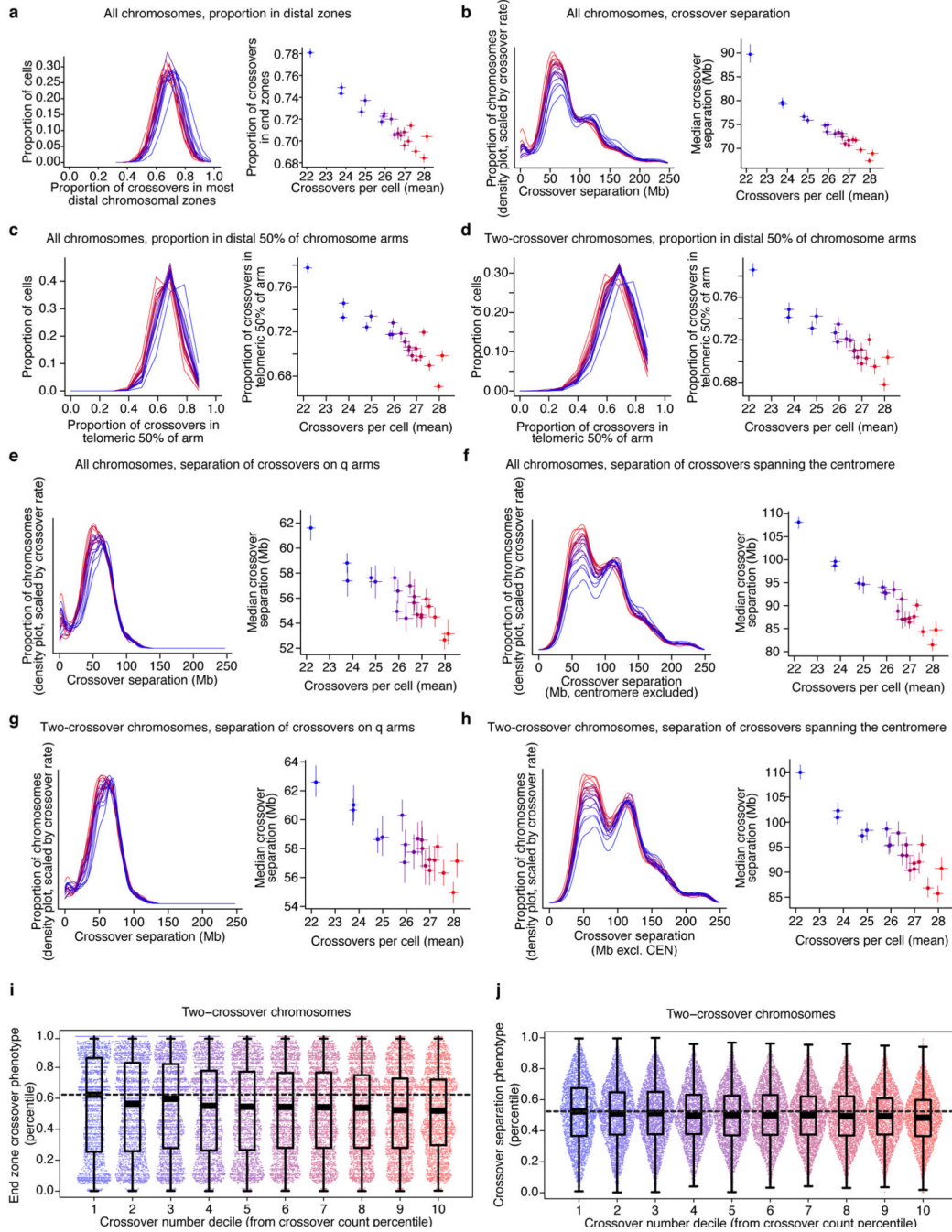
Author Manuscript



Extended Data Fig. 6. Distributions of crossover locations along chromosomes (in “crossover zones”).

a, Each donor’s crossover locations are plotted as a colored line; color indicates the donor’s overall crossover rate (blue: low, red: high); gray boxes show the locations of centromeres (or, for acrocentric chromosomes, centromeres and *p* arms). The midpoint between the SNPs bounding each inferred crossover was used as the position for each crossover in all analyses. To combine data across chromosomes, crossover locations (density plot) are shown on “meta-chromosomes” in which crossover locations are normalized to the length of the

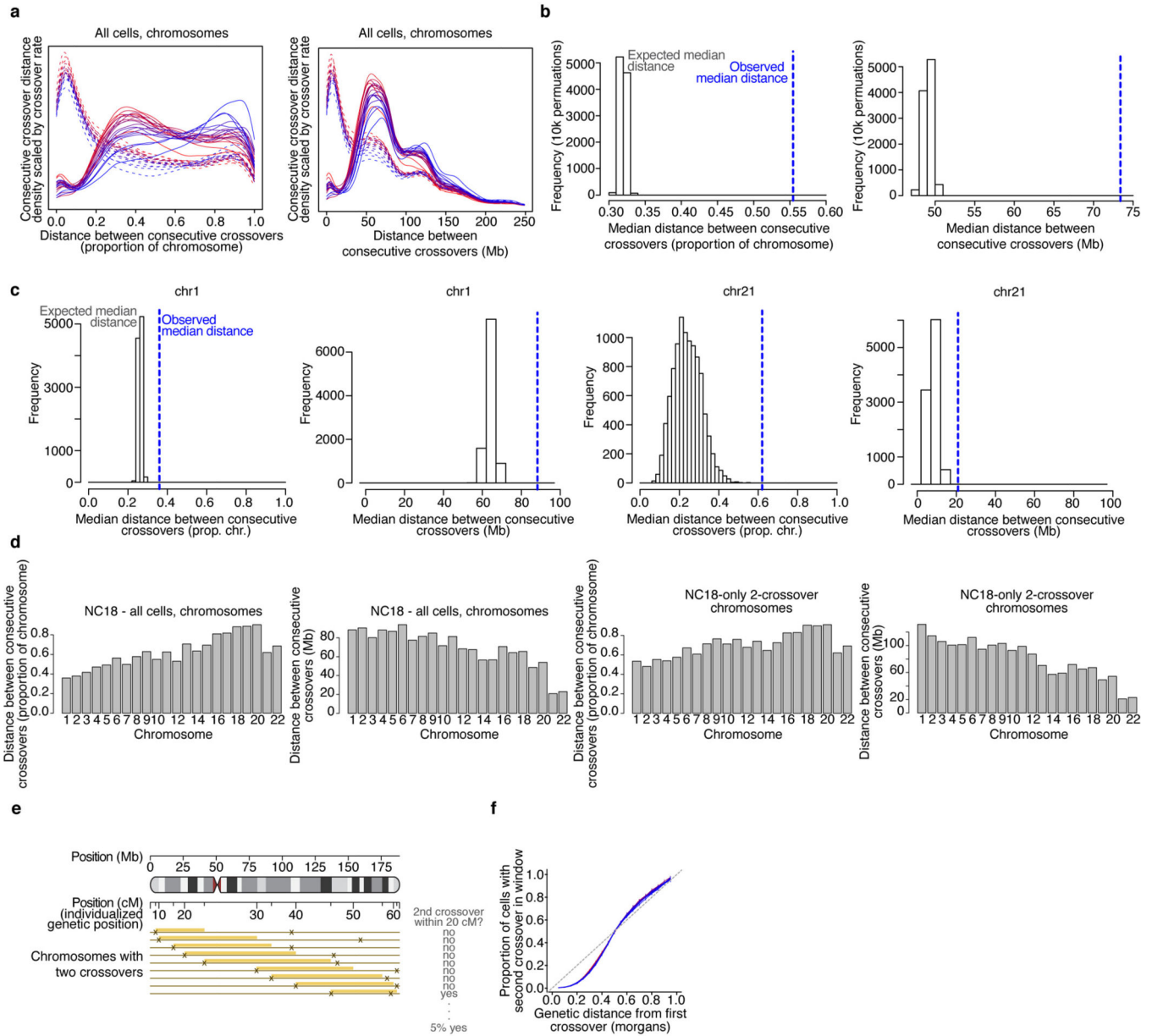
chromosome or arm on which they occurred. For acrocentric chromosomes, only the q arm was considered; for non-acrocentric chromosomes, the p and q arms were afforded space based on the proportion of the non-acrocentric genome (in bp) they comprise, with the centromere placed at the summed p arms' proportion of bp of these chromosomes. Crossover locations were first converted to the proportion of the arm at which they fall, then these positions normalized to the genome-wide p or q arm proportion. **b**, Identification of chromosomal zones of recombination use ("crossover zones") from all donors' crossovers for 22 autosomes. Density plots of crossover location for all sperm donors' total 813,122 crossovers (aneuploid chromosomes excluded; crossover location is the midpoint between SNPs bounding crossovers) along autosomes (hg38) are shown. Crossover zones (bounded by local minima of crossover density) are shown by alternating shades of gray. Diagonally-hatched rectangles indicate centromeres (or centromeres and acrocentric arms).



