

Stability of Rearrangement Measures in the Comparison of Genome Sequences

David Sankoff and Matthew Mazowita

Department of Mathematics and Statistics,
University of Ottawa, 585 King Edward Avenue,
Ottawa, Canada K1N 6N5
{sankoff, mmazo039}@uottawa.ca

Abstract. We present data-analytic and statistical tools for studying rates of rearrangement of whole genomes and to assess the stability of these methods with changes in the level of resolution of the genomic data. We construct data sets on the numbers of conserved syntenies and conserved segments shared by pairs of animal genomes at different levels of resolution. We fit these data to an evolutionary tree and find the rates of rearrangement on various evolutionary lineages. We document the lack of clocklike behaviour of rearrangement processes, the independence of translocation and inversion rates, and the level of resolution beyond which translocations rates are lost in noise due to other processes.

1 Introduction

The goal of this paper is to present data-analytic and statistical tools for studying rates of rearrangement of whole genomes and to assess the stability of these methods with changes in the level of resolution of the genomic data. From secondary data provided by the UCSC Genome Browser, we construct data sets on the number of *conserved syntenies* (pairs of chromosomes, one from each species, containing at least one sufficiently long stretch of homologous sequence) and the number of *conserved segments* (i.e., the total number of such stretches of homologous sequence) shared by pairs of animal genomes at levels of resolution from 30 Kb to 1Mb. For each lineage, we calculate rates of interchromosomal and intrachromosomal rearrangement, with and without the assumption these are due to translocations and inversions of specific kinds.

The key to using whole genome sequences to study evolutionary rearrangements is being able to partition each genome into segments conserved by two genomes since their divergence. In the higher animals and many other eukaryotes this can be extremely difficult, given the high levels of occurrence of transposable elements, retrotranspositions, paralogy and other repetitive and/or inserted sequences, deletions, conversion and uneven sequence divergence. The inference of conserved segments becomes a multistage procedure parametrized

by repeat-masking sensitivities, alignment scores, penalties and thresholds and various numerical criteria for linking smaller segments into large ones. Successful protocols have been developed independently by two research groups [5, 7] using somewhat different strategies to combine short regions of elevated similarity to construct the conserved segments, bridging singly gapped or doubly gapped regions where similarity does not attain a threshold criterion and ignoring short inversions and transpositions that have rearranged one sequence or the other. We develop our data sets on synteny and segments from the continuously updated output of the UCSC protocol made available on the genome browser.

Building on the ideas in [11–13], we derive an estimator for the number of reciprocal translocations responsible for the number of conserved synteny between two genomes, and use simulations to show that the bias and standard deviation of this estimator are less than 5 %, under two models of random translocation, with and without strict conservation of the centromere. By contrasting the number of conserved segments with the number of conserved synteny, we can also estimate the number inversions or other intra-chromosomal events.

Our results include:

- A loss of stability of the data at resolutions starting at about 100 Kb.
- A highly variable proportion of translocations relative to inversions across lineages, for a fixed level of resolution.
- The relative stability of translocation rates compared to inversion rates as resolution is refined.
- The absence of correlation between accumulated rearrangements and chronological time elapsed, especially beyond 20 Myr.

2 The Data

We examined the UCSC browsers for six mammalian species and the chicken and constructed our sets of segments and synteny for the pairs shown in Table 1. We did not use other browsers, or other nets on the four browsers in the table, because either the browser-nets pair are not posted, or because only a build older than the one used in our table was posted.

For each of the pairs in Table 1, four data sets were constructed, one at each of the 1 Mb, 300 Kb, 100 Kb and 30 Kb levels. This contains all segments larger than the resolution level as measured by the segment starting and ending points. We only counted segments in autosomes, except for the comparisons with chicken where the sex chromosomes were also included.

