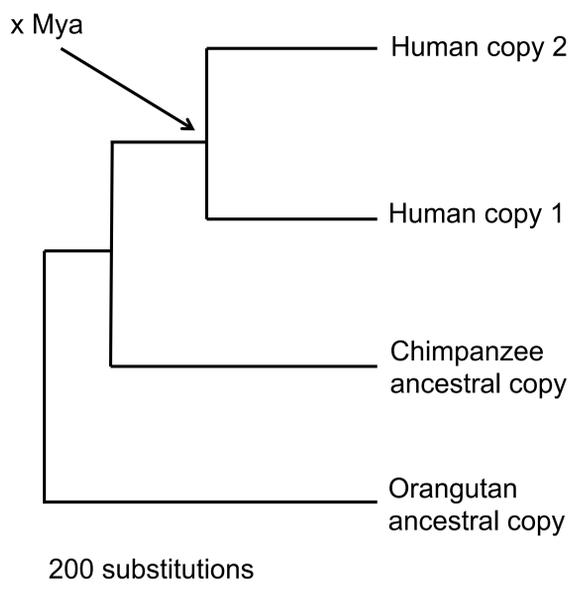


## Introduction

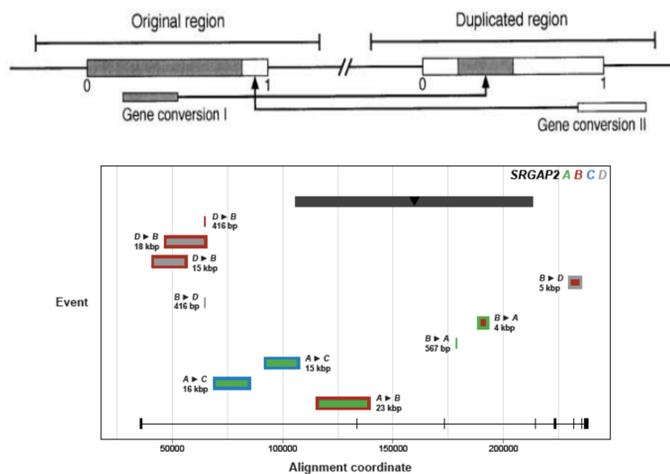
While gene duplications are known to have led to functional novelty in many species, the evolutionary history and functional impact of duplications that arose on the human lineage remains unknown. A subset (notably, *SRGAP2C* and *ARHGAP11B*<sup>2</sup>) has been shown to have important functions in spine maturation and neocortex expansion and may contribute to the risk of neuropsychiatric diseases. Despite their importance for understanding human adaptation and evolution, most of the current research on human-specific duplications has been limited to biomedical applications and the few evolutionary studies to date have not accounted for the confounding effects of paralogous gene conversion that homogenizes the variation between the copies leading to inaccurate estimates of the age of the duplications and potentially of the evidence for selection<sup>3,4</sup>. To estimate the age of duplications, we develop an Approximate Bayesian Computation (ABC) method that models the pairwise divergence between duplicated copies as a function of the time, mutation rate and non-allelic gene conversion rate. Direct estimates for the rates of paralogous gene conversion are limited and hence we use an approach that accounts for the distance between duplicates, local recombination rate, and allelic gene conversion rate to infer these rates. Further we use the revised mutation rate that has recently emerged from studies of human pedigrees<sup>5,6</sup>. Application of our method to *SRGAP2* provides substantially older dates than had previously been reported<sup>4</sup>, indicating that *SRGAP2C* occurred long before the beginning of the neocortex expansion. Together, these analyses allow us to reliably infer the chronology of human-specific duplications and to assess their potential role in critical periods of innovation in the human evolution.

## Using the molecular clock for estimating the age of duplications



The standard approach to estimating the age of duplicates is based on the pairwise sequence divergence among the duplicated copies (in this example, Human copy 1 and Human copy 2). Using an estimate of the direct mutation rate from sequencing human pedigrees or assuming a divergence time to an outgroup such as chimpanzees, the pairwise sequence divergence between duplicated copies can then be translated into an estimate of the time of divergence between the copies. (Regional variation in mutation rate can also be taken into account, e.g., by using a distant outgroup such as orangutan or macaque.) This approach does not account for the impact of paralogous gene conversion that occurs between non-allelic segments that share high sequence similarity and homogenizes the variation between the copies. This omission leads to a systematic downward bias in estimated dates<sup>7,8</sup>.

