

# A Burst of Segmental Duplications in the African Great Ape Ancestor

Tomas Marques-Bonet *et al*, Nature 457: 877-881 (12 Feb  
2009)

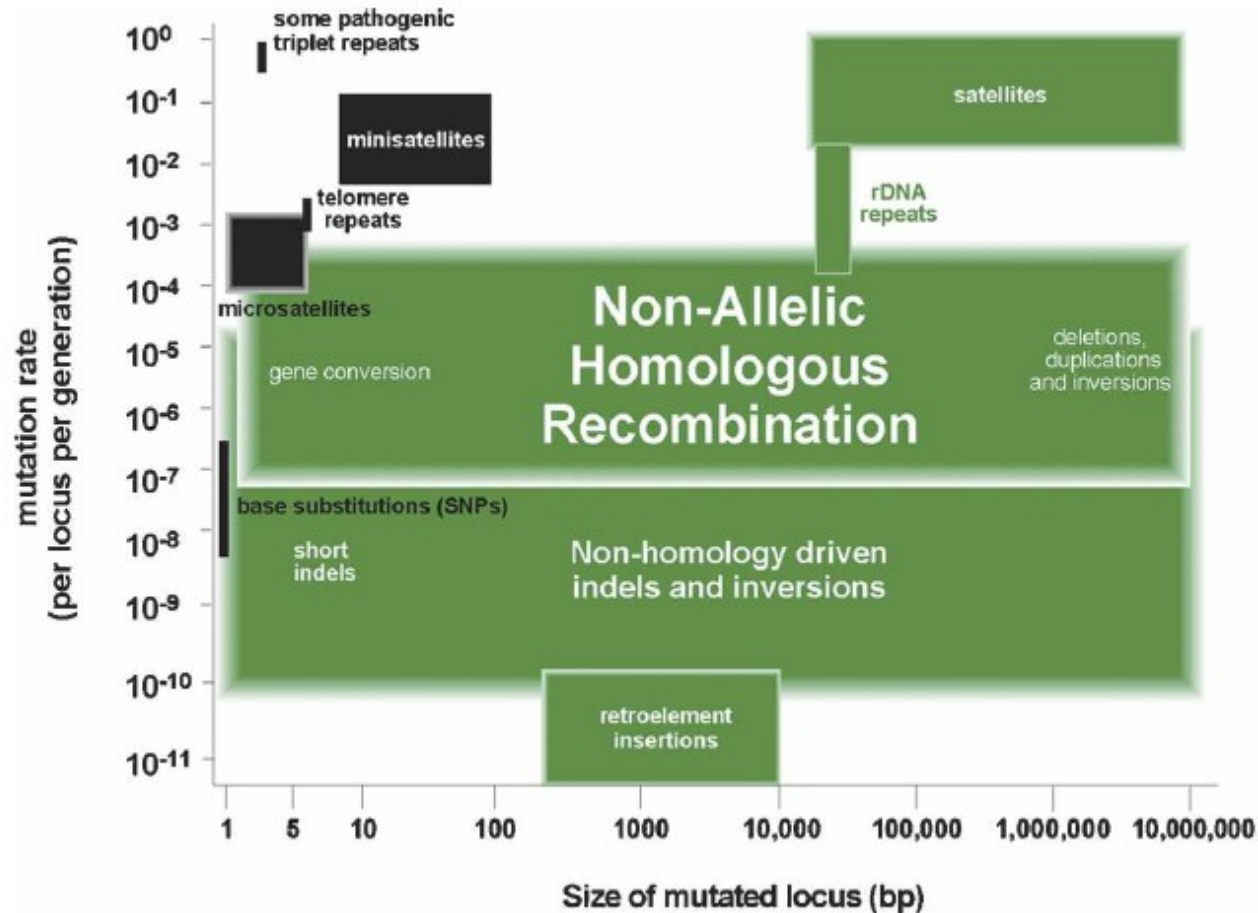
Jclub 31.03.2009  
by Tarmo Puurand

# Classes of human genetic variants

Single nucleotide variant	ATTGGCCTTAACC <b>C</b> CCGATTATCAGGAT ATTGGCCTTAACC <b>T</b> CCGATTATCAGGAT
Insertion–deletion variant	ATTGGCCTTAACCC <b>GAT</b> CCGATTATCAGGAT ATTGGCCTTAACCC <b>---</b> CCGATTATCAGGAT
Block substitution	ATTGGCCTTAAC <b>CCCC</b> GATTATCAGGAT ATTGGCCTTAAC <b>AGTG</b> GATTATCAGGAT
Inversion variant	ATTGGCCTT <b>AACCCCG</b> ATTATCAGGAT ATTGGCCTT <b>CGGGGGT</b> ATTATCAGGAT
Copy number variant	ATT <b>GGCCTTAGGCCTTA</b> ACCCCGATTATCAGGAT ATT <b>GGCCTTA-----</b> ACCTCCGATTATCAGGAT