Extended Data Fig. 7. Crossover placement in end zones, and crossover separation, vary in ways that correlate with crossover rate – among sperm donors and among individual gametes.

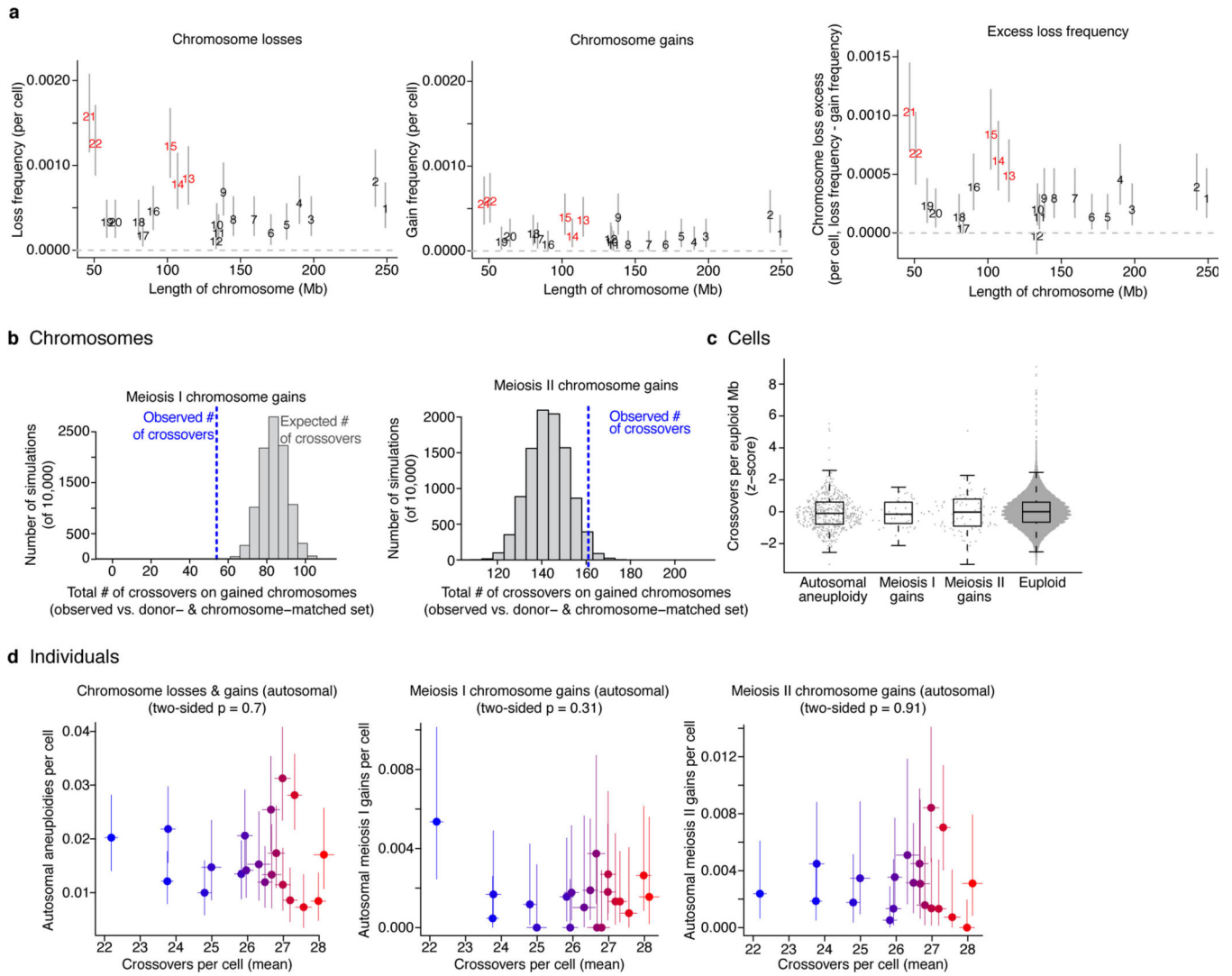
Analyses are shown by donor (a-h, $n = 20$ sperm donors) or by individual gamete (i-j, $n = 31,228$ gametes). In a-h, the left panels show the phenotype distributions for individual donors, and the right panels show the relationship to the donors' crossover rates. To control for the effect of the number of crossovers, the analyses in panels c, d, and g-j use “two-crossover chromosomes” – chromosomes on which exactly two crossovers occurred. For scatter plots (a-h, right), all x axes show mean crossover rate and all error bars are 95%

confidence intervals (y axes are described per panel). **a** and **b**, The proportion of crossovers falling in the most distal chromosome crossover zones (**a**) and crossover separation (**b**) – a readout of crossover interference, the distance between consecutive crossovers (Mb) – vary among 20 sperm donors (left panels; proportion of crossovers in end per cell distributions among-donor Kruskal–Wallis chi-squared = 2,334, $df=19$, $p < 10^{-300}$; all distances between consecutive crossovers among-donor Kruskal–Wallis chi-squared = 3,309, $df=19$, $p < 10^{-300}$). Right panels show both properties (y axes, total proportion of crossovers in distal zones and median crossover separation, respectively) vs. donor's crossover rate (Correlation results for 20 sperm donors: proportion of all crossovers across cells in distal zones Pearson's $r = -0.95$, two-sided $p = 2 \times 10^{-10}$; Pearson's $r = -0.96$, two-sided $p = 1 \times 10^{-11}$). **c**, An alternative method for the proportion of crossovers in the distal regions of chromosomes: proportion of crossovers in the distal 50% of chromosome arms varies across donors (left, among-donor Kruskal–Wallis chi-squared = 2,209, $df=19$, $p < 10^{-300}$) and negatively correlates with recombination rate (right, Pearson's $r = -0.92$, two-sided $p = 2 \times 10^{-8}$; y axis shows actual proportion of crossovers in distal 50%). **d**, As in (**c**), but with proportion of crossovers from two-crossover chromosomes occurring in the distal 50% of chromosome arms. Left, among-donor Kruskal–Wallis chi-squared = 1,058, $df=19$, $p = 2 \times 10^{-212}$; right, correlation with recombination rate Pearson's $r = -0.93$, two-sided $p = 4 \times 10^{-9}$. **e**, as in (**b**) but for consecutive crossovers on the *q* arm of the chromosome. Left, among-donor Kruskal–Wallis chi-squared = 346, $df=19$, $p = 7 \times 10^{-62}$; right, correlation with recombination rate Pearson's $r = -0.90$, two-sided $p = 5 \times 10^{-8}$. **f**, as in (**b**) but for consecutive crossovers on opposite chromosome arms (*i.e.* that span the centromere). Left, among-donor Kruskal–Wallis chi-squared = 1,554, $df=19$, $p = 1 < 10^{-300}$; right, correlation with recombination rate Pearson's $r = -0.96$, two-sided $p = 3 \times 10^{-11}$. **g**, as in (**e**) but for distances between consecutive crossovers on two-crossover chromosomes. Left, among-donor Kruskal–Wallis chi-squared = 181, $df=19$, $p = 2 \times 10^{-28}$; right, correlation with recombination rate Pearson's $r = -0.88$, two-sided $p = 3 \times 10^{-7}$. **h**, as in (**f**) but for distances between consecutive crossovers on two-crossover chromosomes. Left, among-donor Kruskal–Wallis chi-squared = 930, $df=19$, $p = 5 \times 10^{-185}$; right, correlation with recombination rate Pearson's $r = -0.92$, two-sided $p = 1 \times 10^{-8}$. **i, j**, Boxplots show medians and interquartile ranges with whiskers extending to 1.5 times the interquartile range from the box. Each point is a cell. **i**, Within-donor percentile of proportion of crossovers from two-crossover chromosomes falling in distal zones plotted vs. crossover rate decile. Groups are deciles of crossover rate normalized by converting each cell's crossover count to a percentile within-donor (All cells from all donors shown together, n cells in deciles = 3,152, 3,122, 3,276, 3,067, 3,080, 3,073, 3,135, 3,132, 3,090, 3,101, respectively [31,228 total]). Because the initial data is proportions with small denominators, an integer effect is evident as pileups at certain values. **j**, Crossover interference from two-crossover chromosomes (median consecutive crossover separation per cell shown). Each point represents the median of all percentile-expressed distances between crossovers from all two-crossover chromosomes in one cell (percentile taken within-chromosome), groupings and *ns* as in (**i**).



