

Free-living human cells reconfigure their chromosomes in the evolution back to uni-cellularity

Jin Xu^{1,2,6}, Xinxin Peng^{3,6}, Yuxin Chen^{2,6}, Yuezheng Zhang³, Qin Ma^{3,4}, Liang Liang^{3,4}, Ava C. Carter¹, Xuemei Lu^{3*} and Chung-I Wu^{2,5*}

Author affiliation:

¹Center for Personal Dynamic Regulomes and Program in Epithelial Biology, Stanford University School of Medicine, Stanford, California 94305, USA

²State Key Laboratory of Biocontrol, School of Life Sciences, Sun Yat-Sen University, Guangzhou 510275, China.

³Laboratory of Genome Variations and Precision Biomedicine, Beijing Institute of Genomics, Chinese Academy of Sciences, Beijing 100101, P.R. China.

⁴University of Chinese Academy of Sciences, Beijing 100049, China

⁵Department of Ecology and Evolution, University of Chicago, Chicago, Illinois 60637, USA

⁶ These authors contributed equally to this study.

*Corresponding authors: E-mail: ciwu@uchicago.edu or luxm@big.ac.cn

Lead Contact: ciwu@uchicago.edu

Keywords: sex chromosome evolution; dosage compensation; multi-cellularity; uni-cellularity; cancer cell evolution

Abstract

Cells of multi-cellular organisms evolve toward uni-cellularity in the form of cancer and, if humans intervene, continue to evolve in cell culture. During this process, gene dosage relationships may evolve in novel ways to cope with the new environment and may regress back to the ancestral uni-cellular state. In this context, the evolution of sex chromosomes vis-a-vis autosomes is of particular interest. Here, we report the chromosomal evolution in ~600 cancer cell lines. Many of them jettisoned either Y or the inactive X; thus, free-living male and female cells converge by becoming “de-sexualized”. Surprisingly, the active X often doubled, accompanied by the addition of one haploid complement of autosomes, leading to an X:A ratio of 2:3 from the extant ratio of 1:2. Theoretical modeling of the frequency distribution of X:A karyotypes suggests that the 2:3 ratio confers a higher fitness and may reflect aspects of sex chromosome evolution.

Introduction

Genomes of multi-cellular organisms evolve to ensure the survival and reproduction of the whole organisms. With human interventions akin to domestication, hundreds of cell lines survive as free-living cells that are not organized into tissues, organs or individuals[1]. Evolution in such a quasi-unicellular state may be very different from the evolution as multi-cellular entities. Most cell lines are cancerous in origin but a few are derived from normal tissues[2]. Regardless of their origin, they have all evolved characteristics for survival in the unicellular state that is distinct from their natural environments. Cell lines derived from cancer tissues are usually karyotypically less stable than normal cell lines. While this instability may impose a cost, it also permits cancer cell lines to evolve new karyotypes, including polyploidy, more readily than normal cell lines could.

Tumorigenesis has been increasingly viewed as a process of evolution, rather than merely pathological conditions[3][4]. This “ultra-microevolutionary process” is subjected to similar rules including mutation, genetic drift, migration and selection that govern organismal evolution[5]. While this process usually ends when the organism dies, cell lines in the cultured state will continue

51 to evolve. Much like the diversity unleashed by domestication, cultured cell lines, which can be
52 considered “domesticated”, may be informative about the evolutionary potentials at the cellular
53 level.

54

55 In this quasi-unicellular state, gene dosage has been observed to change extensively as
56 polyploidy, aneuploidy (full or partial) and various copy number variations (CNVs) are common in
57 cancer cell lines[6]. Since these cells lines are derived from somatic tissues of man or woman
58 (referred to as male and female cells, for simplicity), they should be different in their sex
59 chromosomes in relation to the autosomes (A’s). Nevertheless, the possibility of separate
60 evolutionary paths has not been raised before. Somatic cells have an inactive X chromosome in
61 females and a Y chromosome in males[7]. Since cell lines presumably do not need sexual characters,
62 we ask how the X:A relationship might have evolved in both male and female cells. More generally,
63 we ask whether the evolution in this relationship may shed light on the emergence of mammalian sex
64 chromosomes and the subsequent evolution.

65

66 In this study, we analyze 620 cancer cell lines that have been genotyped using SNP arrays[8].
67 Among them, 279 are derived from female tissues and 341 from male tissues. We observed the
68 elimination of the Y and the inactive X chromosome, followed by the evolution toward a new
69 equilibrium with 2 active X chromosomes and 3 sets of autosomes (2X:3A). We discuss the
70 implication of these findings for the evolution of sex chromosome, the transition between uni- and
71 multi-cellularity and cancers biology.

72

73 **Results**

74

75 **Convergent sex chromosomes evolution between sexes**

76 A most common form of genomic changes in cell lines is the loss of heterozygosity (LOH) when
77 one of the two homologous chromosomes is eliminated[6]. We therefore examine single nucleotide
78 polymorphisms (SNPs) across the 620 cell lines for occurrences of LOH on each autosome and the X
79 chromosome. Male and female cell lines are separately analyzed.

80

81 **Figure 1A** shows the LOH frequency for each autosome (black dots) and the red dot represents
82 the sex chromosomes (X in female and Y in male). For autosomes, the percentages of LOH are
83 remarkably similar between sexes, with a correlation coefficient of 0.94 among 620 cell lines. There
84 is a slight tendency for the smaller autosomes to have higher LOH rate ($R \sim -0.4$, $p \sim 0.046$, Figure
85 S1). The median percentage of LOH is about 13% for autosomes. However, the losses of X (36% in
86 females) and Y (40% in males) stand out. Given its rank as the 7th largest chromosome, the X is not
87 expected to be lost in more than 15% of cell lines. Since the X expression is not lost, we infer that
88 it’s the inactive X(or Xi) that is eliminated.

89

90 Female lines lose the inactive X (Xi) and male lines lose the Y chromosome at a higher rate than
91 other chromosomes. The two sexes may thus be expected to converge toward having a single sex
92 chromosome. Furthermore, given that spontaneous LOH is not infrequent and the loss cannot be
93 regained, long term cultures might evolve to complete LOH for sex chromosomes as well as
94 autosomes. The genome-wide low rate of LOH suggests selection holding back such changes. The
95 strong correlation between sexes further reflects a balance between the production and elimination of
96 LOH’s, likely involved natural selection.

97

98 A most unexpected finding is that, accompanying the loss of the Y or Xi, an extra X chromosome
99 is often gained. **Figure 1B** shows approximately equal numbers of male cell lines with one or two X
100 chromosomes (partial X aneuploidy not counted). This extra X is active because the inactivating

101 XIST lncRNA is silenced in male cell lines (**Figure 1C**), consistent with previous findings[9]. XIST
102 does not become activated in free-living cells that do not already express this. The expression of X-
103 linked genes is higher in those male lines with two X's than in those with one X and the up-
104 regulation occurs along the length of the X chromosome (**Figure 1D**).

105
106 The pattern is more complex in female lines which, in their original state, contain an X_a and an
107 X_i, the latter expressing XIST[10][11][12]. We focus on female lines that experienced LOH of the
108 X, which should be genetically equivalent to male lines that have lost the Y. These female lines
109 indeed evolve in a manner identical with the male lines. First, female lines that have gained an X are
110 almost as frequent as lines with one X, much like the male lines with one vs. two X_a's (**Figure 1E**).
111 Second, female lines with an additional X do not express Xist and all X's can thus be presumed
112 active (**Figure 1F**)[13]. As in male lines, the X does not switch its state after chromosome
113 duplication.

114
115 Cancer cell lines usually have high rate of aneuploidy and could be heterogeneous within the
116 line, thus making its status difficult to assess. To assess the level of within-line heterogeneity, we
117 chose two representative cell lines to count the X chromosomes in individual cells using fluorescent
118 in situ hybridization (FISH). The two lines are A549 (a male cell line from adenocarcinomic
119 alveolar basal epithelium) and HeLa (a female cervical cancer cell line). Neither line expresses XIST
120 (Table S2), suggesting that all X chromosomes are active. **Figure 2A-B** shows results from
121 individual A549 and HeLa cells with two and three X's. **Figure 2C-D** shows the X karyotype
122 distributions. While there is a modest degree of heterogeneity within each line, almost all cells have
123 two or more active X chromosomes. While labor intensity of assays and cell availability limited our
124 sample size, we nevertheless can conclude that within-cell line heterogeneity does not seem to
125 undermine our conclusions.

126
127 **Evolution toward a new X:A expression ratio ($E_{X/A}$)**

128 With an extra copy of active X, the “expression phenotype” is expected to change. The ratio of
129 the median gene expression on the X to that on the autosomes($E_{X/A}$) is of particular interest. $E_{X/A}$ has
130 been reported to be around 0.5~0.8 for normal mammalian tissues[14–16]. We assayed $E_{X/A}$ by
131 separating lines derived from cancerous and normal tissues. **Figure 3A** shows that $E_{X/A}$ distributions
132 center on ~0.84 in normal cell lines and on 1 in cancerous cell lines. Given the controversy in the
133 assay of $E_{X/A}$, we also varied the threshold for counting expressed transcripts (see Materials and
134 Methods). By varying the threshold (**Figure 3B**), $E_{X/A}$ ranges from 0.78 to 1.05 in normal cell lines
135 but is consistently higher by approximately 15% in cancer cell lines. The same pattern is seen in the
136 RNA-seq data (Figure S2).