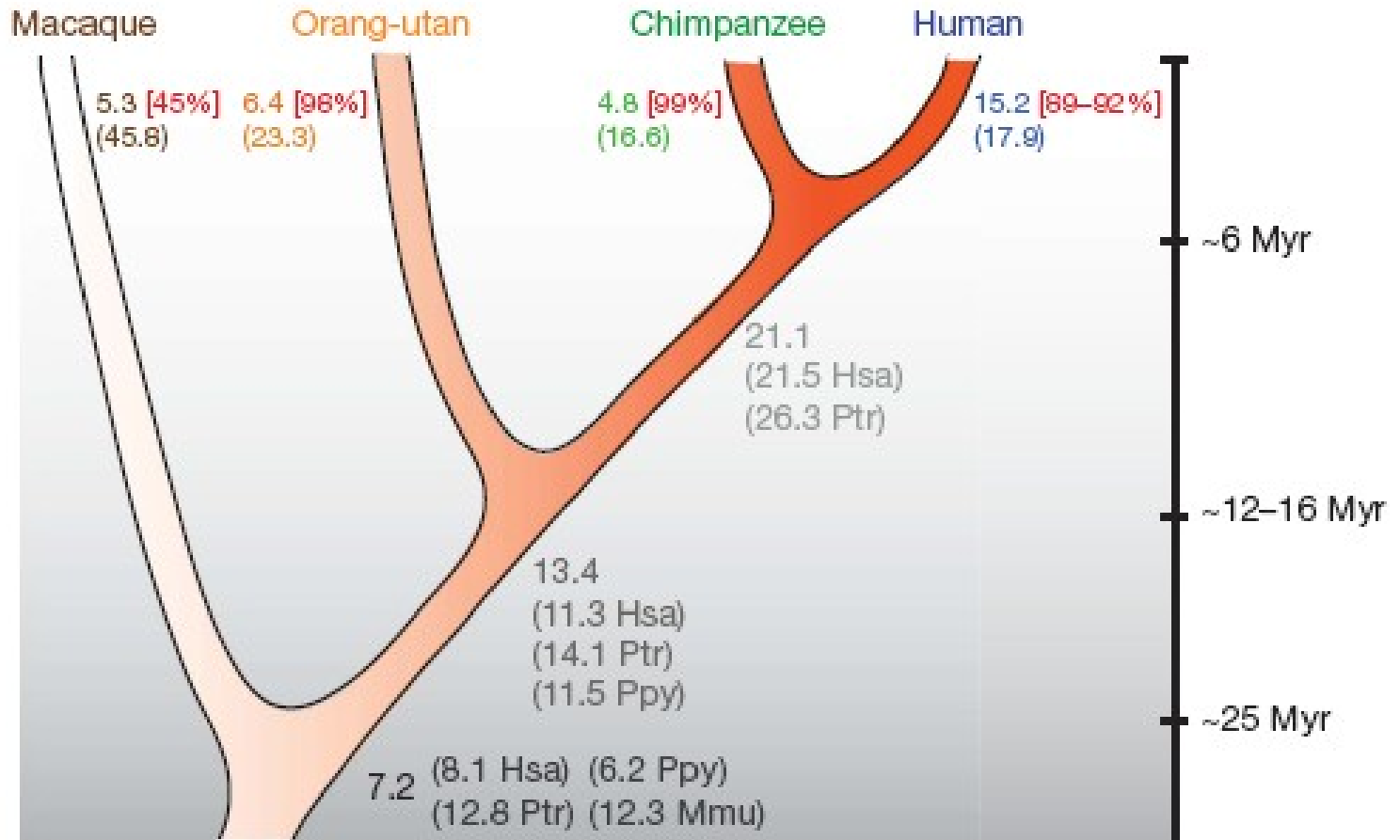
Structural variants

# Markers mutation size and rate in humans



Freeman et al, Genome Research, 2006

# Generally accepted human/great-ape phylogeny



# Detection of segmental duplications

*Supplementary Note Table 1. Primate genome datasets.*

Species	Sample ID	Source	# WGS sequence reads	Phred Quality threshold	# of reads/alignments	# Non-redundant reads placed with quality threshold	# WGS reads required (> 3 standard deviation (Autosomes/X chrom))
Human	#N/A	70% one male human being + 30% pool individuals	27,449,655	20	24,577,141	22,402,464	81/51
Chimpanzee	NS06006	male chimpanzee (Clint)	31,366,275	27	25,493,514	23,393,800	105/59
Orangutan	PR01109	female orangutan (Susie)	25,514,441	30	19,297,789	17,764,564	78
Macaque	IDI17573	female Macaque	22,590,543	27	16,769,443	13,380,372	75

*The total number of reads, the number of reads mapped against the human reference genome, and the non-redundant number of reads mapped in the reference genome are shown. PHRED quality threshold used for every species and the ID of the samples are also reported.*

# Detection of duplications (I)

- Repeatmasked sequences (build35)
- Megablast
- WGS read-depth for 6/7 consecutive 5 kbp window
- Thresholds determination for every primate separately. Some kind of calibration: WGS sequence dataset vs. sequences obtained from BAC-based clones.

# Detection of duplications (II)

- Reads: > 200 bp (20 for human, 27 for chimp, 30 for orangutan and 27 for macaque)
- >94% human and great-ape alignments and >88% for human-macaque alignment
- Regions not in human were analyzed with help of sequence contigs from certain primate

# Classes of primate segmental

**Table 1 | Classes of primate segmental duplication**

Category	Segmental duplications	Segmental duplications >20 kb	Validation (%)	Copy-number-corrected duplicated base pairs			
				Hsa	Ptr	Ppy	Mmu
Hsa	51,458,805	15,236,422	89–92	17,847,869	–	–	–
Ptr	11,239,390	4,789,874	99	–	16,583,946	–	–
Ppy	30,553,228	6,417,679	98	–	–	23,327,737	–
Mmu	24,962,092	5,360,646	45	–	–	–	45,810,964
Mmu*	35,493,466	7,715,410	85	–	–	–	18,266,656
Hsa/Ptr	32,392,480	21,061,194	NA	21,524,417	26,304,286	–	–
Hsa/Ptr/Ppy	25,450,827	13,402,545	NA	11,259,061	14,012,351	11,541,148	–
Hsa/Ptr/Ppy/Mmu	14,094,156	7,156,616	NA	8,092,997	12,820,607	6,176,876	12,542,691
Total	190,150,978	73,424,976	–	58,724,344	69,721,190	41,045,761	30,809,347

Duplications were divided into eight categories based on the WSSD analysis of each primate genome (subsequent analyses were restricted to segmental duplications >20 kb in length). Lineage-specific and shared duplication content are indicated. Percentage validation indicates the fraction of species-specific duplications confirmed by cross-species array comparative genomic hybridization. Because the human genome was used, we corrected for copy number and examined sequence contigs not aligned to the human genome (see Methods). Segmental duplications assigned to the Y chromosome were not considered.