Extended Data Fig. 8. Crossover interference in individual sperm donors and on chromosomes.
a, Solid lines show density plots (scaled by donor’s crossover rate) of the observed distance (separation) between consecutive crossovers as measured in the proportion of the chromosome separating them (left) and in genomic (Mb) distance (right), one line per donor ($n = 20$). Dashed lines show the distance between consecutive crossovers when crossover locations are permuted randomly across cells to remove the effect of crossover interference.
b, The median of observed distances between consecutive crossovers for one donor (NC18, 10th lowest recombination rate of 20 donors; blue dashed line) is shown with a histogram of the medians of $n = 10,000$ among-cell crossover permutations (both permutation one-sided p s < 0.0001). Units, proportion of the chromosome (left) and genomic (Mb) distance (right).
c, Crossover separation on example chromosomes; plots and n s are as in (b). (Permutation

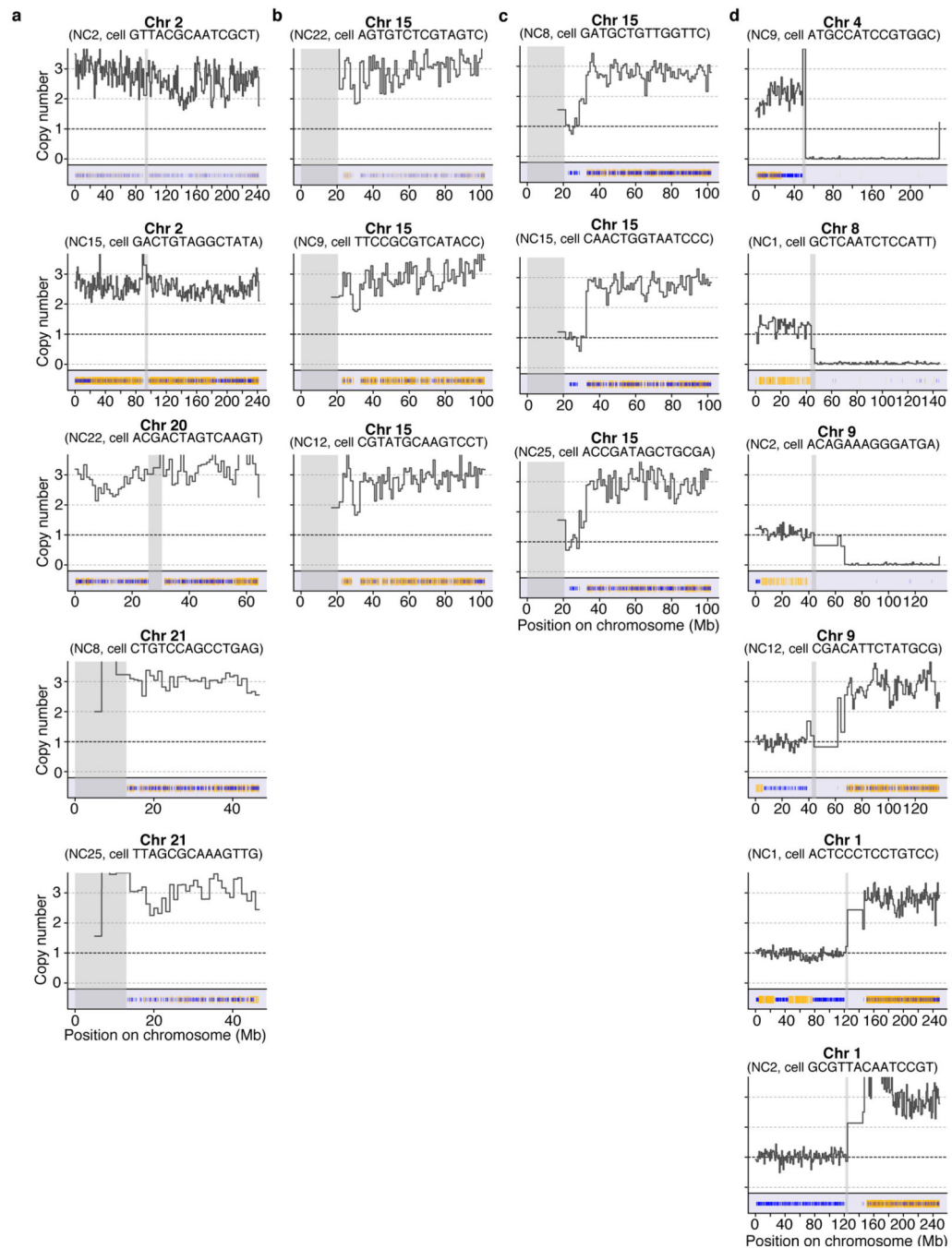
one-sided $p < 0.0001$ for all chromosomes in all sperm donors except occasionally chromosome 21, where especially few double crossovers occur). **d**, Median distances between donor NC18's consecutive crossovers for each autosome for all inter-crossover distances (top) and inter-crossover distances only from chromosomes with two crossovers (bottom). Units are proportion of the chromosome (left) and genomic (Mb) distance (right). **e**, Schematic: analyzing crossover interference in individualized genetic distance (one 20 cM window shown) using a donor's own recombination map. **f**, When parameterized using each donor's own genetic map, sperm donors' crossover interference profiles across multiple genetic distance windows (as shown in **e**) do not differ ($n = 20$ sperm donors, Kruskal–Wallis chi-squared = 0.22, $df = 19$, $p = 1$ using 20 estimates [cM distances] for each of 20 donors). Error bars, binomial 95% confidence intervals on proportion of cells with a second crossover in the window given. This suggests that inter-individual variation in crossover interference, while substantial when measured in base pairs, is negligible when measured in donor-specific genetic distance, pointing to a shared influence upon crossover interference and crossover rate.



Extended Data Fig. 9. Relationships of aneuploidy frequency to chromosome size and recombination.

a. The across-donor per-cell frequency of chromosome losses (left) and gains (center), plotted against the length of the chromosome (hg38; for losses across $n = 22$ chromosomes, Pearson's $r = -0.29$, two-sided $p = 0.19$ and for gains across $n = 22$ chromosomes, Pearson's $r = -0.23$, two-sided $p = 0.30$). Right, the per-chromosome rate of losses exceeding gains (number of losses minus number of gains divided by number of cells) is plotted against the length of the chromosomes (across $n = 22$ chromosomes, Pearson's $r = -0.29$, two-sided $p = 0.19$). Red labels, acrocentric chromosomes. Error bars, 95% binomial confidence intervals on per-cell frequency (number of events / number of cells, all 31,228 cells included). **b-d,** Relationship between aneuploidy frequency and recombination. Only autosomal whole-chromosome aneuploidies are included. **b,** Left, Total number of crossovers on MI nondisjoined chromosomes (blue line; chromosomes analyzed, called as transitions between the presence of one haplotype and both haplotypes on the gained chromosome) compared to $n = 10,000$ donor- and chromosome-matched sets (35×2 chromosomes per set) of properly segregated chromosomes (gray histogram; permutation). (54 total crossovers on MI gains vs.

84.2 mean total crossovers on sets of matched chromosomes, one-sided permutation $p < 0.0001$, for the hypothesis that gained chromosomes have fewer crossovers). Right, as left but for gains occurring during MII (71 MII-derived gained chromosomes of one whole copy from all individuals with fewer than 5 crossovers called on gained chromosome). (One-sided permutation $p = 0.98$ for MII from $n = 10,000$ permutations, for the hypothesis that gained chromosomes have fewer crossovers; sister chromatids nondisjoined in MII capture all crossovers whereas matched chromosomes do not: matched simulations and homologs nondisjoined in MI capture only a random half of crossovers occurring on that chromosome in the parent spermatocyte). **c**, Crossovers per non-aneuploid megabase from each cell from each donor, split by aneuploidy status (n cells = 498, 50, 92, 30,609, left-to-right; “euploid” excludes cells with any autosomal whole- or partial-chromosomal loss or gain and “gains” includes gains of one or more than one chromosome copy; Mann–Whitney test $W = 7,264,117, 722,191, 1,370,376$; two-sided $p = 0.07, 0.49, 0.66$ for all autosomal aneuploidies, meiosis I (MI) gains, and meiosis II (MII) gains, respectively, all compared against euploid). Each cell is one point; boxplots show medians and interquartile ranges with whiskers extending to 1.5 times the interquartile range from the box. **d**, Per-cell crossover rates vs. per-cell aneuploidy (loss and gain) rates, $n = 20$ donors (colored by crossover rate). p values shown in subtitles are for two-sided Pearson’s correlation tests. Error bars are 95% confidence intervals on mean crossover rate (x axis) and on observed aneuploidy frequency (y axis).