137
138
139 **The concerted evolution of autosomes as a set**

140 While sex chromosomes evolve, autosomes should also evolve. Since the generation of
141 aneuploidy may happen independently for each autosome, a key question is whether selection
142 operates on the autosomes as a set. Does natural selection favor cells that have full sets of
143 autosomes?

144 **Figure 4A** shows the distribution of chromosome number across the 620 cell lines we studied.
145 Apparently, cancerous cell lines acquire autosomes during evolution. The distribution of ploidy
146 ($n=22$) number shows peaks at 2 and 3, indicates many cell lines appear to be in transition between
147 full diploidy and triploidy of 44 and 66 autosomes. Similarly, the majority of sublines of HeLa cells
148 we examined have 55-75 chromosomes centering about the triploid count of 69 (Figure S4A).
149 Indeed, autosomes appear to exist as a full complement with $n=22$. Although autosomes may evolve
150 as a set, cells most likely add one autosome at a time. It is hence desirable to track each chromosome

151 individually. Single cells were individually isolated from a HeLa cell line and subsequently grown to
152 a sub-line of 10^6 cells. We subjected 6 such sub-lines to whole genome sequencing such that each
153 chromosome can be tracked individually. Smaller chromosomes are indeed more erratic in their
154 numbers in cell lines. Only the largest 14 chromosomes (13 autosomes and X), which together
155 account for ~75% of the genome, are used to test the convergence of autosomes. The cutoff is based
156 on the observation that chromosome 13 is the largest autosome yielding viable trisomic new-
157 borns[17–19]. We reason that, if whole organisms can survive trisomy, the fitness consequence of
158 the particular aneuploidy would probably be very small at the cellular level.

159

160

161 In all 6 lines, each of the 13 autosomes has 2 – 4 copies, ranging from an average of 2.62 to 3.23
162 (Table S1). If each autosome behaves independently, the number of autosomes that increase by x
163 copies ($x = 0, 1, 2$ etc.) should follow a Poisson distribution with a mean of λ . Two different lines,
164 with $\lambda = 10/13$ and $\lambda = 16/13$, are shown in **Figure 4B and C**. In the former, all cells have $x = 0$ or
165 $x=1$ and, in the latter, all cells have $x = 1$ or $x=2$ (Table S1). The data suggest that each autosome
166 increases by one copy and only after all of the 13 autosomes have gained an extra copy do further
167 increases continue. Figure S4B shows the composite distribution of the five lines with $\lambda < 1$. The
168 pattern, like that of **Figure 4B**, is statistically significant ($P = 0.0021$ by the χ^2 test) with an excess at
169 $x = 1$. These results suggest that the larger autosomes evolve cohesively as a set. With autosomes
170 evolving as a cohesive unit, X:A can be represented by whole numbers of 1:2, 2:3 etc.

171

172 **Evolution of the C(Xa:A) ratio underlying $E_{X/A}$**

173 We now summarize the evolution of cell lines by their C(Xa:A) genotypes. C(Xa:A) is the
174 number of active X chromosomes and the ploidy number of autosomes (in multiples of 22) and is
175 equal to C(1,2) in normal cells. For the purpose of counting on active Xa's, data from most male
176 lines are usable. For female lines, only data from the LOH lines of the X can be used. Between the
177 two sexes, C(Xa:A) distributions are very similar and the combined distribution is used in the
178 analysis (Figure S4C).

179 Shown in **Figure 4D**, most lines have the C(1:2) or C(2:3) genotype which together account for
180 2/3 of the lines. Given that C(1:2) is the starting genotype, its common occurrence at 37.4% is not
181 surprising. The high frequency of C(2:3), however, is unexpected. To reach C(2:3) from the starting
182 point of C(1:2), cells should evolve to either C(2:2) or C(1:3) first, but neither genotype is commonly
183 seen in these cells lines. In contrast, C(2:3) at 29.2% is the second most common genotype. If we
184 include the two genotypes, C(2:4) and C(3:3), that are derivatives of C(2:3), this inclusive C(2:3)
185 cluster is the most common genotype. The model of the next section helps to interpret the
186 observation.

187

188 **A model for the evolution of free-living cells**

189 The pathways of chromosomal evolution can be diagrammed as a series steps in **Figure 5A**.
190 Each node represents a C(Xa:A) genotype, the abundance of which is reflected in the size of the
191 node. Thicker arrows indicate faster transitions which add/delete one X while the thinner arrow
192 denotes the slower transition of adding/deleting the whole set of autosomes. The fitness of each
193 genotype, W , is assumed to be determined by the Xa/A ratio. In general, one would expect the wild
194 type (W_1) to be the fittest genotype and we particularly wish to know whether that is indeed the case
195 here.

196

197 We first model the evolution under strict neutrality where all nodes have the same fitness. For
198 simplicity, genotypes are grouped into 3 clusters centering around the 3 dominant genotypes, C(1:2),
199 C(2:2) and C(2:3), the frequencies of which are x_1 , x_2 and x_3 , respectively. Each cluster consists of
200 the dominant genotype as well as the less common ones adjacent to it (see **Figure 5A**). For instance,

201 x2 is the sum of the frequencies of C(2:2) and C(3:2) and x1 is those of C(1:2), C(1:1) and half of
 202 C(1:3). The frequency of the last one, being adjacent to both C(1:2) and C(2:3), is split between the
 203 two clusters. Tallying up the numbers in **Figure 4D**, we obtain $x_1 = 0.41$, $x_2 = 0.092$ and $x_3 = 0.482$
 204 with a total of 0.984, excluding the marginal genotypes. The analysis below can be expanded to
 205 account for each genotype separately. The transitions between clusters are defined as follows:
 206

$$x_1(T) \begin{matrix} u \\ au \end{matrix} \rightleftharpoons x_2(T) \begin{matrix} v \\ bv \end{matrix} \rightleftharpoons x_3(T)$$

207 where u and v are the transition rates and $x_i(T)$ is the frequency of cluster i at time T . Let $X(T)$ be the
 208 vector of $[x_1(T), x_2(T), x_3(T)]$, expressed as
 209
 210

$$211 \quad X(T) = X(0) \begin{bmatrix} 1-u & u & 0 \\ au & 1-au-v & v \\ 0 & bv & 1-bv \end{bmatrix}^T \quad (1)$$

212
 213 When $T \gg 0$,

$$214 \quad [x_1(T), x_2(T), x_3(T)] \sim [ab, b, 1] / z \quad (2)$$

215 where $z = ab + b + 1$. The genotype frequencies evolve toward the equilibrium, $[ab, b, 1] / z$, which
 216 depends on a and b , but not u and v . We posit that $a > 1$ and $b > 1$ because, as the chromosome
 217 number increases, the probability of chromosome gain/loss increases as well. By Eq. 2, $x_1(T) >$
 218 $x_2(T) > x_3(T)$ when $T \gg 0$. In short, the relative frequency should be in the descending order of
 219 C(1:2), C(2:2) and C(2:3) if there is no fitness difference among genotypes. This predicted
 220 inequality at $T \gg 0$ is very different from the observed trend.
 221
 222

223 Eq. 2 assumes that cell lines have been evolving long enough to approach this equilibrium. A
 224 more appropriate representation should be $X(T)$ where T can reflect the time a cell line has been in
 225 culture. It is algebraically simpler if T is measured by the rate of chromosomal changes, u or v , rather
 226 than by the actual cell generation (Eq. 1, **Figure 5B** and legends). We also assume $u > v$ as u
 227 involves only the X but v involves the whole set of autosomes. With the initial condition of $X(0) =$
 228 $[1, 0, 0]$, **Figure 5B** shows that the C(2:3) cluster approaches the equilibrium more slowly than the
 229 other two clusters. Therefore, the observed high frequency of the C(2:3) cluster ($x_3 = 0.482$ vs. $x_1 =$
 230 0.41 and $x_2 = 0.092$) is incompatible with a neutrally evolving model of chromosome numbers. The
 231 discrepancy is true at all time points and is more pronounced at smaller T 's.
 232

233 Rejecting the neutral evolution model, we now incorporate fitness differences into **Figure 4A**
 234 with $W_1 = 1$ [for C(1:2) and C(2:4)], $W_2 = 1+s$ [for C(2:2)] and $W_3 = 1+t$ [for C(2:3)] where s and t
 235 can either be positive or negative. Here, we add a fourth genotype, C(2:4). In the supplement, we
 236 model 4 genotypes with $x_1 - x_4$ for the frequencies of C(1:2), C(2:2), C(2:3) and C(2:4)
 237 respectively. An expanded transition matrix is used to model selection, followed by a normalization
 238 step (Eq. S1). The solution in the form of $X(T) = X(0) M^T$ is given in Eq. S2 and the equilibrium
 239 $X(T)$ is given in Eq. S3.
 240

241 We are particularly interested in whether $t > 0$ in the 4-genotype model, i.e., whether C(2:3) has a
 242 higher fitness than the wild type, C(1:2). We observe that $[x_1, x_2, x_3, x_4] = [0.374, 0.087, 0.292,$
 243 $0.128]$ where $x_3 = 0.292$ is more than 3 times higher than $x_2 = 0.087$ and is close to $x_1 = 0.374$. Eq.
 244 S3 shows that $s < 0$ is necessary for x_2 to be smaller than x_3 , and $t > 0$ is necessary for x_3 to be close to
 245 x_1 (see Supplement). **Figure 5C** is an example in which $s = -0.5$ and $t = 0.5$. The equilibrium at T
 246 $\gg 0$ is indeed close to the observed values.