One complication in identifying the segment stems from the key technique of the net construction in [5], which allows very long double gaps in alignments. These gaps are often so long that they contain nested large alignments with chromosomes other than that of the main alignment. Whenever a gap in an alignment contained a nested alignment larger than the level of resolution, we

Table 1. Browsers and nets providing segments

| Browser Species↓ | Net Species and Build Number | | | | | |
|---------------------|------------------------------|------------------|--------------|------------|----------------|--------------------|
| | human Hg17 | chimp PanTro1 | mouse Mm5 | rat rn3 | dog canFam1 | chicken GalGal2 |
| human | | ✓ | ✓ | ✓ | ✓ | ✓ |
| mouse | ✓ | | | ✓ | ✓ | ✓ |
| rat | ✓ | | ✓ | | | |
| dog | ✓ | | ✓ | | | |

Table 2. Data on conserved syntenies c' , conserved segments n , inferred translocations \hat{t} and inferred inversions \hat{i}

| Net | Human Browser | | | | Mouse Browser | | | | Rat Browser | | | | Dog Browser | | | | |
|---------|---------------|-----|-----------|-----------|---------------|------|-----------|-----------|-------------|------|-----------|-----------|-------------|-----|-----------|-----------|-------|
| | c' | n | \hat{t} | \hat{i} | c' | n | \hat{t} | \hat{i} | c' | n | \hat{t} | \hat{i} | c' | n | \hat{t} | \hat{i} | |
| Human | | | | | 143 | 1099 | 77.0 | 463.0 | 134 | 1484 | 71.2 | 660.8 | 106 | 636 | 34.5 | 264.5 | |
| | | | | | 120 | 496 | 59.4 | 179.1 | 114 | 711 | 55.8 | 289.7 | 89 | 356 | 25.7 | 133.3 | |
| | | | | | 107 | 328 | 50.3 | 104.2 | 102 | 358 | 47.3 | 121.7 | 83 | 256 | 22.6 | 86.4 | |
| | | | | | 100 | 241 | 45.6 | 65.4 | 94 | 241 | 41.9 | 68.6 | 79 | 180 | 20.6 | 50.4 | |
| Mouse | 30Kb | 181 | 1186 | 103.9 | 478.1 | | | | | 98 | 1833 | 44.6 | 861.9 | 189 | 731 | 79.8 | 266.7 |
| | 100Kb | 121 | 533 | 57.0 | 198.5 | | | | | 65 | 741 | 23.9 | 336.6 | 172 | 428 | 70.2 | 124.8 |
| | 300Kb | 108 | 351 | 48.3 | 116.2 | | | | | 48 | 189 | 14.3 | 70.2 | 158 | 321 | 62.5 | 79.0 |
| | 1Mb | 103 | 256 | 45.1 | 71.9 | | | | | 42 | 88 | 11.1 | 22.9 | 144 | 246 | 54.8 | 49.2 |
| Rat | 30Kb | 179 | 1806 | 102.1 | 789.9 | 97 | 2008 | 43.7 | 950.8 | | | | | | | | |
| | 100Kb | 119 | 822 | 55.7 | 344.3 | 60 | 699 | 21.4 | 318.6 | | | | | | | | |
| | 300Kb | 104 | 393 | 45.7 | 139.8 | 43 | 172 | 12.1 | 64.4 | | | | | | | | |
| | 1Mb | 95 | 249 | 40.1 | 73.4 | 41 | 87 | 11.1 | 22.9 | | | | | | | | |
| Dog | 30Kb | 210 | 1206 | 131.6 | 460.4 | 197 | 1028 | 128.0 | 376.5 | | | | | | | | |
| | 100Kb | 108 | 465 | 48.3 | 173.2 | 178 | 502 | 108.1 | 133.4 | | | | | | | | |
| | 300Kb | 85 | 274 | 34.0 | 92.0 | 162 | 341 | 93.1 | 67.9 | | | | | | | | |
| | 1Mb | 80 | 194 | 31.0 | 55.0 | 148 | 251 | 81.0 | 35.0 | | | | | | | | |
| Chicken | 30Kb | 183 | 1834 | 101.2 | 804.3 | 199 | 1773 | 122.5 | 754.0 | | | | | | | | |
| | 100Kb | 114 | 992 | 50.8 | 433.7 | 159 | 1025 | 86.6 | 415.9 | | | | | | | | |
| | 300Kb | 84 | 598 | 32.5 | 255.0 | 134 | 601 | 67.4 | 223.1 | | | | | | | | |
| | 1Mb | 71 | 330 | 25.0 | 128.5 | 108 | 335 | 49.5 | 108.0 | | | | | | | | |
| Chimp | 30Kb | 79 | 2824 | 30.4 | 1370.6 | | | | | | | | | | | | |
| | 100Kb | 33 | 543 | 5.4 | 255.1 | | | | | | | | | | | | |
| | 300Kb | 25 | 143 | 1.5 | 59.0 | | | | | | | | | | | | |
| | 1Mb | 24 | 65 | 1.0 | 20.5 | | | | | | | | | | | | |