## Non-allelic homologous gene conversion



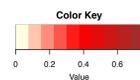
Figures from top: Teshima and Inman, 2004 and bottom: Nuttle et al., 2013

Non-allelic gene conversion acts as a copy paste mechanism leading to unidirectional transfer of short stretches of sequence between paralogous genes<sup>7,8</sup>. Tracks of shared sequences have been empirically documented in *SRGAP2* and other human-specific duplications<sup>9</sup>, and thus this phenomenon needs to be taken into account in dating duplications<sup>8</sup>.

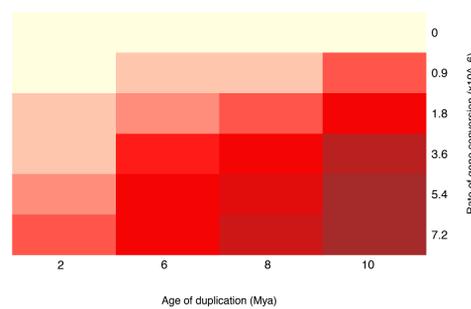
## Impact of non-allelic gene conversion on the estimated age of duplications

To characterize the bias in estimated ages that arises when ignoring gene conversion, we performed coalescent simulations using *ms\_dup*<sup>8</sup>. We assumed that the mutation rate  $\mu = 1.2 \times 10^{-8}$  per bp per gen.<sup>5,6</sup>, the rate of recombination  $r = 10^{-8}$  per bp per gen., the rate at which a site is involved in a gene conversion track  $g$  is shown on the Y-axis, the gene conversion track length  $l = 75$  bp<sup>10</sup> and the length of duplicated region is 200 Kb. Shown in the heat plot is the bias, computed as:

$$\text{Bias} = (T_{\text{truth}} - T_{\text{estimated}}) / T_{\text{truth}}$$

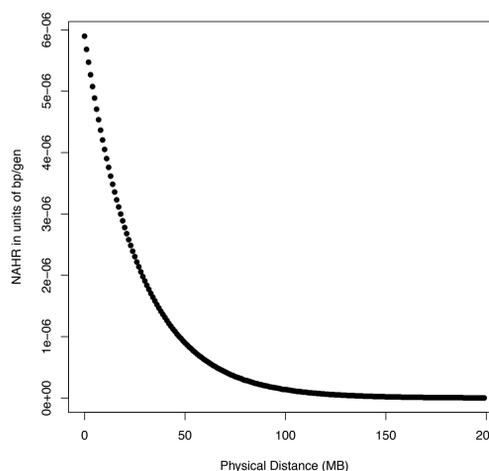


Bias in estimated dates



## Model to infer regional estimates of non-allelic gene conversion

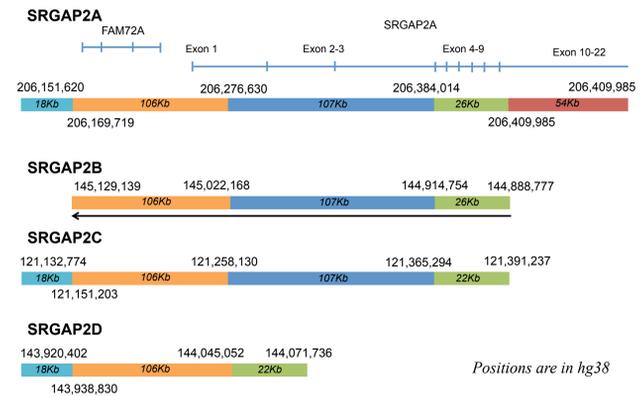
Modeling NAHR as an exponential with distance