247 In conclusion, it appears that the extant state in multicellular organisms of C(1:2) is not the fittest
248 genotype for free-living mammalian cells. The observed genotypic distributions suggest that C(2:3)
249 may have a higher fitness than the wild type, C(1:2).

250

251 Discussion

252 Free-living mammalian cells like all living things speed up the evolution when the environment
253 changes. The practice of cell culturing, however, is to slow down the evolution in order to preserve
254 cell lines' usefulness as proxies for the source tissues. Nevertheless, changes are inevitable and the
255 evolution of sex chromosomes is but one example. It should be noted that cell lines derived from
256 cancerous tissues and normal tissues are different in one important aspect. Cell lines derived from
257 normal tissues generally do not undergo karyotypic changes at an appreciable rate[20–22]. They are
258 therefore much less responsive to selection in cultured conditions that favor new karyotypes. Cancer
259 cell lines, having been through more rounds of passages, have generally experienced stronger
260 selection more frequently than normal cell lines.

261

262 Our observations suggest that the extant X:A relationship (C(1:2)) may not be optimal for free-
263 living mammalian cells. The highest fitness peak, instead, appears to be closer to the karyotype of
264 C(2:3) as free-living cells reproducibly evolve toward this new karyotype. The fitness peaks in free-
265 living cells being different from that of the multi-cellular organisms is not unexpected. With many
266 possible conflicts between individual cells and the community of cells (i.e., the organism), the
267 interest of the community may lie in its ability to regulate the growth potential of its constituents.
268 Free-living cells, on the other hand, are driven by selection to realize their individual proliferative
269 capacity relative to other cells.

270

271 The convergence among these many cell lines to C(2:3) is unexpected in the context of cancer
272 evolution. The TCGA project (reference) has shown that cancer evolution is a process of divergence,
273 not convergence. Indeed, only 2 genes have been mutated in more than 10% of all cancer cases and
274 tumors of the same tissue origin from two different patients may often share no mutated genes at
275 all[5][23]. Therefore, the karyotypic convergence reported here is rather unusual.

276

277

278 We note that C(2:3) toward which cultured cells evolved happens to be the smallest possible
279 increase in the X/A ratio from C(1:2). The higher fitness of C(2:3) than C(1:2) in free-living cells
280 may lend new clues to the debate about the evolution of mammalian sex chromosomes[16][24]. With
281 X-inactivation, it has been suggested that $E_{X/A}$ could have been reduced, or even halved[14][24]. The
282 debate is about whether, and by how much, $E_{X/A}$ might have increased in evolution. Our observation
283 that free-living cells continue to evolve toward C(2:3) raised the possibility that the evolutionary
284 increase in $E_{X/A}$ has not been complete, in comparison with the ancestral $E_{X/A}$.

285

286 Finally, this study of cancerous cell lines may also have medical implications. The common
287 view that tumorigenesis is an evolutionary phenomenon posits that individual cells in tumors evolve
288 to enhance self-interest[3][4][25][26]. A corollary would be that tumorigenesis may have taken the
289 first few steps toward uni-cellularity. This extended view is supported by many expression studies
290 as well as the higher likelihood of obtaining cell lines from tumors than from normal tissues[2]. An
291 alternative view, posits that tumors remain multi-cellular in organization[27]. These different views
292 have been critically examined recently[5]. It is possible that cancer cells in vivo may have been
293 gradually evolving toward a new optimum. In that case, cancer cells in men and women are
294 converging in their sex chromosome evolution and become more efficient in proliferation in this new
295 de-sexualized state.

296 Figure legends

297
298 **Figure 1.** Convergence in sex chromosomes in culture human cells.

299 (A) Percentage of lines with LOH (loss of heterozygosity). Each black dot represents an autosome
300 and the red dot represents X and Y. LOH in male and female lines are separately displayed on the X
301 and Y-axes. (B) Percentage of cell lines with either one or two Xa's in male lines. (C) Expression
302 level of Xist in male cell lines, [Y] means with or without Y chromosome. The number of cell lines
303 is on the top of each bar. (D) Expression ratios of X-linked genes between Xa[Y] and XaXa[Y] cell
304 lines. Each grey dot represents a gene, and significant differences are indicated by black dots (t-test,
305 $p < 0.05$). (E) Percentage of cell lines with either one or two Xa's in female lines with whole X
306 chromosome LOH. Female lines with partial X's or non-LOH are not included because it's
307 ambiguous to assign the activation of X's. (F) Expression level of Xist in female cell lines of
308 different X karyotypes. XaO (female lines with a single X), XaXa (female lines with isodisomy of
309 X), XaXb (non-LOH female lines with heterozygous X's). The number of cell lines is on the top of
310 each bar. All lines except XaXb have very low levels of XIST, suggesting active X's. In XaXb lines,
311 the degree of X-inactivation is variable.

312
313 **Figure 2.** (A-B) Representative images of X chromosome FISH in the A549 cell line (A) with two
314 Xs and HeLa (B) with three Xs. DNA is stained with DAPI (blue), and the X chromosome is labeled
315 with Cy3 (red). (C-D) The distribution of the copy number of X's among cells from A549 ($n = 343$)
316 and HeLa ($n = 170$).

317
318 **Figure 3.** Increasing of expression ratio of X versus autosome (EX/A)
319 (A) EX/A distributions among normal (N) and cancer (C) cell lines. NF and NM (or CF and CM)
320 designate normal (or cancer) female and male lines. EX/A in cancer cell lines become larger than
321 those of the normal cell lines. Note that the expression in normal cell lines is narrowly distributed
322 and is close to that of the normal tissue when compared. Although the numbers of NF and NM lines
323 are much smaller than CF and CM lines (17 and 24 vs. 279 and 341), their EX/A distributions are
324 much tighter than in cancer cell lines. The actual counts correspond to kernel density are given in
325 Fig. S4. (B) EX/A ratio in CF, CM, NF and NM lines with filtering with three different cutoffs (see
326 methods). EX/A ratios are consistently higher in CF and CM lines than in NF and NM lines.

327
328 **Figure 4.** Autosomes change in a cohesive manner and coevolution of X and A.
329 (A) The density plot of autosome copy number among 620 cell lines shows peaks at 2 and 3 per
330 autosome. (B-C) The observed distributions of gain in copy number among autosomes in two HeLa
331 sublines. The expected Poisson distributions are also given for sublines with different means ($\lambda =$
332 $10/13, 16/13$; see text). (D) The percentages of C(Xa:A) types among the 620 cell lines.

333
334 **Figure 5** – A model of karyotypic evolution driven by fitness differences.
335 (A) Evolutionary pathways of chromosomal changes. Each node represents a karyotype C(Xa:A) and
336 the size roughly corresponds to its frequency. Cell fitness is assumed to be a function of the Xa/A
337 ratio, which is represented by the Y-axis. The four abundant karyotypes are shown by solid black
338 circles. Red arrows indicate faster changes in X and black arrows indicate slower changes in
339 autosome. Main transitions between the common karyotypes are indicated by thicker arrows. (B)
340 Changes in the frequencies of the three key genotypes as a function of time (T, expressed in units of
341 $1/v$) under fitness neutrality with all W_i 's = 1. The parameters for Eqs. 1 and 2 are $u = 10v$, $a = 2$ and
342 $b = 1.5$. Both the theoretical trajectories and the observed values are given. The C(2:3) cluster (x3)
343 is far more common in the observation than in the neutral model. (C) Changes in 4 karyotypic
344 frequencies under selection according to Eq. S3 with $s = -0.5$ and $t = 0.5$. All other conditions are

345 the same as above. Under selection, a reasonable agreement between the model and the observation
346 can be obtained.
347

348 **Materials and Methods**

349 **Chromosome number estimation of HeLa sub-lines.**

350 The processing of clonal expansion and whole genome sequencing of HeLa lines are
351 described at Zhang et. al. (unpublished data). For each line, the copies of each chromosome are
352 estimated according to the average sequencing depth by Control-FREEC, a tool for assessing copy
353 number using next generation sequencing data[28].

354

355 **Data collection**

356 Three large-scale datasets were used in this study[8,29,30].

357 Genome-wide SNP array data on cancer cell lines and a normal training set were downloaded
358 from The Wellcome Trust Sanger Institute under the data transfer agreement. Among the 755 cancer
359 cell lines, 645 (from 288 females and 357 males) with available gender information were used for
360 genotype information analysis in the present study. The processed data are in PICNIC output file
361 format, which includes information on genotype, loss of heterozygosity and absolute allelic copy
362 number segmentation[8]. Greenman et. al. developed the algorithm, PICNIC (Predicting Integral
363 Copy Number In Cancer), to predict absolute allelic copy number variation in cancer[8]. This
364 algorithm improved the normalization of the data and the determination of the underlying copy
365 number of each segment. It has been used for Affymetrix genome-wide SNP6.0 data from 755 cancer
366 cell lines, which were derived from 32 tissues. The Affymetrix Genome-Wide SNP Array 6.0 has 1.8
367 million genetic markers, including more than 900,000 single nucleotide polymorphism probes (SNP
368 probes) and more than 900,000 probes for the detection of copy number variation (CN probes).