broke the main alignment in two and counted the segments before and after the gap separately, assuming they remained long enough, as well as the segment in the gap.

Table 2 shows the results of our data extraction procedure. Entries for \hat{t} and \hat{i} calculated according to equations (7) and (1), respectively.

3 Models of Translocation

In order to derive and validate our estimator of translocation rates, we model the autosomes of a genome as c linear segments with lengths $p(1), \dots, p(c)$, proportional to the number of base pairs they contain, where $\sum_{i=1}^c p(i) = 1$. We assume the two breakpoints of a translocation are chosen independently according to a uniform distribution over all autosomes, conditioned on their not being on the same chromosome. (There is no statistical evidence [9] that translocational breakpoints cluster in a non-random way on chromosomes, except in a small region immediately proximal – within 50-300Kb – to the telomere in a wide spectrum of eukaryote lineages [6].)

A reciprocal translocation between two chromosomes h and k consists of breaking each one, at some interior point, into two segments, and rejoining the four resulting segments such that two new chromosomes are produced.

In one version of our model, we impose a left-right orientation on each chromosome, such that a left-hand fragment must always rejoin a right-hand fragment. This ensures that each chromosome always retains a segment, however small it may become, containing its original left-hand extremity. This restriction models the conservation of the centromere without introducing complications such as trends towards or away from acrocentricity. With further translocations, if a breakpoint falls into a previously created segment on chromosome i , it divides that segment into two new segments, the left-hand one remaining in chromosome i , while the right-hand one, and all the other segments to the right of the breakpoint, are transferred to the other chromosome involved in the translocation. It is for this version of the model that we will derive an estimator of the number of translocations, and that we will simulate to test the estimator.

In another version of the model, an inverted left-hand fragment may rejoin another left-hand fragment and similarly for right-hand fragments. This models a high level of neocentromeric activity. We will also simulate this model to see how our estimator (derived from the previous model) fares.

We do not consider chromosome fusion and fission, so that the number of chromosomes is constant throughout the time period governed by the model. Later, in our analysis of animal genomes, we simply assume that the case where fusions or fissions occur will be well approximated by interpolating two models (with fixed chromosome number) corresponding to the two genomes being compared.

Moreover, in our simulations, we do not consider the effects of inversions on the accuracy of estimator. In previous work [12], we showed that high rates of long inversions would severely bias the estimator upwards, but that the rates and distribution of inversion lengths documented for mammalian genomes [5, 7] had no perceptible biasing effect.

In our simulations we impose a threshold and a cap on chromosome size, rejecting any translocation that results in a chromosome too small or too large. Theories about meiosis, e.g. [15], can be adduced for these constraints, though there are clear exceptions, such as the “dot” chromosomes of avian and some reptilian and other vertebrate genomes [1, 3].

The total number of segments on a human chromosome i is

$$n^{(i)} = t^{(i)} + 2u^{(i)} + 1, \tag{1}$$

where $t^{(i)}$ is the number of translocational breakpoints on the chromosome, and $2u^{(i)}$ is the number of inversion breakpoints.

4 Prediction and Estimation

We assume that our random translocation process is temporally reversible, and to this effect we show in Figure 1 and Section 5.1 that the equilibrium state of our process well approximates the observed distribution of chromosome lengths in the human genome. In comparing two genomes, this assumption allows us to treat either one as ancestral and the other as derived, instead of having to consider them as diverging independently from a common ancestor.