Extended Data Fig. 10. Additional examples of non-canonical aneuploidy events detected with Sperm-seq, including those shown in Fig. 3f.

Copy number, SNPs, haplotypes, and centromeres are plotted as in Fig. 3a. Donor and cell identity are noted in the panel subtitles. Coordinates are in hg38. Chromosomes 2, 20, 21 (a) and 15 (b) are sometimes present in 3 copies in an otherwise haploid sperm cell. c, A distinct, recurring triplication of much of chromosome 15, from ~33 Mb onwards but not including the proximal part of the *q* arm, also recurs in cells from 3 donors. d, Chromosome

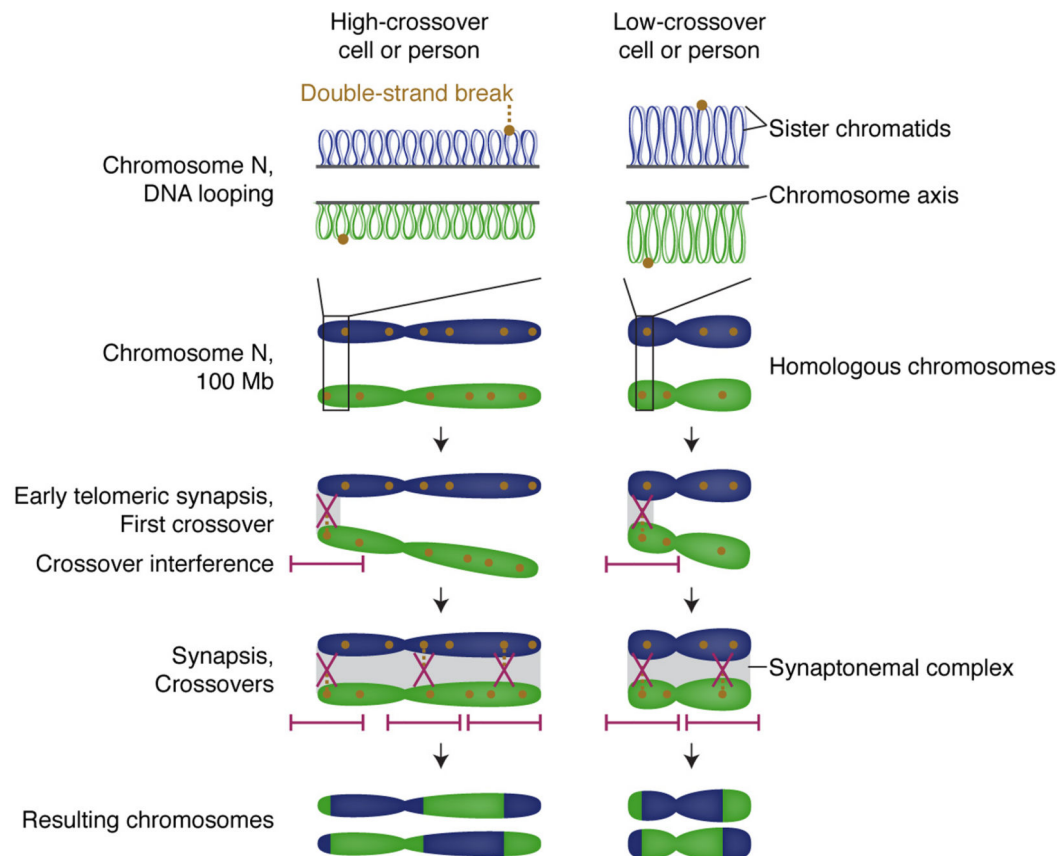
arm-level losses (top) and gains (including in more than one copy, bottom three panels, and a compound gain of the p arm and loss of the q arm, top panel).

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript



In the example depicted:

Chromosome length			
Genomic	100 Mb	=	100 Mb
Physical	10 microns	>	5 microns
Crossover interference			
Genomic	40 Mb	<	90 Mb
Physical	4 microns	=	4 microns
Genomic proportion	0.4	<	0.9
Crossover characteristics			
Crossover number	3	>	2
Fraction of crossovers in distal regions of chromosome arms	2/3	<	1
Crossovers near centromere	Yes	>	No

Extended Data Fig. 11. Single-cell and person-to-person variation in diverse meiotic phenotypes may be governed by variation in the physical compaction of chromosomes during meiosis.

Previous work shows that the physical length of the same chromosome varies among spermatocytes at the pachytene stage of meiosis, likely by differential looping of DNA along the meiotic chromosome axis (e.g. left column shows smaller loops, resulting in more loops total and in greater total axis length compared to the right column with larger loops)^{15,72–75}. This physical chromosome length is correlated across chromosomes among cells from the same individual^{21,76} and correlates with crossover number^{15,20,21,42,73,76}. This length – measured as the length of the chromosome axis or of the synaptonemal complex (the

connector of homologous chromosomes) – can vary two or more-fold among a human’s spermatocytes²¹. We propose that the same process differs on average across individuals and may substantially explain inter-individual variation in recombination rate. On average, individual 1 (left) would have meiotic chromosomes that are physically longer (less compacted) in an average cell than individual 2 (right); one example chromosome is shown in the figure. After the first crossover on a chromosome (likely in a distal region of a chromosome, where synapsis typically begins in male human meiosis before spreading across the whole chromosome^{13–15}), crossover interference prevents nearby double-strand breaks (DSBs) from becoming crossovers; DSBs far away can become crossovers (which themselves also cause interference). More DSBs are likely created on physically longer chromosomes, and crossover interference occurs among non-crossover as well as crossover DSBs⁷⁷. Crossover interference occurs over relatively fixed physical (micron) distances^{43–45,76}; these distances encompass different genomic (Mb) lengths of DNA in different cells or on average in different people due to variable compaction. Thus, crossover interference tends to lead to different total number of crossovers as a function of degree of compaction, resulting in the observed negative correlation (Fig. 2c,e) of crossover rate with crossover spacing (as measured in base pairs). Given that the first crossover likely occurs in a distal region of the chromosome, this model can also explain the negative correlation (Fig. 2b,d) of crossover rate with the proportion of crossovers in chromosome ends. Note: this figure shows the total number of crossovers, crossover interference extent, and crossover locations for both sister chromatids of each homolog combined; in reality, these crossovers are distributed among the sister chromatids, making these relationships harder to detect in daughter sperm cells and requiring large numbers of observations to make relationships among these phenotypes clear.

Extended Data Table 1.

Sperm donor and single-sperm sequencing characteristics and results.

Donor	Ancestry*	Cells (number excluding cell and bead doublets)	Reads per cell (median, thousands)	Genome covered per cell (median, percent)	Heterozygous SNPs in genome (millions)	Unique heterozygous SNP alleles observed per cell (median, thousands)	Crossovers observed (total, thousands)	Crossovers per cell (mean)	Resolution of crossovers (kb, median)	Autosomal aneuploidy events (percent of cells)	Sex chromosome aneuploidy events (percent of cells)
Overall	--	31,228 [‡]	211 [§]	1.0 [¶]	--	24.6 [§]	813 [‡]	26.11 [¶]	240 [¶]	1.6 [§]	0.9 [§]
NC1	Eur.	982	284	1.4	1.95	31.6	26	26.31	189	1.5	0.6
NC2	Eur.	1,680	163	0.8	1.98	18.2	37	22.19	307	2.0	0.7
NC3	Eur.	1,289	190	0.9	1.94	21.5	36	28.13	260	1.7	0.7
NC4	Eur.	1,482	243	1.1	1.98	26.8	40	26.98	243	1.1	0.5
NC6	Afr. Am.	1,370	154	0.8	2.53	23.8	38	27.57	253	0.7	0.3
NC8	As.	1,663	304	1.5	1.81	30.9	45	26.98	229	3.1	0.5
NC9	As.	1,894	245	1.2	1.79	25.6	53	27.98	231	0.8	1.5
NC10	As.	1,154	224	1.1	1.82	23.3	29	24.99	257	1.5	0.3
NC11	Eur.	1,930	202	1.0	1.92	22.8	50	25.82	242	1.3	0.4
NC12	Eur.	2,145	179	0.9	1.91	20.6	51	23.76	270	1.2	1.7
NC13	Eur.	1,514	259	1.2	1.92	28.3	41	27.19	202	0.9	1.0
NC14	Eur.	1,336	296	1.4	1.92	32.4	36	26.65	175	2.5	1.2
NC15	Eur.	1,702	211	1.0	1.93	23.2	42	24.80	268	1.0	0.9
NC16	Eur.	1,785	241	1.2	1.92	26.9	42	23.78	227	2.2	1.3
NC17	Eur.	1,504	220	1.0	1.94	23.8	39	25.92	250	2.1	0.7
NC18	Eur.	1,589	170	0.8	1.93	18.4	42	26.48	317	1.2	0.6
NC22	Afr. Am.	1,693	195	0.9	2.53	29.7	44	25.96	205	1.4	0.7
NC25	Afr. Am.	2,274	175	0.8	2.47	25.8	62	27.31	211	2.8	1.8
NC26	Afr. Am., As.	974	120	0.6	2.55	18.0	26	26.67	355	1.3	0.4
NC27	As. (?)	1,268	267	1.3	1.96	29.2	34	26.80	199	1.7	0.6

* As provided by sperm bank. Afr. Am., of African American ancestry; Eur., of European ancestry; As., of Asian ancestry; (?), conflicting ancestry information given.