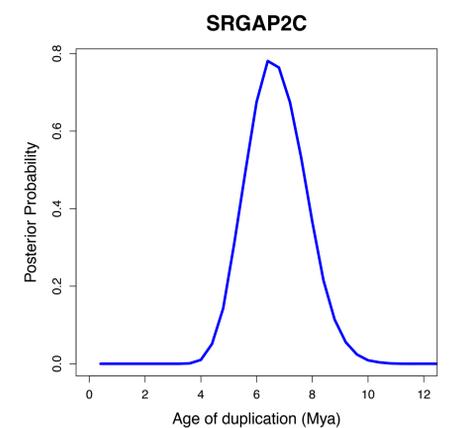
$$g = Ae^{-\lambda d}$$

Here,  $g$  is the rate at which a site is involved in a non-allelic gene conversion event,  $A$  is a regional estimate of the rate of allelic gene conversion<sup>10,11</sup>,  $\lambda$  is the rate of decay of the rate of non-allelic gene conversion with physical distance. This model is motivated from observations in yeast<sup>7</sup> that the rate of ectopic exchange decreases as distance increases, and is highest when duplications are in tandem.

## Estimating the age of human-specific duplication, *SRGAP2*



The *SRGAP2* gene has duplicated three times in human history. Most nonhuman primates only have the ancestral copy, *SRGAP2A*, but contemporary humans have between two to four copies. *SRGAP2C* is fixed in humans and has been functionally shown to play a role in neocortex expansion, by delaying spine maturation and increasing spine density<sup>1</sup>. In contrast, the *SRGAP2B* and *SRGAP2D* are present between zero and four copies<sup>3</sup>. To reconstruct the evolutionary history of human-specific duplications, we developed an Approximate Bayesian Computation (ABC) method that models the pairwise divergence between duplicated copies as a function of the time, regional mutation and gene conversion rates. We implement the approach as follows: (i) We sample parameters from prior distributions:  $g$  is the per bp per gen. rate of non-allelic gene conversion  $\sim U(0, 5.9e-6)$  and  $T$  is the time the duplication arose  $\sim U(1, 12 \text{ Mya})$ . The upper bound for  $T$  is based on the duplicate being human-specific and for  $g$  on the rate of allelic gene conversion in the region<sup>10,11</sup>. (ii) We simulate data using *ms\_dup*<sup>8</sup> given a draw from the priors, and estimate the pairwise sequence divergence in the simulations (iii) We accept parameter values for those simulations in which the divergence in the simulated data is close to observed value (in practice  $< 0.1\%$ ). The parameter values that we accept provide an estimate of posterior distribution on the parameters given the data.



Figures above shows the posterior probability of *SRGAP2C* for non-allelic gene conversion rate of  $1.2e-9$  per bp per generation.

## Conclusions

Gene	% Pairwise divergence between ancestral and duplicated copy	Dates in Mya (based on the ABC method) <sup>§</sup>
<i>SRGAP2B</i>	$0.567 \pm 0.009$	4.0 – 7.6
<i>SRGAP2C</i>	$0.635 \pm 0.009$	4.4 – 8.4
<i>ARHGAP11B</i>	$1.073 \pm 0.007$	8.8 – 12.8

<sup>§</sup> assumes the revised pedigree based mutation rate of  $0.5 \times 10^{-9}$  per bp per year<sup>5,6</sup> and accounts for paralogous gene conversion and mutation rate in the region.

## References

- Charrier, C. et al. *Cell* 149, 923-935 (2012).
- Florio, M. et al. *Science* 347, 1465-1470 (2015).
- Dennis, M. Y. et al. *Cell* 149, 912-922 (2012).
- Antonacci, F. et al. *Nature genetics* 46, 1293-1302 (2014).
- Moorjani P, Gao Z, Przeworski M. *PLoS Biology*; in press <http://www.biorxiv.org/content/early/2016/08/05/058024> (2016).
- Scally A *Current Opinion in Genetics & Development* 41, 36-43 (2016).
- Sasaki, M., Lange, J. & Keeney, S. *Nature reviews Molecular cell biology* 11, 182-195 (2010).
- Teshima, K. M. & Inman, H. *Genetics* 166, 1553-1560 (2004).
- Nuttle, X. et al. *Nature Methods* 10, 903-909 (2013).
- Williams, A. L. et al. *Elife* 4, e04637 (2015).
- Halldorsson, B. V. et al. *Nature Genetics* (2016).