At the outset, assume the first translocation on the lineage from genome A to genome B involves chromosome i . The assumption of a uniform density of breakpoints across the genome implies that the “partner” of i in the translocation will be chromosome j with probability $p_i(j) = \frac{p(j)}{1-p^{(i)}}$. Thus the probability that the new chromosome labelled i contains no fragment of genome A chromosome j , where $j \neq i$, is $1 - p_i(j)$. For small $t^{(i)}$, after chromosome i has undergone $t^{(i)}$ translocations, the probability that it contains no fragment of the genome A chromosome j is approximately $(1 - p_i(j))^{t^{(i)}}$, neglecting second-order events, for example, the event that j previously translocated with one or more of the $t^{(i)}$ chromosomes that then translocated with i , and that a secondary transfer to i of material originally from j thereby occurred.

Then the probability that the genome B chromosome i contains at least one fragment from j is approximately $1 - (1 - p_i(j))^{t^{(i)}}$ and the expected number of genome A chromosomes with at least one fragment showing up on genome B chromosome i is

$$E(c^{(i)}) \approx 1 + \sum_{j \neq i} [1 - (1 - p_i(j))^{t^{(i)}}] \tag{2}$$

so that

$$c - E(c^{(i)}) \approx \sum_{j \neq i} (1 - p_i(j))^{t^{(i)}}, \tag{3}$$

where the leading 1 in (2) counts the fragment containing the left-hand endpoint of the genome A chromosome i itself. We term $c^{(i)}$ the number of *conserved syntenies* on chromosome i .

Suppose there have been a total of t translocations in the evolutionary history. Then

$$\sum_i t^{(i)} = 2t. \tag{4}$$

We can expect these to have been distributed among the chromosomes approximately as

$$t^{(i)} = 2tp(i), \tag{5}$$

so that

$$c^2 - \sum_i E(c^{(i)}) \approx \sum_i \sum_{j \neq i} (1 - p_i(j))^{2tp^{(i)}}. \quad (6)$$

Substituting the $c^{(i)}$ for the $E(c^{(i)})$ in eqn (6) suggests solving

$$c^2 - \sum_i c^{(i)} = \sum_i \sum_{j \neq i} [1 - p_i(j)]^{2\hat{t}p^{(i)}}, \quad (7)$$

for \hat{t} to provide an estimator of t . Newton's method converges rapidly for the range of parameters used in our studies, as long as not all $c^{(i)} = c$. (We know of no comparative map where even one chromosome of one genome shares a significant syntenic segment with every autosome of the other genome, much less a map where every chromosome is thus scrambled.)

5 Simulations

5.1 Equilibrium Distribution of Chromosome Size

Models of accumulated reciprocal translocations for explaining the observed range of chromosome sizes in a genome date from the 1996 study of Sankoff and Ferretti [10]. They proposed a lower threshold on chromosome size in order to reproduce the appropriate size range in plant and animal genomes containing from two to 22 autosomes. A cap on largest chromosome size has also been proposed [15] and shown to be effective [4]. Economy and elegance in explaining chromosome size being less important in the present context than simulating a realistic equilibrium distribution of these sizes, we imposed both a threshold of 50 Mb and a cap of 250 Mb on the process described in Section 3, simply rejecting any translocations that produced chromosomes out of the range. These values were inspired by the relative stability across primates and rodents evident in the data in Table 3, though they are less pertinent for the dog, with a larger number of correspondingly smaller chromosomes, and chicken, which has several very small chromosomes.

Table 3. Shortest and longest chromosome, in Mb

| genome | shortest | longest |
|---------|----------|---------|
| mouse | 61 | 199 |
| human | 47 | 246 |
| rat | 47 | 268 |
| chimp | 47 | 230 |
| dog | 26 | 125 |
| chicken | 0.24 | 188 |

Simulating the translocation process 100 times up to 10,000 translocations each produced the equilibrium distribution of chromosome sizes in Figure 1. The superimposed distribution of human autosome sizes is very close to the equilibrium distribution.