[‡] These numbers are the total number of aneuploidy events divided by the total number of cells multiplied by 100; cells can have more than one event.

[¶] Sum across all cells from all sperm donors.

[§] Median or mean across all individual cells from all sperm donors (31,228 measurements summarized).

Median or mean of aggregate metrics across samples (20 measurements summarized).

Median across all crossovers (813,122 measurements summarized).

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgments

We thank Giulio Genovese for suggestions on analyses, Evan Macosko for advice on technology development, and other members of the McCarroll lab, including Chris Whelan, Steven Burger, and Bob Handsaker for advice. We thank Mark Daly, Joel Hirschhorn, Stephen Elledge, and Samantha Schilit for their insights, 10X Genomics for discussions about reagents, and Christina L. Usher and Christopher K. Patil for contributions to the manuscript text and figures. We thank the anonymous reviewers and others who commented on the preprint version of this article for their input. This work was supported by R01 HG006855 to S.A.M., by a Broad Institute NextGen award to S.A.M., and by a Harvard Medical School Program in Genetics and Genomics NIH Ruth L. Kirchstein training grant to A.D.B.

References

1. Broman KW & Weber JL Characterization of human crossover interference. *American journal of human genetics* 66, 1911–1926, doi:10.1086/302923 (2000). [PubMed: 10801387]
2. Coop G, Wen X, Ober C, Pritchard JK & Przeworski M. High-resolution mapping of crossovers reveals extensive variation in fine-scale recombination patterns among humans. *Science* 319, 1395–1398, doi:10.1126/science.1151851 (2008). [PubMed: 18239090]
3. Halldorsson BV et al. Characterizing mutagenic effects of recombination through a sequence-level genetic map. *Science* 363, doi:10.1126/science.aau1043 (2019).
4. Kong A et al. A high-resolution recombination map of the human genome. *Nature genetics* 31, 241–247, doi:10.1038/ng917 (2002). [PubMed: 12053178]
5. Kong A et al. Common and low-frequency variants associated with genome-wide recombination rate. *Nature genetics* 46, 11–16, doi:10.1038/ng.2833 (2014). [PubMed: 24270358]
6. Kong A et al. Fine-scale recombination rate differences between sexes, populations and individuals. *Nature* 467, 1099–1103, doi:10.1038/nature09525 (2010). [PubMed: 20981099]

7. Myers S, Bottolo L, Freeman C, McVean G & Donnelly P. A fine-scale map of recombination rates and hotspots across the human genome. *Science* 310, 321–324, doi:10.1126/science.1117196 (2005). [PubMed: 16224025]
8. Nagaoka SI, Hassold TJ & Hunt PA Human aneuploidy: mechanisms and new insights into an age-old problem. *Nature reviews. Genetics* 13, 493–504, doi:10.1038/nrg3245 (2012).
9. Broman KW, Murray JC, Sheffield VC, White RL & Weber JL Comprehensive human genetic maps: individual and sex-specific variation in recombination. *American journal of human genetics* 63, 861–869, doi:10.1086/302011 (1998). [PubMed: 9718341]
10. Cheung VG, Burdick JT, Hirschmann D & Morley M. Polymorphic variation in human meiotic recombination. *American journal of human genetics* 80, 526–530, doi:10.1086/512131 (2007). [PubMed: 17273974]
11. Chowdhury R, Bois PR, Feingold E, Sherman SL & Cheung VG Genetic analysis of variation in human meiotic recombination. *PLoS genetics* 5, e1000648, doi:10.1371/journal.pgen.1000648 (2009).
12. Fledel-Alon A et al. Variation in human recombination rates and its genetic determinants. *PloS one* 6, e20321, doi:10.1371/journal.pone.0020321 (2011).
13. Brown PW et al. Meiotic synapsis proceeds from a limited number of subtelomeric sites in the human male. *American journal of human genetics* 77, 556–566, doi:10.1086/468188 (2005). [PubMed: 16175502]
14. Gruhn JR et al. Correlations between Synaptic Initiation and Meiotic Recombination: A Study of Humans and Mice. *American journal of human genetics* 98, 102–115, doi:10.1016/j.ajhg.2015.11.019 (2016). [PubMed: 26749305]
15. Gruhn JR, Rubio C, Broman KW, Hunt PA & Hassold T. Cytological studies of human meiosis: sex-specific differences in recombination originate at, or prior to, establishment of double-strand breaks. *PloS one* 8, e85075, doi:10.1371/journal.pone.0085075 (2013). [PubMed: 24376867]
16. Baudat F & de Massy B. Regulating double-stranded DNA break repair towards crossover or non-crossover during mammalian meiosis. *Chromosome Res* 15, 565–577, doi:10.1007/s10577-007-1140-3 (2007). [PubMed: 17674146]
17. Plug AW, Xu J, Reddy G, Golub EI & Ashley T. Presynaptic association of Rad51 protein with selected sites in meiotic chromatin. *Proceedings of the National Academy of Sciences of the United States of America* 93, 5920–5924 (1996). [PubMed: 8650194]
18. Ioannou D, Fortun J & Tempest HG Meiotic nondisjunction and sperm aneuploidy in humans. *Reproduction*, doi:10.1530/REP-18-0318 (2018).
19. Templado C, Uroz L & Estop A. New insights on the origin and relevance of aneuploidy in human spermatozoa. *Molecular human reproduction* 19, 634–643, doi:10.1093/molehr/gat039 (2013). [PubMed: 23720770]
20. Lynn A et al. Covariation of synaptonemal complex length and mammalian meiotic exchange rates. *Science* 296, 2222–2225, doi:10.1126/science.1071220 (2002). [PubMed: 12052900]
21. Wang S et al. Per-Nucleus Crossover Covariation and Implications for Evolution. *Cell*, doi:10.1016/j.cell.2019.02.021 (2019).
22. Hou Y et al. Genome analyses of single human oocytes. *Cell* 155, 1492–1506, doi:10.1016/j.cell.2013.11.040 (2013). [PubMed: 24360273]
23. Kirkness EF et al. Sequencing of isolated sperm cells for direct haplotyping of a human genome. *Genome research* 23, 826–832, doi:10.1101/gr.144600.112 (2013). [PubMed: 23282328]
24. Lu S et al. Probing meiotic recombination and aneuploidy of single sperm cells by whole-genome sequencing. *Science* 338, 1627–1630, doi:10.1126/science.1229112 (2012). [PubMed: 23258895]
25. Ottolini CS et al. Genome-wide maps of recombination and chromosome segregation in human oocytes and embryos show selection for maternal recombination rates. *Nature genetics* 47, 727–735, doi:10.1038/ng.3306 (2015). [PubMed: 25985139]
26. Wang J, Fan HC, Behr B & Quake SR Genome-wide single-cell analysis of recombination activity and de novo mutation rates in human sperm. *Cell* 150, 402–412, doi:10.1016/j.cell.2012.06.030 (2012). [PubMed: 22817899]

