# ASJP  World Language Tree of Lexical Similarity: Version 3 (July 2010)

## by

## André Müller, Søren Wichmann, Viveka Velupillai, Cecil H. Brown, Pamela Brown, Sebastian Sauppe, Eric W. Holman, Dik Bakker, Johann-Mattis List, Dmitri Egorov, Oleg Belyaev, Robert Mailhammer, Matthias Urban, Helen Geyer, and Anthony Grant

The World Language Tree graphically illustrates relative degrees of lexical similarity holding among 4350 of the world's languages and dialects (henceforth, languages) currently found in the ASJP database (ASJP stands for Automated Similarity Judgment Program). Languages branched more closely together on the ASJP tree are lexically more similar than those branched less closely together. While most lexical resemblance charted in the tree almost certainly is related to genetic affiliation, closely branched languages cannot routinely be assumed to be closely genetically associated since lexical resemblance among languages can be due to factors other than genetic relatedness (see below).

The tree is generated through use of the neighbour-joining computer algorithm originally designed to depict phylogenetic relationships in biology (Saitou & Nei 1987). This is implemented in MEGA 4 (Kumar et al. 2008),[1] the software that we use. The algorithm is applied to a matrix of lexical similarity scores based on Levenshtein (or edit) distances holding between all possible pairs of the 4350 languages (for details about this, including how we modify the Levenshtein distances for our purposes,  see Bakker et al. 2009: 169). All languages of the database are compared to one another with respect to lexical similarity relating to their words for 40 referents determined statistically in Holman et al. (2008) to be most stable among core vocabulary items commonly used in lexicostatistical analysis. The tree is unrooted, but organized around a midpoint, i.e., the point which is equidistant between the two most lexically dissimilar languages in the network. Finally, the tree is annotated to show how it corresponds to the classification used in the latest version of the online World Atlas of Language Structures (Haspelmath et al. 2008),[2] with some updates from Dryer (personal communication). This annotation is presented for ease of orientation, not necessarily because ASJP agrees with it. The language names used are normally simply those

---

[1] http://www.megasoftware.net/
[2] http://wals.info/

of the sources consulted. The sources, as well as corresponding language names of *Ethnologue*, are provided in a continuously updated wiki.[3]

Four factors influence lexical similarity registered in the tree: (1) genetic or genealogical relationship of languages, (2) diffusion (language borrowing), (3) universal tendencies for lexical similarity such as onomatopoeia, and (4) random variation (chance).

Languages branched closely together on the tree may be so because of strong lexical similarity produced by any one or a combination of the four factors. Genetic relationship would appear to be the most dominant factor accounting for close branching, followed next by diffusion. Universal tendencies and chance are less significant contributors to close branching than either genetic relationship or diffusion, but nonetheless clearly contribute to the overall structure of the tree.

Typically, all languages of non-controversial language families such as Austro-Asiatic, Uralic, or Mayan, are respectively branched together on the tree. When some languages of a non-controversial family are not found branched together, this is because they are substantially lexically different from other members of their family despite unambiguously belonging to that family. Occasionally, a language can be so lexically different from co-members of its family that it is found branched more closely with some language or languages with which it is not genetically related at all, usually because of chance lexical similarity or similarity due to borrowing. (When such languages are geographically remote from one another, chance usually explains close branching.)

Typically, branching accords closely with genetic subgroups recognized by experts within non-controversial language families. When branching is not isomorphic with genealogical subgrouping, this often reflects diffusion among languages of the family promoted by language contact. Thus, when used in conjunction with expert classifications of non-controversial language families, the tree can be helpful in calling attention to historical relationships (contact) among genetically related languages that sometimes might not be otherwise apparent.

The tree may also suggest relationships heretofore not noticed among languages that may be profitably investigated. For example, if two languages not known to be related in any way are found together on a terminal branch, this may indicate a relationship between them entailing either inheritance or contact, especially if they are not geographically remote from one another. If the two languages are geographically distant, their close lexical similarity is

---

[3] `http://lingweb.eva.mpg.de/asjp/index.php/ASJP`

more likely explained by chance than by either inheritance or diffusion. Also, language isolates may join one another on a terminal branch because they have nowhere else to go in the tree, creating the illusion that exciting, new far-flung relations may be in evidence. One should be cautious in the interpretation of these cases.

Earlier versions of the ASJP World Language Tree did not include languages regarded as creoles and pidgins, while this version does. Excluding creoles and pidgins would allow human judgments to intrude into the classification (or, in this case non-classification). It is of general interest to show how such languages pattern in the classification when preconceived notions of how they should be treated are avoided.

**References**

Bakker, Dik, André Müller, Viveka Velupillai, Søren Wichmann, Cecil H. Brown, Pamela Brown, Dmitry Egorov, Robert Mailhammer, Anthony Grant, and Eric W. Holman. 2009. Adding typology to lexicostatistics: a combined approach to language classification. *Linguistic Typology* 13: 167-179.

Holman, Eric W., Søren Wichmann, Cecil H. Brown, Viveka Velupillai, André Müller, and Dik Bakker. 2008. Explorations in automated lexicostatistics. *Folia Linguistica* 42.2: 331-354.

Haspelmath, Martin, Matthew S. Dryer, David Gil and Bernard Comrie. 2008. *The World Atlas of Language Structures Online*. Munich: Max Planck Digital Library.

Kumar S., J. Dudley, M. Nei, and K. Tamura K. 2008. MEGA: A biologist-centric software for evolutionary analysis of DNA and protein sequences. *Briefings in Bioinformatics* 9: 299-306.

Saitou, Naruya and Masatoshi Nei. 1987. The neighbor-joining method: a new method for reconstructing phylogenetic trees. *Molecular Biology and Evolution* 4: 406-425.

## Language family abbreviations

| | | | | |
|---|---|---|---|---|
| AA | Afro-Asiatic | | Chn | Chon |
| Aik | Aikana | | Cho | Choco |
| Ain | Ainu | | Chq | Chiquito |
| Ala | Alacalufan | | Cht | Chitimacha |
| Alg | Algic | | Chu | Chumash |
| Alt | Altaic | | CK | Chukotko-Kamchatkan |
| AM | Amto-Musan | | Cmu | Chimúan |
| An | Austronesian | | CN | Cacua-Nukak |
| AP | Awin-Pare | | Cnd | Candoshi |
| Ara | Arafundi | | Cof | Cofán |
| Arc | Araucanian | | Com | Comecrudan |
| Art | Arutani | | Cre | Creoles & Pidgins |
| Aru | Arauan | | Cui | Cuitlatec |
| Arw | Arawakan | | CW | Chapacura-Wanhan |
| Ata | Atakapa | | Dos | Doso |
| AuA | Austro-Asiatic | | Dra | Dravidian |
| Aus | Australian | | EA | Eskimo-Aleut |
| Aym | Aymaran | | EB | East Bougainville |
| Ban | Bangi Me | | EBH | East Bird's Head |
| Bar | Barbacoan | | EGB | East Geelvink Bay |
| Bas | Basque | | Ele | Eleman |
| Beo | Beothuk | | ES | East Strickland |
| Bil | Bilua | | GA | Great Andamanese |
| Bor | Border | | Gcu | Guaicuruan |
| Bos | Bosavi | | GS | Gogodala-Suki |
| Brs | Burushaski | | Gua | Guahiban |
| Bul | Bulaka River | | Had | Hadza |
| Bur | Burmeso | | Hai | Haida |
| Cad | Caddoan | | Har | Harakmbet |
| Cah | Cahuapanan | | HM | Hmong-Mien |
| Cam | Camsá | | Hok | Hokan |
| Car | Cariban | | Hua | Huavean |
| Cay | Cayuvava | | Hui | Huitotoan |
| Chi | Chibchan | | IE | Indo-European |
| Chk | Chimakuan | | IG | Inland Gulf |
| Chl | Cholón | | Ira | Irantxe |
| Chm | Chimila | | Iro | Iroquoian |

| | | | | |
|---|---|---|---|---|
| Ito | Itonama | | May | Mayan |
| Jab | Jabuti | | MGe | Macro-Ge |
| Jap | Japanese | | Mis | Misumalpan |
| Jiv | Jivaroan | | Mol | Molof |
| Kad | Kadugli | | Mom | Mombum |
| Kam | Kamula | | Mon | Monumbo |
| Kap | Kapixana | | Mos | Mosetenan |
| Kar | Karok | | Mov | Movima |
| Kat | Katukinan | | Mrw | Morwap |
| Kau | Kaure | | MUM | Morehead and Upper Maro Rivers |
| Kay | Kayagar | | Mur | Mura |
| Ken | Kenaboi | | Mus | Muskogean |
| KF | Kwomtari-Fas | | MZ | Mixe-Zoque |
| Kho | Khoisan | | Nah | Nahali |
| Kiw | Kiwaian | | Nam | Nambikuaran |
| Kol | Kolopom | | Nat | Natchez |
| Kor | Korean | | NC | Niger-Congo |
| Krt | Kartvelian | | NDa | Nakh-Daghestanian |
| KT | Kiowa Tanoan | | NDe | Na-Dene |
| Ktn | Kutenai | | Nim | Nimboran |
| Kuj | Kujarge | | Niv | Nivkh |
| Kun | Kunza | | NS | Nilo-Saharan |
| Kus | Kusunda | | NWC | Northwest Caucasian |
| Kut | Kuto | | OC | Oregon Coast |
| Kwa | Kwalean | | Odi | Odiai |
| Kwe | Kwerba | | Oks | Oksapmin |
| Kwz | Kwaza | | OM | Oto-Manguean |
| Lav | Lavukaleve | | Pae | Paezan |
| Lec | Leco | | Pan | Panoan |
| LeM | Left May | | Pat | Pataxo |
| Len | Lencan | | Pau | Pauwasi |
| LMa | Lower Mamberamo | | Pen | Penutian |
| LP | Lakes Plain | | Pui | Puinave |
| LS | Leonhard Schultze | | PY | Peba-Yaguan |
| LSR | Lower Sepik-Ramu | | Que | Quechuan |
| Mai | Mairasi | | Sal | Salishan |
| Mar | Marind | | SAn | South Andamanese |
| Mas | Mascoian | | Sav | Savosavo |
| Mat | Matacoan | | Sen | Senagi |

| | | | | |
|---|---|---|---|---|
| Sep | Sepik | War | Warao |
| Sho | Shom Peng | Was | Wasi |
| Sio | Siouan | WBg | West Bougainville |
| Sko | Sko | WBm | West Bomberai |
| Sln | Salinan | WF | Western Fly |
| Slv | Sáliban | WP | West Papuan |
| Snt | Sentani | Wsh | Washo |
| ST | Sino-Tibetan | WY | Wappo-Yukian |
| Tac | Tacanan | Xin | Xincan |
| Tak | Takelma | Yal | Yale |
| Tar | Tarascan | Yam | Yamana |
| Tau | Taushiro | Yan | Yanomam |
| Teb | Teberan-Pawaian | Yaw | Yawa |
| Teq | Tequistlatecan | Yel | Yele |
| Tic | Ticuna | Yen | Yeniseian |
| Tim | Timucua | Yka | Yukaghir |
| TK | Tai-Kadai | Yrb | Yareban |
| TNG | Trans-New Guinea | Yrr | Yaruro |
| TO | Tor-Orya | Yua | Yuat |
| Tof | Tofanma | Yuc | Yuchi |
| Tol | Tol | Yur | Yuracare |
| Ton | Tonkawa | Yuw | Yuwana |
| Tor | Torricelli | Zam | Zamucoan |
| Tot | Totonacan | Zap | Zaparoan |
| Tou | Touo | Zun | Zuni |
| Tru | Trumai | | |
| Tuc | Tucanoan | | |
| TuK | Turama-Kikorian | | |
| Tup | Tupian | | |
| UA | Uto-Aztecan | | |
| UC | Uru-Chipaya | | |
| Un | Unknown | | |
| Ura | Uralic | | |
| Urr | Urarina | | |
| Usk | Usku | | |
| UY | Upper Yuat | | |
| VJ | Vaupés-Japurá | | |
| Wak | Wakashan | | |
| Wao | Waorani | | |

NORTH EFATE NGUNA
NORTH EFATE PWELE
NORTH EFATE SIVIRI
NORTH EFATE SESAKE
NORTH EFATE WORAVIU
LELEPA
ETON
SOUTH EFATE ERATAP
SOUTH EFATE PANGO
NAMAKURA BONGABONGA
NAMAKURA MATASO
NAMAKURA MAKURA
NAMAKURA TONGARIKI
BIERIA VOVO
MAII
APMA
SURU KAVIAN APMA
SA
SEKE
SKE
SOWA
PAAMA FAULILI
PAAMA LAUL
PAAMA LIRONESA
SOUTHEAST AMBRYM MAAT
SOUTHEAST AMBRYM TOAK
NORTH AMBRYM FONAH
NORTH AMBRYM RANON
LONWOLWOL
ORKON
PORT VATO
DAKAKA BAIAP
DAKAKA SESIVI
BAKI
BIEREBO YEVALI
LEWO MAPREMO
BIEREBO TAVIO
BIEREBO BONKOVIA
BIEREBO BURUPIKA
LAMENU
LEWO MATE
LEWO NUL
LEWO FILAKARA
LEWO NIKAURA
LEWO NUVI
LEWO VISINA
LABO WINDUA
LABO
BIG NAMBAS LEVIAMP
BIG NAMBAS UNMET
MARAGUS
TAPE
AXAMB AVOK
AXAMB MAXBAXO
AXAMB
PORT SANDWICH
MASKELYNES
REREP
UNUA
BURMBAR
BURMBAR LEPAXSIVIR
BURMBAR VARTAVO
AULUA
DIXON REEF 1
DIXON REEF 2
LETEMBOI
REPANBITIP
KATBOL TIMBEMBE
KATBOL
SOUTH WEST BAY BENOUR
SOUTH WEST BAY LEMBINWEN
NATI
MALFAXAL
LITZLITZ
VINMAVIS
LAREVAT
PINALUM
WALA

PINALUM
WALA
RANO
TAUTU
URIPIV
URI
ATCHIN
MAE ORAP
MAE
LINGARAK
MPOTOVORO VOVO
VAO
NESE
MALUA BAY PETARMUR
MALUA BAY
MERLAV MERIG
MERLAV
NUME
WETAMUT DORIG
WETAMUT
KORO
LAKONA
LEHALI
MOTLAV
LEHALURUP
MOSINA VETUMBOSO
MOSINA
VATRATA SASAR
VATRATA
EAST AMBAE LOLOMATUI
EAST AMBAE WAILENGI
RAGA
WEST AMBAE
EAST AMBAE LOLSIWOI
MOTA
CENTRAL MAEWO
MARINO
BAETORA
BAETORA NAROVOROVO
BAETORA NASAWA
BAETORA NAVENEVENE
BAETORA TAM
HIW
TOGA
BUTMAS
TUR
POLONOMBAUK
LOREDIAKARKAR
SHARK BAY 1
SHARK BAY 2
RORIA
PIAMATSINA
TASMATE
VUNAPU
VALPEI HUKUA
VALPEI
NOKUKU
TOLOMAKO
TAMBOTALO
ARAKI
TANGOA
NORTH MALO
SOUTH MALO
MAFEA
AORE
TUTUBA
NARANGO NAMBEL
NARANGO
AMBLONG
MOROUAS BATUNLAMAK
MOROUAS
WUSI MANA
WUSI VALUI
WUSI KEREPUA
AKEI PENANTSIRO
WAILAPA
FORTSENAL
AKEI TASIRIKI
AKEI
LAMETIN

An.OCEANIC

AKEI
LAMETIN
MALMARIV
NAVUT
NAVUT MATAE
WUSI NONONA
PALAUAN ] An.PALAUAN
SOUTH EFATE ERAKOR
SOBEI
TARPIA
An.OCEANIC

KIRIBATI
KUSAIE
PINGELAPESE
POHNPEIAN
MOKILESE
EASTERN MARSHALLESE
MARSHALLESE
PULO ANNA
SONSOROLESE
WOLEAIAN
ULITHIAN
CAROLINIAN
PULUWATESE
TRUKESE
An.OCEANIC

MALEU
MENGEN
MUSSAU EMIRA
LOTE
KAULONG
SENGSENG
LAMOGAI
AMARA
MOUK 1
MOUK 2
An.OCEANIC

SAKETA
WOSI
GIMAN
MAILOA
EAST MAKIAN
NGOFAGITA
NGOFAKIAHA
SOMA
PELERI SAMSUMA
TAHANE
BULI
SAWAI UnnamedInSource
AS
GEBE
MATBAT
AMBEL
BUTLEH
BIGA
MISOOL MAYA
SALAWATI MAAYA
LANGANYAN
KAWE
WAUYAI
An.SOUTH HALMAHERA-WEST NEW GUINEA

ARGUNI ] An.CENTRAL MALAYO-POLYNESIAN
KAMBRAM
YARIK
] An.SOUTH HALMAHERA-WEST NEW GUINEA

SUNGLIK LIHIR
PATPATAR
SIAR
KANDAS
MINIGIR
TOLAI
BILUR
SURSURUNGA
TANGA
NALIK
WEST KARA
TIANG
TIGAK
TUNGAK
An.OCEANIC

IRARUTU ] An.SOUTH HALMAHERA-WEST NEW GUINEA
YAKAMUL ] An.OCEANIC
DOBEL
UJIR
NGAIBOR
TARANGAN BARAT FERUNI
An.CENTRAL MALAYO-POLYNESIAN

NGAIBOR
TARANGAN BARAT FERUNI
KAZUKURU An.OCEANIC
EAST MASELA
SOUTHEAST BABAR
SERILI
IMROING
TELA MASBUAR
CENTRAL MASELA
EMPLAWAS
NORTH BABAR
DAI
DAWERA DAWELOOR
WEST DAMAR

An.CENTRAL MALAYO-POLYNESIAN

BURU MASARETE
HUKUMINA
KAYELI

An.CENTRAL MALAYO-POLYNESIAN

FATAKAI NUAULU SERAM An.CENTRAL MALAYO-POLYNESIAN

ONIN 2
SEKAR 2
ONIN
SEKAR
YAMDENA
FORDATA
KEI
MASIWANG
BOBOT
BONFIA
BONFIAY
KOWIAI

An.CENTRAL MALAYO-POLYNESIAN

ELPAPUTIH SERAM 2
ELPAPUTIH SERAM 3
PAULOHI
ELPAPUTIH SAMASURU PAULOHU
ELPAPUTIH SERAM 1
AMAHAI
HARUKU
PELAUW HARUKU
HITU
TULEHU
SAPARUA OUW
NUSA LAUT
SAPARUA HARIA
SAPARUA IHAMAHU
TALUTI LAIMU
TALUTI TAMILOUW
SEPA
TALUTI
ALUNE
SAPOLEWA SOOW KWELE ULUI SERAM
ASILULU LIMA RUMAHSOSAL NUWETETU
ASILULU LIMA WARAKA

An.CENTRAL MALAYO-POLYNESIAN

ANAKALANG
WANUKAKA
BALILEDO
PONDOK
MAMBORU
LAMBOYA
WEJEWA
GAURA NGGAURA
KAMBERA
SOUTHERN KAMBERA
LEWA KAMBERA
UMBU RATU NGGAI KAMBERA
SIKA
MANGGARAI
PALUE
ENDE
LIO
NGADHA
SOA
ELAT KEI BESAR
WATUBELA
GESER
WARU SERAM