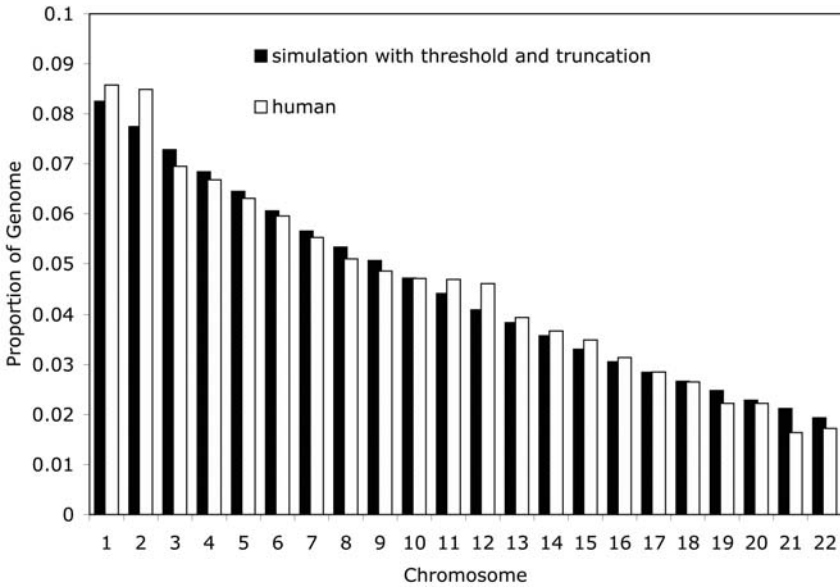


Fig. 1. Comparison of equilibrium distribution of simulated chromosome sizes with human autosome sizes

5.2 Performance of the Estimator

Figure 2 depicts the estimated number of translocations as a function of the true number t in the simulation. The estimator \hat{t} appears very accurate, only lightly biased (less than 5 % for $t < 200$), with small error rates (s.d. less than 5 % for $t < 200$).

6 Fitting the Data to Animal Phylogeny

To infer the rates of rearrangement on evolutionary lineages, we assumed the phylogenetic tree in Figure 3. Because of the limited number of genome pairs for which we have data, we artificially attributed all the chimp-human divergence to the chimp lineage, and could not estimate the translocational divergence during the mammalian radiation, i.e. between the divergence of the dog lineage and the common primate/rodent lineage.

We fit the data in Table 2 to the tree by solving the system of linear equations between the additive path lengths in the tree and the inferred rearrangement distances, namely $c' - c$ the number of new synteny created on a path, \hat{t} the number of translocations inferred to have occurred, $n - c$ the number of segments created, and \hat{i} the inferred number of inversions. Where browser-net pairs were available in both directions, we averaged the two results in Table 2 to produce a single equation. As there are more pairs than edges, we also used an averaging procedure in solving the equations to produce the results in Table 4.