\*Macaque segmental duplications detected in the macaque reference genome using WSSD and WGAC (<94% identity) approaches.

<sup>1</sup>Department of Genome Sciences, University of Washington and the Howard Hughes Medical Institute, Seattle, Washington 98195, USA. <sup>2</sup>Institut de Biologia Evolutiva (UPF-CSIC), 08003 Barcelona, Catalonia, Spain. <sup>3</sup>Sezione di Genetica-Dipartimento di Anatomia Patologica e Genetica, University of Bari, 70125 Bari, Italy. <sup>4</sup>Genome Sequencing Center, Washington University School of Medicine, St Louis, Missouri 63108, USA. <sup>5</sup>Institució Catalana de Recerca i Estudis Avançats (ICREA) and Instituto Nacional de Bioinformática (INB), Dr. Ariguader 88, 08003 Barcelona, Spain.



# Distribution of SDs in kb (all)

Category	Total bp	N	AVG (length)	STD Dev (length)	MAX length	MIN length
HSA specific SDs	<b>51,458,805</b>	5,887	8,741	13,318	292,021	49
PTR specific SDs	<b>11,129,390</b>	1,169	9,520	18,223	341,154	21
PPY specific SDs	<b>30,299,228</b>	3,797	7,980	11,028	275,363	2
MMU specific SDs	<b>24,962,092</b>	2,463	10,135	8,698	149,378	41
HSA / PTR	<b>32,392,480</b>	2,018	16,052	22,340	345,000	36
HSA / PPY	<b>9,787,003</b>	1,586	6,171	6,823	71,000	27
HSA / MMU	<b>3,989,127</b>	740	5,391	6,495	93,000	21
PTR / PPY	<b>1,080,458</b>	244	4,428	4,938	44,610	51
PTR / MMU	<b>577,152</b>	100	5,772	7,094	52,050	545
PPY / MMU	<b>1,650,595</b>	321	5,142	4,158	27,000	26
HSA / PTR / PPY	<b>25,450,827</b>	1,770	14,379	17,384	234,000	21
HSA / PTR / MMU	<b>5,889,226</b>	782	7,531	7,844	63,000	21
HSA / PPY / MMU	<b>3,473,366</b>	529	6,566	6,838	47,260	9
PTR / PPY / MMU	<b>190,558</b>	69	2,762	2,580	15,330	325
HSA / PTR / PPY / MMU	<b>14,094,156</b>	1,011	13,941	13,489	168,780	38
<b>Total</b>	<b>216,424,463</b>	<b>22,486</b>	<b>9,625</b>	<b>13,692</b>	<b>345,000</b>	<b>2</b>

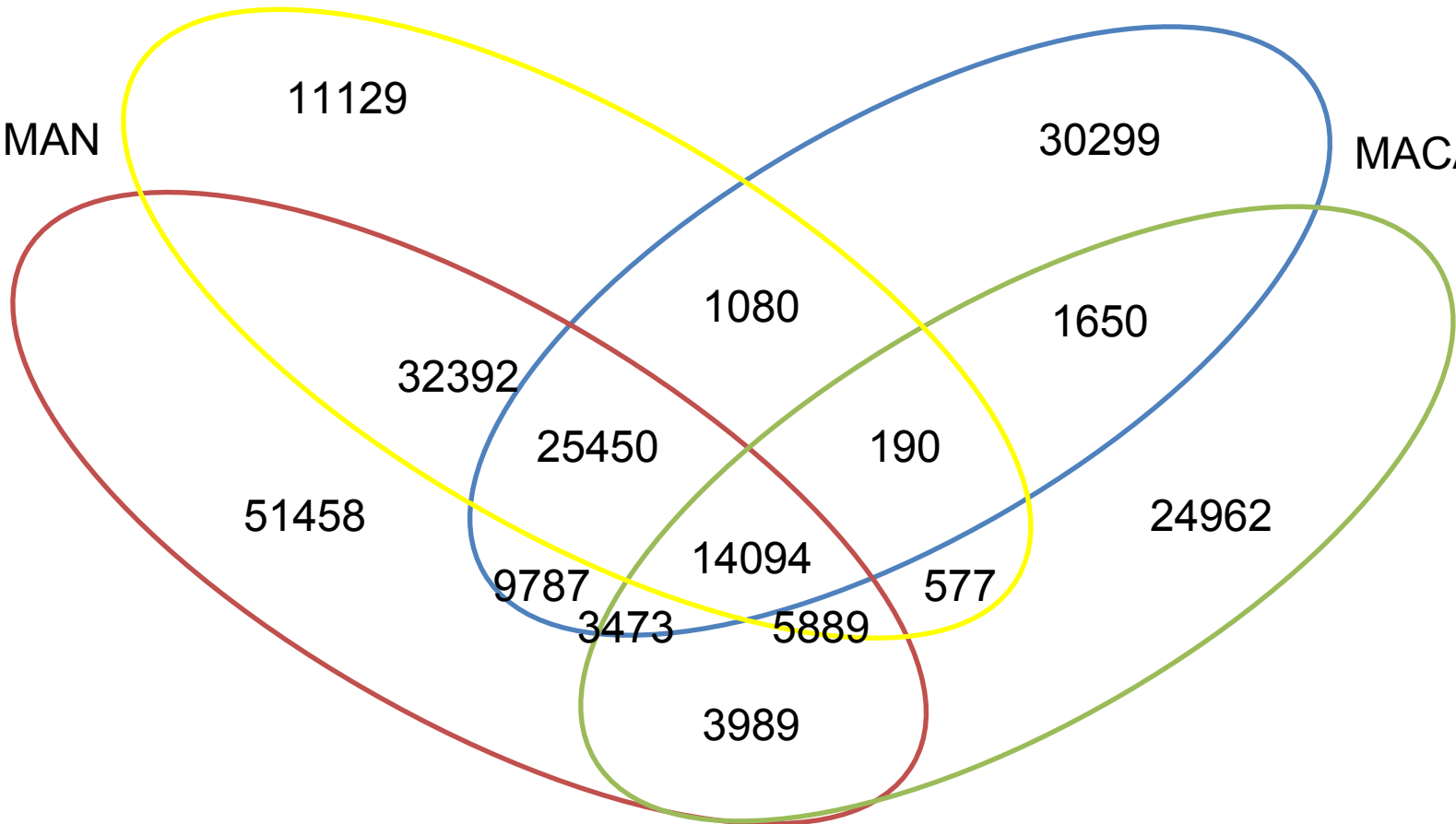
# Distribution of SDs in kb (all)