369

370 The genome-wide gene expression data for 947 human cancer cell lines from 36 tumor types
371 were generated by Barretina et al[29]. as part of the cancer cell line Encyclopedia (CCLE) project
372 using Affymetrix U133 plus 2.0 arrays and are available from the CCLE project website
373 (CCLE_Expression_Entrez_2012-09-29.gct, <http://www.broadinstitute.org/ccle/home>). The
374 expression profiles of 768 cell lines with gender information, representing 337 females and 431
375 males, were used in this study. These cell lines were partially overlapped with the lines used in
376 Greenman et. al. Additionally, RNA-seq data from 41 lymphoblastoid cell lines from 17 females and
377 24 males were downloaded from GEO database (GSE16921)[30].

378

379 **LOH detection and copy number estimation**

380 We used the genotype information and absolute allelic copy number estimation generated
381 from PICNIC to infer LOH, as well as copy number, of a specific chromosome. As for a
382 chromosome, if $\geq 95\%$ of SNP sites were homologous we considered that there was a LOH(loss of
383 heterogeneity) event for this chromosome. Similarly, if $\geq 95\%$ of detected alleles on the chromosome
384 had a constant copy number of 0, 1, 2, 3 or 4, the copy number would be considered as the copy
385 number of the chromosome. The copy number of the Y chromosome was estimated separately. In
386 females, although all sites on Y chromosome should have yielded 0 copies, only $\sim 60\%$ of sites
387 detected by the Y chromosome probes showed a copy number of 0. This result indicated that several
388 X homologous regions on the Y were covered by $\sim 30\%$ of Y probes. Therefore, Y chromosome loss
389 was defined as when more than 60% of SNP probes from the Y chromosome showed a copy number
390 of 0.

391

392 **Sex chromosome genotype inference**

393 The expression level of *XIST* can be used as a proxy to distinguish the active X chromosome
394 from the silent one as this gene was expressed on the inactive X chromosome and functioned *in*
395 *cis*[13]. According to Greenman's and Barretina's studies, 496 cancer cell lines have both copy
396 number and expression data. As expected, *XIST* was silenced in male cell lines, as well as in females
397 with whole X chromosome LOH (Fig. 1C). Based on X chromosome LOH and copy number

398 information, we identified five genotypes, including XaO (female lines with one X-20 lines), XaXa
399 (female lines with isodisomy of X-17lines), XaXb (female lines with heterozygous for the X-28
400 lines), Xa[Y] (male lines with one X-53 lines) and XaXa[Y] (male lines with two X's-69 lines).

401

402 **C(Xa:A)(ratio of active X's to autosomes) calculation**

403 All male (341 lines) and female cell lines with whole X chromosome LOH (103 lines) were
404 employed for C(Xa:A) calculation. C(Xa:A) was defined as the ratio of absolute X copy number to
405 that of all autosomes.

406

407 **$E_{X/A}$ (ratio of X to autosomal expression) calculation**

408 $E_{X/A}$ was defined as the ratio of the expression of X-linked genes to that of autosomal ones.
409 The median values of expressed X-linked and autosomal genes were used to calculate $E_{X/A}$ in both
410 cancerous and normal cell lines. For the datasets from the Affymetrix U133 + 2.0 array, genes with
411 signal intensities ≥ 32 ($\log_2 > 5$) were considered to be expressed. While as for RNA-seq data, genes
412 with RPKM values ≥ 1 were considered to be expressed

413 Previous studies have shown that $E_{X/A}$ value may be affected by gene set used[15]. In addition,
414 several silent genes in normal tissues have been shown to be expressed in tumor tissues[31]. Those
415 genes were dominant on X chromosome, which could result in an increase of $E_{X/A}$. To exclude the
416 possibility that $E_{X/A}$ ratios may be biased in cancerous cell lines, gene sets for $E_{X/A}$ calculation were
417 first selected in normal cell lines by three criteria, with the same sets then selected in cancerous cell
418 lines. The three filtering criteria for gene set selection were RPKM > 0 , 1, and 5 in normal cell lines
419 (Fig. 2C).

420

421 **Differences in X-linked gene expression between Xa[Y] and XaXa[Y] lines**

422 To explore the impact of extra X chromosome on gene expression levels of X-linked genes,
423 118 cell lines with Xa[Y] and 109 cell lines with XaXa[Y] were used. T-test with Benjamini and
424 Hochberg adjusting method was employed to determine genes, the expression of which are
425 significantly changed due to an extra X copy. 648 detected X-linked genes are plotted in Fig. 2A.
426 The free statistical programming language R was used for the statistical analysis (version 3.0.1).

427

428 **X chromosome Fluorescence *in situ* hybridization**

429 HeLa cells (from the Culture Collection of the Chinese Academy of Sciences, Shanghai,
430 China) were cultured in DMEM (Life Technologies) supplemented with 10% fetal bovine serum
431 (FBS), 100 U/ml of penicillin, and 100 μ g/ml of streptomycin. A549 cells (from Mi-lab) were
432 cultured in RPMI-1640 (Life Technologies) with 10% fetal bovine serum (FBS), 100 U/ml of
433 penicillin, and 100 μ g/ml of streptomycin at 37°C with 5% CO₂. Approximately 2×10^6 cells were
434 seeded and cultured in 10 cm dishes with 10 ml growth medium as described above. To synchronize
435 the cells, 200 μ l of thymidine (100 mM) was added to the cells. After incubating for 14 hours, the
436 cells were washed twice with 10 ml PBS and then supplemented with 10 ml growth medium
437 containing deoxycytidine (24 μ M). After incubating for 2 hours, 10 μ l nocodazole (100 μ g/ml) was
438 added to the cells. The cells were incubated for an additional 10 hours.

439

440 After synchronization, cells were harvested and treated with 4 ml hypotonic solution (75 mM,
441 KCl) pre-warmed to 37°C for 30 min. The cells were then fixed via three immersions in fresh
442 fixative solution (3:1 methanol:acetic acid) (15 min each time). The fixed cell suspension was
443 spotted onto a clean microscope slide and allowed to air dry. We used the “XCyting Chromosome
444 Paints” and “Xcyting Centromere Enumeration Probe” (MetaSystems, Germany) for whole X
445 chromosomes and centromere of X chromosome fluorescence *in situ* hybridization (FISH) analysis,
446 respectively. Following the manufacturer's instructions, 10 μ l of probe mixture was added to the
447 prepared slide. The slide was then covered with 22 x 22 mm² cover slip and sealed with rubber

448 cement. Next, the slide was heated at 75°C for 2 min on a hotplate to denature the sample and probes
449 simultaneously, followed by incubation in a humidified chamber at 37°C overnight for hybridization.
450 After hybridization, the slide was washed in 0.4 x SSC (pH 7.0) at 72°C for 2 min, then in 2 x SSC
451 and 0.05% Tween-20 (pH 7.0) at room temperature for 30 seconds, before being rinsed briefly in
452 distilled water to avoid crystal formation. The slide was drained and allowed to air dry. Finally, 5 µl
453 DAPI (MetaSystems) was applied to the hybridization region and covered with a coverslip. The slide
454 was processed and captured using fluorescence microscopy as recommended (Olympus FV1000, 100
455 x objective). The number of Xs were counted for each individual cell. A total of 343 HeLa cells and
456 170 A549 cells were screened.

457

458 **Acknowledgments:** We thank Jian Lu and Xionglei He for comments on an earlier version of the
459 manuscript. Additionally, we thank Yang Shen, Yong E. Zhang, Rui Zhang and Wenfeng Qian for
460 numerous constructive discussions.

461

462 **Competing interests:** The authors have declared that no competing interests exist.