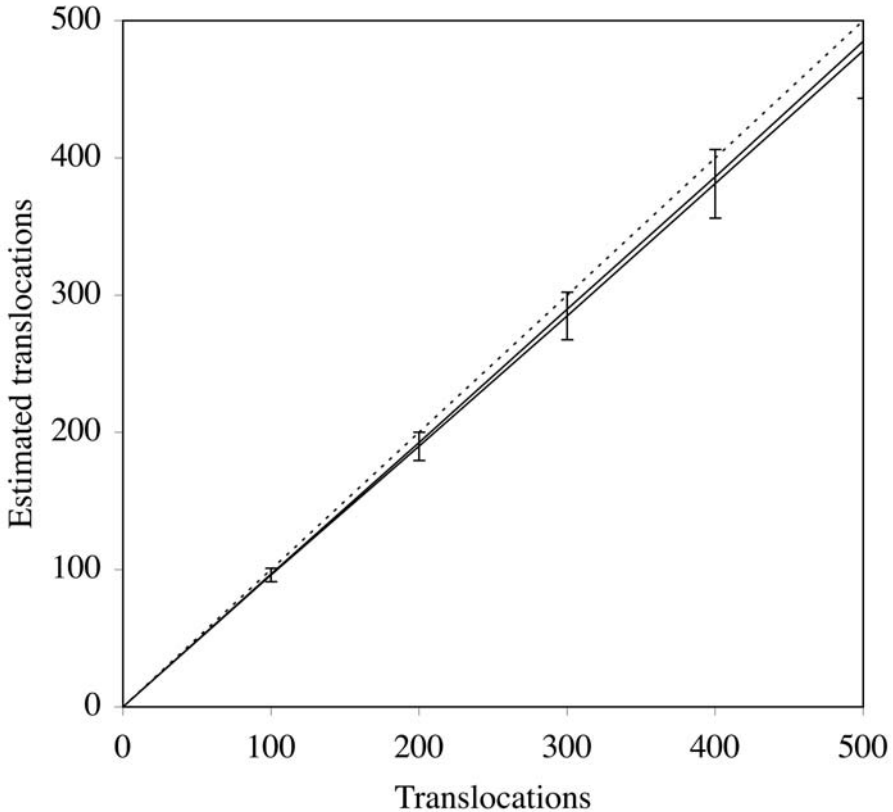


Fig. 2. Mean value, over 100 runs, of \hat{t} as a function of t . Dotted line: $\hat{t} = t$. Lower line with ± 1 s.d. error bars: model with centromere. Upper line: model without centromere

The very approximate temporal edge lengths given in Table 4 were based, for p, on twice the usual estimate (6 Myr) of chimp-human divergence to account for both human and chimp evolution; for m and r the date for the rat-mouse divergence (≈ 20 Myr); for h and d the date for mammalian radiation; for a the same date less the 20 Myr of murid evolution; and for c twice the mammalian-reptile divergence time (310 Myr) less the 85 Myr since the mammalian radiation.

7 Observations

The most striking trend in Table 4 is the dramatic increase in both conserved synteny and conserved segments in almost every lineage (except a) as the level of resolution is refined, starting at 100 Kb but accelerating rapidly at 30 Kb. It seems likely that the increased level of translocations inferred is artificial, the apparent level of conserved synteny reflecting retrotransposition and other interchromosomal process and not reciprocal translocation [16].

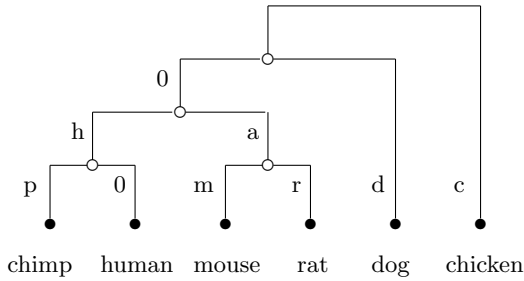


Fig. 3. Unrooted phylogeny for fitting translocation measures. Edges labelled “0” indicate the two endpoints are collapsed; none of the pairwise measures bear on the location of the vertices between chicken and dog, and between human and chimp. These two lengths can be assumed to be short in any case (period of rapid mammalian radiation and human-chimp divergence, respectively)

That this increase does not reflect translocational distance is further evidenced by the loss at 30 Kb of clear trends among the lineages visible at less refined resolutions, such as the very low values for human, mouse and rat compared to the other lineages. Thus we can conclude that below the 100 Kb level, the study of translocational rearrangement by our statistical approach is no longer feasible.

Even at the less refined levels of resolution, any correlation between chronological time and translocational distance breaks down somewhere between 20 and 65 Myr. As has been remarked previously [2], the chicken evidences a low rate of translocation. The dog on the other hand, shows a high rate.

Turning to the results on inversions, the rapid increase in segments and inferred inversions at refined resolutions may, in contrast to translocations, be a real effect. It is known that the inversions of small size are very frequent in these genomes [5], with a mean size less than 1 Kb, so that it can be expected that the number of inversions inferred will continue to accelerate with increased resolution. Indeed, even between 1 Mb and 300 Kb, while translocation rates are relatively stable, inversion rates increase substantially.

The pattern of inversion rates among the lineages is very different from that of translocations. Here, chicken has a very high rate while dog has a low one, the opposite of what was seen with translocations. Perhaps most startling is the high rate recovered for the chimp lineage. This disproportion is likely an artifact; whereas in more distant comparisons the alignments of many inverted segments may not be detected due to sequence divergence, the proximity of the human and chimpanzee genomes allows for a very high recovery rate.

Again with inversions, as with translocations, there is little correlation of lineage-specific rates and the chronological span of the lineage.