CHIMPANZEE

ORANG-UTAN

HUMAN

MACAQUE



# Distribution of SDs in events (>20kb)

Category	Total bp (>20 kb)	N	AVG (length)	STD Dev (length)	MAX length	MIN length
HSA specific SDs	<b>15,236,422</b>	315	48,370	37,561	292,021	20,035
PTR specific SDs	<b>4,789,874</b>	96	49,895	46,215	341,154	20,024
PPY specific SDs	<b>6,417,679</b>	137	46,844	39,283	275,363	20,076
MMU specific SDs	<b>5,360,646</b>	162	33,090	16,531	149,378	20,047
HSA / PTR	<b>21,061,194</b>	479	43,969	31,495	345,000	20,023
HSA / PPY	<b>1,452,735</b>	45	32,283	14,097	71,000	20,161
HSA / MMU	<b>392,712</b>	9	43,635	27,989	93,000	21,698
PTR / PPY	<b>86,700</b>	2	43,350	1,782	44,610	42,090
PTR / MMU	<b>135,794</b>	4	33,949	12,436	52,050	23,748
PPY / MMU	<b>27,000</b>	1	27,000		27,000	27,000
HSA / PTR / PPY	<b>13,402,545</b>	322	41,623	25,497	234,000	20,012
HSA / PTR / MMU	<b>1,545,552</b>	51	30,305	9,943	63,000	20,026
HSA / PPY / MMU	<b>704,864</b>	23	30,646	9,438	47,260	20,065
PTR / PPY / MMU						
HSA / PTR / PPY / MMU	<b>7,156,616</b>	201	35,605	15,546	168,780	20,025
<b>Total</b>	<b>77,770,333</b>	<b>1,847</b>	<b>42,106</b>	<b>30,462</b>	<b>345,000</b>	<b>20,012</b>

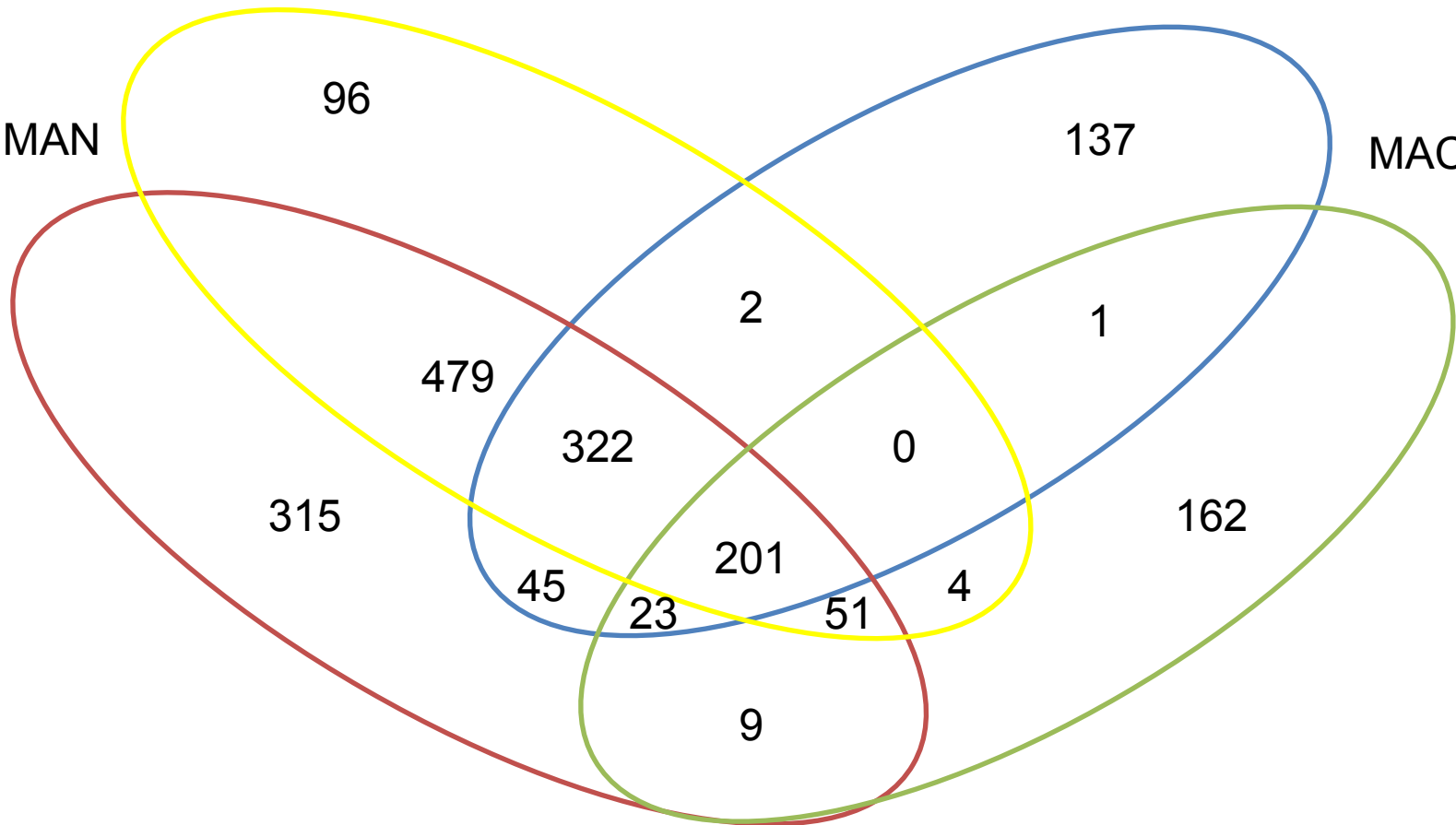
# Distribution of SDs in events (>20kb)