27. Miller D, Brinkworth M & Iles D. Paternal DNA packaging in spermatozoa: more than the sum of its parts? DNA, histones, protamines and epigenetics. *Reproduction* 139, 287–301, doi:10.1530/REP-09-0281 (2010). [PubMed: 19759174]
28. Macosko EZ et al. Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell* 161, 1202–1214, doi:10.1016/j.cell.2015.05.002 (2015). [PubMed: 26000488]
29. Zheng GX et al. Haplotyping germline and cancer genomes with high-throughput linked-read sequencing. *Nat Biotechnol* 34, 303–311, doi:10.1038/nbt.3432 (2016). [PubMed: 26829319]
30. Campbell CL, Furlotte NA, Eriksson N, Hinds D & Auton A. Escape from crossover interference increases with maternal age. *Nat Commun* 6, 6260, doi:10.1038/ncomms7260 (2015). [PubMed: 25695863]
31. Berg IL et al. PRDM9 variation strongly influences recombination hot-spot activity and meiotic instability in humans. *Nature genetics* 42, 859–863, doi:10.1038/ng.658 (2010). [PubMed: 20818382]
32. Myers S, Freeman C, Auton A, Donnelly P & McVean G. A common sequence motif associated with recombination hot spots and genome instability in humans. *Nature genetics* 40, 1124–1129, doi:10.1038/ng.213 (2008). [PubMed: 19165926]
33. Hinch AG et al. Factors influencing meiotic recombination revealed by whole-genome sequencing of single sperm. *Science* 363, doi:10.1126/science.aau8861 (2019).
34. Housworth EA & Stahl FW Crossover interference in humans. *American journal of human genetics* 73, 188–197, doi:10.1086/376610 (2003). [PubMed: 12772089]
35. Sun F et al. Human male recombination maps for individual chromosomes. *American journal of human genetics* 74, 521–531, doi:10.1086/382138 (2004). [PubMed: 14973780]
36. Oliver TR et al. Investigation of factors associated with paternal nondisjunction of chromosome 21. *Am J Med Genet A* 149A, 1685–1690, doi:10.1002/ajmg.a.32942 (2009). [PubMed: 19606484]
37. Page SL & Hawley RS Chromosome choreography: the meiotic ballet. *Science* 301, 785–789, doi:10.1126/science.1086605 (2003). [PubMed: 12907787]
38. Sun F et al. The relationship between meiotic recombination in human spermatocytes and aneuploidy in sperm. *Hum Reprod* 23, 1691–1697, doi:10.1093/humrep/den027 (2008). [PubMed: 18482994]
39. Ferguson KA, Wong EC, Chow V, Nigro M & Ma S. Abnormal meiotic recombination in infertile men and its association with sperm aneuploidy. *Hum Mol Genet* 16, 2870–2879, doi:10.1093/hmg/ddm246 (2007). [PubMed: 17728321]
40. Ma S, Ferguson KA, Arsovska S, Moens P & Chow V. Reduced recombination associated with the production of aneuploid sperm in an infertile man: a case report. *Hum Reprod* 21, 980–985, doi:10.1093/humrep/dei428 (2006). [PubMed: 16373411]
41. Savage AR et al. Elucidating the mechanisms of paternal non-disjunction of chromosome 21 in humans. *Hum Mol Genet* 7, 1221–1227 (1998). [PubMed: 9668162]
42. Tease C & Hulten MA Inter-sex variation in synaptonemal complex lengths largely determine the different recombination rates in male and female germ cells. *Cytogenet Genome Res* 107, 208–215, doi:10.1159/000080599 (2004). [PubMed: 15467366]
43. Wang S, Zickler D, Kleckner N & Zhang L. Meiotic crossover patterns: obligatory crossover, interference and homeostasis in a single process. *Cell Cycle* 14, 305–314, doi:10.4161/15384101.2014.991185 (2015). [PubMed: 25590558]
44. Zhang L, Liang Z, Hutchinson J & Kleckner N. Crossover patterning by the beam-film model: analysis and implications. *PLoS genetics* 10, e1004042, doi:10.1371/journal.pgen.1004042 (2014).
45. Zhang L et al. Topoisomerase II mediates meiotic crossover interference. *Nature* 511, 551–556, doi:10.1038/nature13442 (2014). [PubMed: 25043020]
46. Billings T et al. Patterns of recombination activity on mouse chromosome 11 revealed by high resolution mapping. *PloS one* 5, e15340, doi:10.1371/journal.pone.0015340 (2010).
47. Petkov PM, Broman KW, Szatkiewicz JP & Paigen K. Crossover interference underlies sex differences in recombination rates. *Trends in genetics : TIG* 23, 539–542, doi:10.1016/j.tig.2007.08.015 (2007). [PubMed: 17964681]

48. Bell AD, Mello CJ & McCarroll SA Sperm-seq wet lab protocol: sperm preparation and droplet-based sequencing library generation. Protocol Exchange, doi:10.21203/rs.2.22133/v1 (2020).
49. Bell AD et al. Analysis scripts for: Insights about variation in meiosis from 31,228 human sperm genomes. Zenodo, doi:10.5281/zenodo.2581595 (2019).
50. Bell AD et al. Recombination and aneuploidy data for: Insights about variation in meiosis from 31,228 human sperm genomes. Zenodo, doi:10.5281/zenodo.2581570 (2019).
51. Bell AD, Usher CL & McCarroll SA Analyzing Copy Number Variation with Droplet Digital PCR. *Methods Mol Biol* 1768, 143–160, doi:10.1007/978-1-4939-7778-9_9 (2018). [PubMed: 29717442]
52. Regan JF et al. A rapid molecular approach for chromosomal phasing. *PloS one* 10, e0118270, doi:10.1371/journal.pone.0118270 (2015).
53. Montag M, Tok V, Liow SL, Bongso A & Ng SC In vitro decondensation of mammalian sperm and subsequent formation of pronuclei-like structures for micromanipulation. *Mol Reprod Dev* 33, 338–346, doi:10.1002/mrd.1080330316 (1992). [PubMed: 1449801]
54. Samocha-Bone D et al. In-vitro human spermatozoa nuclear decondensation assessed by flow cytometry. *Molecular human reproduction* 4, 133–137 (1998). [PubMed: 9542970]
55. Taylor AC Titration of heparinase for removal of the PCR-inhibitory effect of heparin in DNA samples. *Mol Ecol* 6, 383–385 (1997). [PubMed: 9131813]
56. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. arXiv, doi:arXiv:1303.3997v2 (2013).
57. McKenna A et al. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome research* 20, 1297–1303, doi:10.1101/gr.107524.110 (2010). [PubMed: 20644199]
58. Van der Auwera GA et al. From FastQ data to high confidence variant calls: the Genome Analysis Toolkit best practices pipeline. *Curr Protoc Bioinformatics* 43, 11 10 11–33, doi:10.1002/0471250953.bi1110s43 (2013). [PubMed: 25431634]
59. Sherry ST et al. dbSNP: the NCBI database of genetic variation. *Nucleic Acids Res* 29, 308–311 (2001). [PubMed: 11125122]
60. Genomes Project C et al. A global reference for human genetic variation. *Nature* 526, 68–74, doi:10.1038/nature15393 (2015). [PubMed: 26432245]
61. Kent WJ et al. The human genome browser at UCSC. *Genome research* 12, 996–1006, doi:10.1101/gr.229102 (2002). [PubMed: 12045153]
62. Tyner C et al. The UCSC Genome Browser database: 2017 update. *Nucleic Acids Res* 45, D626–D634, doi:10.1093/nar/gkw1134 (2017). [PubMed: 27899642]
63. Genovese G et al. Clonal hematopoiesis and blood-cancer risk inferred from blood DNA sequence. *N Engl J Med* 371, 2477–2487, doi:10.1056/NEJMoa1409405 (2014). [PubMed: 25426838]
64. Bansal V & Bafna V. HapCUT: an efficient and accurate algorithm for the haplotype assembly problem. *Bioinformatics* 24, i153–159, doi:10.1093/bioinformatics/btn298 (2008). [PubMed: 18689818]
65. Selvaraj S, J, R. D., Bansal, V. & Ren, B. Whole-genome haplotype reconstruction using proximity-ligation and shotgun sequencing. *Nat Biotechnol* 31, 1111–1118, doi:10.1038/nbt.2728 (2013). [PubMed: 24185094]
66. Loh PR et al. Reference-based phasing using the Haplotype Reference Consortium panel. *Nature genetics* 48, 1443–1448, doi:10.1038/ng.3679 (2016). [PubMed: 27694958]
67. Loh PR, Palamara PF & Price AL Fast and accurate long-range phasing in a UK Biobank cohort. *Nature genetics* 48, 811–816, doi:10.1038/ng.3571 (2016). [PubMed: 27270109]
68. Handsaker RE, Korn JM, Nemesh J & McCarroll SA Discovery and genotyping of genome structural polymorphism by sequencing on a population scale. *Nature genetics* 43, 269–276, doi:10.1038/ng.768 (2011). [PubMed: 21317889]
69. Handsaker RE et al. Large multiallelic copy number variations in humans. *Nature genetics* 47, 296–303, doi:10.1038/ng.3200 (2015). [PubMed: 25621458]
70. Hyndman RJ & Khandakar Y. Automatic time series forecasting: the forecast package for R. *Journal of Statistical Software* 26, 1–22 (2008). [PubMed: 19777145]