Finally, when inversion and translocation rates are compared at a fixed level of resolution, no systematic association can be seen.

Table 4. Tree edge-lengths estimated from pairwise interchromosomal measures in Table 2. Negative entries indicate poor fit of the data to this tree

| | | Time (Myr) | | | | | | |
|-------|--|-----------------------|--------|-------|--------|--------|-------|--------|
| | | h | a | m | r | c | d | p |
| | | 85 | 65 | 20 | 20 | 535 | 85 | 12 |
| | | New syntenies created | | | | | | |
| | | h | a | m | r | c | d | p |
| 1Mb | | 13.5 | 52.8 | 14.8 | 7.3 | 28.5 | 43.0 | 2.0 |
| 300Kb | | 10.9 | 60.6 | 15.5 | 10.5 | 45.0 | 49.3 | 3.0 |
| 100Kb | | 18.5 | 57.8 | 23.8 | 19.3 | 66.0 | 57.5 | 11.0 |
| 30Kb | | 56.9 | 42.6 | 42.0 | 36.0 | 99.8 | 75.5 | 57.0 |
| | | Translocations | | | | | | |
| | | h | a | m | r | c | d | p |
| 1Mb | | 6.0 | 31.6 | 7.8 | 3.3 | 14.6 | 24.2 | 1.0 |
| 300Kb | | 3.6 | 37.8 | 8.0 | 5.2 | 25.3 | 28.4 | 1.5 |
| 100Kb | | 7.1 | 38.5 | 12.6 | 10.1 | 39.6 | 34.0 | 5.4 |
| 30Kb | | 34.7 | 31.8 | 23.9 | 20.2 | 66.7 | 48.3 | 30.4 |
| | | New segments created | | | | | | |
| | | h | a | m | r | c | d | p |
| 1Mb | | 96.3 | 95.8 | 36.0 | 32.0 | 198.0 | 74.5 | 43.0 |
| 300Kb | | 141.1 | 115.1 | 62.8 | 98.3 | 419.5 | 109.3 | 121.0 |
| 100Kb | | 224.0 | 45.5 | 224.5 | 476.0 | 741.0 | 161.5 | 521.0 |
| 30Kb | | 585.5 | -163.0 | 699.5 | 1201.5 | 1222.0 | 310.0 | 2802.0 |
| | | Inversions | | | | | | |
| | | h | a | m | r | c | d | p |
| 1Mb | | 42.1 | 16.3 | 10.2 | 12.7 | 83.9 | 13.1 | 20.5 |
| 300Kb | | 67.0 | 19.8 | 23.4 | 43.9 | 184.0 | 26.3 | 59.0 |
| 100Kb | | 104.9 | -15.8 | 99.7 | 227.9 | 330.4 | 46.8 | 255.1 |
| 30Kb | | 188.2 | -43.5 | 325.8 | 580.6 | 543.8 | -33.0 | 1370.6 |

8 Discussion

We have proposed an estimator of the number of translocations intervening between two rearranged genomes, based only on the numbers of conserved syntenies on each chromosome, the lengths of the chromosomes and a simplified random model of interchromosomal exchange. This estimator proves to be very accurate in simulations, which is remarkable given that it only explicitly takes into account the first-order effects of interchromosomal exchange.

In this paper, we applied our estimator to animal genome comparisons at various levels of resolution. This showed that translocation estimates are stable at coarse resolutions, while inversions increased markedly. This reflects the discovery of high numbers of smaller-scale local arrangements recognizable from

genomic sequence [5]. At very detailed levels of resolution, inferred translocations numbers probably reflect processes other than translocation, though increased inversion inferences are more likely to reflect the inversion process.

Our estimates of the number of translocations and inversions in the evolutionary divergence of animals are only about a half of what has been published by Pevzner and colleagues [7, 8, 2] for corresponding level of resolution. Their estimates are based on an algorithmic reconstruction of the details of evolutionary history. Our model assumes each translocation and inversion creates two new segments, but the algorithms require a number of rearrangements almost equal to the number of segments to account for how the segments are ordered on the chromosomes. This accounts for the difference between the two sets of results.