CHIMPANZEE

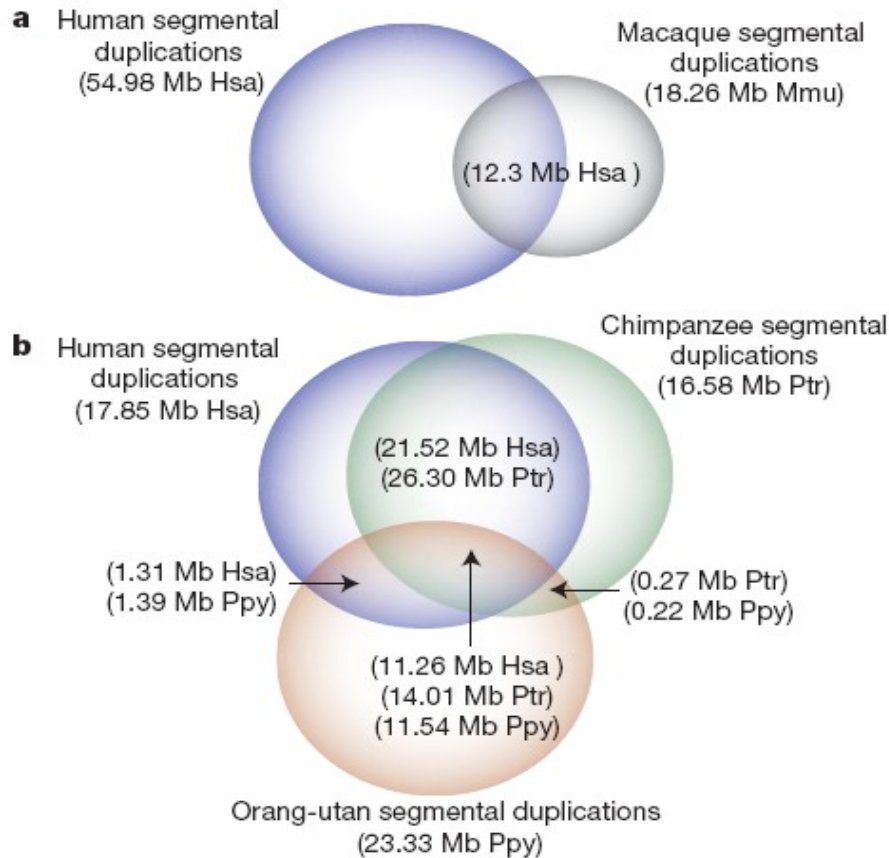
ORANG-UTAN

HUMAN

MACAQUE



# Shared vs. lineage-specific duplications



# Validation

- FISH- 58 lineage specific (14 human, 24 chimp, 20 orangutan) and 38 complex regions > 40 kbp in length
- Interspecific array comparative genomic hybridisation

*Supplementary Note Table 9. Segmental duplications and copy-number polymorphism.*

<b>Human copy-number polymorphisms (n = 8 individuals)</b>					
<b>SD Category</b>	<b># CNV SD intervals</b>	<b>Total Length (bp)</b>	<b>CNV SD Intervals</b>	<b>Total Length CNV SD</b>	<b>% CNV SD</b>
Human specific SDs	199	9,809,268	106	5,018,693	<b>33.9%</b>
Human/chimpanzee shared SDs	300	12,222,058	179	8,839,136	<b>42.0%</b>
Human/chimp/orang shared SDs	235	10,303,447	87	3,099,098	<b>23.1%</b>
Human/chimp/orang/macaque shared SDs	145	5,114,155	56	2,042,461	<b>28.5%</b>
Chimpanzee-specific SDs	91	4,684,302	2	42,140	<b>0.9%</b>
Orangutan-specific SDs	134	6,344,870	2	51,655	<b>0.8%</b>
Macaque-specific SDs	148	4,894,873	12	414,331	<b>7.8%</b>
<b>Total</b>	<b>1,252</b>	<b>53,372,973</b>	<b>444</b>	<b>19,507,514</b>	<b>26.8%</b>