463

464 **Reference**

- 465 1. Bruce Alberts, Alexander Johnson, Julian Lewis, Martin Raff, Keith Roberts, and P.W.
466 (2002). Isolating Cells and Growing Them in Culture. In *Molecular Biology of the Cell* (New
467 York: Garland Science).
- 468 2. Hayflick, L. (1998). A Brief History of the Mortality and Immortality of Cultured Cells. *Keio*
469 *J. Med.* *47*, 174–182.
- 470 3. Nowell, P. (1976). The clonal evolution of tumor cell populations. *Science* (80-). *194*, 23–28.
- 471 4. Merlo, L.M.F., Pepper, J.W., Reid, B.J., and Maley, C.C. (2006). Cancer as an evolutionary
472 and ecological process. *Nat. Rev. Cancer* *6*, 924–935.
- 473 5. Wu, C.-I., Wang, H.-Y., Ling, S., and Lu, X. (2016). The Ecology and Evolution of Cancer:
474 The Ultra-Microevolutionary Process. *Annu. Rev. Genet.* *50*, annurev-genet-112414-054842.
- 475 6. Roschke, A. V., Tonon, G., Gehlhaus, K.S., McTyre, N., Bussey, K.J., Lababidi, S., Scudiero,
476 D. a., Weinstein, J.N., and Kirsch, I.R. (2003). Karyotypic Complexity of the NCI-60 Drug-
477 Screening Panel. *Cancer Res.* *63*, 8634–8647.
- 478 7. Charlesworth, B. (1991). Evoluton of Sex Chromosomes. *Science* (80-). *251*, 1030–3.
- 479 8. Greenman, C.D., Bignell, G., Butler, A., Edkins, S., Hinton, J., Beare, D., Swamy, S.,
480 Santarius, T., Chen, L., Widaa, S., *et al.* (2010). PICNIC: an algorithm to predict absolute
481 allelic copy number variation with microarray cancer data. *Biostatistics* *11*, 164–75.
- 482 9. Guttenbach, M., Koschorz, B., Bernthaler, U., Grimm, T., and Schmid, M. (1995). Sex
483 chromosome loss and aging: in situ hybridization studies on human interphase nuclei. *Am. J.*
484 *Hum. Genet.* *57*, 1143–50.
- 485 10. Chow, J.C., Yen, Z., Ziesche, S.M., and Brown, C.J. (2005). Silencing of the mammalian X
486 chromosome. *Annu. Rev. Genomics Hum. Genet.* *6*, 69–92.
- 487 11. Plath, K., Mlynarczyk-Evans, S., Nusinow, D. a, and Panning, B. (2002). Xist RNA and the
488 mechanism of X chromosome inactivation. *Annu. Rev. Genet.* *36*, 233–78.
- 489 12. Ng, K., Pullirsch, D., Leeb, M., and Wutz, A. (2007). Xist and the order of silencing. *EMBO*
490 *Rep.* *8*, 34–9.
- 491 13. Richardson, A.L., Wang, Z.C., De Nicolo, A., Lu, X., Brown, M., Miron, A., Liao, X.,
492 Iglehart, J.D., Livingston, D.M., and Ganesan, S. (2006). X chromosomal abnormalities in
493 basal-like human breast cancer. *Cancer Cell* *9*, 121–32.
- 494 14. Xiong, Y., Chen, X., Chen, Z., Wang, X., Shi, S., Wang, X., Zhang, J., and He, X. (2010).
495 RNA sequencing shows no dosage compensation of the active X-chromosome. *Nat. Genet.*
496 *42*, 1043–7.
- 497 15. Deng, X., Hiatt, J.B., Nguyen, D.K., Ercan, S., Sturgill, D., Hillier, L.W., Schlesinger, F.,
498 Davis, C. a, Reinke, V.J., Gingeras, T.R., *et al.* (2011). Evidence for compensatory
499 upregulation of expressed X-linked genes in mammals, *Caenorhabditis elegans* and
500 *Drosophila melanogaster*. *Nat. Genet.* *43*, 1179–85.
- 501 16. Kharchenko, P. V, Xi, R., and Park, P.J. (2011). Evidence for dosage compensation between
502 the X chromosome and autosomes in mammals. *Nat. Genet.* *43*, 1167-9–2.
- 503 17. Taylor, a I. (1968). Autosomal trisomy syndromes: a detailed study of 27 cases of Edwards’
504 syndrome and 27 cases of Patau’s syndrome. *J. Med. Genet.* *5*, 227–252.
- 505 18. Patterson, D. (2009). Molecular genetic analysis of Down syndrome. *Hum. Genet.* *126*, 195–
506 214.
- 507 19. Kleijer, W.J., van der Sterre, M.L.T., Garritsen, V.H., Raams, A., and Jaspers, N.G.J. (2006).
508 Prenatal diagnosis of the Cockayne syndrome: survey of 15 years experience. *Prenat. Diagn.*
509 *26*, 980–984.
- 510 20. Shirley, M.D., Baugher, J.D., Stevens, E.L., Tang, Z., Gerry, N., Beiswanger, C.M., Berlin,
511 D.S., and Pevsner, J. (2012). Chromosomal variation in lymphoblastoid cell lines. *Hum.*
512 *Mutat.* *33*, 1075–1086.
- 513 21. Frazer, K.A., Ballinger, D.G., Cox, D.R., Hinds, D.A., Stuve, L.L., Gibbs, R.A., Belmont,

- 514 J.W., Boudreau, A., Hardenbol, P., Leal, S.M., *et al.* (2007). A second generation human
515 haplotype map of over 3.1 million SNPs. *Nature* 449, 851–861.
- 516 22. Pickrell, J.K., Marioni, J.C., Pai, A.A., Degner, J.F., Engelhardt, B.E., Nkadori, E., Veyrieras,
517 J.-B., Stephens, M., Gilad, Y., and Pritchard, J.K. (2010). Understanding mechanisms
518 underlying human gene expression variation with RNA sequencing. *Nature* 464, 768–72.
- 519 23. Kandoth, C., McLellan, M.D., Vandin, F., Ye, K., Niu, B., Lu, C., Xie, M., Zhang, Q.,
520 McMichael, J.F., Wyczalkowski, M.A., *et al.* (2013). Mutational landscape and significance
521 across 12 major cancer types. *Nature* 502, 333–339.
- 522 24. Lin, F., Xing, K., Zhang, J., and He, X. (2012). Expression reduction in mammalian X
523 chromosome evolution refutes Ohno ’ s hypothesis of dosage compensation.
- 524 25. Chen, H., Lin, F., Xing, K., and He, X. (2015). The reverse evolution from multicellularity to
525 unicellularity during carcinogenesis. *Nat. Commun.* 6, 6367.
- 526 26. Chen, H., and He, X. (2016). The Convergent Cancer Evolution toward a Single Cellular
527 Destination. *Mol. Biol. Evol.* 33, 4–12.
- 528 27. Almendro, V., Marusyk, A., and Polyak, K. (2012). Cellular Heterogeneity and Molecular
529 Evolution in Cancer. *Annu. Rev. Pathol. Mech. Dis.* 8, 121023133009008.
- 530 28. Boeva, V., Popova, T., Bleakley, K., Chiche, P., Cappel, J., Schleiermacher, G., Janoueix-
531 Lerosey, I., Delattre, O., and Barillot, E. (2012). Control-FREEC: A tool for assessing copy
532 number and allelic content using next-generation sequencing data. *Bioinformatics* 28, 423–
533 425.
- 534 29. Barretina, J., Caponigro, G., Stransky, N., Venkatesan, K., Margolin, a a, Kim, S., Wilson,
535 C.J., Lehar, J., Kryukov, G. V, Sonkin, D., *et al.* (2012). The Cancer Cell Line Encyclopedia
536 enables predictive modelling of anticancer drug sensitivity. *Nature* 483, 603–607.
- 537 30. Cheung, V.G., Nayak, R.R., Wang, I.X., Elwyn, S., Cousins, S.M., Morley, M., and Spielman,
538 R.S. (2010). Polymorphic Cis- and Trans-Regulation of Human Gene Expression. *PLoS Biol.*
539 8, e1000480.
- 540 31. Hofmann, O., Caballero, O.L., Stevenson, B.J., Chen, Y.-T., Cohen, T., Chua, R., Maher, C. a,
541 Panji, S., Schaefer, U., Kruger, A., *et al.* (2008). Genome-wide analysis of cancer/testis gene
542 expression. *Proc. Natl. Acad. Sci. U. S. A.* 105, 20422–7.
- 543