Acknowledgements

Thanks to Phil Trinh for help with the genome browsers and other tools. Research supported in part by a Discovery grant from the Natural Sciences and Engineering Research Council of Canada (NSERC). DS holds the Canada Research Chair in Mathematical Genomics and is a Fellow in the Evolutionary Biology Program of the Canadian Institute for Advanced Research.

References

1. Bed'hom, B. (2000). Evolution of karyotype organization in *Accipitridae*: A translocation model. In Sankoff, D. and Nadeau, J. H. (eds) *Comparative Genomics: Empirical and Analytical Approaches to Gene Order Dynamics, Map Alignment and Evolution of Gene Families*. Dordrecht, NL, Kluwer, 347–56.
2. Bourque, G., Zdobnov, E., Bork, P., Pevzner, P.A. and Tesler, G. (2005) Comparative architectures of mammalian and chicken genomes reveal highly variable rates of genomic rearrangements across different lineages. *Genome Research*, **15**, 98–110.
3. Burt, D.W. (2002). Origin and evolution of avian microchromosomes. *Cytogenetic and Genome Research*, **96**, 97–112.
4. De, A., Ferguson, M., Sindi, S. and Durrett, R. (2001). The equilibrium distribution for a generalized Sankoff-Ferretti model accurately predicts chromosome size distributions in a wide variety of species. *Journal of Applied Probability*, **38**, 324–34.
5. Kent, W. J., Baertsch, R., Hinrichs, A., Miller, W. and Haussler, D. (2003). Evolution's cauldron: Duplication, deletion, and rearrangement in the mouse and human genomes. *Proceedings of the National Academy of Sciences, USA*, **100**, 11484–9.
6. Mefford, H.C. and Trask, B.J. (2002). The complex structure and dynamic evolution of human subtelomeres. *Nature Reviews in Genetics*, **3**, 91–102; 229.
7. Pevzner, P. A. and Tesler, G. (2003). Genome rearrangements in mammalian genomes: Lessons from human and mouse genomic sequences. *Genome Research*, **13**, 37–45
8. Pevzner, P. A. and Tesler, G. (2003). Human and mouse genomic sequences reveal extensive breakpoint reuse in mammalian evolution. *Proceedings of the National Academy of Sciences, USA*, **100**, 7672–7

9. Sankoff, D., Deneault, M., Turbis, P. and Allen, C.P. (2002) Chromosomal distributions of breakpoints in cancer, infertility and evolution. *Theoretical Population Biology*, **61**, 497–501.
10. Sankoff, D. and Ferretti, V. (1996). Karotype distributions in a stochastic model of reciprocal translocation. *Genome Research*, **6**, 1–9.
11. Sankoff, D., Ferretti, V. and Nadeau, J.H. (1997) Conserved segment identification. *Journal of Computational Biology*, **4**, 559–65.
12. Sankoff, D. and Mazowita, M. (2005) Estimators of translocations and inversions in comparative maps. In Lagergren, J. (ed.) *Proceedings of the RECOMB 2004 Satellite Workshop on Comparative Genomics, RCG 2004*, Lecture Notes in Bioinformatics, **3388**, Springer, Heidelberg, 109–122
13. Sankoff, D., Parent, M.-N. and Bryant, D. (2000). Accuracy and robustness of analyses based on numbers of genes in observed segments. In Sankoff, D. and Nadeau, J. H. (eds) *Comparative Genomics: Empirical and Analytical Approaches to Gene Order Dynamics, Map Alignment and Evolution of Gene Families*. Dordrecht, NL, Kluwer, 299–306.
14. Sankoff, D. and Trinh, P. (2004). Chromosomal breakpoint re-use in the inference of genome sequence rearrangement. *Proceedings of RECOMB 04, Eighth International Conference on Computational Molecular Biology*. New York: ACM Press, 30–5.
15. Schubert, I. and Oud, J.L. (1997). There is an upper limit of chromosome size for normal development of an organism. *Cell*, **88**, 515–20.
16. Trinh, P., McLysaght, A. and Sankoff, D. (2004) Genomic features in the breakpoint regions between syntenic blocks. *Bioinformatics*, **20**, I318–I325.