---

**Chimpanzee copy-number polymorphisms (n = 8 individuals)**

---

<b>SD Category</b>	<b># CNV SD intervals</b>	<b>Total Length (bp)</b>	<b>CNV SD Intervals</b>	<b>Total Length CNV SD</b>	<b>% CNV SD</b>
Human-specific SDs	255	12,842,592	50	1,985,369	<b>13.4%</b>
Human/chimpanzee shared SDs	312	12,224,102	167	8,837,092	<b>42.0 %</b>
Human/chimp/orang shared SDs	204	8,591,738	118	4,810,807	<b>35.9%</b>
Human/chimp/orang/macaque shared SDs	110	3,761,322	91	3,395,294	<b>47.4%</b>
Chimpanzee-specific SDs	35	1,443,956	58	3,282,486	<b>69.4%</b>
Orangutan-specific SDs	135	6,343,149	1	53,376	<b>0.8%</b>
Macaque-specific SDs	149	4,829,875	11	479,329	<b>9.0%</b>
<b>Total</b>	<b>1,200</b>	<b>50,036,734</b>	<b>496</b>	<b>22,843,753</b>	<b>31.3%</b>

---



---

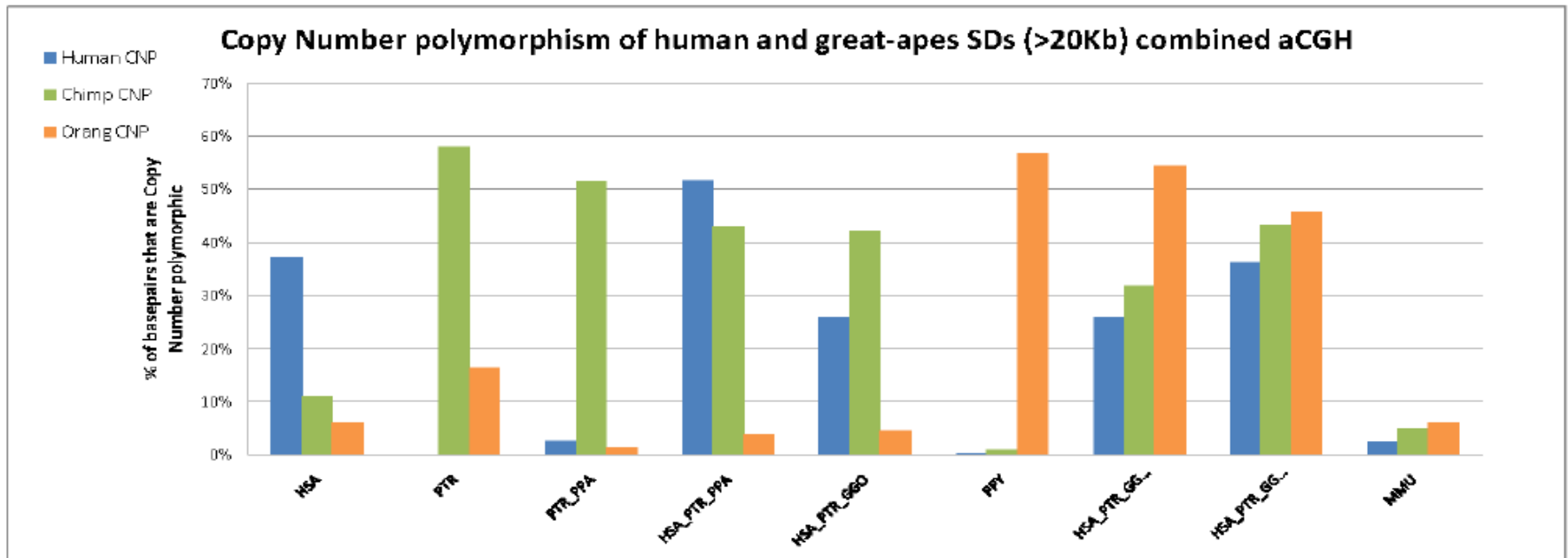
**Orangutan copy-number polymorphisms (n = 8 individuals)**

---

<b>SD Category</b>	<b># CNV SD intervals</b>	<b>Total Length (bp)</b>	<b>CNV SD Intervals</b>	<b>Total Length CNV SD</b>	<b>% CNV SD</b>
Human-specific SDs	276	13,794,990	29	1,032,971	<b>7.0%</b>
Human/chimpanzee shared SDs	452	20,060,117	27	1,001,077	<b>4.8%</b>
Human/chimp/orang shared SDs	146	5,667,445	176	7,735,100	<b>57.7%</b>
Human/chimp/orang/macaque shared SDs	107	3,688,249	94	3,468,367	<b>48.5%</b>
Chimpanzee specific SDs	88	4,560,303	5	166,139	<b>3.5%</b>
Orangutan specific SDs	72	2,934,820	64	3,461,705	<b>54.1%</b>
Macaque specific SDs	146	4,826,751	14	482,453	<b>9.1%</b>
<b>Total</b>	<b>1,287</b>	<b>55,532,675</b>	<b>409</b>	<b>17,347,812</b>	<b>23.8%</b>