Figure 1

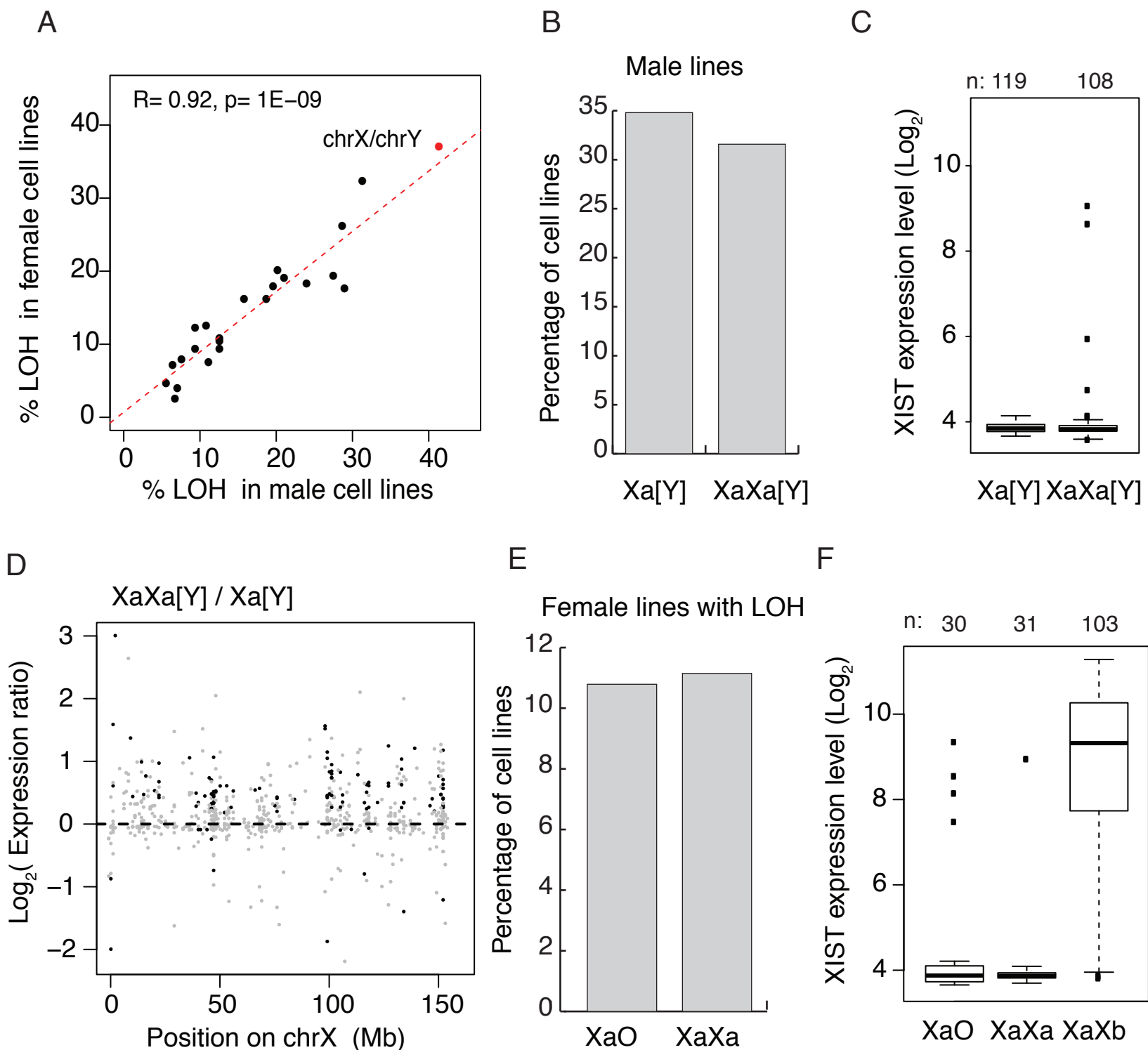
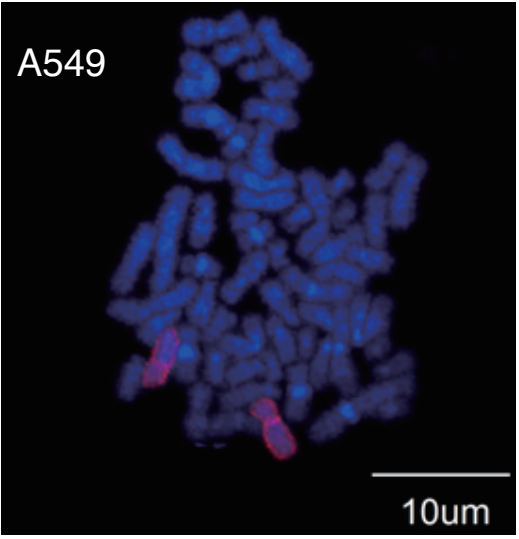
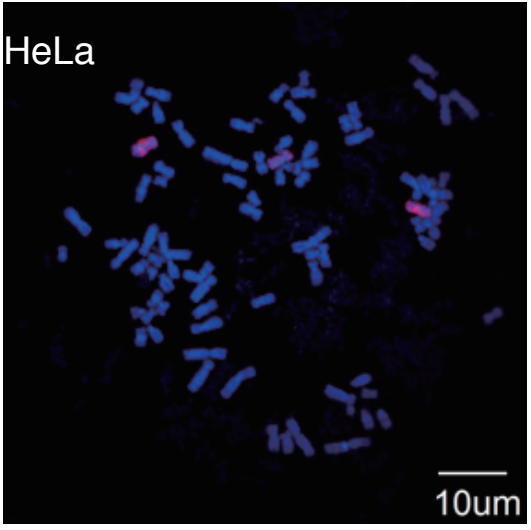


Figure 2

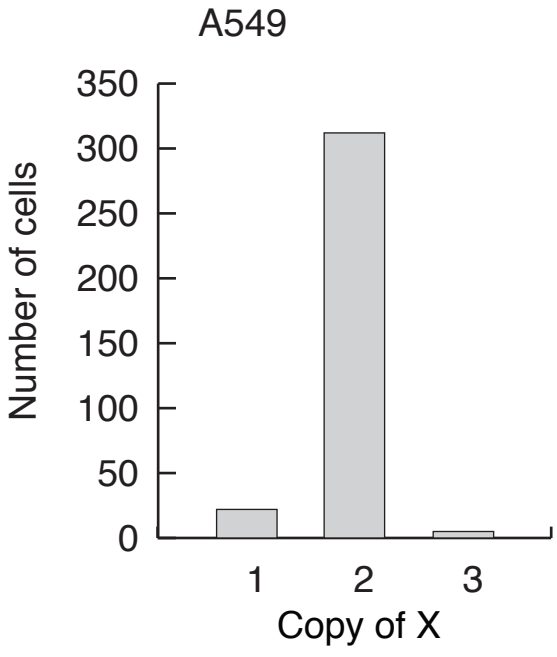
A



B



C



D

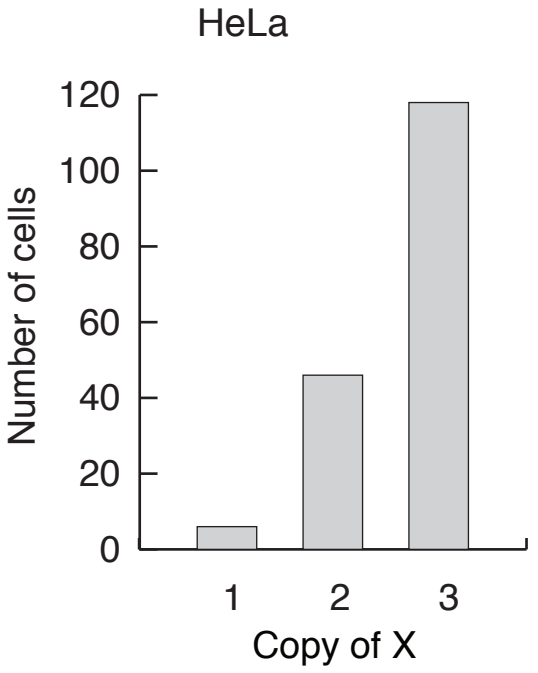
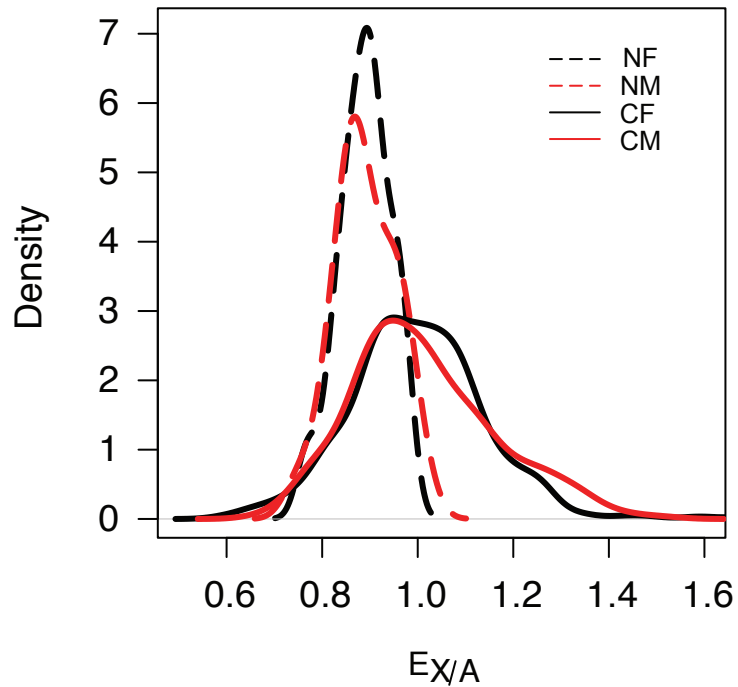


