# CLASSIFYING LALO LANGUAGES: SUBGROUPING, PHONETIC DISTANCE, AND INTELLIGIBILITY[*]

**Cathryn Yang**
SIL International

**Abstract**: Lalo is a Central Ngwi (Loloish) language cluster spoken in western Yunnan, China by fewer than 300,000 speakers. Previously, most Lalo varieties were undocumented, and Lalo was thought to have only two dialects. This paper presents the first rigorous classification of Lalo languages, based on a dialectological study conducted in eighteen Lalo villages, which included the recording of 1,001-item vocabularies, comprehension tests, and sociolinguistic interviews. Synthesizing results from diachronic subgrouping, intelligibility, and phonetic distance measured by Levenshtein distance, this paper argues that the Lalo cluster actually comprises at least seven closely related languages. Phonetic distance correlates strongly with comprehension, and the NeighborNet network based on phonetic distance concurs with the diachronic subgrouping at shallow time depths, providing further validation of the dialectometric techniques. This paper suggests that such techniques are useful in the classification of lesser-known, endangered, indigenous languages, which often have an urgent need for language maintenance efforts.

**Keywords**: Lalo, classification, subgrouping, phonetic distance, intelligibility

## 1. INTRODUCTION

Lalo, spoken in western Yunnan, China, is one of the least studied of the Central Ngwi (Loloish) language clusters, with no documentation of most Lalo varieties until this study. Chen, Bian & Li (1985: 195), in their wide-sweeping survey of many Ngwi languages, claimed that Lalo had only two dialects, without specifying details of how the dialects differed. To address this gap, a dialectological study was conducted from February 2008 to January 2009 in Dali, Baoshan, Lincang, and Pu'er prefectures. The results of this study suggest that the Lalo cluster actually has at least three major dialect clusters, all mutually non-intelligible, and four smaller languages located on the periphery of Lalo geographic distribution.

Evidence for this classification comes from three methodologies with complementary perspectives: a diachronic, qualitative perspective from the comparative method; a synchronic, quantitative perspective from dialectometry, and a perceptual perspective from comprehension testing. Pelkey (2011) developed this multi-strand approach to classification, which he terms "integrational dialectology." The Proto-Lalo phonological system and over 900 lexical items were reconstructed, and shared innovations in morphology, phonology, and lexicon were used for subgrouping. Aggregate phonetic distance between varieties was measured quantitatively with a string edit distance algorithm known as Levenshtein distance and then analysed with NeighborNet network analysis and multi-dimensional scaling (MDS). To investigate the perceptual impact of dialect differences, speakers of most varieties were tested on their comprehension of short narratives recorded in a selection of varieties. Each methodology serves as a control on the interpretation of the other methods' results, an important triangulation when documenting and classifying varieties for the first time.

This study applies recently developed dialectometric techniques to an endangered, under-documented language cluster of China. The new field of dialectometry, pioneered by Séguy (1971) and Goebl (1982, 2006), uses quantitative, aggregate analysis of linguistic variation in search of new ways of understanding the relation between linguistic variation and explanatory factors (Nerbonne & Kretzschmar 2006). Levenshtein distance is one of a growing number of dialectometric techniques, including analysis of variation at various linguistic levels (e.g., syntax, phonology, etc.)  Only a handful of studies have applied Levenshtein distance to languages of the Sinosphere: Tang 2009: 117-136 on Chinese, Yang & Castro 2008 on Bai and Zhuang, Yang 2009 on Nisu, and Stanford in prep on Sui. Pelkey (2011: 281-285, 353-355) uses NeighborNet network analysis based on lexical distance in his classification of the Phula languages, with mixed results. Ben Hamed (2005) and Ben Hamed & Wang (2006) find that NeighborNet networks for Chinese languages correlate with Chinese history and societal trends. This study finds that Levenshtein distance feeding into NeighborNet and MDS renders classifications that mostly agree with the historical subgrouping. These results, along with the strong correlation found between Levenshtein distance and intelligibility, are in line with previous studies that validate Levenshtein distance as a dialectometric technique (e.g. Heeringa and Gooskens 2004; Heeringa 2004; Gooskens 2006). Many indigenous languages of East Asia are under-documented and endangered, and such tools are especially appropriate because of the urgent need for language maintenance and planning efforts. Overall phonetic similarity, even if the similarity is a result of contact, drift, or retentions, is helpful information when making language planning decisions.