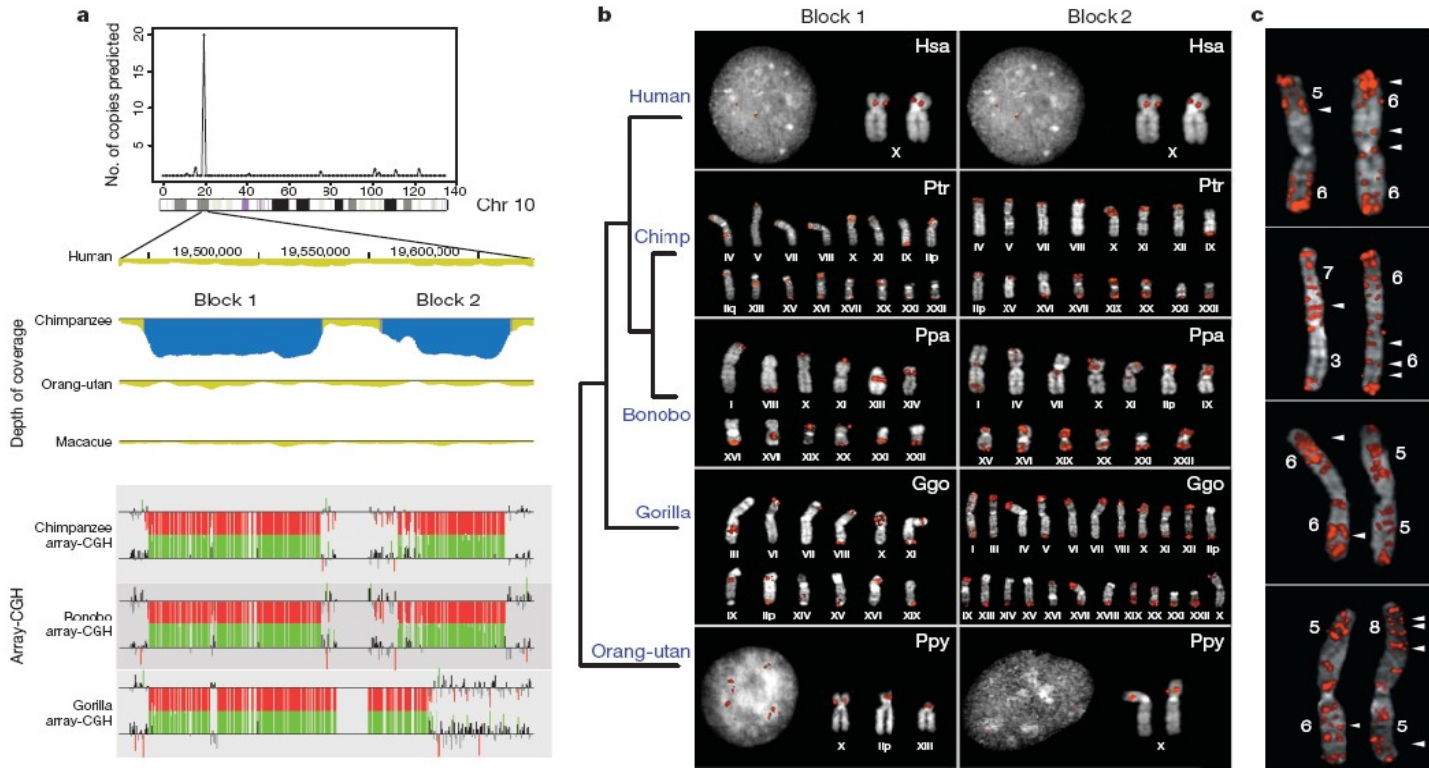
---

# CNV-s in SD regions



*Supplementary Note Fig. 5. Copy-number polymorphism of human and great-ape SDs. In this figure, SDs were further categorized (see Fig. 2c) using arrayCGH information from gorilla and bonobo. The same trends reported in the text are observed.*

# Extreme duplications



**Figure 3 | Convergent gene duplication expansion in African great apes but not humans.** **a**, Two regions on chromosome 10 have expanded in chimpanzee, gorilla and bonobo when compared to human based on computational and interspecific array-CGH experiments (see Fig. 1 legend). **b**, FISH confirms 23–50 copies in chimpanzee and bonobo (*Ppa*, *Pan paniscus*), and >100 copies in gorilla (*Ggo*, *Gorilla gorilla*) (Methods).

End-sequence pair analysis using gorilla and chimpanzee WGS sequences reveals that all but the ancestral location are non-orthologous, indicating independent expansions in both lineages. **c**, Detailed analysis of eight homologues of gorilla chromosome 1 reveals interstitial locations of the block 2 duplication that show variation both in copy number and in terms of location.

# Rates of duplications (Mbp)

*Supplementary Note Table 14. Great-ape comparisons.*

	subt/1000bp	SDs (Mb)	Rate Mb/subt per 1000 bp	Million Year per branch	SDs (Mb)	Rate Mb/Mya
Human terminal branch	5.4	13.6	<b>2.519</b>	6	13.6	<b>2.267</b>
Chimpanzee terminal branch	5.56	6.1	<b>1.097</b>	6	6.1	<b>1.017</b>
Human/chimpanzee shared branch	1.07	9.32	<b>8.710</b>	2	9.32	<b>4.660</b>
Gorilla terminal branch	7.19			8		
Human/chimpanzee/gorilla shared branch	7.62	16.82	<b>2.207</b>	4	16.82	<b>4.205</b>
Orangutan terminal branch	15.03	20.33	<b>1.353</b>	12	20.33	<b>1.694</b>
Human/chimpanzee/orangutan shared branch	14.7213	15.97	<b>1.085</b>	13	15.97	<b>1.228</b>