Figure 3

A



B

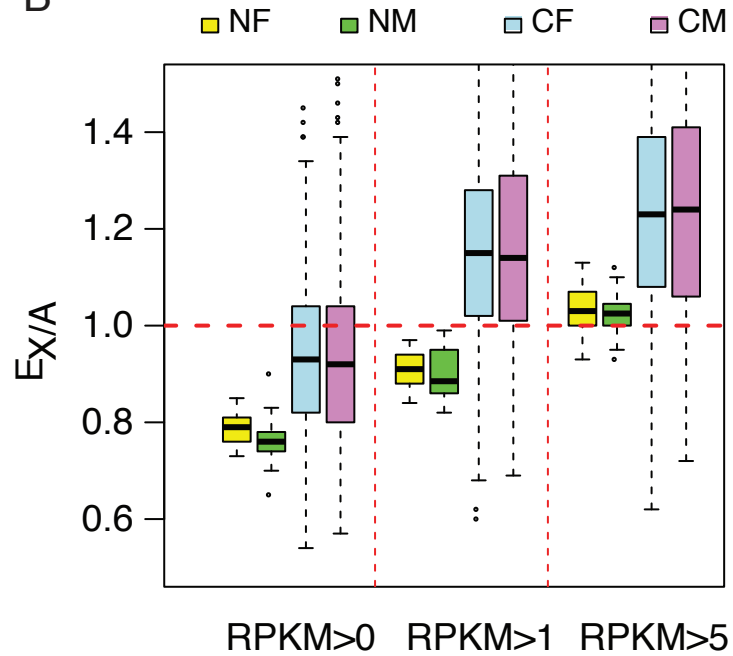
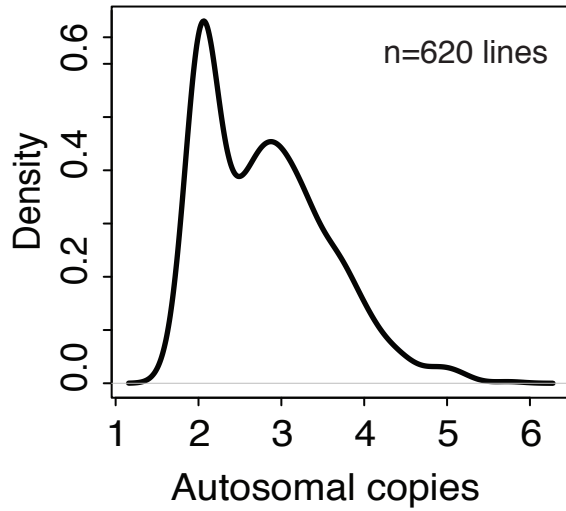
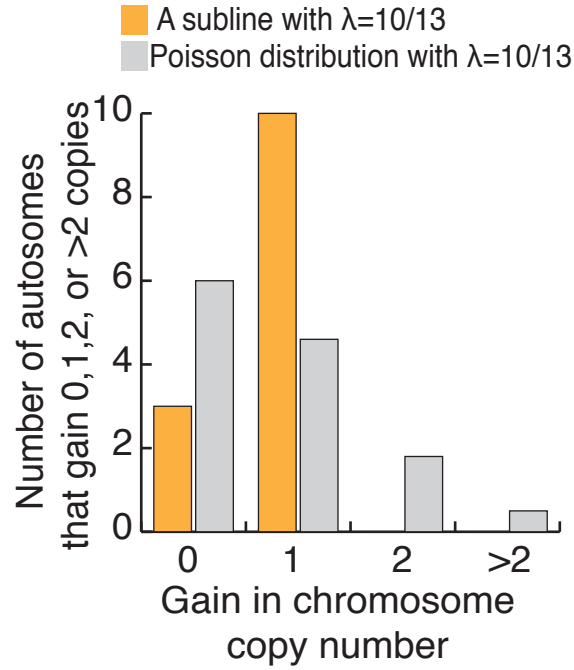


Figure 4

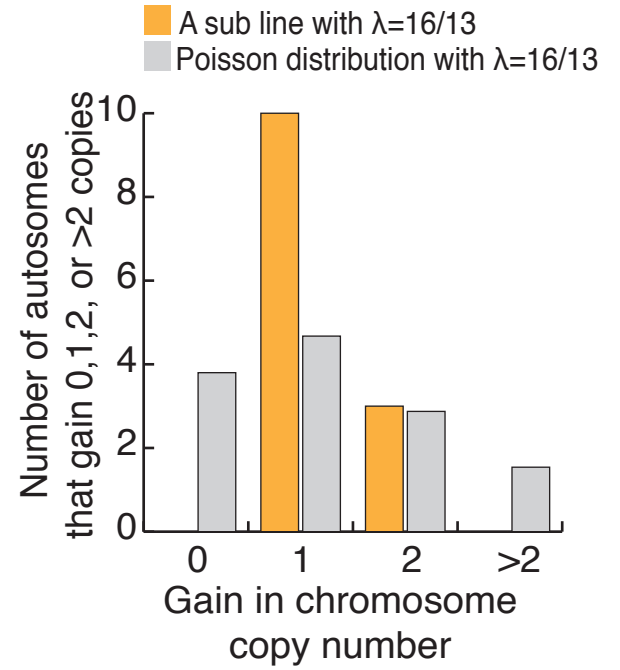
A



B



C

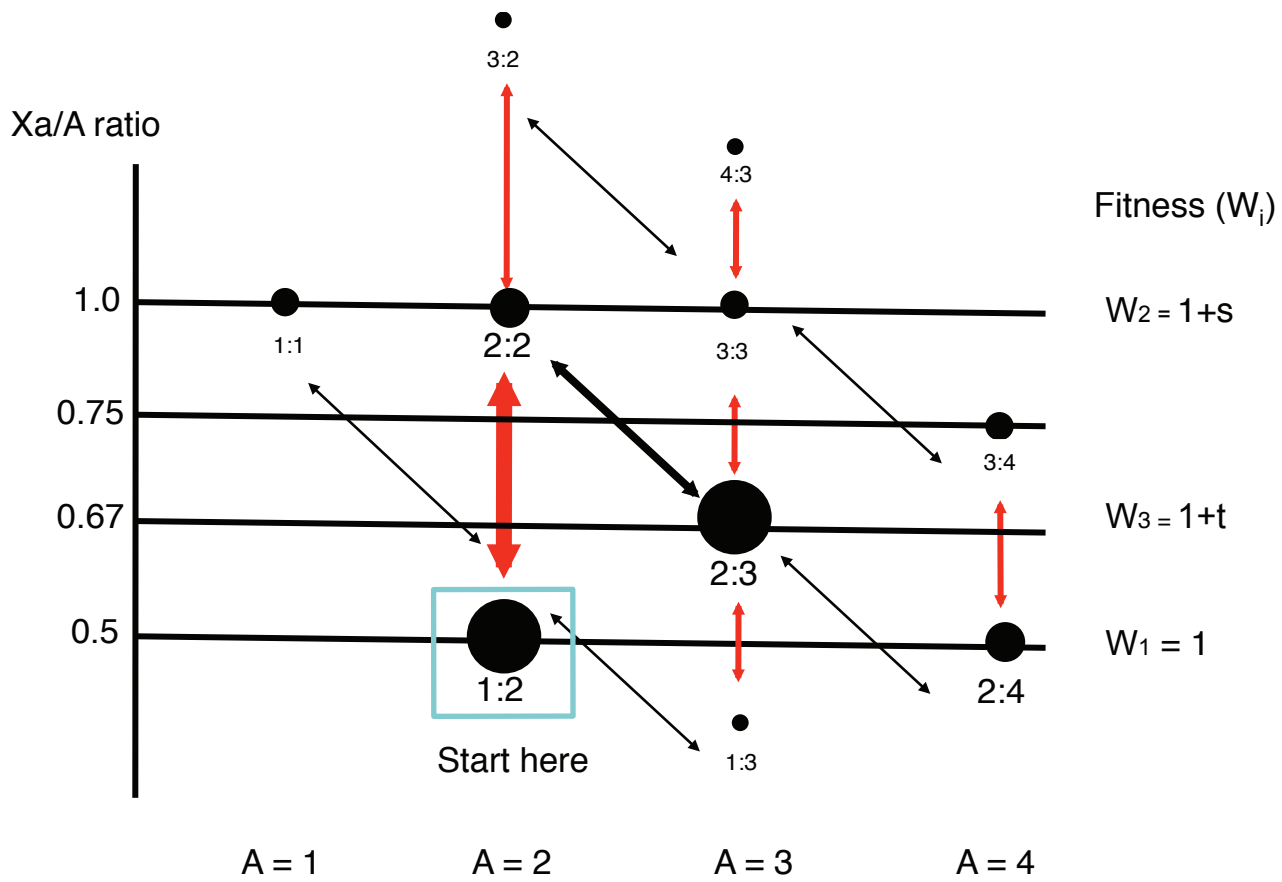


D

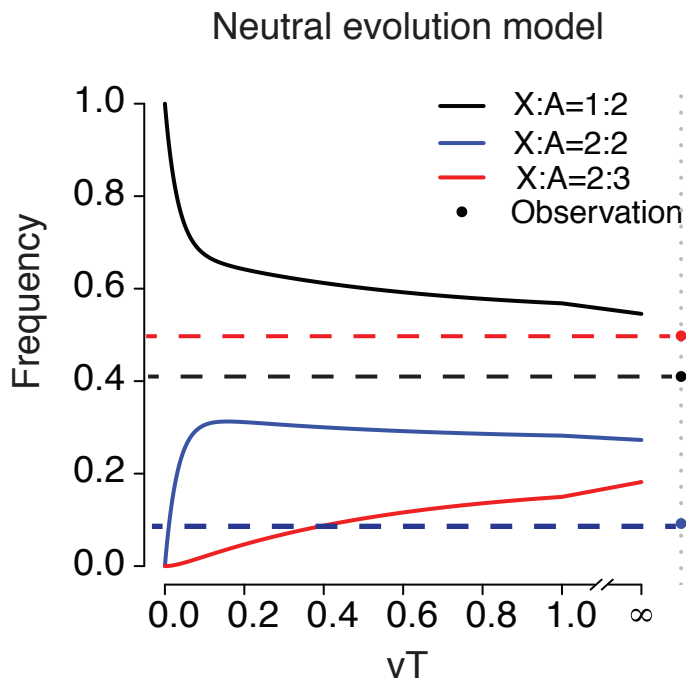
		Autosome			
		1	2	3	4
ChrX	1	0.0	37.4	7.3	0.0
	2	0.0	8.7	29.2	12.8
	3	0.0	0.5	2.5	1.1
	4	0.0	0.0	0.0	0.5