71. R H et al. forecast: Forecasting functions for time series and linear models. R package version 8.4, <<http://pkg.robjhyndman.com/forecast>> (2018).
72. Kauppi L et al. Distinct properties of the XY pseudoautosomal region crucial for male meiosis. *Science* 331, 916–920, doi:10.1126/science.1195774 (2011). [PubMed: 21330546]
73. Kleckner N, Storlazzi A & Zickler D. Coordinate variation in meiotic pachytene SC length and total crossover/chiasma frequency under conditions of constant DNA length. *Trends in genetics : TIG* 19, 623–628, doi:10.1016/j.tig.2003.09.004 (2003). [PubMed: 14585614]
74. Revenkova E et al. Cohesin SMC1 beta is required for meiotic chromosome dynamics, sister chromatid cohesion and DNA recombination. *Nat Cell Biol* 6, 555–562, doi:10.1038/ncb1135 (2004). [PubMed: 15146193]
75. Zickler D & Kleckner N. Meiotic chromosomes: integrating structure and function. *Annu Rev Genet* 33, 603–754, doi:10.1146/annurev.genet.33.1.603 (1999). [PubMed: 10690419]
76. Wang S et al. Inefficient Crossover Maturation Underlies Elevated Aneuploidy in Human Female Meiosis. *Cell* 168, 977–989 e917, doi:10.1016/j.cell.2017.02.002 (2017). [PubMed: 28262352]
77. Blat Y, Protacio RU, Hunter N & Kleckner N. Physical and functional interactions among basic chromosome organizational features govern early steps of meiotic chiasma formation. *Cell* 111, 791–802 (2002). [PubMed: 12526806]

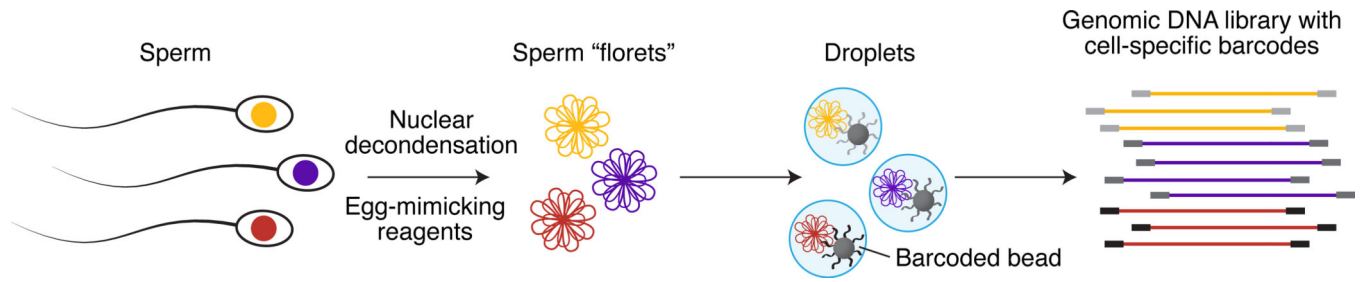


Fig. 1. "Sperm-seq" overview.
Schematic of our droplet-based single-sperm sequencing method.

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

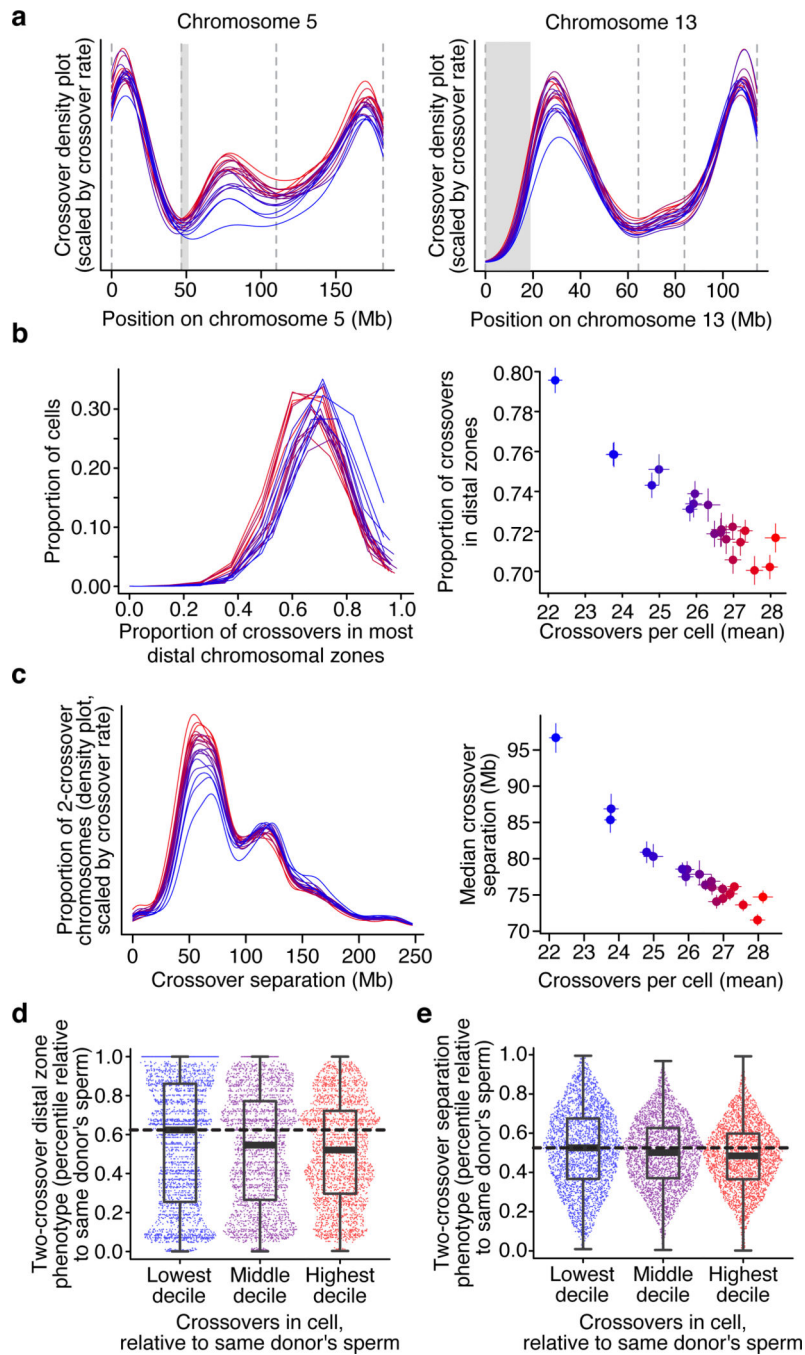


Fig 2. Variation in crossover positioning and crossover separation (interference). Color indicates crossover rate of donor or cell (blue: low, red: high). **a**, Crossover location density plots for each donor ($n = 20$). Dashed gray vertical lines: crossover zone boundaries. **b-e**, Crossover positioning and separation (interference) on chromosomes with two crossovers. **b-c**, Inter-individual variation among $n = 20$ sperm donors. Error bars: 95% confidence intervals. **b**, Left, per-cell proportion of crossovers in the most distal crossover zones (Kruskal–Wallis chi-squared = 1,034, $df = 19$, $p = 2 \times 10^{-207}$). Right, mean crossover rate (x axis) vs. the proportion of all crossovers (on two-crossover chromosomes) occurring

in distal zones (y axis, total proportion) (Pearson's $r = -0.95$, two-sided $p = 8 \times 10^{-11}$). **c**, Left, density plot of separation between consecutive crossovers (Kruskal–Wallis chi-squared = 1,792, $df = 19$, $p < 10^{-300}$). Right, mean crossover rate (x axis) vs. median crossover separation (y axis) on two-crossover chromosomes (Pearson's $r = -0.95$, two-sided $p = 7 \times 10^{-11}$). **d-e**, Among-cell covariation of crossover rate with distal zone use (**d**) or crossover interference (**e**). Phenotypes are analyzed as percentiles relative to sperm from the same donor. Boxplots: midpoints, medians; boxes, 25th and 75th percentiles; whiskers, minima and maxima. **d**, Single-cell distal-zone use (the proportion of crossovers on two-crossover chromosomes that are in the most distal zones) vs. crossover rate (n cells per decile = 3,152, 3,080, 3,101 for first, fifth, and tenth deciles, respectively; Mann–Whitney $W = 5,271,934.5$, two-sided $p = 2 \times 10^{-9}$ between first and tenth deciles.) **e**, Single-cell crossover-separation (the median of all fractions of a chromosome separating consecutive two-crossover chromosome crossovers in each cell) vs. crossover rate (Mann–Whitney $W = 148,548,161$, two-sided $p = 3 \times 10^{-53}$ between first [$n = 11,658$] and tenth [$n = 23,154$] deciles; all inter-crossover separations used in test).