An.CENTRAL MALAYO-POLYNESIAN

ROTE UnnamedInSource
TERMANU
UAB METO

An.CENTRAL MALAYO-POLYNESIAN

KEMAK
TUKUDEDE An.CENTRAL MALAYO-POLYNESIAN

KEMAK
TUKUDEDE — An.CENTRAL MALAYO-POLYNESIAN
MAMBAE
SELARU ] An.CENTRAL MALAYO-POLYNESIAN
APUTAI
PERAI
TUGUN — An.CENTRAL MALAYO-POLYNESIAN
ERAI
ILIUN
TALUR ] An.CENTRAL MALAYO-POLYNESIAN
TETUN ] An.CENTRAL MALAYO-POLYNESIAN
TETUN DILI ] Cre.TETUN BASED
NILA
TEUN
SERUA
EAST DAMAR
KISAR
ROMA — An.CENTRAL MALAYO-POLYNESIAN
LETINESE
LUANG
LUNGGA
SIMBO
DUKE
ROVIANA
UGHELE
MAROVO
MBAREKE
HOAVA
KUSAGHE
VAGHUA
VARISI
RIRIO
SISIQA
AVASO BABATANA — An.OCEANIC
LOMAUMBI BABATANA
KATAZI BABATANA
TUNOE BABATANA
CHEKE HOLO
GHOVE BLABLANGA
ZABANA
KOKOTA
BLABLANGA
ZAZAO
ALU MONO
FAURO MONO
MONO
TORAU
TEOP
BANONI — An.OCEANIC
TAIOF
HAKU
NEHAN
SOLOS
DAMI
MAISIN ] An.OCEANIC
SALIBA PAPUA NEWGUINEA
SUAU
BUHUTU
SEWA BAY
DOBU
MOLIMA
UBIR
GAPAPAIWA
TAWALA
WEDAU
BWAIDOKA — An.OCEANIC
DIODIO
IAMALELE
MINAVEHA
GUMAWANA
KILIVILA
MUYUW
BILBIL
GEDAGED
MATUKAR
TAKIA
AROP LOKEP
MANAM
MBULA PAPUA NEW GUINEA
KIS
KAIRIRU

KIS
KAIRIRU
KAIEP
WOGEO
YABEM
NUMBAMI
BARIAI
KOVE
SEIMAT
WUVULU
LEVEI
LIKUM
NYINDROU
LOU
NAUNA
SIVISA TITAN
LEIPON
LONIU

An.OCEANIC

KUNI
LALA
DOURA
GABADI

An.OCEANIC

MOTU ] An.OCEANIC
MOTU HIRI ] Cre.MOTU BASED
MEKEO ] An.OCEANIC
RORO
SINAUGORO ] An.OCEANIC

PENRHYN
TIKOPIA
RAROTONGAN
MAORI
MANIHIKI
EMAE
RENNELLESE
TONGAN
NIUE
NIUAFOOU
WALLISIAN
PUKAPUKA
VAEAKAU TAUMAKO
EAST FUTUNA
NANUMEA
TOKELAU
ANIWA
FUTUNA
FILA
MELE
RAPA NUI
NORTH MARQUESAN
KAPINGAMARANGI
NUKUORO
LUANGIUA
HAWAIIAN
HAWAIIAN 2
RURUTUAN
TAHITIAN
SAMOAN
TUAMOTUAN

An.OCEANIC

BILEKI NAKANAI
NAKANAI
MAUTUTU

An.OCEANIC

FIJIAN(WAYAN) UnnamedInSource
WAYAN FIJIAN
FIJIAN
BOLA
VITU

An.OCEANIC

ROTUMAN ] An.OCEANIC
URUAVA ] An.OCEANIC
NIAS NORTHERN
NIAS SOUTHERN ] An.SUMATRA

NGGERI GHARI
TANDAI GHARI
NDI GHARI
NGGAE GHARI
GHARI
NGINIA GHARI
POLEO TALISE
KOO TALISE
MALANGO
MBIRAO

MALANGO
MBIRAO
MALAGHETI TALISE
MOLI TALISE
TOLO
NGGELA
LENGO
PARIPAO LENGO
BUGHOTU
MANDARA
TABAR
BAROO BAURO
SANTA ANA OWA
BAURO
HAUNUNU BAURO
FAGANI
KAHUA
TAWAROGA KAHUA
UKI NI MASI SAA
ULAWA SAA
AULU SAA
SAA
AROSI TAWATANA
ONEIBIA AROSI
LONGGU
WAIAHAA AREARE
MAASUPA AREARE
MARAU
OROHA
DORIO
LANGALANGA
KWAIO
KWARAAE
MBAELELEA
MBAENGGUU
FATALEKA
TOABAITA
LAU
KWAI
WALADE LAU

An.OCEANIC

BUANG
BUANG MAPOS
SUDEST
TEANU

An.OCEANIC

MATO
MISIMA PANAEATI
NIMOA

An.OCEANIC

YAPESE ]An.YAPESE
NGULUWAN ]Cre.YAPESE BASED
MOR ]An.SOUTH HALMAHERA-WEST NEW GUINEA
WAREMBORI ]LMa.LOWER MAMBERANO
WAROPEN
BIAK
NUMFOR
DUSNER
RON
KURUDU
WABO
AMBAI
WADAPI LAUT
SERUI-LAUT
WANDAMEN
MUNGGUI
POM
PAPUMA
ANSUS
WOI

An.SOUTH HALMAHERA-WEST NEW GUINEA

ADZERA
WAMPAR
ARIBWATSA
MUSOM

An.OCEANIC

BAROK
LAMASONG
MADAK
SAKAO
SIE
URA
ASUMBOA
WHITESANDS IARKEI
WHITESANDS LONIEL

WHITESANDS IARKEI
WHITESANDS LONIEL
NORTH TANNA
LENAKEL LENAUKAS
LENAKEL LONASILIAN
KWAMERA PORT RESOLUTION
KWAMERA YATUKWEY
KWAMERA ISIAI
SOUTHWEST TANNA IKITI
SOUTHWEST TANNA ENFITANA
SOUTHWEST TANNA IMREANG
SOUTHWEST TANNA IKIYAU
SOUTHWEST TANNA LAPWANGTOAI

An.OCEANIC

KANAKANABU
SAAROA
TSOU

An.TSOUIC

ATAYAL
SEDIQ

An.ATAYALIC

CHAM EASTERN
CHAM WESTERN
CHRU
NORTHERN ROGLAI
RADE
TSAT

An.MALAYIC

MODANG
SEGAI

An.KAYAN-MURIK

CHAMORRO

An.CHAMORRO

FAVORLANG
PAPORA
PAPORA 2

An.PAIWANIC

MAMUJU
MANDAR
SADAN
TAE
BUGINESE
SOPPENG BUGINESE
EMBALOH
MAKASAR
MAKASSARESE
SELAYAR
COASTAL KONJO
KONJO
BANTIK
SANGIL
SANGIL SARANGANI ISLANDS
SANGIR
SANGIR 2
TABUKANG SANGIR

An.SULAWESI

SASAK

An.BALI-SASAK

TOBA BATAK

An.SUMATRA

TONSAWANG
TOMBULU
TONSEA
TONDANO
TONTEMBOAN

An.SULAWESI

GAYO

An.GAYO

BONGGI
CENTRAL DUSUN
TAMBUNAN DUSUN
IDAAN
TIMUGON MURUT

An.NORTHWEST MALAYO-POLYNESIAN

SAMAL
SIASI SAMA
BALANGINGI SAMA
MAPUN
YAKAN
INABAKNON
LEMO BAJAU
LUWUK BAJAU
INDONESIAN BAJAU
KAYUADI BAJAU
KOLO BAWAH BAJAU
KALEROANG BAJAU
LANGARA LAUT BAJAU
LAKONEA BAJAU
PADEI LAUT BAJAU
ANAIWOI BAJAU
PITULUA BAJAU
BAJOE BAJAU
LAPULU BAJAU

An.SAMA-BAJAW

BAJOE BAJAU
LAPULU BAJAU
LAKARAMBA BAJAU
BOEPINANG BAJAU
LAURU BAJAU
MORAMO BAJAU
DONDO
TOMINI
LAUJE
LAUJE AMPIBABO
TAJIO
PENDAU
TAJE PETAPA
TAJE TANAMPEDAGI
BALAESANG
DAMPELAS
BOANO
TOTOLI
DAA
BAREE
UMA
MATO NO UWE
WABHULA
DESA WALI
BATU ATAS
MASIRI
KUMBEWAHA
PASARWAJO
BONERATE
TUKANG BESI NORTHERN
TUKANG BESI SOUTHERN
KAMARU
LASALIMU
BANGGAI
WOLIO
TOLAKI LAIWUI
TOLAKI WIWIRANO
TOLAKI ASERA
TOLAKI KONAWE
TOLAKI MEKONGGA
RAHAMBUU
WARU
WARU LALOMERUI
KODEOHA
MORI ATAS
PADOE
TOMADINO
MORI BAWAH
MORI BAWAH WATU
TOLAKI
MORONENE
MORONENE TOKOTUA
WAWONII
WAWONII MENUI
KULISUSU
KORONI
TALOKI

An.SULAWESI

GORONTALO
KAIDIPANG
BUOL
MONGONDOW

An.SULAWESI

SOBOYO    An.CENTRAL MALAYO-POLYNESIAN

BUSOA
KAIMBULAWA
KADATUA
MUNA
LIABUKA
WASUAMBA
KAMBOA
LAWELE
KAPONTORI
TODANGA

An.SULAWESI

DHAO
SAVU
BIMA
KEDANG
LAMALERA
LAMAHOLOT ILE MANDIRI
ALOR
ALOR/BARAHUSA/KABIR

An.CENTRAL MALAYO-POLYNESIAN

ALOR
ALOR/BARAHUSA/KABIR
ALOR/KALABAHI
MALAGASY MAHAFALY
MALAGASY VEZO
MALAGASY TANDROY 1
MALAGASY TANDROY 2
MALAGASY SAKALAVA 1
MALAGASY SAKALAVA 2
MALAGASY ANTANKARANA
MALAGASY TSIMIHETY
MALAGASY ANTAISAKA
MALAGASY ZAFISORO
MALAGASY BARA
MALAGASY TAIMORO
MALAGASY ANTAMBAHOAKA
MALAGASY BETSIMISARAKA
MALAGASY FIANARANTSOA
MALAGASY SIHANAKA
MALAGASY AMBOSITRA
MALAGASY MERINA
An.BARITO
INDONESIAN
MALAY
SEKOLA LONTHOIR BANDA
SEKOLA NEIRA BANDA
INDONESIAN 2
INDONESIAN JAKARTA
OGAN
BESEMAH
PALEMBANG MALAY
AI BANDA
AMBON MALEIS
KUPANG MALAY
AMBONESE MALAY
MANADONESE
TERNATE PASAR
IBAN
MINANGKABAU
BANJARESE MALAY
DELANG
TAMUAN
LOM
KERINCI
KERINCI 2
An.MALAYIC
MADURESE  ] An.MADURESE
ACEH  ] An.MALAYIC
BALI  ] An.BALI-SASAK
SUNDANESE  ] An.SUNDANESE
TENGGER NGADAS
MALANG
YOGYAKARTA
An.JAVANESE
LAMPUNG
LAMPUNG NYO ABUNG/KOTABUMI
LAMPUNG NYO MELINTING
ABUNG SUKADANA LAMPUNG NYO
MENGGALA TULANG BAWANG LAMPUNG NYO
KOMERING
KAYU AGUNG ASLI KOMERING
KAYU AGUNG PENDATANG KOMERING
RANAU LAMPUNG API
WAY LIMA LAMPUNG API
JABUNG LAMPUNG API
PUBIAN LAMPUNG API
SUNGKAI LAMPUNG API
BELALAU LAMPUNG API
TALANG PADANG LAMPUNG API
KOTA AGUNG LAMPUNG API
KALIANDA LAMPUNG API
KRUI LAMPUNG API
SUKAU LAMPUNG API
DAYA LAMPUNG API
PERJAYA ULU KOMERING
WAY KANAN LAMPUNG API
ADUMANIS ULU KOMERING
ILIR KOMERING
An.LAMPUNGIC
KUALAN
SINGHI
BEKATI
An.LAND DAYAK
MOKEN  ] An.MOKLEN
REJANG  ] An.REJANG

REJANG ] An.REJANG
BELAIT
LONG TERAWAN BERAWAN } An.NORTHWEST MALAYO-POLYNESIAN
BINTULU
MUKAH MELANAU } An.NORTHWEST MALAYO-POLYNESIAN
LAHANAN
BUKAT ] An.KAYAN-MURIK
LONG ANAP KENYAH ] An.NORTHWEST MALAYO-POLYNESIAN
REJANG KAYAN ] An.KAYAN-MURIK
BARIO KELABIT ] An.NORTHWEST MALAYO-POLYNESIAN
TUNJUNG
DUSUN WITU
SAMIHIM
MAANYAN
PAKU
DUSUN DEYAH
DUSUN MALANG
LAWANGAN
TAWOYAN
MURUNG SIANG
SIANG
KADORIH
NGAJU BAAMANG
NGAJU OLOH MANGTANGAI
NGAJU PULOPETAK
KAPUAS KAHAYAN
KATINGAN
An.BARITO
TIRURAY
TAGABILI 2
TBOLI
TAGABILI
BILAAN KORONADAL
BILAAN SARANGANI
KORONADAL BILAAN
SARANGANI BILAAN
An.SOUTH MINDANAO
AMIS
C AMIS
PAIWAN
PUYUMA
An.PAIWANIC
BUNUN ] An.ATAYALIC
KAVALAN ] An.PAIWANIC
SIRAYA
SAISIYAT
PAZEH
THAO
An.PAIWANIC
UMIRAY DUMAGET AGTA ] An.NORTHERN PHILIPPINES
KANKANAY NORTHERN
NORTHERN KANKANAY
BONTOC GUINAANG
CENTRAL BONTOC
BALANGAO
BALANGAW
ITNEG BINONGAN
ITNEG BINONGAN 2
KALINGA LIMOS
KALINGA GUINAANG
LUBUAGAN KALINGA
BATAD IFUGAO
IFUGAO BATAD 2
AMGANAD IFUGAO
IFUGAO AMGANAD 2
BAYNINAN IFUGAO
IFUGAO BAYNINAN 2
INIBALOI
INIBALOI 2
KALLAHAN KAYAPA PROPER
KAYAPA KALLAHAN
KALLAHAN KELEYQIQ
KALLAHAN KELEYQIQ IFUGAO
KAPAMPANGAN
ILOKANO
PANGASINAN
SAMBAL BOTOLAN
SAMBAL BOTOLAN 2
CASIGURAN DUMAGAT AGTA
DUMAGAT CASIGURAN
CASIGURAN NEGRITO
DUPANINGAN AGTA
AGTA
CENTRAL AGTA
An.NORTHERN PHILIPPINES

AGTA
CENTRAL AGTA
YOGAD
GADDANG
GADDANG 2
ISNAG
ISNEG
IBANAG
ATTA PAMPLONA 2
PAMPLONA ATTA
ILONGOT KAKIDUGEN
ILONGOT
KAKIDUGEEN ILONGOT
ISAMORONG IVATAN
IVASAY IVATAN
IVATAN
ITBAYAT IVATAN
IMOROD
IRARALAY IVATAN
IVATAN BATANES ISLANDS 2
IVATAN BASCO
ITBAYATEN BATANES ISLANDS
ITBAYATEN IVATAN

An.NORTHERN PHILIPPINES

TAGBANWA KALAMIAN
TAGBANWA KALAMIAN 2          An.MESO-PHILIPPINE

MAMANWA
MAMANWA 2                    An.MESO-PHILIPPINE
CENTRAL TAGBANWA

MANOBO ILIANEN
MANOBO ILIANEN 2
MANOBO MATIGSALUG
MANOBO TIGWA
MANOBO WESTERN BUKIDNON 2
MANOBO WESTERN BUKIDNON
WESTERN BUKIDNON MANOBO
BINUKID
BINUKID 2
MANOBO ATA
MANOBO ATA 2
KAGAYANEN
MANOBO DIBABAWON
MANOBO DIBABAWON 2
MANOBO COTABATO
MANOBO KALAMANSIG COTABATO
MANOBO SARANGANI
MANOBO SARANGANI 2

An.SOUTHERN PHILIPPINES

TAUSUG
TAUSUG 2                     An.MESO-PHILIPPINE

BATAK PALAWAN
PALAWAN BATAK
PALAWANO BATAK
TAGBANWA ABORLAN
TAGBANWA ABORLAN 2           An.MESO-PHILIPPINE
MOLBOG
SOUTHWEST PALAWANO

MARANAO
IRANUN                       An.SOUTHERN PHILIPPINES
MORO MAGINDANAU

HANUNOO                      An.MESO-PHILIPPINE

KALAGAN KAAGAN
KALAGAN TAGAKAOLO
MANDAYAN ISLAM PISO
MANDAYAN CARAGA
MANDAYAN BOSO
MANSAKA 3
MANSAKA
MANSAKA 2
KALAGAN
KALAGAN 2

An.MESO-PHILIPPINE

CENTRAL SUBANEN
SUBANUN SINDANGAN
SUBANON SIOCON
WESTERN SUBANON             An.SOUTHERN PHILIPPINES

KINARAY-A
WARAY WARAY
HILIGAYNON
BUTUANON
SURIGAONON
CAPIZNON
AKLANON

CAPIZNON
AKLANON
CEBUANO
TAGALOG
CENTRAL BICOLANO
DARAGA
OAS
BUHI
IRIGA
LIBON
NORTHERN CATANDUANES
LEGAZPI
NAGA
SOUTHERN CATANDUANES
SOUTHERN SORSOGON
MASBATE
NORTHERN SORSOGON

An.MESO-PHILIPPINE

NEA/NEMBOI
NEA
NEA/NOOLI
BANUA
NAMBAKAENGO/MALO
NANGGU

An.OCEANIC

AYIWO ] An.OCEANIC
POROME ] Kiw.KIWAIAN

DEHU
NENGONE
IAAI
AJIE
GRAND COULI
XARACUU
CEMUHI
NELEMWA
JAWE
NEMI

An.OCEANIC

KAYUPULAU ] An.OCEANIC
URARINA ] Urr.URARINA
PELE ATA WASI ] Was.WASI
BILUA
NDOVELE BILUA ] Bil.BILUA
SAVOSAVO ] Sav.SAVOSAVO
LAVUKALEVE ] Lav.LAVUKALEVE
MBANIATA ] Tou.TOUO
MAIRASI/FARANJAO
MAIRASI
SEMIMI ETNA BAY
SEMIMI

Mai.MAIRASI

ANGKAMUTHI
YADHAYKENU
ATAMPAYA
URADHI ATAMPAYA
URADHI ANGKAMUTHI
URADHI YADHAYKENU
LINNGITHIGH
MPAKWITHI ANGUTHIMRI
YIR YORONT
KUKU UWANH
WIK MUNGKAN