*The rates of segmental duplication (>20 kbp) accumulation on different branches were compared as a function of millions of years since divergence and as a function of the genetic distance (single basepair substitutions)<sup>38</sup> between the species.*

# Rates of duplications (events)

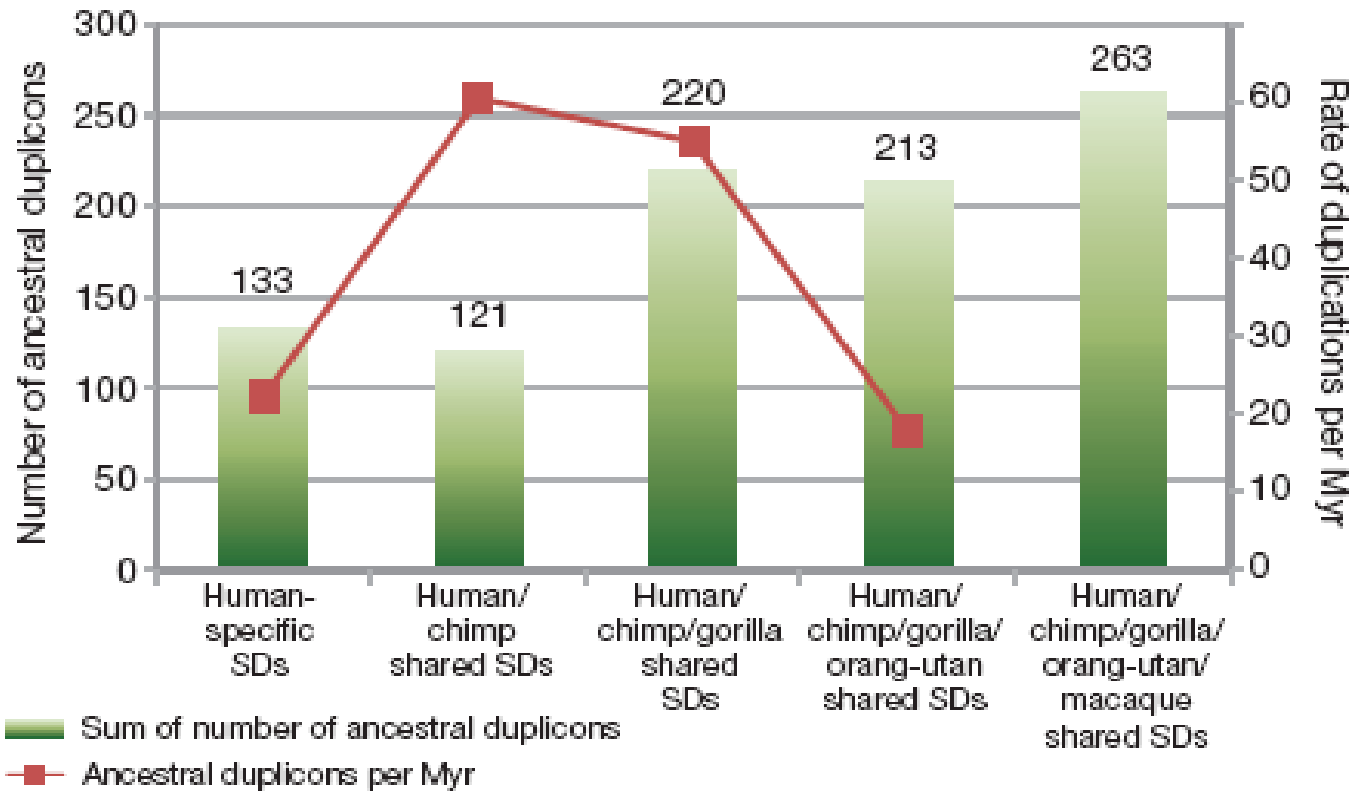
*Supplementary Note Table 15. Hominid rates of duplication (events >20 kbp).*

	<b>Number of chained sub-units human Duplicons</b>	<b>Rate of Duplications (Events /Myr)</b>	<b>subt/1000bp</b>	<b>Rate of Duplications (Events /substitutions)</b>
HSA specific	133	<b>22.17</b>	5.40	<b>24.63</b>
PTR specific			5.56	
HSA/PTR shared	121	<b>60.50</b>	1.07	<b>113.08</b>
HSA/PTR/GGO shared	220	<b>55.00</b>	7.62	<b>28.87</b>
PPY specific			15.03	
HSA/PTR/PPY Shared	213	<b>16.38</b>	14.72	<b>14.47</b>

*950 duplicons detected previously<sup>11</sup> were used as a surrogate for duplication events. Two measures of time were applied to calculate the rates: a) million years of divergence and b) genetic distance estimates<sup>38</sup>.*

# Rates of SD

c



# Other analysis

- Nonrandom distribution of primate SD
- Gene duplication analyses
- SD and disease susceptibility loci
- Duplication status vs. copy number
- Lineage specific deletions

# Venter vs. Watson (CNV-s)

*Supplementary Note Table 3. Copy-number variation of shared and individual-specific human SDs.*

<b>Cat_SD</b>	<b>Invariant</b>	<b>%</b>	<b>CNVs</b>	<b>%</b>	<b>Grand Total</b>	<b>Fisher exact test P-value (vs. Shared)</b>
<b>SHARED</b>	109	<i>18.3</i>	486	<i>81.6</i>	595	
<b>VENTER</b>	18	<i>41.9</i>	25	<i>58.1</i>	43	0.0001989417
<b>WATSON</b>	17	<i>60.7</i>	11	<i>39.3</i>	28	0.0000002341
<b>Grand Total</b>	<i>144</i>		<i>522</i>		<i>666</i>	