Figure 5

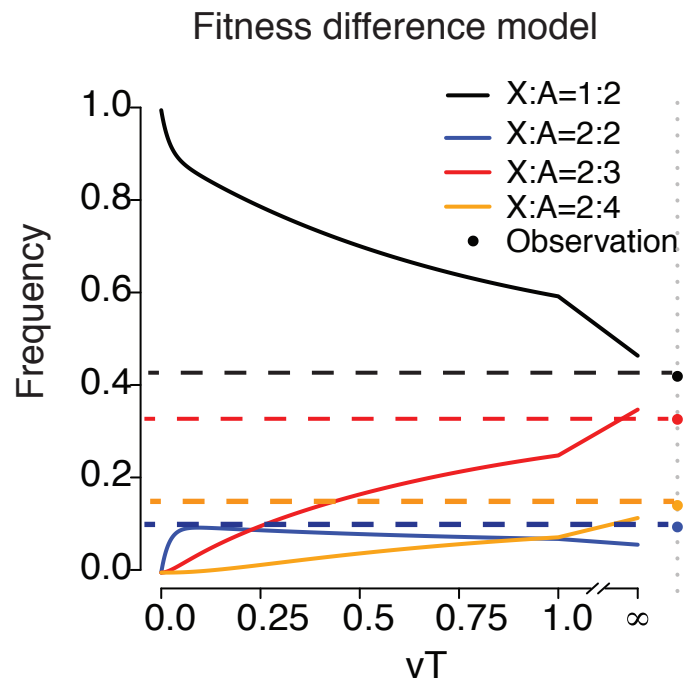
A



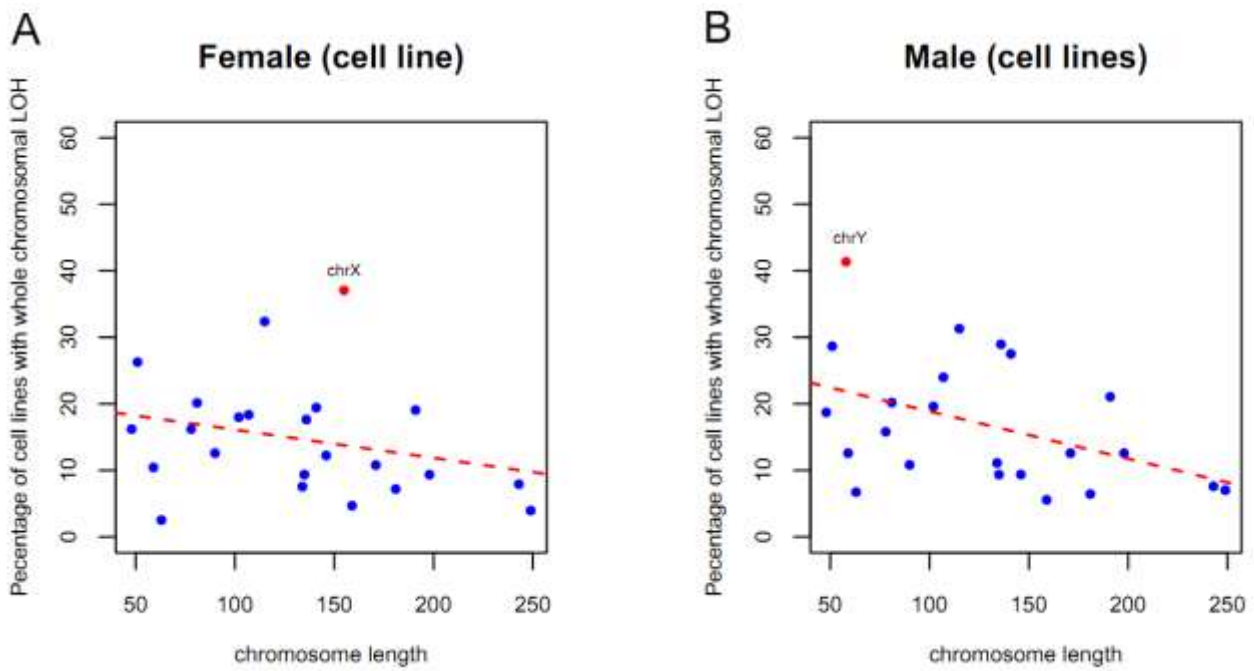
B



C



- 1 **Figure S1:** The frequency of chromosomes loss show negative correlation to their length. There is a
- 2 slight tendency for the smaller autosomes to have higher LOH's than for the larger ones ($R \sim -0.4$,
- 3 $p \sim 0.046$). X chromosome shows significant deviated from the regression line.



4

Figure S2: $E_{X/A}$ ratio in cancerous and normal cell lines by RNA-seq. The gene expression information (gtf files) by RNA-seq was downloaded from UCSC (<ftp://hgdownload.cse.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeCshlLongRnaSeq/>). There are 7 cancerous and 11 normal cell lines respectively. $E_{X/A}$ was calculated by the median value of expressed X-linked and autosomal genes. Expressed genes were selected as RPKM >1.

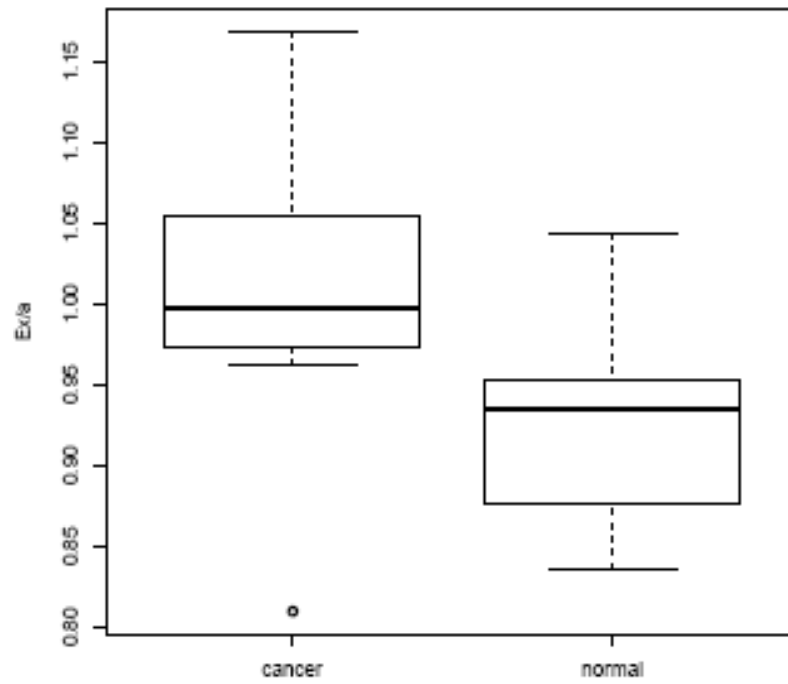


Figure S3: The frequency spectrum of $E_{X/A}$ in male and female cancerous cell lines compared to normal male and female cell lines. The median expression values of X and autosomes genes were used to compute $E_{X/A}$ for each cell lines. The proportion of cell lines within in a bin (0.1) was plotted as Y-axis. The $E_{X/A}$ of cancerous cell lines show a strong right shift compared to that of normal cell lines.

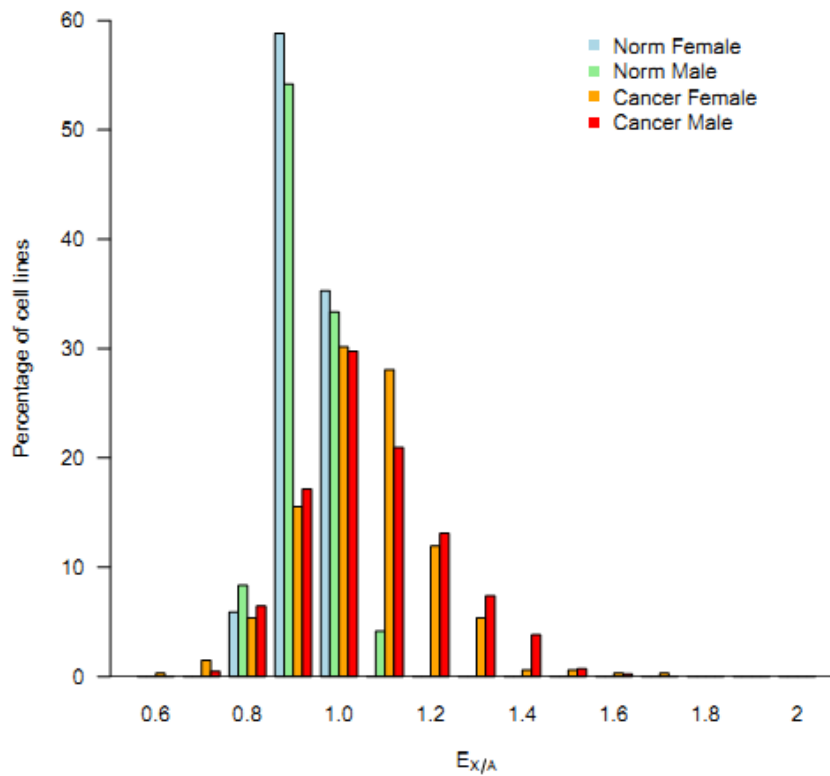


Figure S4: (A) Both the ancestral and sub-clonal HeLa population have 55-75 chromosomes centering around the triploid count of 69. (B) The composite distribution of the five lines with $\lambda < 1$. And the comparison to Poisson distribution. (C) The frequency spectrum of $C(Xa:A)$ in male and female cancerous cell lines. The median values for copy numbers on autosomes and X chromosome were used to calculate $C(Xa:A)$. The proportion of cell lines within a bin (0.1) was plotted as the Y-axis. The discrete peaks denote the four major genotypes (X:3A; X:2A; 2X:3A; 2X:2A).