Aus.PAMA-NYUNGAN

KALA LAGGAW YA ] Aus.PAMA-NYUNGAN

TADAKSAHAK
TASAWAQ
DENDI
ZARMA
TONDI SONGWAY KIINI
KOROBORO SENNI
DJENNE CHIINI
KOYRA CHIINI

NS.SONGHAY

MARIDAN
MARITHIEL
MARITYABEN
MARENGAR
MARAMANADJI
MARANUNGGU
AMI
MANDA

Aus.WESTERN DALY

MURRINHI PATHA ] Aus.MURRINH-PATHA
NGANGIKURRUNGGURR
NGENGOMERI ] Aus.SOUTHERN DALY
KAMOR
YUNGGOR ] Aus.EASTERN DALY

KAMOR
YUNGGOR ] Aus.EASTERN DALY
PUNGUPUNGU
WADYGINY ] Aus.ANSON BAY
TYARAITY ] Aus.NORTHERN DALY
MATNGALA ] Aus.EASTERN DALY
MALAKMALAK ] Aus.NORTHERN DALY
GANGGALIDA ] Aus.TANGKIC
TIWI ] Aus.TIWIAN
MANGARAYI ] Aus.MANGARRAYI
UMBUGARLA ] Aus.UMBUGARLA-NGUMBUR
INGURA ] Aus.ANINDILYAKWA
YAGUA ] PY.PEBA-YAGUAN
GARAWA ] Aus.GARRWAN
LIMILNGAN ] Aus.LIMILNGAN
GUNWINGGU MANYALLALUK MAYALI
GUNWINGGU GUN DJEIHMI
GUNWINGGU KUNWINJKU
GUNWINGGU KUNINJKU
GUNWINGGU KUNE
GUNBALANG WARLANG
BUAN ] Aus.GUNWINYGIC
RAINBARNGO ] Aus.REMBARNGA
NGALAKAN ] Aus.NGALAKAN
NGANDI ] Aus.NGANDI
BURARRA ] Aus.BURARRAN
DJEEBBANA ] Aus.NDJEBBANA
GOROGONE ] Aus.BURARRAN
DJAUAN ] Aus.DJAUAN
MANGERR
ERRE        Aus.GIIMBIYU
URNINGANGG
GAAGUDJU ] Aus.GAAGUDJU
YANYUWA ] Aus.PAMA-NYUNGAN
MARGU
AMURDAK
IWAIDJA        Aus.IWAIDJAN
MAWNG
MARA
WANDARANG        Aus.MARAN
ALAWA
NUNGGUBUYU ] Aus.NUNGGUBUYU
DJINANG
DJINBA
DHUWAL
YOLNGUMATHA
DJAPU
DALWONGO
YANANGO
WANGURI
GALBU
RIRAIDJANGO
WARAMERI        Aus.PAMA-NYUNGAN
RIDARNGO
GOBABINGO
MARARBA
GOMAIDJ
MANGGALILI
GUNIN/KWINI
WUNAMBAL
NGARINYIN
WURLA        Aus.WORORAN
UNGGUMI
WORRORA
KITJA ] Aus.DJERAGAN
BUNABA
GOONIYANDI        Aus.BUNUBAN
MIRIWUNG ] Aus.DJERAGAN
DJAMINDJUNG ] Aus.JAMINJUNGAN
WAGIMAN ] Aus.WAGIMAN
WARDAMAN
YANGMAN        Aus.YANGMANIC
DJINGILI
GUDANJI        Aus.WEST BARKLY
WAMBAYA
MANTJILTJARA
PINTUPI
PITJANTJATJARA YANKUNTJATJARA
MARTU WANGKA
YULPARIJA

MARTU WANGKA
YULPARIJA
NGAANYATJARRA
WAJARRI
NGADJUNMAYA
NHANDA
MARTUTHUNIRA
YINDJIBARNDI
NGALOOMA
PANYTYIMA
NYUNGA NORTHERN
NYUNGA EASTERN
NYUNGA SOUTH WESTERN
WARLPIRI
WALMAJARRI
DJARU
GURINDJI
DHURGA
THURAWAL
NGUNAWAL
KAURNA
WIRANGU
ADNYAMATHANHA
ARABANA
PITTA PITTA
DIYARI
YANDRUWANDHA
DARLING
MURUWARI
MALYANGAPA
NGURA
MAYKULAN
NGAWUN
MAYI YAPI
MAYI THAKURTI
MAYAGUDUNA
KALKATUNGU
YALARNNGA
COLAC
WATHAWURRUNG
WOIWURRUNG
MADIMADI
DJADJALA
WEMBAWEMBA
BUNGANDITJ
WARRNAMBOOL
DYAABUGAY
YIDINY
GUUGU YALANDYI
GUUGU YIMIDHIR
GUMBAYNGGIR
YAYGIR
BIDJARA
KUNGGARI
MARGANY
GUNYA
GUWAMU
GANGULU
BIRRI
WIRRI
DYIRBAL
WARUNGU
WULGURU
NYAWAYGI
WARGAMAY
DYANGADI
WORIMI
DARRKINYUNG
AWABAKAL
GIDABAL
WAALUBAL
YUGAMBAL
DUUNGIDJAWU
WULIWULI
GURENG GURENG
BAYALI
WIRADHURI
GAMILARAAY
WANGAAYBUWAN NGIYAMBAA
YUWANA

Aus.PAMA-NYUNGAN

Yuw.YUWANA

YUWANA
YUWANA 2 ] Yuw.YUWANA
SALIBA COLUMBIA ] Slv.SALIBAN
PAPI PAUPE ] LS.LEONHARD SCHULTZE
MOCHICA ] Cmu.CHIMUAN
SANDAWE ] Kho.SANDAWE
MAI BRAT ] WP.NORTH-CENTRAL BIRDS HEAD
ATAKAPA ] Ata.ATAKAPA
ALFENDIO ] Ara.ARAFUNDI
TONKAWA ] Ton.TONKAWA
MEPHAA ACATEPEC
TEOCUITLAPA MEPHAA
TLACOAPA MEPHAA
TLAPANEC MALINALTEPEC
AZOYU TLAPANEC
SUBTIABA
OM.SUBTIABA-TLAPANEC
ABUN JI/MADIK
ABUN/KARON
ABUN
WP.NORTH-CENTRAL BIRDS HEAD
MPUR ] WP.KEBAR
KALABRA
MORAID
TEHIT/TEHIDYIT
TEHIT
SEGET/WALIEM
SEGET
MOI/WAIPU
MOI/STOKHOF FLASSY
MOI
WP.WEST BIRDS HEAD
GEBUSI
HONIBO
OIBAE
KUBO
SAMO
ODOODEE
AGALA
ES.EAST STRICKLAND
LAMA ] NC.GUR
DOSO ] Dos.DOSO
TARASCAN ] Tar.TARASCAN
ESKAYAN ] UNCLASSIFIED.UNCLASSIFIED
PARIMANKUTINMA ] Aus.PAMA-NYUNGAN
COMECRUDO ] Com.COMECRUDAN
TAURAP/BORUMESSU ] Bur.BURMESO
TAUSHIRO ] Tau.TAUSHIRO
LEKO ] Lec.LECO
FULNIO ] MGe.YATE
ABUI TAKALELANG ] TNG.WEST TIMOR-ALOR-PANTAR
KADIWEU ] Gcu.GUAICURUAN
TOBA
MOCOVI
PILAGA
Gcu.GUAICURUAN
LAMALAMA COASTAL
LAMALAMA INLAND
UMBUYKAMU
Aus.PAMA-NYUNGAN
NAURUAN ] An.OCEANIC
YORTA YORTA ] Aus.PAMA-NYUNGAN
KAYTETYE
ALYAWARRA
WESTERN ARRERNTE
ARRERNTE CENTRAL
ARRERNTE EASTERN
Aus.PAMA-NYUNGAN
APINAYE
KAYAPO
APANIEKRA
KRAHO
KREYE
PYKOBJE
PANARA
SUYA
XERENTE
XAVANTE
XAVANTE 2
MGe.GE-KAINGANG
GUAHIBO
PLAYERO
CUIBA
JITNU
GUAYABERO
Gua.GUAHIBAN
TIKUNA ] Tic.TICUNA
ONA
TEHUELCHE
Chn.CHON PROPER

ONA
TEHUELCHE ] Chn.CHON PROPER
PAYAGUA ] Un2.UNKNOWN2
MACA
NIVACLE
CHOROTE
WICHI LHAMTES GUISNAY
WICHI LHAMTES VEJOZ
MAWES/DAI ] Kwe.MAWES
MAWES/WARES
TLINGIT ] NDe.TLINGIT
KUSUNDA ] Kus.KUSUNDA
CADDO ] Cad.CADDOAN
KARITIANA ] Tup.ARIKEM
BUGUN ] ST.MIRISH
SULUNG
CHENG
OI
BRAO
SAPUAN
JRU
LOVEN
NHAHEUN
NHAHEUN 2
JEH
JEH 2
RENGAO
BAHNAR
BAHNAR 2
CUA
TAMPUAN
SRE
CHRAU
STIENG
BUNOR
PREH
KUAN
GAR
ROLOM
SEDANG
TODRAH UnnamedInSource
KATU
KATU EASTERN
NGE
PACOH
BRU
BRU(WESTERN) UnnamedInSource
KUY
KUI THAILAND
KUI THAILAND 2
MON 2
MON NONG DUU
MON
NYAKUR
NYAKUR 2
SEMAQ BERI UnnamedInSource
SEMELAI
JAHAI
KENSIW
CEQ WONG
JAH HUT
TEMIAR
SEMAI
SEMAI 2
KHMER
KHMER 2 ] AuA.KHMER
SURIN
PEAR B
CHONG H
KASONG
PEAR
THAVUNG SO
RUC
RUC 2
MALIENG
MUNG KOI
VIETNAMESE
VIETNAMESE 2
VIETNAMESE 3
CAR

Mat.MATACOAN

AuA.BAHNARIC

AuA.KATUIC

AuA.MONIC

AuA.ASLIAN

AuA.PEARIC

AuA.VIET-MUONG

MAZHELONG BAI
GONGXING BAI
TUOLUO BAI
EGA BAI
ENQI BAI
JINMAN BAI
LUOBENZHUO BAI
AMOY MINNAN CHINESE
HAINAN MINNAN CHINESE ] ST.CHINESE
SUZHOU WU ] ST.CHINESE
HAKKA ] ST.CHINESE
THONG BOI ] Cre.HAKKA BASED
CANTONESE ] ST.CHINESE
WUTUN ] Cre.CHINESE BASED
MANDARIN
MANDARIN 2 ] ST.CHINESE
NISI
NISI TAGIN
BENGNI
NISHING
BOKAR
MISING
APATANI
IDU MISHMI
TARAON
PYU ] KF.PYU
NAGA MZIEME
NAGA ZEME
NAGA LIANGMAI
NAGA MARAM
KARBI
DHAMMAI
MIJI
NAGA POCHURI
NAGA SUMI
NAGA CHOKRI
NAGA RENGMA
NAGA LOTHA
NAGA MAO
BIYUE
LAHU
LISU
LISU 2
DAFANG
XIDE
LALO
NANJIANG
AKHA
BURMESE
NUSU
ZAIWA
ACHANG
XIANDAO
PHUNOI
BISU
MPI
JINO
JINO 2
NAXI ] ST.NAXI
JIARONG ] ST.RGYALRONG
KAREN GEBA
KAREN SGAW
BWE KAREN
KAREN PAO
KAREN YINBAW
KAYAN
KAREN MANUMANAW
KAREN YINTALE
KAYAH LI EASTERN
KAYAH WESTERN
A TONG
WANANG
KACHARI
GARO
TIWA
DIMASA
KOK BOROK 2
RABHA
CHANG
NAGA CHANG

ST.MIRISH

ST.KUKI-CHIN-NAGA

ST.DHAMMAI UNCLASSIFIED

ST.KUKI-CHIN-NAGA

ST.BURMESE-LOLO

ST.KAREN

ST.BARIC

CHANG
NAGA CHANG
NAGA NOCTE
NAGA WANCHO
NAGA KONYAK
NAGA PHOM
DEURI ] ST.BARIC
DHIMAL ] ST.BODIC
JERUNG
WAMBULE
BAHING
SUNWAR
SUNWAR 2
THULUNG
KHALING
DUMI
DUMI 2
LEPCHA ] ST.LEPCHA
KAMAN ] ST.MIRISH
MARPHATAN THAKALI
THAKALI
TAMANG
MANANGE
GHACHOK
GURUNG
CHANTYAL
KINNAURI
BYANSI
DARMIYA
DOLAKHA NEWAR
THANGMI
NAGA AO
NAGA SANGTAM
NAGA TANGKHUL
LIMBU
MEWA KHOLA LIMBU
YAKHA
KULUNG
PUMA
TSHANGLA
NAGA TARAO
MEITEI
CHIN HAKA
LUSHAI
CHIN TEDIM
CHIN THADO
CHEPANG
MAGAR
KHAM TAKALE
KHAM
TAKA KHAM
DRUNG
NUNG CHINA
RAWANG
KAIKE
RGYALTHANG
KYIRONG
LOWA
JIREL
SHERPA
TIBETAN CENTRAL
TIBETAN LHASA
TIBETAN WRITTEN
EASTERN BALTI
WESTERN BALTI
PURIK
KARGIL BALTI
CHORBAT BALTI
KHARMANG BALTI
SKARDU BALTI
KHAPALU BALTI
RONDU BALTI
SHIGAR BALTI
BOZE
GINGAREDE
BINE/BOZE GIRINGAREDE
BINE/SOGAL
SOGAL
BINE/SEBE
SEBE

ST.BODIC
ST.KUKI-CHIN-NAGA
ST.BODIC
ST.KUKI-CHIN-NAGA
ST.BODIC
ST.NUNGISH
ST.BODIC

SEBE
BINE/TATI
TATI
BINE/MASINGLE
MASINGLE
BINE/KUNINI
KUNINI
BINE/IREPI DRAGELI
DRAGELI
IRUPI
MERIAM
GIZRA/GIJARA
GIZRA/KUPERE
GIZRA/TOGO
GIZRA/WAIDORO
WAIDORO
KUPERE
TOGO
GIDRA/JIBU
YUTA
GUIAM
IAMEGA
ZIM
KAPAL
WIPIM
PEAWA
UME
ABAM
PODARI
WONIE
GAMAEWE
DOROGORI
KURU

WF.WESTERN FLY

LANI
WANO
ANGGURUK YALI
KINIAGEIMA
PYRAMID WODO
UPPER PYRAMID DANI
MID GRAND VALLEY DANI
HITIGIMA DANI
TANGMA DANI

TNG.DANI

BERIA ] NS.EASTERN SAHARAN
KARAS ] WBm.WEST BOMBERAI
IHA
MBAHAM ] WBm.WEST BOMBERAI

KAMANG/LANGKURU
KAMANG/PIDO II
KAMANG/PIDO I
KAMANG
KAMANG/LETLEY
KAMANG/KOLOMANE
WERSING ] TNG.KOLANA-TANGLAPUI
SAWILA ] TNG.KOLANA-TANGLAPUI

TNG.WEST TIMOR-ALOR-PANTAR

MAKASAE
FATALUKU
OIRATA

TNG.WEST TIMOR-ALOR-PANTAR

HAMAP
KABOLA/HAMAP
KABOLA/AIMOLI
KABOLA/PITUMBANG
ADANG PITUNG
KABOLA
KELON/HALERMAN
KELON/PROBUR
KLON
KUI INDONESIA
KAFOA
ABUI/ATIMELANG
ABUI/MAKADAI
LAMMA/BIANGWALA
LAMMA/MAUTA TUBAL
LAMMA
KAWA
LAMMA/KALONDAMA
NEDEBANG
BLAGAR/APURI PURA
BLAGAR/LIMARAHING
BLAGAR/TEREWENG
BLAGAR/BAKALANG