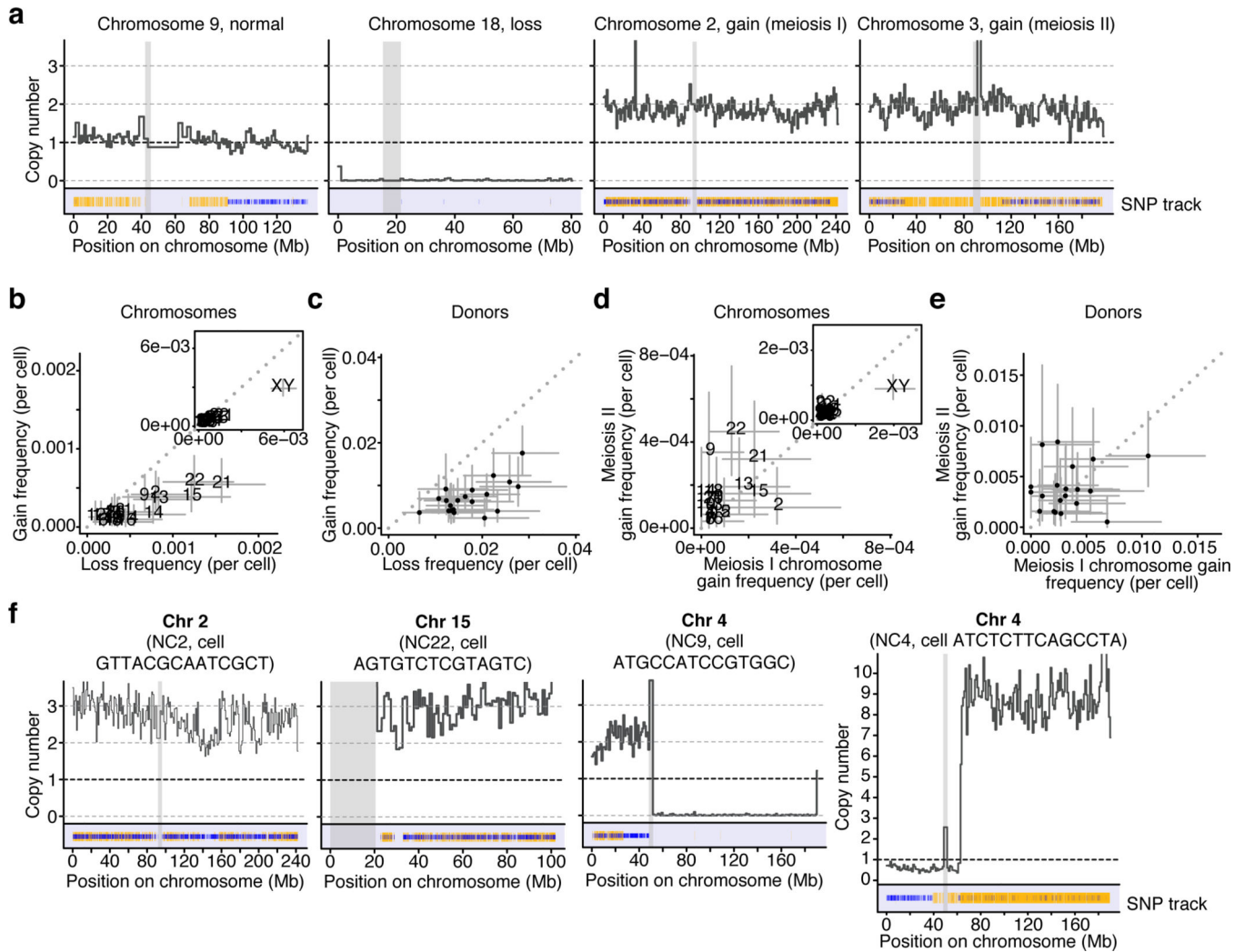


Fig. 3. Aneuploidy in sperm from 20 sperm donors.

a, Example chromosomal ploidy analyses. Thick dark gray line: DNA copy number measurement (normalized sequence coverage in 1 Mb bins); blue (haplotype 1) and yellow (haplotype 2) vertical lines: observed heterozygous SNP alleles, plotted with 90% transparency; gray vertical boxes: centromeres (hg38). **b-e**, Frequencies (number of events divided by number of cells) of various aneuploidy categories. $n = 23$ chromosomes (**b, d**) and $n = 20$ donors (**c, e**). Error bars are 95% binomial confidence intervals. **b**, Frequencies of whole-chromosome losses (x axis) vs. gains (y axis) for each chromosome (excluding XY Pearson's $r = 0.88$, two-sided $p = 7 \times 10^{-8}$; including XY [inset] Pearson's $r = 0.99$, two-sided $p < 10^{-300}$). **c**, Per-sperm-donor aneuploidy rates (axes as in **b**) (excluding XY [not shown] Pearson's $r = 0.51$, two-sided $p = 0.02$; including XY Pearson's $r = 0.62$, two-sided $p = 0.003$). **d**, Frequencies of whole-chromosome gains occurring during MI (x axis) and MII (y axis) for each chromosome (excluding XY Pearson's $r = 0.32$, two-sided $p = 0.15$; including XY [inset] Pearson's $r = 0.85$, two-sided $p = 3 \times 10^{-7}$). **e**, Frequencies of whole-chromosome gains occurring during MI (x axis) and MII (y axis) for each donor (axes as in **d**) (excluding XY [not shown] Pearson's $r = 0.06$, two-sided $p = 0.80$; including XY

Pearson's $r = 0.17$, two-sided $p = 0.47$). **f**, Example genomic anomalies detected in sperm cells, plotted as in **(a)**.

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

Extended Data Table 1.

Sperm donor and single-sperm sequencing characteristics and results.

Donor	Ancestry [*]	Cells (number excluding cell and bead doublets)	Reads per cell (median, thousands)	Genome covered per cell (median, percent)	Heterozygous SNPs in genome (millions)	Unique heterozygous SNP alleles observed per cell (median, thousands)	Crossovers observed (total, thousands)	Crossovers per cell (mean)	Resolution of crossovers (kb, median)	Autosomal aneuploidy events (percent of cells) [†]	Sex chromosome aneuploidy events (percent of cells) [‡]
Overall	--	31,228 [‡]	211 [§]	1.0	--	24.6 [§]	813 [‡]	26.1	240	1.6 [§]	0.9 [§]
NC1	Eur.	982	284	1.4	1.95	31.6	26	26.31	189	1.5	0.6
NC2	Eur.	1,680	163	0.8	1.98	18.2	37	22.19	307	2.0	0.7
NC3	Eur.	1,289	190	0.9	1.94	21.5	36	28.13	260	1.7	0.7
NC4	Eur.	1,482	243	1.1	1.98	26.8	40	26.98	243	1.1	0.5
NC6	Afr. Am.	1,370	154	0.8	2.53	23.8	38	27.57	253	0.7	0.3
NC8	As.	1,663	304	1.5	1.81	30.9	45	26.98	229	3.1	0.5
NC9	As.	1,894	245	1.2	1.79	25.6	53	27.98	231	0.8	1.5
NC10	As.	1,154	224	1.1	1.82	23.3	29	24.99	257	1.5	0.3
NC11	Eur.	1,930	202	1.0	1.92	22.8	50	25.82	242	1.3	0.4
NC12	Eur.	2,145	179	0.9	1.91	20.6	51	23.76	270	1.2	1.7
NC13	Eur.	1,514	259	1.2	1.92	28.3	41	27.19	202	0.9	1.0
NC14	Eur.	1,336	296	1.4	1.92	32.4	36	26.65	175	2.5	1.2
NC15	Eur.	1,702	211	1.0	1.93	23.2	42	24.80	268	1.0	0.9
NC16	Eur.	1,785	241	1.2	1.92	26.9	42	23.78	227	2.2	1.3
NC17	Eur.	1,504	220	1.0	1.94	23.8	39	25.92	250	2.1	0.7
NC18	Eur.	1,589	170	0.8	1.93	18.4	42	26.48	317	1.2	0.6
NC22	Afr. Am.	1,693	195	0.9	2.53	29.7	44	25.96	205	1.4	0.7
NC25	Afr. Am.	2,274	175	0.8	2.47	25.8	62	27.31	211	2.8	1.8
NC26	Afr. Am., As.	974	120	0.6	2.55	18.0	26	26.67	355	1.3	0.4
NC27	As. (?)	1,268	267	1.3	1.96	29.2	34	26.80	199	1.7	0.6

^{*} As provided by sperm bank. Afr. Am., of African American ancestry; Eur., of European ancestry; As., of Asian ancestry; (?), conflicting ancestry information given.

[†] These numbers are the total number of aneuploidy events divided by the total number of cells multiplied by 100; cells can have more than one event.

[‡] Sum across all cells from all sperm donors.

[§] Median or mean across all individual cells from all sperm donors (31,228 measurements summarized).

Median or mean of aggregate metrics across samples (20 measurements summarized).
Median across all crossovers (813,122 measurements summarized).

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript