

Supplemental Data

Ancient Substructure in Early mtDNA

Lineages of Southern Africa

Chiara Barbieri, Mário Vicente, Jorge Rocha, Sununguko W. Mpoloka, Mark Stoneking, and Brigitte Pakendorf

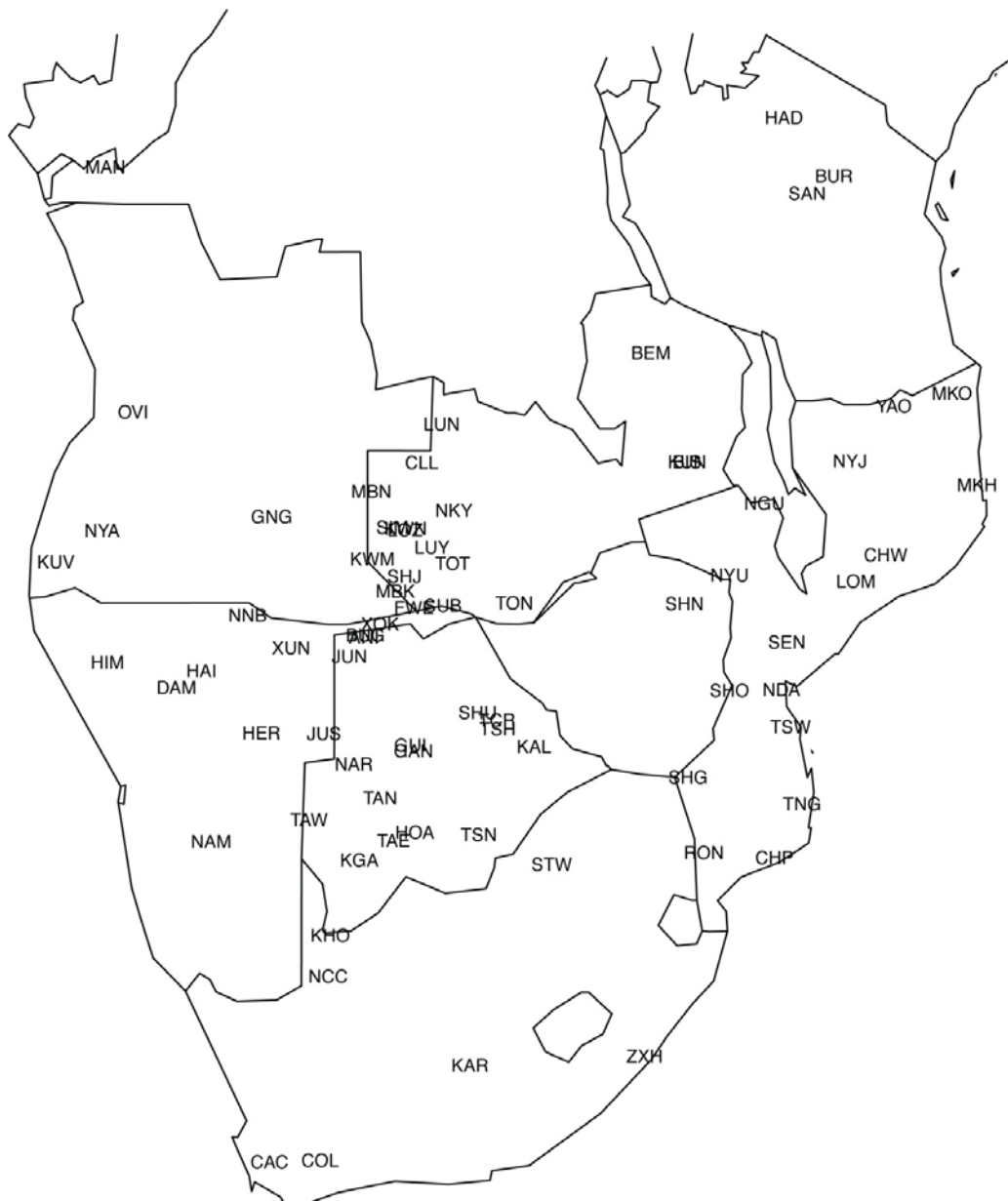


Figure S1: Map of approximate locations of the populations included in the surfer map (Figure 1 in the main text).

Population codes as indicated in Supplementary Table 1.

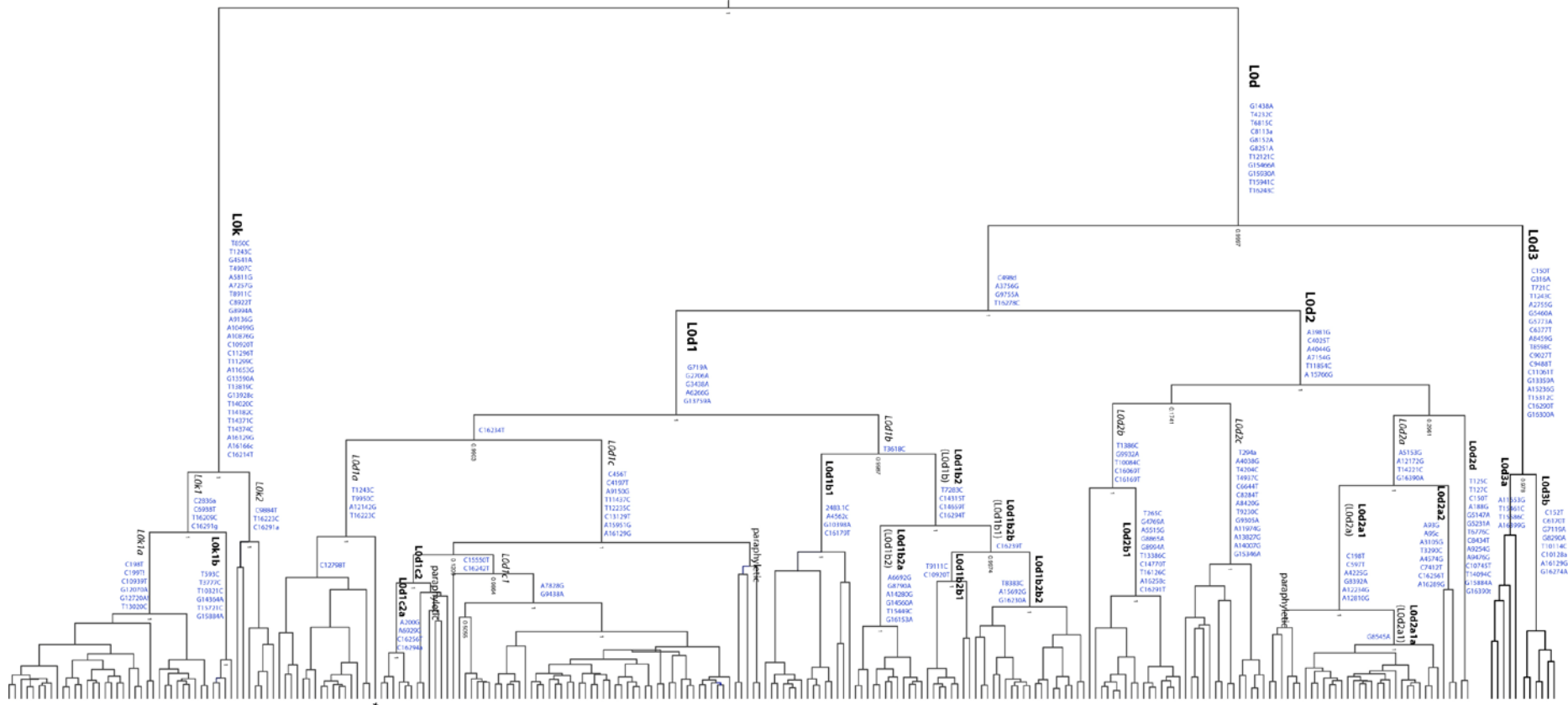


Figure S2: Tree of the 254 unique haplotypes of the dataset.

The tree is based on full sequences, without a time scale, and includes the major branches L0k, L0d1, L0d2 and L0d3. Sub-branches that have not changed are labeled in italic font; new branches defined here are labeled in bold font. When previously defined branches have to be renamed, the older label is indicated in brackets. The posterior probabilities associated with major nodes are shown. Mutations defining branches are shown in blue font: transversions are indicated with lowercase, and back mutations to an ancestral state are indicated with an exclamation mark (!). The individual marked with an asterisk is mentioned in the note about positions 199 and 16266 in Supplementary Table 3.

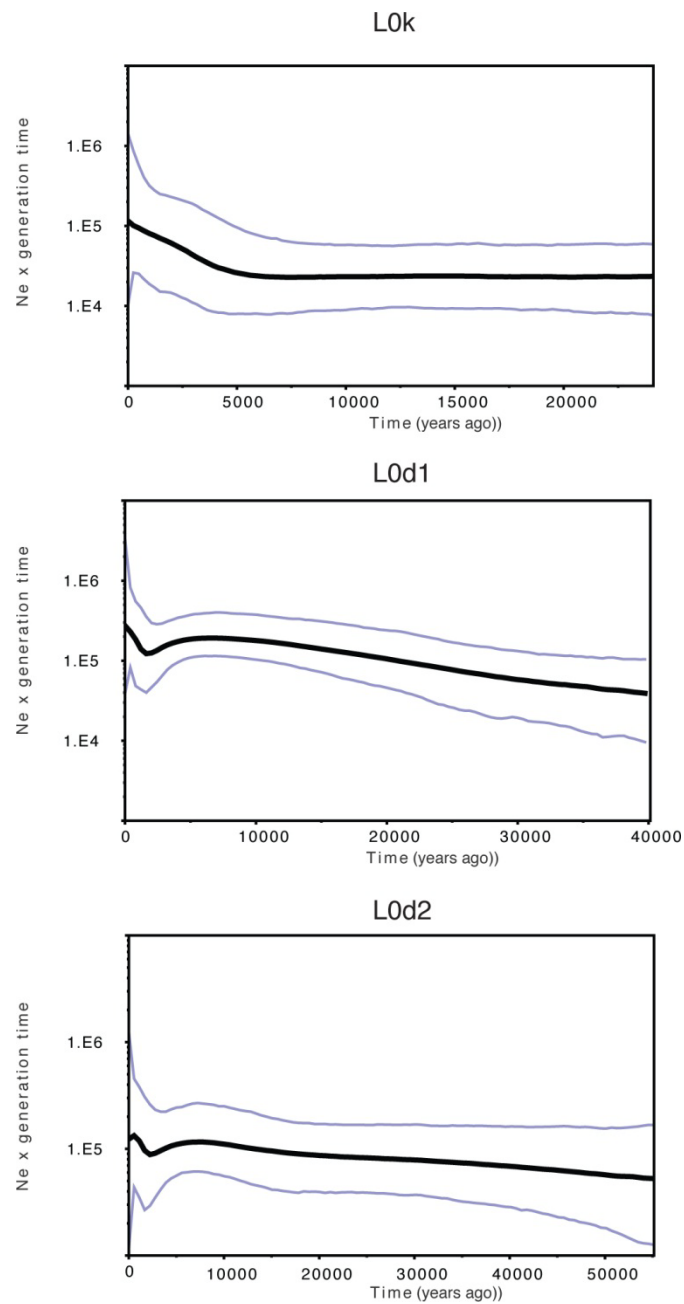


Figure S3: Bayesian Skyline Plots (BSP) of the L0k, L0d1 and L0d2 haplogroups.

The BSPs are based on the mtDNA coding region, estimated with 10 million iterations. The y axis for each plot is the product of the effective population size and the generation time and the x axis shows time using a linear relaxed clock with the substitution rate of 1.26×10^{-8} per site per year.

Ovimbundu	OVI	Angola	Bantu	14.915771	-11.888853	60	3.3	1.6	1.7	1.7	0.0	1.7	0.0	0.0	PRESENT STUDY
Anikhoë	ANI	Botswana	Khoë	21.8850954	-18.3734521	18	44.4	22.2	44.4	0.0	0.0	22.2	0.0	0.0	PRESENT STUDY
Xokhwe	XOK	Botswana	Khoë	22.3761	-17.9957	17	17.6	11.8	17.6	0.0	0.0	11.8	0.0	0.0	PRESENT STUDY
Bugakhoe	BUG	Botswana	Khoë	21.9367	-18.3219	14	42.9	28.6	42.9	0.0	0.0	28.6	0.0	0.0	PRESENT STUDY
Naro	NAR	Botswana	Khoë	21.5840541	-22.0320817	35	77.1	17.1	51.4	25.7	0.0	17.1	0.0	0.0	PRESENT STUDY
G ana	GAN	Botswana	Khoë	23.3889	-21.6523	15	93.3	6.7	80.0	13.3	0.0	6.7	0.0	0.0	PRESENT STUDY
G ui	GUI	Botswana	Khoë	23.2946698	-21.486584	31	93.5	3.2	51.6	41.9	0.0	3.2	0.0	0.0	PRESENT STUDY
Hai om	HAI	Namibia	Khoë	16.9694944	-19.3450768	51	68.6	13.7	39.2	27.5	2.0	13.7	0.0	0.0	PRESENT STUDY
Nama	NAM	Namibia	Khoë	17.2608889	-24.2660935	29	79.3	3.4	37.9	34.5	6.9	0.0	3.4	0.0	PRESENT STUDY
Damara	DAM	Namibia	Khoë	16.2257392	-19.8301838	38	13.2	0	10.5	2.6	0.0	0.0	0.0	0.0	PRESENT STUDY
Shua	SHU	Botswana	Khoë	25.3321307	-20.5502369	42	35.7	2.4	35.7	0.0	0.0	2.4	0.0	0.0	PRESENT STUDY
TcireTcire	TCR	Botswana	Khoë	25.9166477	-20.7658488	12	50	16.7	41.7	8.3	0.0	16.7	0.0	0.0	PRESENT STUDY
Tshwa	TSH	Botswana	Khoë	25.9365757	-21.0249347	22	54.5	0	50.0	4.5	0.0	0.0	0.0	0.0	PRESENT STUDY
†Hoan	HOA	Botswana	K'xa	23.4351167	-23.9989176	13	100	0	92.3	7.7	0.0	0.0	0.0	0.0	PRESENT STUDY
!Xuun	XUN	Botswana	K'xa	19.6826306	-18.6907202	27	55.5	33.3	44.4	11.1	0.0	33.3	0.0	0.0	PRESENT STUDY
Ju 'hoan North	JUN	Botswana	K'xa	21.4524476	-18.9372569	40	72.5	22.5	50.0	22.5	0.0	22.5	0.0	0.0	PRESENT STUDY
Ju 'hoan South	JUS	Botswana	K'xa	20.6815392	-21.151918	44	70.5	25	50.0	20.5	0.0	25.0	0.0	0.0	PRESENT STUDY
Taa East	TAE	Botswana	Tuu	22.8206545	-24.2365162	30	100	0	46.7	53.3	0.0	0.0	0.0	0.0	PRESENT STUDY
Taa North	TAN	Botswana	Tuu	22.4158579	-23.0145647	25	84	16	68.0	16.0	0.0	16.0	0.0	0.0	PRESENT STUDY
Taa West	TAW	Botswana	Tuu	20.2727412	-23.639938	31	74.2	22.6	51.6	22.6	0.0	22.6	0.0	0.0	PRESENT STUDY
Shona	SHN	Zimbabwe	Bantu	31.593017	-17.413546	59	1.7	1.7							Castrì et al. 2009
Kunda	KUN	Zambia	Bantu	31.671753	-13.325485	36	2.8	0							De Filippo et al. 2010
Bisa	BIS	Zambia	Bantu	31.67175	-13.325483	46	0	0							De Filippo et al. 2010, present study
SA Coloured	COL	South Africa	Indoeuropean	20.562744	-33.449777	563	60	0							Quintana-Murci et al. 2010
Chopi	CHP	Mozambique	Bantu	34.317627	-24.726875	27	0	0							Salas et al. 2002
Chwabo	CHW	Mozambique	Bantu	37.679443	-16.003576	20	0	0							Salas et al. 2002
Lomwe	LOM	Mozambique	Bantu	36.778564	-16.762468	20	0	0							Salas et al. 2002
Makhwa	MKH	Mozambique	Bantu	40.447998	-13.987376	20	0	0							Salas et al. 2002
Makonde	MKO	Mozambique	Bantu	39.700927	-11.350797	19	5.3	0							Salas et al. 2002
Ndau	NDA	Mozambique	Bantu	34.537353	-19.890723	19	30	0							Salas et al. 2002

Nguni	NGU	Mozambique	Bantu	34.010009	-14.51978	11	0	0							Salas et al. 2002
Nyanja	NYJ	Mozambique	Bantu	36.602783	-13.304103	20	0	0							Salas et al. 2002
Nyungwe	NYU	Mozambique	Bantu	32.955322	-16.594081	20	0	0							Salas et al. 2002
Ronga	RON	Mozambique	Bantu	32.186279	-24.58709	21	19	0							Salas et al. 2002
Sena	SEN	Mozambique	Bantu	34.691162	-18.521283	21	0	0							Salas et al. 2002
Shangaan	SHG	Mozambique	Bantu	31.70288	-22.411029	22	4.5	0							Salas et al. 2002
Shona	SHO	Mozambique	Bantu	32.955322	-19.911384	18	0	0							Salas et al. 2002
Tonga	TNG	Mozambique	Bantu	35.152587	-23.180764	20	5	0							Salas et al. 2002
Tswa	TSW	Mozambique	Bantu	34.801025	-20.96144	19	15.8	0							Salas et al. 2002
Yao	YAO	Mozambique	Bantu	37.965087	-11.716788	10	0	0							Salas et al. 2002
Karretjie Mense	KAR	South Africa	Indoeuropean ^a	25.101013	-30.712638	30	100	0							Schlebusch et al. 2011
Cape Colured	CAC	South Africa	Indoeuropean	19.037475	-33.495598	20	45	0							Schlebusch 2010
Khomani	KHO	South Africa	Tuu	20.872192	-26.971038	57	98.2	0							Schlebusch 2010
Manyanga	MAN	DRC	Bantu	14.058837	-4.82826	14	0	0							Schlebusch2010
Northern Cape Coloured	NCC	South Africa	Indoeuropean	20.804443	-28.149503	40	92.5	0							Schlebusch 2010
Sotho Tswana	STW	South Africa	Bantu	27.572021	-24.926295	22	22.7	0							Schlebusch 2010
Zulu Xhosa	ZXH	South Africa	Bantu	30.384521	-30.448674	36	44.4	2.8							Schlebusch 2010
Burunge	BUR	Tanzania	Cushitic	36.119384	-5.090944	38	3	0							Tishkoff et al. 2007
Hadza	HAD	Tanzania	Khoisan (isolated)	34.603271	-3.403758	79	0	0							Tishkoff et al. 2007
Sandawe	SAN	Tanzania	Khoisan (isolated)	35.306396	-5.594118	82	5	0							Tishkoff et al. 2007

^a this population used to speak a Tuu language but has shifted to Afrikaans.

Table S2: List of positions (numbered in accordance with the RSRs/rCRS) with missing data that were excluded from the analysis. Polymorphic sites are underlined.

<u>316</u>
<u>1243</u>
3106
3492
3516
<u>3981</u>
<u>4232</u>
<u>5515</u>
<u>5936</u>
6716
<u>6938</u>
7412
<u>8563</u>
<u>10550</u>
10589
11854
<u>13020</u>
13198
<u>13386</u>
<u>14770</u>
<u>15530</u>
<u>15930</u>
<u>15941</u>
<u>16069</u>
<u>16093</u>
<u>16169</u>
<u>16212</u>
16215
<u>16230</u>
<u>16242</u>
<u>16243</u>
16474

Table S3: Notes on some of the haplogroup-defining mutations.

Mutation	Remarks
C152T	This mutation defines L0d3b, but is also present in an individual belonging to L0d3a as well as being found sporadically in other branches of L0d and L0k.
A188G	In Supplementary Figure 2, this is shown only for L0d2d; however, this mutation also occurs in nearly all the individuals belonging to L0d1b1, with only 2 exceptions.
C198T	In Supplementary Figure 2, this is shown for L0d2a1 and L0k1a; however, this mutation also defines a minor subbranch of L0d1c1 (rather than defining L0d1c1 as a whole, as previously thought).
199	<p>The evolutionary pathway involving L0k cannot be resolved, since L0k2 and L0k1b carry a C at this position, while L0k1a carries a T, which is the state reconstructed for the RSRS. In Supplementary Figure 2 we show the C199T back mutation as defining L0k1a; however, with our dataset it is equally likely that two independent T-C transitions occurred on the branches leading to L0k2 and L0k1b, with L0k1a retaining the ancestral T.</p> <p>In addition, L0d1a carries a C at this position with the exception of three lineages not forming a clade. One of these is a deeply divergent lineage represented by only one individual from Botswana (indicated by an asterisk in Supplementary Figure 2); thus, one could postulate either three back mutations from the mutation defining L0d1a as a whole, or consider T199C a defining mutation only for the subclade L0d1a1, with two back mutations having occurred subsequently. Since C16266a is also missing in the divergent lineage (see below), one should perhaps consider both T199C and C16266a as mutations defining the subclade L0d1a1, with subsequent back mutations (C199T) or novel mutations (A16266G) in some individuals.</p>
294	In Supplementary Figure 2, we show the T-A transversion defining L0d2c; in addition, a T-C transition defines a subbranch of L0d1c1.
A7828G	Rather than defining branch L0d1c1 as a whole, as previously suggested, this is missing from one individual and thus defines only a subbranch, as shown in Supplementary Figure 2.
C8922T	This is found in L0k2, L0k1b, and several branches of L0k1a, but is missing from one subbranch of L0k1a. The most plausible reconstruction is that the transition occurred on the branch leading to L0k, as previously assumed,
G8994A	In Supplementary Figure 2, this is shown only for L0d2b1 and L0k; however, this mutation also defines a small subbranch of the paraphyletic branch of L0d1c.
A9136G	This mutation defining L0k mutates back to A in a subbranch of L0k2.
A9347G	This mutation is at the root of haplogroup L0, but almost all of the L0k2 individuals present a back mutation at this site, with the exception of the sample from Yemen.

G9438A	Rather than defining branch L0d1c1 as a whole, as previously suggested, this is missing from one individual and thus defines only a subbranch, with a further back mutation to A9438G found in one sequence.
A11653G	This mutation defines L0k as well as L0d3a.
C15550T	Together with C16242T, this is the only mutation defining branch L0d1c1; A7828G and G9438A are missing from one divergent lineage and thus define only a subset of L0d1c1 (with a further back mutation to A9438G found in one sequence), while C198T defines an even smaller branch within L0d1c1.
T15586C	This mutation defining L0d3a mutates back to T in one individual of the same subbranch.
A16129G	This mutation defines L0d1c, L0k and L0d3b, as well as a subbranch of L0d1b2a. Given the hypervariability of this position, it is not surprising that several back mutations to A occur in the tree – the most notable being a back mutation in the individual from Yemen whose sequence up to now was the only lineage known for L0k2. Therefore, A16129G was previously considered a mutation defining only L0k1; with our extended dataset we show that it defines all of L0k.
T16209C	This mutation, which defines L0k1, also appears in a subbranch of L0d1a.
C16242T	Together with C15550T, this is the only mutation defining branch L0d1c1; A7828G and G9438A are missing from one divergent lineage and thus define only a subset of L0d1c1 (with a further back mutation to A9438G found in one sequence), while C198T defines an even smaller branch within L0d1c1.
16266	Like the T-C transition at 199, C16266a is not found in all the sequences belonging to L0d1a; rather, four sequences carry a G at this position. Since one of these is the divergent lineage represented by an asterisk in Supplementary Figure 2 (as mentioned for position 199 above), one should perhaps consider both T199C and C16266a as mutations defining the subclade L0d1a1, with subsequent back mutations (C199T) or novel mutations (A16266G) in some individuals.
16291	While a C-T transition defines branch L0d2b1, it also defines a subbranch of the paraphyletic sister clade of L0d1c2. Furthermore, L0k1 is defined by a G at this position, with a subsequent G to A transition on a subbranch of L0k1a; L0k2 carries an A at this position. While Phylotree (http://www.phylotree.org/tree/subtree_L.htm , Build 15) reconstructs a C-G transversion for L0k as a whole and a G-A transition for L0k2, from the data available to us it appears impossible to decide whether a C-G or C-A transversion took place on the branch leading to L0k. Therefore, in Supplementary Figure 2 the mutations defining L0k1 and L0k2 are both listed as transversions, even though the actual evolutionary path would have involved just one transversion (on the branch leading to L0k) and one transition (on either L0k1 or L0k2).
16294	While a C-T transition defines branch L0d1b2, and a C-A transversion defines branch L0d1c2a, the paraphyletic sister branch of L0d1c2a is defined by a G at this position, with the exception of one individual who carries an A.
A16300G	This mutation, which defines L0d3, mutates back to A in two individuals of branch L0d3b

The mutations are numbered in accordance with the RSRs/rCRS sequence.

Supplemental References:

Barbieri, C., Butthof, A., Bostoen, K., and Pakendorf, B. (2012). Genetic perspectives on the origin of clicks in Bantu languages from southwestern Zambia. *Eur J Hum Genet*. doi: 10.1038/ejhg.2012.192. Aug 29. [Epub ahead of print]

Castrì, L., Tofanelli, S., Garagnani, P., Bini, C., Fosella, X., Pelotti, S., Paoli, G., Pettener, D., and Luiselli, D. (2009). mtDNA variability in two Bantu-speaking populations (Shona and Hutu) from Eastern Africa: implications for peopling and migration patterns in sub-Saharan Africa. *Am J Phys Anthropol* 140, 302-311.

de Filippo, C., Heyn, P., Barham, L., Stoneking, M., and Pakendorf, B. (2010). Genetic perspectives on forager-farmer interaction in the Luangwa valley of Zambia. *Am J Phys Anthropol* 141, 382-394.

Quintana-Murci, L., Harmant, C., Quach, H., Balanovsky, O., Zaporozhchenko, V., Bormans, C., van Helden, P.D., Hoal, E.G., and Behar, D.M. (2010). Strong Maternal Khoisan Contribution to the South African Coloured Population: A Case of Gender-Biased Admixture. *Am J Hum Genet* 86, 611-620.

Schlebusch, C.M., de Jongh, M., and Soodyall, H. (2011). Different contributions of ancient mitochondrial and Y-chromosomal lineages in 'Karretjie people' of the Great Karoo in South Africa. *J Hum Genet* 56, 623-630.

Schlebusch, C.M. (2010). Genetic variation in Khoisan-speaking populations from southern Africa. PhD thesis, University of the Witwatersrand, Johannesburg.

Tishkoff, S.A., Gonder, M.K., Henn, B.M., Mortensen, H., Knight, A., Gignoux, C., Fernandopulle, N., Lema, G., Nyambo, T.B., Ramakrishnan, U., et al. (2007). History of click-speaking Populations of Africa inferred from mtDNA and Y chromosome genetic variation. *Mol Biol Evol* 24, 2180-2195