TNG.WEST TIMOR-ALOR-PANTAR

BLAGAR/TEREWENG
BLAGAR/BAKALANG
BLAGAR/RETTA PURA
BLAGAR/RETTA TERNATE
KAERA
SAR INDONESIA
TEWA
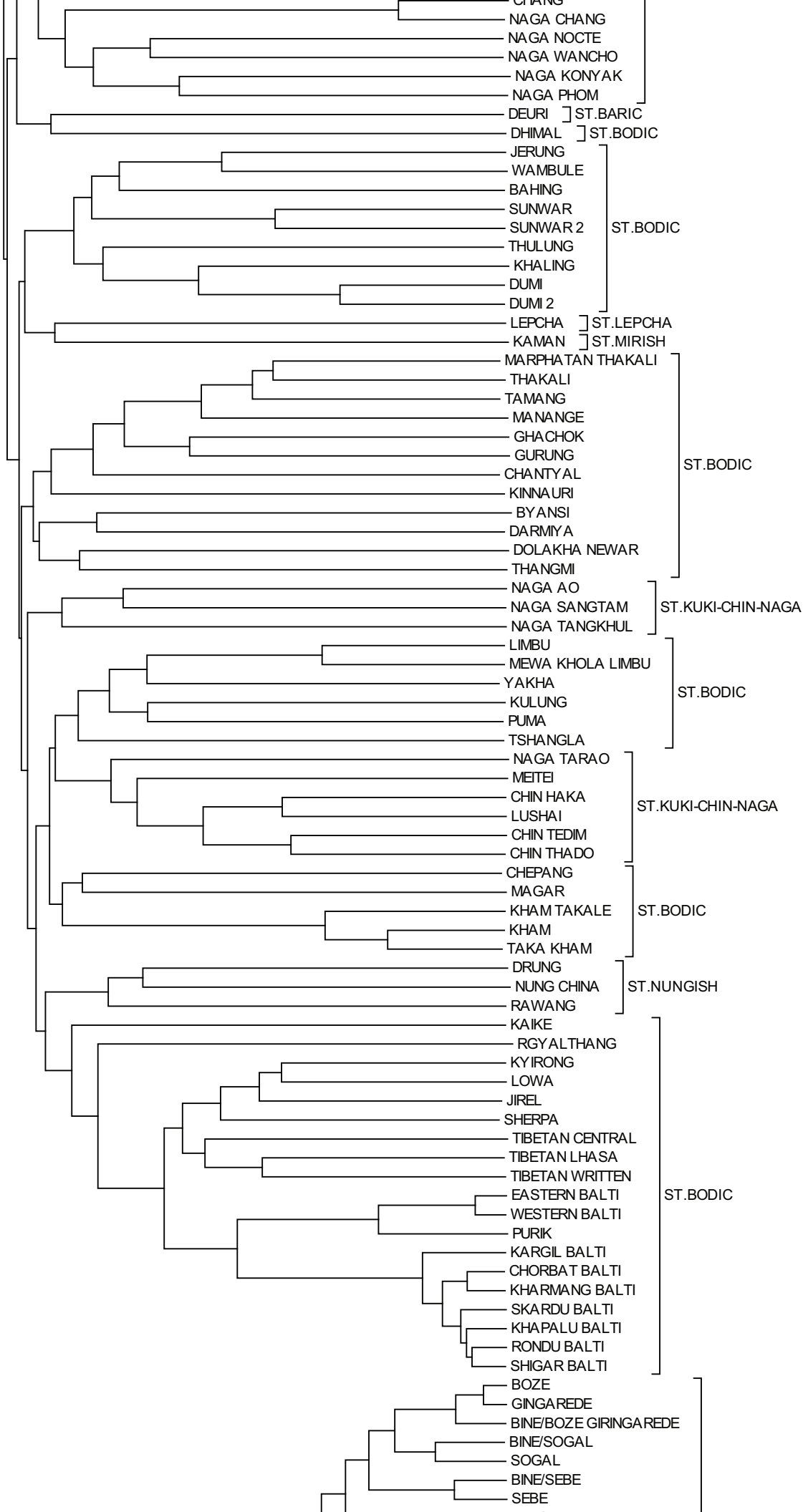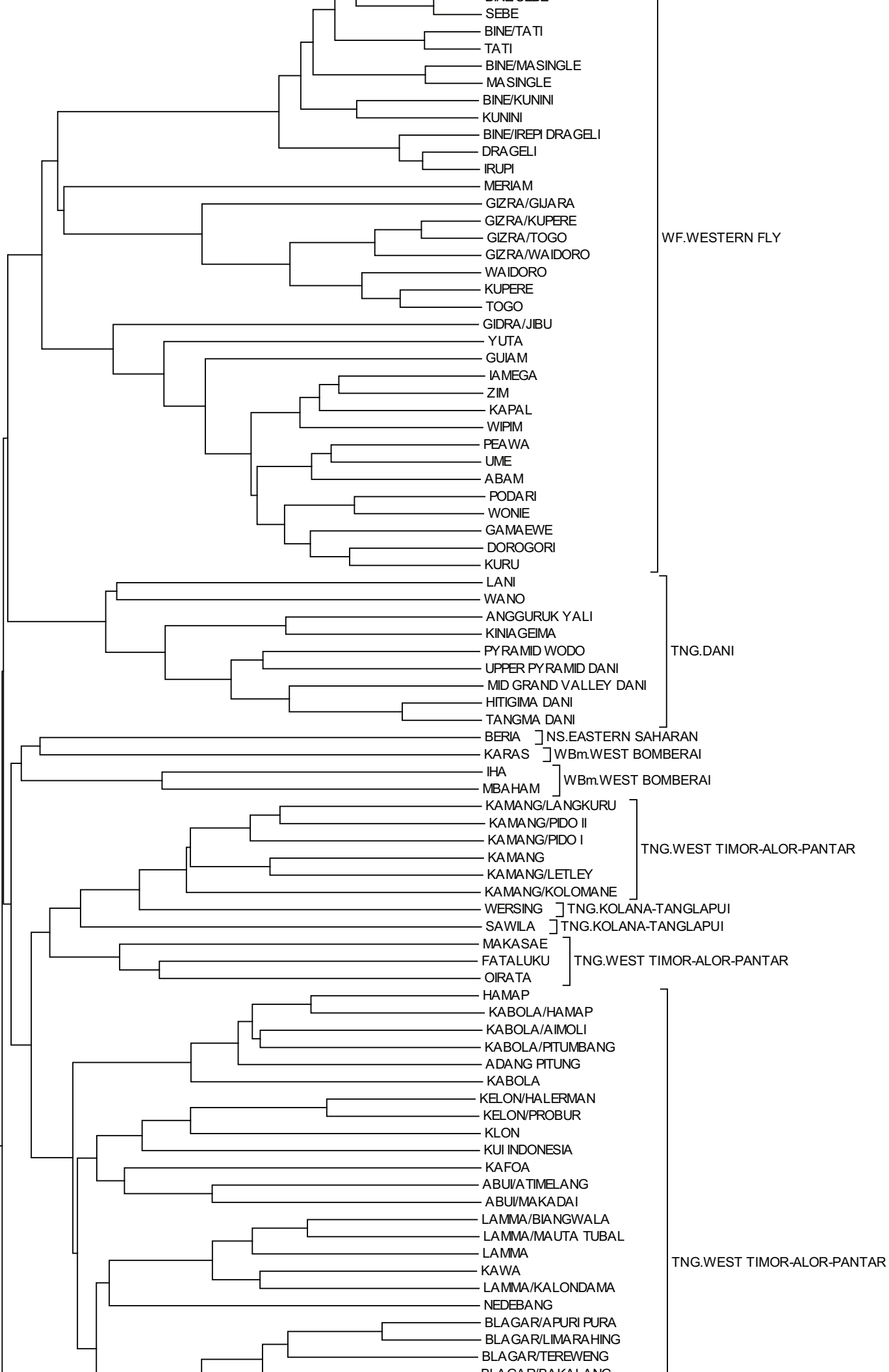TEWA/DEING
TEWA/LEBANG
TEWA/MADAR
TEWA/SARGANG
YURACARE
YURACARE YURUJURE ] Yur.YURACARE
GUNUNA KUNE ] Chn.PUELCHE
HAIDA ] Hai.HAIDA
MORWAP ] Mrw.MORWAP
SHASTA ] Hok.SHASTA
MELEKU JAIKA ] Chi.RAMA
CAHUARANO
ARABELA Zap.ZAPAROAN
ZAPARO
NAHUATL ACAXOCHITLAN
NAHUATL JALATLACO
NAHUATL SAN JERONINO AMANALCO
NAHUATL CUACUILA HUAUCHINANGO
NAHUATL SAN PEDRO TLALCUAPAN SANTA ANA C
NAHUATL SANTA ANA TLACOTENCO
NAHUATL XALATZALA TLAPA
NAHUATL XALPATLAHUAC
NAHUATL JICOCINGO ZACATLAN
NAHUATL SAN MIGUEL AYOTLA
NAHUATL XILOCUAUTLA HUAUCHINANGO
NAHUATL HUATLATLAUCA
NAHUATL XAALITLA TEPECUAUILCO
NAHUATL HUITZILTEPEC ZUMPANGO DEL RIO
NAHUATL SAN AGUSTIN OAPAN
NAHUATL SAN FRANCISCO TLALNEPANTLA
NAHUATL STA MA COAPAN
NAHUATL RAFAEL DELGADO
NAHUATL SAN JOSE MIAHUATLAN
NAHUATL SAN PABLO ZOQUITLAN
TETELCINGO NAHUATL
NAHUATL MONTEGRANDE PLATON SANCHEZ
NAHUATL CHINANCAHUATL ZACUALTIPAN
NAHUATL ZAHUASTIPAN SAN AUGUSTIN METZQUI
NAHUATL BUENOS AIRES ALAMO
NAHUATL CHICONTEPEC
NAHUATL CUAMELC0 TIANGUISTENGO
NAHUATL LAS BALSAS
NAHUATL IXHUATLAN DE MADERO
NAHUATL LA REFORMA TEPEHUACAN DE GUERRER   UA.AZTECAN
NAHUATL HUEYATI YAHUALICA
NAHUATL COXCATLAN
NAHUATL XOCHIATIPAN
NAHUATL CUATLAMAYAN ANTONIO SANTOS
NAHUATL TLALNEPANTLA TAMAZUNCHALE
NAHUATL XILITLA
NAHUATL ACATLAN
NAHUATL ZITLALA
NAHUATL ATLIACA TIXTLA
NAHUATL CHICHIQUILA
NAHUATL COYOTEPEC
NAHUATL CHILACACHAPA CUETZALA DEL PROGRE
NAHUATL COATEPEC COSTALES
NAHUATL SAN AGUSTIN DE BUENAVENTURA
NAHUATL SAN PEDRO JICORA
NAHUATL POMARO AQUILA
NAHUATL QUETZALAPA AZOYU
PIPIL
TABASCO NAHUATL CUPILCO
HIGHLAND PUEBLA NAHUATL
NAHUATL CHILOCOYO HUEHUETLA
NAHUATL IXTACAMAXTITLAN
NAHUATL CHIGNAUTLA
NAHUATL AYOTOXCO
NAHUATL REYES DE VALLARTA TUZAMAPAN
NAHUATL MECAYAPAN
NAHUATL PAJAPAN
POCHUTLA NAHUATL
NORTHERN PAIUTE ]

POCHUTLA NAHUATL
NORTHERN PAIUTE
SOUTHERN PAIUTE
UTE 1
UTE 2
KAWAIISU
COMANCHE
SHOSHONI
HOPI    UA.HOPI
TUBATULABAL    UA.TUBATULABAL
CAHUILLA
LUISENO    UA.TAKIC
PIMA BAJO
TOHONO OODHAM    UA.TEPIMAM
TEPECANO
EL NAYAR CORA    UA.CORACHOL
HUICHOL
GUARIJIO
HUARIJIO    UA.TARAHUMARAN
CENTRAL TARAHUMARA
OPATA
YAQUI
MAYO
MAYO LOS CAPOMOS    UA.CAHITA
YAQUI 2
CANDOSHI    Cnd.CANDOSHI
KUNZA    Kun.KUNZA
WAGARABAI
MIAN
ANGIYAKMIN FAIWOL
BIMIN
TIFAL
TELEFOL
NINGGIRUM KAWOMA    TNG.OK
NINATIE MUYU
NORTH KATI
DIGUL MUYU
METOMKA MUYU
SOUTH KATI
MORAORI    MUM.MOREHEAD AND UPPER MARO RIVERS
WARAO    War.WARAO
QAWASQAR    Ala.ALACALUFAN
KORANA
NAMA
KXOE
NARO    Kho.CENTRAL KHOISAN
GWI
GXANA
KWADI
LEHAR
SAFEN
NON    NC.NORTHERN ATLANTIC
NDUT FALOR
PALOR
NIMBORAN/BESUM
NIMBORAN
MEKWEI/MARIBU    Nim.NIMBORAN
MEKWEI/KENDATE
MEKWEI/WABRON
QUICHUA BOLIVAR CACHISAGUA
QUICHUA TUNGURAHUA SALASACA
QUICHUA CHIMBORAZO TROJE
QUICHUA AZUAY
QUICHUA COTOPAXI PAPAURCO
QUICHUA COTOPAXI COMPANIA GRANDE
QUICHUA TUNGURAHUA GUAPANTE
QUICHUA COTOPAXI TIGUA
QUICHUA LOJA
QUICHUA PICHANCHA    Que.QUECHUAN
QUICHUA CANYAR
QUICHUA CHIMBORAZO NIZAG
QUICHUA IMBABURA 2
QUECHUA IMBABURA
QUECHUA CHACHAPOYAS
INGA PUTUMAYO
QUECHUA PASTAZA
QUECHUA AYACUCHO
QUECHUA HUAYLAS ANCASH
CENTRAL AYMARA
JAQARU    Aym.AYMARAN

CENTRAL AYMARA
JAQARU
JAQARU 2
] Aym.AYMARAN

CALLAWALLA ] Cre.PUQUINA BASED
TUWARI ] LS.LEONHARD SCHULTZE
AWAKE ] Art.ARUTANI

CHIPAYA
CHIPAYA 2
UCHUMATAQU
] UC.URU-CHIPAYA

YANA ] Hok.YANA

POMO CENTRAL
POMO NORTHERN
POMO NORTHEASTERN
POMO EASTERN
SOUTHEASTERN POMO
KASHAYA
SOUTHERN POMO
] Hok.POMOAN

KANOE ] Kap.KAPIXANA
SERI ] Hok.SERI

HAVASUPAI
WALAPAI
PAIPAI
YAVAPAI
KILIWA
COCOPA
TIPAI
DIEGUENO
MOJAVE
MARICOPA
YUMA
] Hok.YUMAN

HADZA ] Had.HADZA
XINCA ] Xin.XINCAN
BLACKFOOT ] Alg.ALGONQUIAN
YAMANA ] Yam.YAMANA
HIGHLAND TEQUISTLATEC
OAXACA CHONTAL
] Teq.TEQUISTLATECAN

ACHUAR
HUAMBISA
JIVARO SHUAR
AGUARUNA
] Jiv.JIVAROAN

GELAO
GELAO WANZI
GELAO LAOZHAI
GELAO QIAOSHANG
LACHI
] TK.KADAI

SUI JUNG CHIANG
SUI LINGAM
SUI PYO
SUI
TEN
MAK
MAONAN
KAM ZHANGLU
SOUTHERN DONG
MULAO
] TK.KAM-TAI

LAKKJA ] TK.KADAI
BUYANG
PUBIAO
LAHA
PAHA
] TK.KADAI

LI BAODING
LI TONGSHI
] TK.HLAI
ONG BE ] TK.KAM-TAI
JIAMAO ] TK.HLAI

THAI 2
TAY
DIOI
YAY
YAY 2
PO AI
TAI PO AI
TAI WUMING
WU MING
SAEK 2
SAEK
SAEK 3
ZHUANG NORTHERN
TAI
THAI
DEHONG

THAI
DEHONG
NUNG
ZHUANG SOUTHERN
LU
YANG
TAI LUNGCHOW
TAI PING SIANG
TAI LEI PING
TAI NING MING
LUNG MING
TAI LUNGMING
TAI WESTERN NUNG
NUNG 2
THO
LAOTIEN
SHAN
SHAN 2
LUNG CHOW
WHITE THAI 3
SIAMESE 2
AHOM
TAI CHIENGMAI
TAI NONG KHAI
SIAMESE
BLACK THAI
WHITE THAI
BLACK TAI 2
WHITE THAI 2
LUE CHIENG HUNG
LUE MUONG YONG

TK.KAM-TAI

HILDI
WAMDIU
MARGI
HUBA
KILBA
PUTAI
CIBAK
BURA
NGGWAHYI
FALI GILI
FALI KIRIYA
BAZZA
FUTU
GHYE
KAMALE
NKAFA
DABA
BACAMA MULYEN
NZANYI
GUDU
SHARWA
TSUVAN
ZIZILIVAKAN
GUDE
FALI MUCHELLA
FALI BAGIRA
JIMI CAMEROON
CINENI
GLAVDA
GAVA
DGHWEDE
WANDALA
VAME
ZULGO
MAFA
MOFU GUDUR
MOLOKO
MAKARY KOTOKO

AA.BIU-MANDARA

BIDIYA
MIGAMA
KAJAKSE
MUBI
MOKILKO

AA.EAST CHADIC

MASANA
MUSEY
MASANA POGO
PEVE PALA
MARBA
MESME

AA.MASA

MESME
HERDE
PEVE LAME
POLCI BULI
POLCI ZUL
GEJI
DASS DWAT
POLCI
SAYA
DYARIM 1
DYARIM 2
BOGHOM
JIMI
HAUSA
HAUSA 2
KANO
NIMBIA
GITATA
ARABISHI
GARAKU
GWAGWA
KARSHI

AA.WEST CHADIC

MIYA
WARJI

AA.WEST CHADIC

GAANDA GABIN
GAANDA
BOGA
HWANA
TERA PIDLIMDI
TERA

AA.BIU-MANDARA

ANGAS
MUPUN
MWAGHAVUL
GOEMAI
KOFYAR
MISHIP
PERO
TANGALE
DERA
BOLE
NGAMO
KAREKARE
GERA
GERA 2
GALAMBU
GERUMA
KIRFI

AA.WEST CHADIC

BADE
NGIZIM

AA.WEST CHADIC

ZENAGA
AHAGGAR TUAREG
TAMAHAQ TAHAGGART
TAMASHEQ
TAMAJEQ TAYART/AIR
TETSERRET
OUARGLA BERBER
TAGARGRENT
TUMZABT
FIGUIG
NAFUSI
FOQAHA
SIWA BERBER
SIWI
GHADAMES
AWJILAH
GHOMARA
SENHAJA DE SRAIR
TAMAZIGHT CENTRAL ATLAS/AYT IZDEG
TAMAZIGHT CENTRAL ATLAS/AYT NDHIR
TAMAZIGHT CENTRAL ATLAS/NTIFA
TASHELHIT/IDA USEMLAL
KABYLE/GREATER KABYLIA AT MANGELLAT
METMATA
NAFUSI MATMATA
TARIFIT GUELAIA
BENI SNOUS WESTERN ALGERIAN BERBER
TARIFIT/BENI IZNASSEN

AA.BERBER

CHAHA
INNEMOR

INNEMOR
GETO
MESQAN
SODDO
ARGOBBA
AMHARIC
AMHARIC 2
HARARI
WALANI SILTE
ZWAY
GAFAT
MESMES
MESMES 2
AA.SEMITIC

HARSUSI
MEHRI
SAURI
SOQOTRI
AA.SEMITIC

ARAMAIC
HERTEVIN
MLAHSO
TUROYO
MIDOB UnnamedInSource
MODERN MANDAIC
AA.SEMITIC

TIGRE
TIGRE 2
TIGRINYA
TIGRINYA 2
HEBREW
AA.SEMITIC

KIBERA KENYA    Cre.ARABIC BASED

MALTESE
DELLYS
MOROCCAN ARABIC
NORTH LEVANTINE ARABIC
SYRIAN ARABIC
STANDARD ARABIC
PALESTINIAN ARABIC
TUNISIAN ARABIC MAGHRIB
EGYPT ARABIC
ARABIC GULF SPOKEN
OGADEN ARABIC
AA.SEMITIC

ARANDAI/BARAU
BARAU
ARANDAI/WERIAGAR
WERIAGAR
ARANDAI/SEBYAR
ARANDAI
ARANDAI/TAROF
TAROF
ARANDAI/NAJARAGO
ARANDAI/KASUWERI
KASUWERI
KAMPONG BARU
PURAGI
INANWATAN
INANWATAN/BIRA
INANWATAN/ITIGO
INANWATAN/SOLOWAT
Mar.SOUTH BIRDS HEAD

GOLA    NC.SOUTHERN ATLANTIC

KPELLE
KPELLE GUINEA
LOMA LIBERIA
BANDI
MENDE
NC.WESTERN MANDE

HUITOTO MINICA
HUITOTO NIPODE
HUITOTO MURUI
OCAINA
Hui.HUITOTO

MANTAURAN
RUKAI
An.TSOUIC

MATLATZINCA SAN FRANCISCO OXTOTILPAN
MATLATZINCA SAN FRANCISCO OXTOTILPAN 2
MAZAHUA
MAZAHUA CENTRAL
OTOMI TOLUCA
OTOMI MEZQUITAL
OTOMI QUERETARO
OM.OTOMIAN

APALAI
TIRIYO
WAYANA
KALINA

WAYANA
KALINA
MAQUIRITARI
HIXKARYANA
KAXUYANA
WAIWAI
CARIJONA
BAKAIRI                    Car.CARIBAN
ARARA DO PARA
IKPENG
AKAWAIO
MAKUSHI
TAUREPANG
PANARE
YUKPA
WAIMIRI ATROARI
KUIKURO
MAMAINDE
NAMBIKWARA        Nam.NAMBIKUARAN
SABANE
IRANTXE        ] Ira.IRANTXE
TOARIPI
UARIPI
AHEAVE
KEURU                  Ele.ELEMAN PROPER
OPAO
OROKOLO
AMUZGO        ] OM.AMUZGOAN
JICAQUE        ] Tol.TOL
CHATINO WESTERN HIGHLAND
ZACATEPEC CHATINO
CHATINO TATALTEPEC
ZAPOTEC SAN LUCAS QUIAVINI
ALOAPAM ZAPOTEC
YARENI ZAPOTEC
CHOAPAN ZAPOTEC
ISTHMUS ZAPOTEC
OCOTLAN ZAPOTEC
SANTA INES YATZECHI ZAPOTEC
TILQUIAPAN ZAPOTEC
ZAPOTEC SIERRA DE JUAREZ
LACHIXIO ZAPOTEC
ZAPOTEC YATZACHI
ZAPOTEC ZOOGOCHO           OM.ZAPOTECAN
CAJONOS ZAPOTEC
ZAPOTEC YALALAG
COATLAN LOXICHA ZAPOTEC
ZAPOTEC LOXICHA
TEXMELUCAN ZAPOTEC
SANTA MARIA QUIEGOLANI ZAPOTEC
SAN PEDRO QUIATONI ZAPOTEC
ZAPOTEC MITLA
AMATLAN ZAPOTEC
AYOQUESCO ZAPOTEC
ZAPOTEC MIXTEPEC
QUIOQUITANI QUIERI ZAPOTEC
XANAGUIA ZAPOTEC
BARA
WAIMAYA
TUYUCA
YURUTI
TUCANO
PIRATAPUYO
WANANO
MACUNA
CARAPANA
TATUYO                     Tuc.TUCANOAN
DESANO
SIRIANO
BARASANO
TANIMUCA
CUBEO
KOREGUAYE
OREJON
SECOYA
SIONA
BORORO
UMOTINA        ] MGe.BORORO
NATCHEZ        ] Nat.NATCHEZ
IK        ] NS.KULIAK

NATCHEZ ] Nat.NATCHEZ
IK ] NS.KULIAK
RANQUELCHE
MAPUDUNGUN ] Arc.ARAUCANIAN
MAPUDUNGUN 2
AINU SARU ] Ain.AINU
CHIMARIKO ] Hok.CHIMARIKO
KUTENAI ] Ktn.KUTENAI
TUSCARORA
MOHAWK
ONEIDA
SENECA ] Iro.NORTHERN IROQUOIAN
CAYUGA
ONONDAGA
EVENKI
EWENKI 1
NEGIDAL
KUR URMI
EWENKI 2
OROQEN
SOLON
HEZHE
NANAI CHINA
OROCHI
UDEHE
UILTA ] Alt.TUNGUSIC
NANAI 2
NANAI
ULCHA
MANCHU
MANCHU 2
XIBE
EVEN
EWEN
MONGUOR
SANTA MONGOLIAN
MOGHOL
BURIAT MONGOLIA
DAUR ] Alt.MONGOLIC
TACHENG DAGUR
KALMYK
MONGOLIAN
JARAWA ] SAn.SOUTH ANDAMANESE
ONGE
AKA BEA
SOUTH ANDAMAN
AKA BALE
AKA CARI
AKA KEDE ] GA.GREAT ANDAMANESE
AKA BO
AKA KOL
OKO JUWOI
MENYA ] TNG.ANGAN
BARUYA ] TNG.ANGAN
WINTU ] Pen.WINTUAN
DUMO
DUSUR
LEITRE
WUTUNG ] Sko.WESTERN SKO
SANGKE
SKOU
KUJARKE ] Kuj.KUJARGE
ABIPON ] Gcu.GUAICURUAN
CAVENENA
TACANA
ARAONA ] Tac.TACANAN
ESE EJJA
KONIBO
KULINA PANO
MATIS
MATSES
KAXARARI
POYANAWA
CAPANAHUA
AMAHUACA
SHARANAHUA
CASHIBO ] Pan.PANOAN
SHIPIBO
CHACOBO
SHANENAWA

CHACOBO
SHANENAWA
ARARA PANO
YAWANAWA
YAMINAWA
KAXINAWA
KATUKINA PANO
MARUBO
FOX
KICKAPOO
MENOMINEE
MIAMI
MAHICAN
DELAWARE MUNSEE
UNAMI UnnamedInSource
ATIKAMEKW
MONTAGNAIS
NASKAPI
NANTICOKE UnnamedInSource
WAMPANOAG NATICK
ABNAKI WESTERN
MOHEGAN
POTAWATOMI
CHIPPEWA
ALGONQUIN
OJIBWA EASTERN
OJIBWA (SEVERN) UnnamedInSource
OJIBWE MINNESOTA
MICMAC
PASSAMAQUODDY MALISEET UnnamedInSource
KASKASKIA ILLINOIS
Alg.ALGONQUIAN

NEZ PERCE
UMATILLA SAHAPTIN ] Pen.SAHAPTIAN
MOLALA ] Pen.MOLALA
HUAVE
HUAVE 2 ] Hua.HUAVEAN
CHITIMACHA ] Cht.CHITIMACHA

TOTONAC COYUTLA
TOTONAC HIGHLAND
PAPANTLA TOTONAC
TOTONAC OZELONACAXTLA
TOTONAC COATEPEC
TOTONAC OLINTLA
TOTONAC FILOMENA MATA
TOTONAC UPPER NEXACA
TOTONAC TEJERIA
XICOTEPEC TOTONAC
TOTONAC MISANTLA
TEPEHUA TLACHICHILCO
TEPEHUA HUEHUETLA
TEPEHUA PISA FLORES
Tot.TOTONACAN

TAPACHULTEC
NORTH HIGHLAND MIXE
ULTERIOR MIXE C
LOWLAND MIXE
SAYULA POPOLUCA
OLUTA POPOLUCA
AYUTLA MIXE
SOUTH HIGHLAND MIXE
SOTEAPAN ZOQUE
TEXISTEPEC ZOQUE
MARIA CHIMALAPA
MIGUEL CHIMALAPA
ZOQUE RAYON
CHIAPAS ZOQUE
ZOQUE FRANCISCO LEON
MZ.MIXE-ZOQUE

CHINANTEC OJITLAN
CHINANTEC PALANTLA
CHINANTEC SAN FELIPE USILA
LEALAO CHINANTEC
OM.CHINANTECAN
ALBANIAN TOSK ] IE.ALBANIAN
GILYAK ] Niv.NIVKH
WIYOT ] Alg.WIYOT
YUROK ] Alg.YUROK
BEOTHUK ] Beo.BEOTHUK
KAMSA ] Cam.CAMSA
TIMUCUA ] Tim.TIMUCUA
HANIS COOS ] OC.COOSAN
LENGUA
SANAPANA ANGAITE | Mas.MASCOIAN

LENGUA
SANAPANA ANGAITE
SANAPANA ENLHET — Mas.MASCOIAN

KYAIMBARANG
MIYAK — Yua.YUAT

HARUAI/WAIBUK
WIYAW
ARAMO
HAGAHAI/ARAMO II — UY.UPPER YUAT
NANGENUWETAN
PINAI/I
PINAI 1

GABIANO
PAKA
PIAME
SANIO
HEWA — Sep.SEPIK HILL
ALAMBLAK
BAHINEMO
KAPRIMAN
KWOMA
MENDE PNG — Sep.MIDDLE SEPIK
IWAM/MAY — Sep.UPPER SEPIK
NAMIA — Sep.YELLOW RIVER
AWTUW
POUYE — Sep.RAM

BOIKIN
YENGORU
KWUSAUN
NGALA
NYAURA — Sep.MIDDLE SEPIK
MANAMBU
YELOGU
MAPRIK
WOSERA

ALSEA — OC.ALSEA
SIUSLAW — OC.SIUSLAWAN
BRAHUI — Dra.NORTHERN DRAVIDIAN
KURUKH
SAURIA PAHARIA — Dra.NORTHERN DRAVIDIAN
TODA — Dra.SOUTHERN DRAVIDIAN
MALAYALAM
TAMIL
RAVULA
BADAGA — Dra.SOUTHERN DRAVIDIAN
KANNADA
TULU
TELUGU — Dra.SOUTH-CENTRAL DRAVIDIAN
KUI
KUVI — Dra.SOUTH-CENTRAL DRAVIDIAN
PENGO — Dra.SOUTH-CENTRAL DRAVIDIAN
GONDI UnnamedInSource
KONDA 1 UnnamedInSource — Dra.SOUTH-CENTRAL DRAVIDIAN
SOUTHERN GONDI
KODAVA UnnamedInSource
KOROMFE UnnamedInSource — Dra.SOUTHERN DRAVIDIAN
PARJI — Dra.CENTRAL DRAVIDIAN
GADABA POTTANGI OLLAR UnnamedInSource — Dra.CENTRAL DRAVIDIAN
NORTHWESTERN KOLAMI — Dra.CENTRAL DRAVIDIAN
KOTA UnnamedInSource — Dra.SOUTHERN DRAVIDIAN

AWING
PINYIN
MANKON
BAFUT
BAMBILI
MUNDANI
YEMBA 1
YEMBA 2
BAMILEKE
MEDUMBA
FEFE 1
FEFE 2
DZODINKA
LIMBUM — NC.BANTOID
MUNGAKA
SHUPAMEM
YAMBA
NDE BUKWOK
NDE YULANA
DE WUNGTSE

DE WUNGTSE
MFUMTE
NDANDA SUD
GHOMALA 1
GHOMALA 2
WUSHI
TIKAR AKUEN
TWUMWU
YUCHI ] Yuc.YUCHI
WAORANI ] Wao.WAORANI
CENTRAL SIERRA MIWOK
NORTHERN SIERRA MIWOK
S SIERRA MIWOK
PLAINS MIWOK
OLAMENTKE
MIWOK BODEGA
MIWOK LAKE
MUTSUN
RUMSEN
YOKUTS TINLINNEH
YOKUTS UNKNOWN
MAIDU KONKAU
MAIDU NORTHWEST NAKUM
NE MAIDU
NISENAN
WASHO ] Wsh.WASHO
MOSETEN ] Mos.MOSETENAN
WAPPO ] WY.WAPPO
YUKI ] WY.YUKIAN
DOROMU/ARAMAIKA
DOROMU/BAREIKA
DOROMU/LOFAIKA
DOROMU
MARIA/MARANOMU 1
MARIA
MISKITO
CACAOPERA
ULWA
GUAYMI
NGABERE MOVE
PECH ] Chi.PAYA
BORUCA ] Chi.TALAMANCA
CHIMILA ] Chm.CHIMILA
DAMANA
IKA
KOGUI
CHIBCHA ] Chi.CHIBCHAN PROPER
RAMA ] Chi.RAMA
CUNA
KUNA YALA SAN BLAS
TUNEBO ] Chi.CHIBCHAN PROPER
BUGLERE MURIRE ] Chi.GUAYMI
BARI COLUMBIA ] Chi.MOTILON
DORASQUE ] Chi.DORASKE
TERIBE ] Chi.TALAMANCA
BRIBRI
CABECAR CHIRIPO
CABECAR
WESTERN QUICHE MOMOSTENANGO
WESTERN QUICHE TOTONICAPAN
CENTRAL QUICHE SANTA MARIA CHIQUIMULA
WESTERN QUICHE SANTA CATARINA IXTAHUACAN
SIPAKAPENSE
EASTERN QUICHE RABINAL
SACAPULTECO SACAPULAS CENTRO
CENTRAL QUICHE
USPANTEKO
TZUTUJIL WESTERN
TZUTUJIL SAN JUAN LA LAGUNA
TZUTUJIL SANTIAGO ATITLAN
SOUTHERN CAKCHIQUEL SAN ANDRES ITZAPA
WESTERN CAKCHIQUEL PATZUN
NORTHERN CAKCHIQUEL SAN MARTIN JILOTEPEQ
NORTHERN CAKCHIQUEL TECPAN
POCOMAM EASTERN
POCOMAM SAN LUIS JILOTEPEQUE
POQOMCHI WESTERN
WESTERN POCOMAM SAN CRISTOBAL VERAPAZ
AGUACATEC
AGUACATECO AGUACATAN

Pen.MIWOK

Pen.COSTANOAN

Pen.YOKUTS

Pen.MAIDUAN

TNG.MANUBARAN

Mis.MISUMALPAN

Chi.GUAYMI

Chi.ARUAK

Chi.KUNA

Chi.TALAMANCA
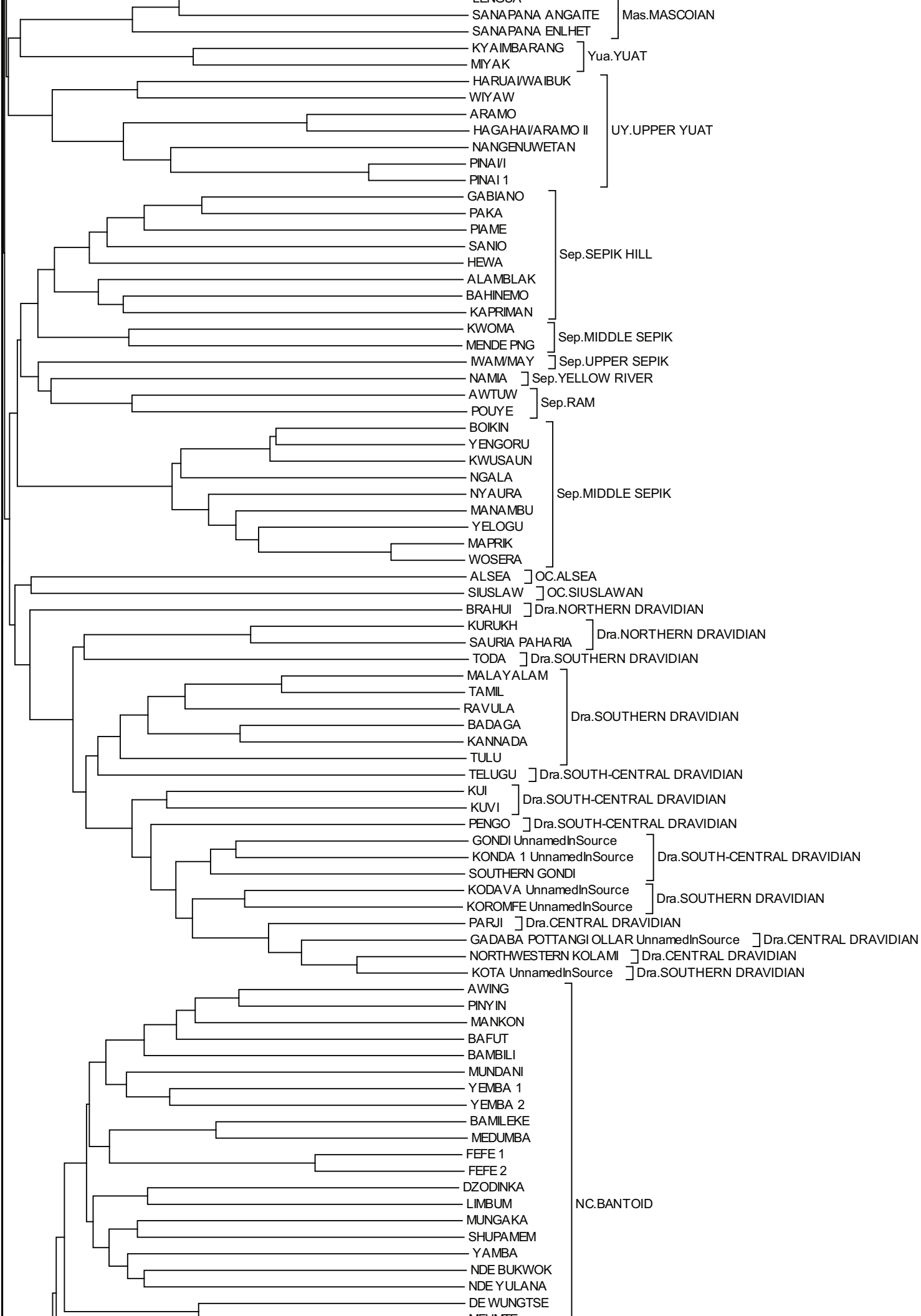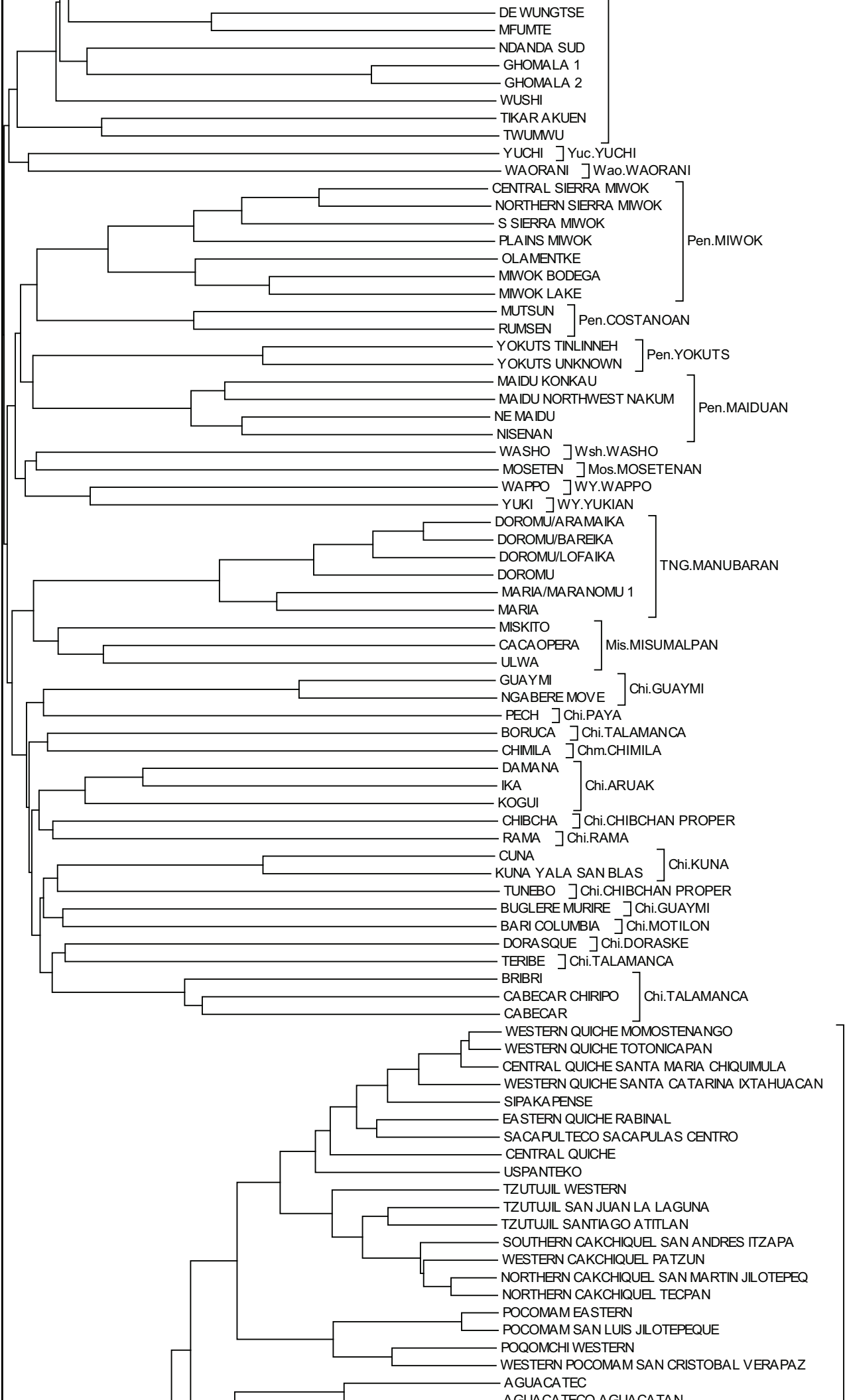
AGUACATEC
AGUACATECO AGUACATAN
IXIL CHAJUL
IXIL NEBAJ
MOCHO
JACALTEC
JACALTECO WESTERN
ACATECO SAN MIGUEL ACATAN
QANJOBAL SANTA EULALIA
TECO TECTITAN
MAM NORTHERN
MAM TODOS SANTOS CUCHUMATAN 1
MAM TODOS SANTOS CUCHUMATAN 2
MAM TACANA 1
MAM TACANA 2
MAM TECTITAN
MAM TACANA 3
MAM CAJOLA
MAM SAN MIGUEL SIGUILA
MAM CONSEPCION CHIQUIRICHAPA
MAM SAN JUAN OSTUNCALCO
MAM SAN MARTIN SACATEPEQUEZ
SOUTHERN MAM SAN JUAN OSTUNCALCO
MAM CABRICAN
NORTHERN MAM SAN ILDEFONSO IXTAHUACAN
MAM SAN SEBASTIAN H
MAM SANTIAGO CHIMALTENANGO
MAM SANTA BARBARA
MAM SAN PEDRO NECTA
MAM SAN RAFAEL PETZAL
MAM TODOS SANTOS CUCHUMATAN 3
MAM SAN ILDEFONSO IXTAHUACAN
MAM SAN GASPAR IXCHIL
MAM SAN JUAN ATITLAN
KEKCHI
EASTERN KEKCHI CAHABON
WESTERN KEKCHI SAN JUAN CHAMELCO
TZELTAL OXCHUC
TZELTAL BACHAJON
TZOTZIL SAN ANDRES
ZINACANTAN TZOTZIL
TOJOLABAL
CHUJ
CHUJ SAN MATEO IXTATAN
CHOL TILA
CHOL TUMBALA
CHONTAL TABASCO
CHORTI
CHORTI JOCOTAN
MOPAN
MOPAN SAN LUIS PETEN
LACANDON
ITZAJ
MAYA YUCATAN
CHICOMUCELTEC
HUASTEC

May.MAYAN

GAVIAO DO JIPARANA
GAVIAO DO RONDONIA
CINTA LARGA
SURUI DO RONDONIA
MONDE

Tup.MONDE

MEKENS
TUPARI
MAKURAP

Tup.TUPARI

AWETI    ] Tup.AWETI
SATERE MAWE    ] Tup.MAWE-SATERE
KURUAYA
MUNDURUKU

Tup.MUNDURUKU

JURUNA
XIPAYA

Tup.TUPI-GUARANI

SIRIONO
GUARANI
ACHE
ARAWETE
XETA
GUARANI KAIWA
MBYA
CHIRIGUANO
TAPIETE
WAYAMPI

TAPIETE
WAYAMPI
AMONDAVA
URUEWAUWAU
TENHARIM
APIAKA
PARINTINTIN
KOKAMA
OMAGUA
NHEENGATU
ASURINI
PARAKANA
TAPIRAPE
GUAJA
URUBU KAAPOR
ANAMBE
GUAJAJARA
TEMBE
SURUI DO PARA
ASURINI XINGU
AVA CANOEIRO
KAMAYURA
KAYABI
ZOE

Tup.TUPI-GUARANI

GANZA
NORTHERN MAO
HOZO
MAJI
NAO
SHEKO
BENCHO
SHE
KEFA
MOCHA 2
MOCHA
SOUTHERN MAO
SHINASSHA
SHINASSHA 2
JANJERO
GIDICHO
KOYRA
ZAYSE
ZERGULLA
KACHAMA
MALE ETHIOPIA
BASKETO
DACHE
DORZE
OYDA
WELAMO
KULLO
WOLAYTTA

AA.NORTH OMOTIC

BAUZI/NOIADI
BAUZI
TARUNGGAREH
TURUNGGARE UNKNOWN DIAL

EGB.EAST GEELVINK BAY

LENCA EL SALVADOR
LENCA HONDURAS

Len.LENCAN

IGBO ONITSHA
IZI
OGBA
ECHIE
EKPEYE

NC.IGBOID

EAST LIMBA
YELETNYE

NC.SOUTHERN ATLANTIC
Yel.YELE

BIAFADA
PAJADE
BALANTA
BALANTA GANJA
NALU
DIOLA
JOLA
PAPEL
MANKAN
MANJACA CHURO
MANJACA COSTA DE BAIXO

NC.NORTHERN ATLANTIC

DANARU
URIGINA
SUMAU
USINO

SUMAU
USINO
USU
BIYOM
TAUYA
KWATO
DUDUELA
ERIMA
ASAS
SINSAURU
KESAWAI
SAUSI
KOLOM
SUROI
DUMPU
ARAWUM
LEMIO
PULABU
GANGLAU
SAEP
YABONG
JILIM
RERAU
YANGULAM
BOM
BONGU
MALE PAPUANG
SONGUM
KARE
ISEBE
PANIM
AMELE
GUMALU
SIHAN
BAU
BEMAL
MUNIT
GIRAWA
BAGUPI
NAKE
MURUPI
SARUGA
SAMOSA
MOSIMO
WAMAS
GARUS
RAPTING
REMPI
YOIDIK
GARUH
KAMBA
BAIMAK
MAWAN
GAL
MATEPI
SILOPI
UTU

TNG.MADANG

KAURE
KOSARE
Kau.KAURE

SUMARARU
WOMO
POKO RAWO
Sko.SERRA HILLS

SUMO
RAMO
BARUPU
POO
Sko.WARUPU

KATIATI
SILEIBI
EMERUM
FAITA
MUSAK
OSUM
WADAGINAM
MORESADA
IKUNDUN
PONDOMA
PAYNAMAR
ANGAUA
ATEMPLE

TNG.MADANG

BARGAM
TNG.MADANG

A TEMPLE

BARGAM  ⎤ TNG.MADANG
MUGIL   ⎦

MOLOF   ⎤ Mol.MOLOF

KOBON   ⎤ TNG.MADANG

DIMIR
ABASAKUR
KOGUMAN
PILA
SAKI
TANI
PAY
HINIHON
MAWAK
KOWAKI
MUSAR
WANAMBRE
ULINGAN
BEPOUR
MOERE
BUNABUN
MALAS
AMAIMON
WANUMA
YABEN
BILAKURA
UKURIGUMA
PARAWEN
YARAWATA          TNG.MADANG

NIHALI   ⎤ Nah.NAHALI
NAHALI   ⎦

WOLOF    ⎤ NC.NORTHERN ATLANTIC

FOE   ⎤ TNG.KUTUBUAN

FASU     ⎤ TNG.FASU
NAMUMI   ⎦

GUAMBIANO
TOTORO
AWA PIT         Bar.BARBACOAN
CAYAPA
COLORADO

MOR 2   ⎤ TNG.MOR

KENABOI 1   ⎤ Ken.KENABOI

HATAM   ⎤ WP.HATAM

MAFULU   ⎤ TNG.GOILALAN

KARAJA   ⎤ MGe.KARAJA

KAKIAE   ⎤ Ele.TATE

AFOA   ⎤ TNG.GOILALAN

KOVIO   ⎤ TNG.GOILALAN

MAIPUA   ⎤ Ele.PURARI
PURARI   ⎦

ASARO
GAHUKU/ASARO
GAFUKU
GAHUKU
SIANE
YABIYUFA
GENDE
ISABI
FORE
GIMI
KAMANO KAFE
BENABENA
YAGARIA
YATE
BINUMARIEN
TAIRORA/BINUMARIEN
N TAIRORA
WAFFA
AWA
AWA 2
AUYANA
GADSUP
AGARIBI
GADSUP/AGARABI          TNG.EASTERN HIGHLANDS

LONGUDA   ⎤ NC.ADAMAWA

LORHON   ⎤ NC.GUR

DEM   ⎤ TNG.DEM

YAREBA   ⎤ Yrb.YAREBAN

DAMAL   ⎤ TNG.DAMAL

CAYUBABA   ⎤ Cay.CAYUVAVA

CAYUBABA ] Cay.CAYUVAVA

KARKAR YURI
YURI
YAFI        Pau.PAUWASI
BIKSI ] Sep.BIKSI

YALE KOSAREK
UNA
BIME        TNG.MEK
EIPOMEK

GWAMA
SOUTHERN KOMA
NORTHERN KOMA    NS.KOMAN
CENTRAL KOMA
UDUK

KISSI
KRIM
MMANI        NC.SOUTHERN ATLANTIC
SHERBRO

KIOWA
TEWA ARIZONA        KT.KIOWA-TANOAN
TEWA SAN JUAN PUEBLO

BINANDERE
MAMBARE RIVER
KORAFE YEGHA        TNG.BINANDEREAN
SUENA
TAFOTA BARUGA

KORAK
WASKIA        TNG.MADANG
HOAN ] Kho.SOUTHERN KHOISAN
NORTHERN TUJIA ] ST.TUJIA

PAWAIAN
PAWAIA        Teb.PAWAIAN
DARIBI
FOLOPA        Teb.TEBERAN

AUA
GAWIL
MELPA
KANDAWO
NARAK
GOLIN        TNG.CHIMBU
KUMAN
SINASINA
BOUMAI
DOM

S KIWAI/SC/TURETURE
TURETURE
S KIWAI/SC/MAWATA
DOMORI
KIWAI
WABUDA
BAMU
BAMU 2
ANIGIBI        Kiw.KIWAIAN
GIBAIO
GOPE
URAMA
KEREWO
MORIGI

TABO/WAIA
WAIA        MUM.MOREHEAD AND UPPER MARO RIVERS
TAPAPI
SUKI ] GS.SUKI

GOGODALA/ARI
GOGODALA/GIRARA
GOGODALA
GOGODALA/GAIMA        GS.GOGODALA
GOGODARA
ADIBA
GOGODALA/ADIBA

BUIN
MOTUNA        EB.EAST BOUGAINVILLE
NASIOI

DUBU
DUBU UNKNOWN DIAL
DUBU/AFI        Pau.PAUWASI
TOWEI
WIRU ] TNG.WIRU
TOFAMNA ] Tof.TOFAMNA
RUMU ]
OMATI

RUMU
OMATI
IKOBI
MENA
TuK.TURAMA-KIKORIAN

ORIG ] NC.RASHAD

KAROK ] Kar.KAROK

KONERAW
MOMBUN
Mom.MOMBUM

IRIA/ASIENARA
IRIA
KAMORO
SEMPAN
ASMATH NORTH
ASMAT CENTRAL
CASUARINA COAST ASMAT
CITAK
TNG.ASMAT-KAMORO

ATOHWAIM/KAUGAT
KAUGAT
KAYGIR
TAMAGARIO
Kay.KAYAGAR

KLAMATH ] Pen.KLAMATH-MODOC

YELMEK/JAB
YELMEK
MAKLEW
MEKLEW
Bul.BULAKA RIVER

KOIARI
KOIARI 2
KOITA
MOUNTAIN KOIARI
AOMIE
BARAI
ESE MANAGALASI
TNG.KOIARIAN

KIMAGHAMA
RIANTANA
NDOM
Kol.KOLOPOM

MONI
KAPAUKU
WODANI
TNG.WISSEL LAKES-KEMANDOGA

HUBE
TOBO
DEDUA
BORONG
BURUM MINDIK
MOMOLILI
NABAK
KOMBA
SELEPET
TIMBE
ONO
MIGABAC
KATE
MAPE
NEK
NUKNA
YOPNO
AWARA
WANTOAT
TNG.FINISTERRE-HUON

KOROWAI
SAWUY
KAETI DUMUT
KAETI
WAMBON
AGHU
PISA
SIAGHA
YENIMU
TNG.AWJU-DUMUT

BOAZI/SOUTH
SOUTH BOAZI
BOAZI
BOAZI/BOAZI
BOAZI/KUINI
KUINI
BEGUA
ZIMAKANI
MARIND/MARIND
MARIND
MARIND/TUGERI
WARKAY
YAKAY
YAQAY
Mar.MARIND PROPER

YAKAY
YAQAY
MAIBI
YARIBA
LEMBENA
ENGA
KYAKA ENGA
INIAI
BISORIO
PIKARU
HULI HOLE
HULI
SAU
KEWA
KEWA/S/POLE
POLE

TNG.ENGAN

BOGAYA
DUNA

TNG.DUNA-BOGAYA

KAMULA     Kam.KAMULA
PARE     AP.AWIN-PARE

DIBIYASO
BEAMI
ETORO
EDOLO
BEDAMINI
BIAMI
KASUA
KASUA 2
ONABASULU
AIMELE
SONIA
SUNIA
KALULI
BOSAVI
KALULI 2

Bos.BOSAVI

BANIVA
BANIWA
CURRIPACO
TARIANA
TARIANO
ACHAGUA
CABIYARI
PIAPOCO
WAREKENA
MAIPURE
YAVITERO
YUCUNA
INAPARI
PALIKUR
MAWAYANA
WAPIXANA
BAHUANA
BARE
GUINAU
PARESI
SARAVEKA
YAWALAPITI
MEHINAKU
WAURA
MASHCO PIRO
PIRO
MAXINERI
APURINA
CHAMICURO
BAURE
IGNACIANO
TRINITARIO
KINIKINAU
TERENA
PARAUJANO
WAYUU
LOKONO
TAINO
GARIFUNA
ISLAND CARIB
AMUESHA 1
AMUESHA 2
CAQUINTE
NOMATSIGUENGA
MACHIGUENGA

Arw.ARAWAKAN

NOMATSIGUENGA
MACHIGUENGA
ASHENINKA
CAMPA DE PERENE
RESIGARO  Arw.ARAWAKAN
MUINANE
BORA  Hui.BORAN
MIRANA
KONUA  WBg.WEST BOUGAINVILLE
ROTOKAS
MORE
PAKAANOVA  CW.CHAPACURA-WANHAN
DUNGERWAB TSI
IAUGA/DUNGERWAB
IAUGA/PARB  MUM.MOREHEAD AND UPPER MARO RIVERS
PARB
CHOLON  Chl.CHOLON
OBISPENYO
CRUZENYO
VENTURENYO  Chu.CHUMASH
CHUMASH BARBARENO
INESENYO
YAMBES
YAMPES
TORRICELLI
KOMBIO
ARAPESH2  Tor.KOMBIO-ARAPESH
ARAPESH
WAM
WOM
BASQUE  Bas.BASQUE
KARIRI XOCO  MGe.KARIRI
AU
NINGIL  Tor.WAPEI-PALEI
FUR  NS.FUR
KONDA 2  Mar.SOUTH BIRDS HEAD
SHABO  NS.SHABO
SALKA
TSUVADI  NC.KAINJI
CICIPU
OKSAPMIN  Oks.OKSAPMIN
IPIKO
MINANIBAI  IG.INLAND GULF
TAO SUAMATO
LANGA  NS.KOMAN
GREEK  IE.GREEK
TANAH MERAH  TNG.TANAHMERAH
WICHITA
ARIKARA  Cad.CADDOAN
PAWNEE
ARIEPI
SARAWANDORI
AMBADAIRU
TINDARET  Yaw.YAWA
KONTI UNAI
MARIADEI
CHEYENNE  Alg.ALGONQUIAN
L MOREHEAD/PEREMKA
PEREMKA  MUM.MOREHEAD AND UPPER MARO RIVERS
DISOHA
SESE  NS.GUMUZ
GUMUZ
MULAHA/IAIBU
MULAHA  Kwa.KWALEAN
MULAHA/MULAHA
MAGORI  An.OCEANIC
LAUA
DOMU  TNG.MAILUAN
MAILU
SALISH STRAITS
SAMISH
SONGISH
CLALLAM
COWICHAN
MUSQUEAM  Sal.CENTRAL SALISH
SQUAMISH
SLIAMMON
LUSHOOTSEED
TWANA  Sal.CENTRAL SALISH
CHEHALIS UPPER
Sal.TSAMOSAN

TWANA
CHEHALIS UPPER
COWLITZ — Sal.TSAMOSAN
TILLAMOOK — Sal.TILLAMOOK
LILLOOET
THOMPSON
OKANAGAN COLVILLE
COEUR DALENE — Sal.INTERIOR SALISH
MONTANA SALISH
SPOKANE
BELLA COOLA — Sal.BELLA COOLA
CATAWBA — Sio.SIOUAN
QUILEUTE — Chk.CHIMAKUAN
XOON MASARWA
XOON NUEN
XOON — Kho.SOUTHERN KHOISAN
NU
LILAU
MONOMBO — Mon.MONUMBO
BUNGAIN
BUNA
KAKARA BUNA
MANDI PAPUANG
MUNIWARA — Tor.MARIENBERG
URIMO
KAMASAU
ELEPI/SAMAP
ELIPI
JICALTEPEC MIXTEC
MIXTEC CHAYUCO
MIXTEC PENYOLES
YOSONDUA MIXTEC
MIXTECO DE SAN JUAN COLORADO — OM.MIXTECAN
MIXTEC ALCOZAUCA
CUICATEC
TRIQUI CHICAHUAXTLA
TRIQUI COPALA
EYAK — NDe.EYAK
KANAMARI
KATUKINA — Kat.KATUKINAN
KATAWIXI
WALIO — LS.LEONHARD SCHULTZE
AMARAKAERI — Har.HARAKMBET
IXCATEC
POPOLOCA DE SAN VICENTE COYOTEPEC
POPOLOCA METZONTLA
POPOLOCA SAN JUAN ATZINGO
CHOCHO OCOTLAN
CHOCHOTEC
SAN LORENZO CUAUNECUILTITLA
MAZATEC CHIQUIHUITLAN
HUAUTLA DE JIMENEZ
SAN JERONIMO TECOATL
MAZATLAN DE FLORES — OM.POPOLOCAN
SAN MIGUEL HUAUTLA
SAN JUAN CHIQUIHUITLAN
JALAPA DE DIAZ
SAN PEDRO IXCATLAN
SAN BARTOLOME AYAUTLA
SAN MIGUEL SOYOLTEPEC
DAW
NADEB
JUPDA — VJ.VAUPES-JAPURA
YUHUP
SHOM PENG — Sho.SHOM PENG
PIRAHA — Mur.MURA
ITONAMA — Ito.ITONAMA
MOVIMA — Mov.MOVIMA
BANARO — LSR.GRASS
ANOR
RAO — LSR.ANNABERG
GAMEI
KAIAN — LSR.LOWER RAMU
MIKAREW MAKARUB
GIRI
GIRI KIRE — LSR.MIKAREW
KIRE
AWYI/NJAO
AWYI
AWYI UNKNOWN DIAL

AWYI
AWYI UNKNOWN DIAL
AWYI/KEMBRE
AWYI/KIMBRIMORO
MANEM/IMOM
TAIKAT/GIRERE
ARZO/TAMI
TAIKAT
MANEM/SKOFRO
MANEM/MEREM
SENGGI
WARIS UNKNOWN DIAL
WAINA
IMONDA
WARIS

Bor.BORDER

ORYA UNKNOWN DIAL
ORYA
SAWE

TO.ORYA

BERIK
BERRIK PAPUA

TO.TOR

SAMAROKENA/TOMAYO
TOMAJO
SABERI
KWERBA/AIRMATI
KWERBA/NAIDJBEDJ
KWERBA/KAUWERAWET II
KWERBA/KAUWERAWET I

Kwe.KWERBA

NAGATIMAN
NAGATMAN

Yal.YALE

CHAMACOCO
AYOREO
AYOREO 2

Zam.ZAMUCOAN

DARA GAZ KHORASANI 1
LOTF ABAD KHORASANI
DARA GAZ KHORASANI 2
SHIRWAN KHORASANI
ZEYARAT KHORASANI
SHURAK KHORASANI
ASADLI KHORASANI
DOUGHAI KHORASANI
QUCHAN KHORASANI
JONK KHORASANI
GUJGI KHORASANI
MARESHK KHORASANI
CHARAM SARJAM KHORASANI
QARA BAGH KHORASANI
PIR KOMAJ KHORASANI
HARW E OLYA KHORASANI
JOGHATAY KHORASANI
HOKM ABAD KHORASANI
SOLTAN ABAD KHORASANI
RUH ABAD KHORASANI
SHEYH TEYMUR KHORASANI
LANGAR KHORASANI
AZERBAIJANI NORTH
TURKMEN
GAGAUZ UnnamedInSource
TURKISH 2
SALAR
UYGHUR
KARAKALPAK
KAZAKH
KYRGYZ
KHALAJ
UZBEK
CRIMEAN TATAR
KARACHAY BALKAR
KUMYK
NOGAI
BASHKIR
TATAR
KAZAN TATAR
MISHER TATAR
DOLGAN
SAKHA
KHAKAS
ALTAI
SHOR
TOFA
TUVAN
KARAIM

Alt.TURKIC

TUVAN
KARAIM
TURKISH
CHUVASH
KUOT ] Kut.KUOT
DITIDAHT
MAKAH
NOOTKA ] Wak.SOUTHERN WAKASHAN
HEILTSUK
KWAKWALA ] Wak.NORTHERN WAKASHAN
INGASSANA ] NS.EASTERN JEBEL
ASSAN
KOTT
ARIN
PUMPOKOL
KET
YUGH RECENT ] Yen.YENISEIAN
TABLA UNKNOWN DIAL
TABLA/W
TABLA
TABLA/C
SENTANI ] Snt.SENTANI
QIANG LONGXI
QIANG MIANCHI ] ST.QIANGIC
TIRIO ] TNG.TIRIO
BARDI ] Aus.NYULNYULAN
MOLMO ONE ] Tor.WEST WAPEI
TEMEIN ] NS.TEMEIN
PURI ] MGe.PURI
WALMAN ] Tor.WAPEI-PALEI
KUKWO ] Tor.URIM
SRENGE ] Tor.WAPEI-PALEI
GANISH HUNZA BURUSHASKI
HUSSAINABAD HUNZA BURUSHASKI
MURTAZABAD HUNZA BURUSHASKI
HOPER NAGAR BURUSHASKI
UYUM NAGAR BURUSHASKI
HAIDERABAD HUNZA BURUSHASKI
NAZIMABAD HUNZA BURUSHASKI
CENTRAL YASIN BURUSHASKI
NORTHERN YASIN BURUSHASKI
BURUSHASKI ] Brs.BURUSHASKI
KOREAN ] Kor.KOREAN
BUNAK ] TNG.WEST TIMOR-ALOR-PANTAR
BANAWA
JAMAMADI
JARAWARA
MADIJA
DENI
KULINA ARAUA
PAUMARI ] Aru.ARAUAN
ALEUT ] EA.ALEUT
TANACROSS
HARE
KUTCHIN
C CARRIER
JICARILLA
SAN CARLOS
NAVAHO
CHIRICAHUA
LIPAN
NAVAJO
JICARILLA APACHE
CARRIER
CHIPEWYAN
SARCEE
GALICE
HUPA 2
BEAVER
KATO
HUPA
MATTOLE ] NDe.ATHAPASKAN
LODA
LOLODA
TOBELO
TABARU
MADOLE
MODOLE
PAGU
GALELA

PAGU
GALELA
TIDORE
SAHU
WEST MAKIAN
SABALE
BOBAWA
MALAPA
NGOFABOBAWA
TALAPAO
TAGONO
TAFASOHO

WP.NORTH HALMAHERAN

DIME
ARI
BANNA

AA.SOUTH OMOTIC

LARAGIYA
LARRAKIA

Aus.LARAGIYAN

CHAYAHUITA
JEBERO

Cah.CAHUAPANAN

N ITELMEN
S ITELMEN

CK.SOUTHERN CHUKOTKO-KAMCHATKAN

CHUKCHEE
ALUTOR
KORYAK

CK.NORTHERN CHUKOTKO-KAMCHATKAN

EASTERN FRISIAN
NORTHERN LOW SAXON
LIMBURGISH
LUXEMBOURGISH
STANDARD GERMAN
YIDDISH EASTERN
ALSATIAN
BERNESE GERMAN
SAXON UPPER
SWABIAN
PLAUTDIETSCH
FRISIAN WESTERN
NORTH FRISIAN AMRUM
BRABANTIC
AFRIKAANS
DUTCH
STELLINGWERFS
WESTVLAAMS
FRANS VLAAMS
ZEEUWS
FAROESE
ICELANDIC
DANISH
NORWEGIAN BOKMAAL
SWEDISH
JAMTLANDIC
NORWEGIAN NYNORSK TOTEN

IE.GERMANIC

BERBICE DUTCH CREOLE
NEGERHOLLANDS

Cre.DUTCH BASED

AUKAN
SRANAN TONGO
SARAMACCAN
SARAMACCAN 2

Cre.ENGLISH BASED

SCOTS    ] IE.GERMANIC

TOK PISIN
TORRES STRAIT CREOLE
BISLAMA
KRIOL
NGUKURR BAMYILI CREOLE

Cre.ENGLISH BASED

KRIO SL
KRIO

Cre.ENGLISH BASED

GULLAH    ] Cre.ENGLISH BASED

HAWAI CREOLE ENGLISH    ] Cre.ENGLISH BASED

LIMONESE CREOLE
VINCENTIAN CREOLE
JAMAICAN CREOLE

Cre.ENGLISH BASED

ENGLISH    ] IE.GERMANIC

GEECHEE    ] Cre.ENGLISH BASED

NIGERIAN PIDGIN
GHANAIAN PIDGIN ENGLISH
KAMTOK
PICHI

Cre.ENGLISH BASED

CIMBRIAN    ] IE.GERMANIC

MICHIF    ] Cre.FRENCH BASED

HAITIAN CREOLE
ST LUCIAN CREOLE FRENCH

HAITIAN CREOLE
ST LUCIAN CREOLE FRENCH
KARIPUNA CREOLE
GUADELOUPE CREOLE
REUNIONNAIS
SEYCHELLES CREOLE 2
MAURITIAN
SEYCHELLES CREOLE
FRENCH ] IE.ROMANCE
ARPITAN ] IE.ROMANCE

Cre.FRENCH BASED

ANGOLAR
SANTOMENSE      Cre.PORTUGUESE BASED
PRINCIPENSE

ARAGONESE
GALICIAN
JUDEO ESPAGNOL
SPANISH
ZAMBOANGUENO ] Cre.SPANISH BASED
PAPIAMENTO ] Cre.SPANISH BASED
ROMANSH SURSILVAN ] IE.ROMANCE

IE.ROMANCE

CAPE VERDEAN CREOLE
PAPIA KRISTANG
KORLAI
PORTUGUESE ] IE.ROMANCE

Cre.PORTUGUESE BASED

CATALAN
FRIULIAN
ITALIAN
SICILIAN UnnamedInSource
ROMANIAN
ROMANIAN 2

IE.ROMANCE

GAELIC SCOTTISH
IRISH GAELIC
MANX
BRETON
WELSH

IE.CELTIC

BOSNIAN
CROATIAN
SERBOCROATIAN
SLOVENIAN
BULGARIAN
MACEDONIAN
CZECH
SLOVAK
POLISH
UPPER SORBIAN
LOWER SORBIAN
LOWER SORBIAN 2
RUSSIAN
BELARUSIAN
NINILCHIK RUSSIAN
UKRAINIAN

IE.SLAVIC

LATVIAN
LITHUANIAN
] IE.BALTIC

EASTERN ARMENIAN
WESTERN ARMENIAN
] IE.ARMENIAN

BAFFA PASHTO
OGHI PASHTO
BATAGRAM PASHTO
CHARSADDA PASHTO
SWABI PASHTO
DIR PASHTO
MADYAN PASHTO
MINGORA PASHTO
MARDAN PASHTO
PESHAWAR PASHTO
NORTHERN PASHTO
BAJAUR PASHTO
MOHMAND PASHTO
CHERAT PASHTO
JALLOZAI PASHTO
BAR PASHTO
MALLAGORI PASHTO
PARACHINAR PASHTO
THAL PASHTO
HANGU PASHTO
ZAKHA KHEL AFRIDI PASHTO
JAMRUD AFRIDI PASHTO
TIRAH AFRIDI PASHTO
NINGRAHAR PASHTO
SHINWARI PASHTO

NINGRAHAR PASHTO
SHINWARI PASHTO
MIRAN SHAH PASHTO
WANA PASHTO
BANNU PASHTO
KARAK PASHTO
LAKKI MARWAT PASHTO
CHAMAN PASHTO
QUETTA PASHTO
PASHIN KAKARI PASHTO
KANDAHAR PASHTO
PISHIN PASHTO
WANECI
ORMURI
EASTERN FARSI
TAJIK
PERSIAN
JUDEO TAT
SIVEREK
KURDISH KURMANJI
MUNJANI
YIDGHA
SHUGHNI
YAGHNOBI
ISHKOMAN WAKHI
YASIN WAKHI
CHAPURSAN WAKHI
CENTRAL GOJAL WAKHI
SHIMSAL WAKHI
DIGOR OSSETIAN
IRON OSSETIAN

IE.IRANIAN

SINHALA
VEDDA
KASHMIRI

IE.INDIC

GARAM CHISHMA KHOWAR
KESU KHOWAR
ODIR KHOWAR
USHU KHOWAR
PARGAM NISAR KHOWAR
CHATORKHAND KHOWAR
EASTERN KATIVIRI
SHEKHANI
LADAKHI

IE.INDIC

ST.BODIC

HUNGARIAN VEND ROMANI
ROMUNGRO ROMANI
GURVARI ROMANI
LATVIAN ROMANI
EAST SLOVAK ROMANI
CRIMEAN ROMANI
URSARI ROMANI
LOVARA ROMANI
BANATISKI GURBET ROMANI
SREMSKI GURBET ROMANI
GURBET ROMANI
MACEDONIAN DZAMBAZI ROMANIAN
SINTI ROMANI
VLAX ROMANI
SEPECIDES ROMANI
BUGURDZI ROMANI
SOFIA ERLI ROMANI
MACEDONIAN ARLI ROMANI
KALDERAS ROMANI
KOSOVO ARLI ROMANI
WELSH ROMANI
LITHUANIAN ROMANI
NORTH RUSSIAN ROMANI
BURGENLAND ROMANI
SELICE ROMANI
DOLENJSKI ROMANI
FINNISH ROMANI
ANGLOROMANI
NORWEGIAN ROMANI
SWEDISH ROMANI

IE.INDIC

Cre.ROMANI BASED

Cre.ROMANI BASED

ASHRET PHALURA
BIORI PHALURA
PURIGAL PHALURA
SAVI
KOHISTANI INDUS
KALAMI
KALKOTI

KALAMI
KALKOTI
TORWALI
USHOJO
DOMAAKI
GAWAR BATI
DAMELI
DAMELI 2
ZUGUNUK KALASHA
GURU KALASHA
KRAKAL KALASHA
BENGALI
CHAKMA UnnamedInSource
MARATHI
GUJARATI
MAITHILI
HINDI
NEPALI
LAMANI
AGRA GUJARI
URDU
PESHAWAR CITY HINDKO
PAKHA GOLAM
TALAGANG
WAD PAGGA
ATTOCK
JAMMUN
SINGO DI GARI
ATTOCK CITY HINDKO
TALAGANG HINDKO
PAKHA GOLAM HINDKO
WAD PAGGA HINDKO
KOHAT CITY HINDKO
SHERPUR HINDKO
BALAKOT HINDKO
MANSEHRA HINDKO
SINGO DI GARHI HINDKO
SERAIKI
MENDHAR GUJARI
SOUTHERN AZAD KASHMIR GUJARI
CENTRAL AZAD KASHMIR GUJARI
NORTHERN AZAD KASHMIR GUJARI
SOUTHERN HAZARA GUJARI
KAGHAN GUJARI
SETTLED SWAT GUJARI
TRANSHUMANT SWAT GUJARI
GILGIT GUJARI
DIR GUJARI
CHITRAL GUJARI
KUNAR GUJARI

IE.INDIC

BORANA OROMO
ORMA
W OROMO
MECHA OROMO
BORANA OROMO 2
EASTERN OROMO
BUSSA 2
KOMSO
KOMSO 2
BUSSA
GIDOLE
GIDOLE 2
BAISO
BAISO 2
RENDILLE
RENDILLE 2
SOMALI
SOMALI 2
AFAR
AFAR 2
SAHO
SAHO 2
BURJI
BURJI 2
HADIYYA 2
LIBIDO
HADIYYA
KAMBAATA
KAMBAATA 2
ALABA
GEDEO

AA.EASTERN CUSHITIC

ALABA
GEDEO
GEDEO 2
SIDAMO
SIDAMO 2
TSAMAI
TSAMAI 2
GAWWADA
GOBEZE
GAWWADA 2
WERIZE
BIRAYLE
ONGOTA ] AA.BIRAYLE
DAHALO ] AA.SOUTHERN CUSHITIC
DAASANACH
DAASANACH 2
EL MOLO
ARBORE
ARBORE 2 ] AA.EASTERN CUSHITIC
AWNGI
AWNGI 2
KEMANT 2
XAMTANGA
XAMTANGA 2
BILIN 2
BILIN
KEMANT ] AA.CENTRAL CUSHITIC
KWADZA ] AA.SOUTHERN CUSHITIC
MBABARAM ] Aus.PAMA-NYUNGAN
KAMBOT/KAMBARAMBA
KAMBOT ] LSR.BOTIN
MOGOGODO ] AA.EASTERN CUSHITIC
INNER MBUGU BUMBULI ] Cre.CUSHITIC BASED
ALAGWA
BURUNGE
IRAQW
IRAQW 2 ] AA.SOUTHERN CUSHITIC
DUVLE
FOAU ] LP.DUVLE-FOAU
OFAYE ] MGe.OPAYE
WETAWIT ] NS.BERTA
BEJA ] AA.BEJA
RASAWA
SAPONI
AWERA ] LP.COASTAL
FAIA
KIRIKIRI/FAIA
KIRIKIRI
EDOPI
IAU
FAYU/SEHUDATE
FAYU
TAUSE/WEIRATE
WEIRATE
TAUSE
DEIRATE
TAUSE/DEIRATE
AIKWAKAI/SIKARITAI
SIKARITAI
PAPASENA
WARITAI
DOUTAI
BIRITAI
ERITAI
OBOKUITAI
OBUKUITAI ] LP.TARITU
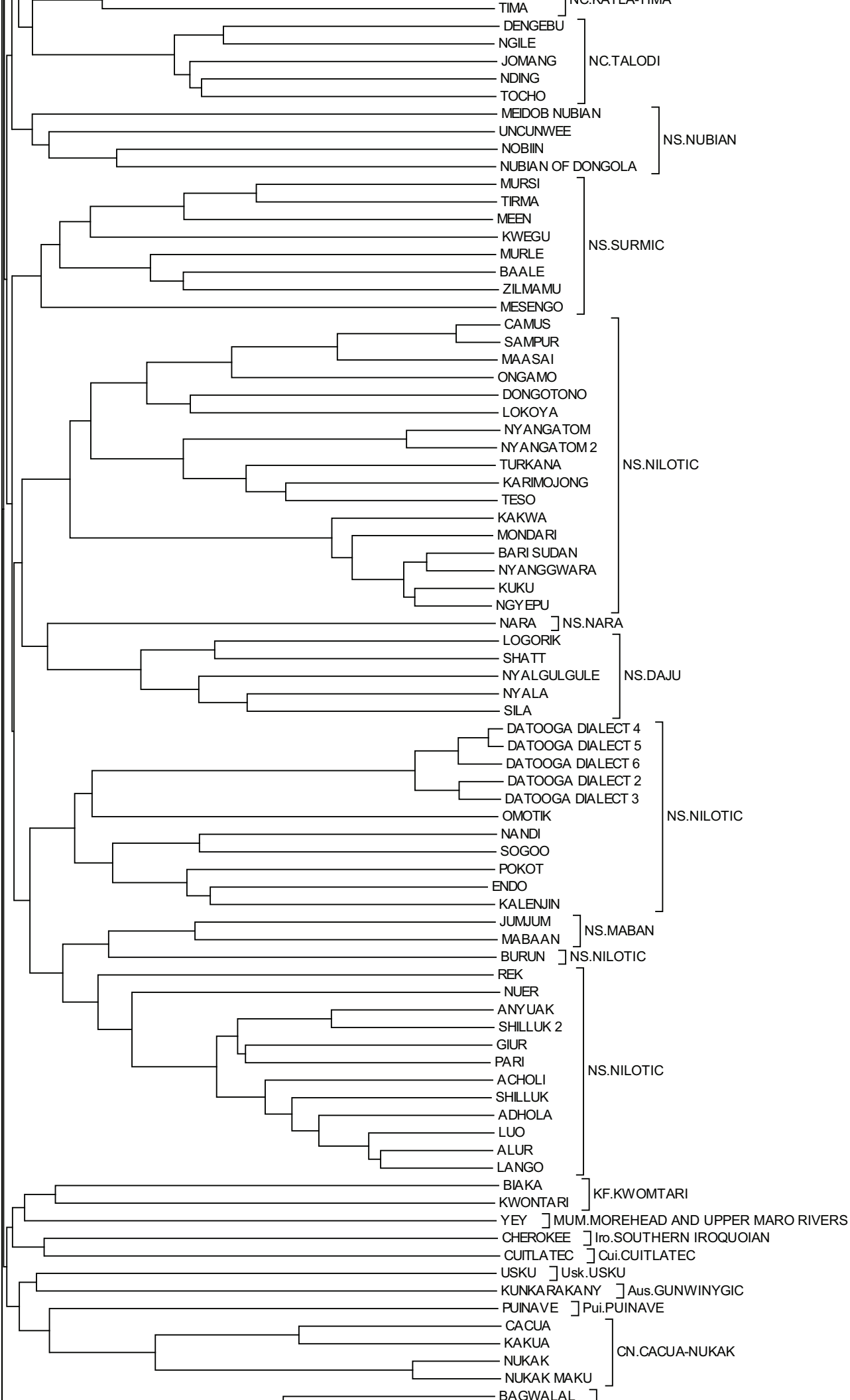ABUL
EBANG
UTORO
LARU SUDAN
RERE
SHIRUMBA
LOGOL
MORO
TIRO
KO1
WARNANG ] NC.HEIBAN
KATLA
TIMA ] NC.KATLA-TIMA
DENGEBU

TIMA　　　　NC.KATLA-TIMA

DENGEBU
NGILE
JOMANG　　　NC.TALODI
NDING
TOCHO

MEIDOB NUBIAN
UNCUNWEE
NOBIIN　　　　NS.NUBIAN
NUBIAN OF DONGOLA

MURSI
TIRMA
MEEN
KWEGU
MURLE　　　　NS.SURMIC
BAALE
ZILMAMU
MESENGO

CAMUS
SAMPUR
MAASAI
ONGAMO
DONGOTONO
LOKOYA
NYANGATOM
NYANGATOM 2
TURKANA　　　NS.NILOTIC
KARIMOJONG
TESO
KAKWA
MONDARI
BARI SUDAN
NYANGGWARA
KUKU
NGYEPU

NARA　　　NS.NARA

LOGORIK
SHATT
NYALGULGULE　　NS.DAJU
NYALA
SILA

DATOOGA DIALECT 4
DATOOGA DIALECT 5
DATOOGA DIALECT 6
DATOOGA DIALECT 2
DATOOGA DIALECT 3
OMOTIK
NANDI　　　　NS.NILOTIC
SOGOO
POKOT
ENDO
KALENJIN

JUMJUM
MABAAN　　　NS.MABAN
BURUN　　　NS.NILOTIC

REK
NUER
ANYUAK
SHILLUK 2
GIUR
PARI
ACHOLI
SHILLUK　　　NS.NILOTIC
ADHOLA
LUO
ALUR
LANGO

BIAKA
KWONTARI　　　KF.KWOMTARI

YEY　　　MUM.MOREHEAD AND UPPER MARO RIVERS

CHEROKEE　　　Iro.SOUTHERN IROQUOIAN

CUITLATEC　　　Cui.CUITLATEC

USKU　　　Usk.USKU

KUNKARAKANY　　　Aus.GUNWINYGIC

PUINAVE　　　Pui.PUINAVE

CACUA
KAKUA
NUKAK　　　CN.CACUA-NUKAK
NUKAK MAKU

BAGWALAL

NUKAK MAKU

BAGWALAL
TINDI
KARATA
BOTLIKH
GODOBERI
ANDI                    NDa.AVAR-ANDIC-TSEZIC
CHAMALAL
AKHVAKH
AVAR
DARGWA          NDa.LAK-DARGWA
LAK
ARCHI
ARCHI 2
UDI
KHINALUG
RUTUL
TSAKHUR                 NDa.LEZGIC
BUDUKH
KRYZ
LEZGI
AGUL
TABASARAN
BATS
CHECHEN          NDa.NAKH
INGUSH
BEZHTA
BEZHTA 2
HUNZIB
HINUKH                  NDa.AVAR-ANDIC-TSEZIC
TSEZ
INKHOKWARI
KHWARSHI
BAIBAI           KF.BAIBAI
FAS
GUATO           MGe.GUATO
DUKA EKOR/MONGGOWAR
MONGGOR
DERA/MUNGGUAFI            Sen.SENAGI
AMGOTRO
DERA/AMGOTRO
ALABAMA
KOASATI
MIKASUKI                Mus.MUSKOGEAN
CREEK
MOBILIAN JARGON       Cre.MUSKOGEAN BASED
CHICKASAW
CHOCTAW              Mus.MUSKOGEAN
CHIQUITANO         Chq.CHIQUITO
BIKARU          Sep.SEPIK HILL
URAT           Tor.WAPEI-PALEI
SOUGB
MENINGGO           EBH.EAST BIRDS HEAD
MEYAH
ABAZA
ABKHAZ
UBYKH               NWC.NORTHWEST CAUCASIAN
ADYGHE
KABARDIAN
COAST TSIMSHIAN        Pen.TSIMSHIANIC
KRENAK          MGe.BOTOCUDO
MAXAKALI          MGe.MAXAKALI
PATAXO          Pat.PATAXO
SHIRISHANA
YANAM 2
YANAM
SANIMA
SANUMA               Yan.YANOMAM
YANOMAMI
YANOMAM
YANOMAME
NAKWI
NIMO/NAKWI
NIMO
AMA                 LeM.LEFT MAY
BO
ROCKY PEAK
OWINIGA
AMTO
MUSAN           AM.AMTO-MUSAN

AMTO
MUSAN — AM.AMTO-MUSAN
MUSIAN

BUSA PAPUANG ] Odi.ODIAI
DEMTA/MURIS
DEMTA UNKNOWNDIAL
DEMTA/AMBORA — Snt.SENTANI
DEMTA

KAINGANG
XOKLENG — MGe.GE-KAINGANG
ARIKAPU
JABUTI — Jab.JABUTI

KOLYMA YUKAGIR
YUKAGHIR TUNDRA — Yka.YUKAGHIR
ZUNI ] Zun.ZUNI
SALINAN ] Sln.SALINAN
ACHUMAWI
ATSUGEWI — Hok.PALAIHNIHAN
GUACHI ] Un.UNKNOWN
TRUMAI ] Tru.TRUMAI

NOWOLOKAKAN
SIENKOKAKAN
WOJENEKAKAN
BODUGUKAKAN
FOLOKAKAN
GBELEBANKAKAN
JULA VEHICULAIRE
TUDUGUKAKAN
VANDUGUKAKAN
SIAKAKAN
KANIKAKAN
KARANJANKAN
JULA DE KONG
KOYAGAKAN
SAGAKAKAN
KOROKAN
NIGBIKAN
BARALAKAKAN
FINANGAKAN
MAUKAKAN
WORODUGUKAKAN
KOROKAKAN
TENENGAKAN
KONO
VAI
KURANKO
LELE GUINEA
MARKA
BAMBARA
MALINKE
MANDINKA
XAASONGAXANGO
LIGBI
SUSU
YALUNKA
SEEKU
BANKA
DUUNGOMA
SONINKE
BOZO HAINYAHO
BOZO JENAAMA
BOZO TIEYAHO — NC.WESTERN MANDE

WAN
WAN 2
BISA
BOKOBARU
BUSA NIGERIA
BOKO
ILLO BUSA
KYENGA
SHANGA
BENG
BENG 2
TOURA
YAKOUBA
GURO
YAURE
MANN
MWAN — NC.EASTERN MANDE

BANGI ME ] Ban.BANGI ME

MMAN
BANGI ME ⎤ Ban.BANGI ME
FULA
PULAR
FULFULDE MAASINA
SERER SINE
NC.NORTHERN ATLANTIC

NYAMBEENGGE 1
NYAMBEENGGE 2
BUNOGE
AMPARI PA
ANA TINGA
YANDA
TEBUL URE
WALO KUMBE
DOGON
DOGON JAMSAY
NC.DOGON

SVAN
GEORGIAN
LAZ
MINGRELIAN
Krt.KARTVELIAN

EMBERA CHAMI
EMBERA TADO
EPENA BASURUDO
EPENA SAIJA
CATIO
EMBERA DARIEN
NORTHERN EMBERA
WOUNAAN
Cho.CHOCO

ENGGANO
BANJAR SARI ENGGANO
MALAKONI ENGGANO
An.SUMATRA

CROW
HIDATSA
MANDAN
BILOXI
OFO
TUTELO
ASSINIBOINE
LAKHOTA
LAKOTA
WINNEBAGO
IOWA-OTO
OMAHA-PONCA
KANSA
OSAGE
QUAPAW
Sio.SIOUAN

FINNISH
KARELIAN
ESTONIAN VORO
ESTONIAN
VEPS
LIVE
Ura.FINNIC

HELSINKI STADIN SLANGI ⎤ Cre.FINNISH BASED

MEADOW MARI
UDMURT
KOMI PERMYAK
KOMI ZYRIAN
Ura.FINNIC

ERZYA
MOKSCHA MORDWINISCH
SOUTH SAAMI
KILDIN SAAMI
NORTH SAAMI
LULE SAAMI
INARI SAAMI
SKOLT SAAMI
Ura.FINNIC

CSANGO
HUNGARIAN
KHANTY
MANSI
Ura.UGRIC

NENETS
SELKUP
Ura.SAMOYEDIC

TEDAGA
KANURI
MANGA
NS.WESTERN SAHARAN

KIBET
RUNGA
MASALIT
MABA CHAD
MIMI
NS.MABAN

CHIAPANEC
CHOROTEGA
OM.CHIAPANEC-MANGUE

CHIAPANEC
CHOROTEGA ] OM.CHIAPANEC-MANGUE

KOPAR/SINGARIN
KOPAR
MURIK/KARAU
MURIK
ANGORAM/KAMBRINDO } LSR.LOWER SEPIK
ANGORAM
YIMAS
CHAMBRI/KILIMBIT
CHAMBRI

JUHOAN
KUNG EKOKA } Kho.NORTHERN KHOISAN
OUNG

BEJOND
NGAMBAY
MBAY
SAR CHAD
NA
NDOKA
GULA MERE
GULA SARA
GULA ZURA
KABA DEME SARA
KABBA
BAGIRMI } NS.BONGO-BAGIRMI
KENGA
FER
KARA 2
YULU
BAKA SUDAN
BONGO
FORMONA
SINYAR

BALESE
MAMVU ] NS.MANGBUTU-EFE

MANGBETU
MANGBETU 2 ] NS.MANGBETU

LENDU DJADHA
LENDU TADHA
LENDU NJAWDHA
LENDU DDRADHA } NS.LENDU
LENDU PIDHA
LENDU
NGITI

LOKAI
MAADI
LULUBA
PANDIKERI
KELIKO
LUGBARA
LOGO
OJIGA
OJILA } NS.MORU-MAADI
AGI
WADI
ANDRI
BALIMBA
MIZA
KEDIRU
LAKAMADI

KALI
NDOK MBALI
NGOUMI
NJAK MBAI
KO
PAM } NC.ADAMAWA
DAMA
GALKE
MBUM
MUNDANG
TUPURI

RIKBAKTSA ] MGe.RIKBAKTSA

JAPANESE KYOTO
JAPANESE
JAPANESE 2
NORTHERN AMAMI OSHIMA } Jap.JAPANESE
YONAMINE
NAHA
SHURI

NAHA
SHURI
AIKANA ] Aik.AIKANA
KWAZA ] Kwz.KWAZA
COFAN ] Cof.COFAN
HUMENE/MANUGORO
HUMENE
KWALE
Kwa.KWALEAN
TULISHI
KURONDI
FAMA
KRONGO
KEIGA
KAMDANG
KANGA
CHIRORO
KATCHA
KADUGLI 2
MIRI
Kad.KADUGLI
AGOB/DABU
DABU
AGOB/AGOB
KAWAM
DIBOLUG
AGOB/BUGI
MUM.MOREHEAD AND UPPER MARO RIVERS
BACAMA ] AA.BIU-MANDARA
TAKELMA ] Tak.TAKELMA
BAGA BINARI
LANDOMA
TEMNE
BAGA MADURI
NC.SOUTHERN ATLANTIC
YUPIK SIRENIK
ST LAWRENCE YUPIK
YUPIK CENTRAL SIBERIAN CHAPLINO
C YUPIK
WEST GREENLANDIC
INUPIAQ
KANGIRYUARMIUTUN
INUKTITUT EASTERN CANADIAN
EA.ESKIMO
APOI
BASAN
EAST OLODIAMA
BUMO
OYAKIRI
GBARAIN
EKPETIAMA
KOLOKUMA
OPOROMO
OGBOIN
EAST TARAKIRI
IKIBIRI
OGBE IJO
OPEREMO
WEST TARAKIRI
KABOU
KUNBO
IDUWINI
OGULAGHA
GBARANMATU
AROGBO
FURUPAGHA
IZON
MEIN
BISENI
AKITA
ORUMA
NKORO
AKAHA
NEMBE
IBANI
KALABARI
OKRIKA
DEFAKA
NC.IJOID
PAEZ ] Pae.PAEZAN
YARURO ] Yrr.YARURO
TEGEM ] NC.TEGEM
TIBEA ] NC.BANTOID
ISAKA ] Sko.KRISA
PAASAAL
VAGALA

KRUMAN
GODIE
BIJOGO ] NC.BIJAGO
MAMARA SENOUFO
SENOUFO SUPYIRE
NAFAARA
CEFO
NC.GUR
MBRE ] NC.unclassified
KONKOMBA
KONKOMBA 2
BASSARI
MOBA
YOM
DAGBANI
MAMPRULI
HANGA
DAGAARE
KUSAL
MOORE
FRAFRA
NINKARE
NC.GUR
IIGAU
ISHEU
IYINNO
IKAAN
IKAKUMO
UKAAN
NC.UKAAN
AKUNNU ] NC.AKPES
BEKWARRA ] NC.CROSS RIVER
IDOMA ] NC.IDOMOID
OKO OSANYE ENI ] NC.OKO
IGEDE ] NC.IDOMOID
IGBIRRA ] NC.NUPOID
GBARI
NUPE
NC.NUPOID
GADE ] NC.NUPOID
ELOYI ] NC.PLATOID
JIJILI
JILI
NC.PLATOID
KANA
TEE
BAN OGOI
ABUA
OGBRONUAGUM
KOHUMONO
MBEMBE
NC.CROSS RIVER
OKA
YORUBA
AYERE
ARIGIDI
NC.DEFOID
DEGEMA
ENGENNI
DEGEMA 2
DEGEMA 3
EGENE
EPIE 2
EPIE
ISOKO
URHOBO 2
ERUWA
OKPE
UVBIE
URHOBO
EDO
EMAI
OKPAMHERI
EHUEUN
UKUE
EMHALHE
UHAMI
IBILO
AOMA
EDO 2
GHOTUO
UNEME
AUCHI
AVBIANWU
NC.EDOID
CINDA 5
KUKI
REGI

KURI
REGI
CINDA 1
CINDA 2
ROGO
SEGEMUK
SHAMA
MADAKA
PONGU 2
PONGU
LELA
DUKA
FAKAI A
ROR
LARU NIGERIA
LOPA
DOOHWAAYAAYO
MOMI
YANDANG
YOTI
BALI NIGERIA
KPASHAM
ADIOUKROU
OBOLO
USAKADE
EFIK
IBIBIO 1
ENWANG
UDA
EBUGHU
ORO
EFAI
ILUE
OKOBO
IBINO
IKO
EKIT
ETEBI
ANAANG
IBIBIO 2
ITUMBUSO
IBUORO
UKWA
AGHEM WEH
AGHEM WUM
BU CAMEROON
AGHEM ISU
MMEN
BABANKI
BUM B CAMEROON
KUO
KOM MBIZINAKU
KOM
BABUNGO
KENSWEINSEI
LAMNSOQ
NDEMLI
AMASI
AMBELE
BALEP
EKPARABONG
BENDEGHE
EJAGHAM
NDE
NTA
NSELLE
EFUTOP
EKAJUK
NNAM
NKIM
NKUMM
MODELE
ESIMBI
IPULO
TIV 1
TIV 2
ALINGA
TUNEN
BONEK
MANDI CAMEROON

NC.KAINJI

NC.ADAMAWA

NC.KWA

NC.CROSS RIVER

NC.BANTOID

MANDI CAMEROON
NYOKON
YAMBETA
NUGUNU
NUBACA
LIBIE
KALONG
MMALA
ABAR
MISSONG
MBU
BU
MUNDABLI
KOSHIN
FANG
CUNG
KEMEZUNG
MASHI
NSARI
MUNGONG
NCANE
NONI
DONG
MBE
BUJEBA
BAS KENYANG
HAUT KENYANG
CENTRAL KENYANG
KENYANG/KITWII
KENDEM
DENYA/BASHO
DENYA/BAJWO
DENYA/BITIEKU
DENYA/TAKAMANDA
KONJA NDUNG
KONJA SUNDANI
MAMBILA ATTA
MAMBILA
NIZAA
WAWA                    NC.BANTOID
VUTE MBANJO
VUTE YOKO
AFI AMANDA
KAMINO
BATU ANWE
BURU
GBUGYAR
RIJA
KEJA
NDEYWAN
NCO
NJIGBAN
MADA
NCEKPE
RINZE
NUNKU
NINGYE
NINKA
ANIB
NINKYOP
BU NIGERIA
CE
KWANKA
SHALL
HASHA
SAMBE
TESU
TORO
JIBU
AKE                     NC.PLATOID
EGGON
KUTEB
GWARA
ITOO
IDU
NYANKPA
HYAM
KULU
PE
YANGKAM

PE
YANGKAM
TAPSHIN
AYU
NDUN
FIRAN
GANANG
IZERE
ITEN
BEROM F
TAHOSS
CARA
RUKUL
FYEM
HOROM
KULUNG NIGERIA
MBULA NIGERIA
BILE
LABIR
ZAAMBO
BANKALA
BADA
DUGURI

NC.BANTOID

LIBOBI
LIFONGA
LIBOBI 2
LIFONGA 2
BOMBOMA
BOMBOMA 2
BOBANGI
BOBANGI 2
MABALE
MABALE 2
LIBINZA
LIBINZA 2
ZAMBA
ZAMBA 2
DOKO
DOKO 2
BUJA
BUJA 2
LINGOMBE
LINGOMBE 2
BANGALA
LINGALA
EGBUTA
LOMONGO
POVE
BAKUERI
DUALA
LEGA
TUKI
EWONDO
MENGISA
POL
FANG GABON
BASAA
BUM A CAMEROON
BAFIA
LEFA
BANKON
NKONGHO
NLA MBO
NLE MBO
BAFAW
EHO MBO
MIENGE
BAKOSSI
NNINONG
ELUNG
MANEHAS
BABONG
BANEKA
BAKAKA
BALONDO

NC.BANTOID

OBELI
OYABI
NTSIAMI
NDOUBA
NKOMO KELLE

NDOUBA
NKOMO KELLE
NKOMO OLOLI
YABA MBETI
OYUOMI TCHERRE
OBAA
OYUOMI MBAMA
TEGE EWO
MBAMBA SIBITI
MBAMBA LIWEME
MBEDE
MBEDE 2
NDUMU
NDUMU 2
BANTOU DU GABON NDJABI
DUMA
GISIRA
YIPUNU
LAADI
YAKA
KIHOLU
MBALA
URUUND

NC.BANTOID

XHOSA
ZULU
SWAZI
LESOTHO
SILOZI

NC.BANTOID

ICIBEMBA
SHIYEYI
LHUKONZO
LUGANDA
KINYARWANDA
NKORE KIGA

NC.BANTOID

MBALANHU
NDONGA
KWANYAMA
KWANGALI
NKOYA

NC.BANTOID

HEHE    NC.BANTOID

CIYAWO
MWERA
KIMATENGO
NYANJA
SHONA
KESUKUMA
KINYAMWEZI
MBUGWE
KAGULU
VALANGI

NC.BANTOID

RUNYANKORE
RUTOORO
HUNDE
HAYA
BENDE
GWERE
KIKEREWE
JITA
NILAMBA
NYAKYUSA
LUBA
KILUBA
KITABWA

NC.BANTOID

PARE TANZANIA    NC.BANTOID
NORMAL MBUGU    Cre.CUSHITIC BASED
ILWANA    NC.BANTOID
KOTI
SHIMAORE
CHONYI
GIRYAMA
RABAI
DIGO
DURUMA
CHWAKA
BUU
SWAHILI MWANI
NGHWELE
SWAHILI
SWAHILI PATE
SWAHILI TIKUU

NC.BANTOID

SWAHILI PATE
SWAHILI TIKUU
SWAHILI MAKUNDUCHI
SWAHILI PEMBA
SWAHILI MVITA
SWAHILI CHIRAZI
SWAHILI VUMBA

1000