Based on the shared innovations given in Section 4, Lalo is divided into three distinct clusters, Eastern (E), Western (W), and Central (C), each of which fulfills criteria for consideration as an independent language. These clusters comprise the Core Lalo group and are mostly located in the traditional Lalo homeland of

southern Dali Prefecture. Four peripherally located Lalo languages, Eka, Mangdi, Yangliu and Xuzhang, represent different waves of migration out of the Lalo homeland. Comprehension tests show that comprehension between the Core languages, and between peripheral and Core languages, is low, while intra-cluster comprehension for E and C clusters is high. Groupings based on phonetic distance corroborate the historical subgrouping for the most part, while also displaying the effects of significant contact between certain varieties.

## 2. BACKGROUND

Lalo's phylogenetic pedigree is Tibeto-Burman, Ngwi-Burmese (Lolo-Burmese), Ngwi (Loloish), Central Ngwi (Bradley 2002). Lalo belongs to the Yi ethnic minority, and Chinese linguists classify Lalo as the "Western Yi dialect" due to Lalo's distribution in western Yunnan. Lalo has less than 300,000 speakers, an estimate based on Chinese demographic sources and this study's ethnolinguistic vitality research. The Central cluster has approximately 213,000 speakers; West, 44,000; East, 15,000; Yangliu, 7,000; Eka, 3,000; Mangdi, 3,000; and Xuzhang, 2,000. All Lalo languages show a reflex of the Proto-Lalo autonym $*la^2lo^Hpa^L$, e.g. W-YL [$la^{21}lo^{33}pa^{53}$] and C-LJ [$la̠^{21}lo̠^{33}pa̠^{21}$]. Eka speakers' autonym is now [$o^{21}k^ha^{24}$], but elder speakers remember a time when they called themselves [$la^{21}lu̠^{33}po^{21}$]. Other diagnostic criteria for membership in the Lalo cluster are given in Section 4. Lalo speakers are mostly located in southern Dali Prefecture, especially Weishan County, considered the traditional homeland of the Lalo. Historically, this area is the home of the Meng clan, who ruled the Nanzhao Kingdom (737-902 A.D.). Many Core Lalo claim to be descendants of the Meng clan. Figure 1 below is a map of the area where most Lalo are located.

Other Central Ngwi languages such as Lahu and Lisu have been the focus of extensive documentation and dialectological research, e.g. Bradley 1979a, 1997; Matisoff 1973 [1982], 1988. For Lalo, however, there has been no rigorous attempt at classification, as few Lalo varieties have been documented, with the notable exception of Björverud's (1998) grammar of the Central Lalo variety in Longjie Village, Weishan County. Chen et al. (1985: 198), based on linguistic data from the 1950's, describe Lalo as having two dialects, East Mountain and West Mountain, divided by the valley that bisects Weishan County into eastern and western halves. However, they give no linguistic justification for the division. Later large scale surveys of Yi languages (Wang 2003; Zhu 2005) echo Chen et al. in applying the East-West dichotomy to Lalo as a whole. "East Mountain" served as a category for any variety that was not West Mountain, including groups ranging from Midu County in the east to as far west as Baoshan (Zhu 2005).

However, the labels "East Mountain" and "West Mountain" only have relevance within Weishan County. The East Mountain label, as used by previous researchers, conflated the true East Mountain in northeast Weishan with the Eastern Lalo dialect cluster spoken in Dali Municipality. The incorrect perception that Lalo speakers in those areas belonged to the same dialect group was

probably due to their geographic proximity and similar traditional costume, as well as the lack of linguistic data. This article demonstrates that East and West Mountain Lalo both belong to the Central Lalo subgroup, which in turn is only one of several Lalo subgroups.



*Figure 1. Map of Yunnan, China. Most Lalo live in Dali, Baoshan, Lincang and Pu'er prefectures.*

## 3. DATA COLLECTION

Fieldwork conducted during 2008 entailed data collection in eighteen Lalo villages in nine counties. Field sites were chosen in consultation with Lalo speakers with the goal of documenting the widest range of dialect diversity. Table 1 gives locations and abbreviations for the chosen field sites, grouped together based on the proposed classification. East Mountain and West Mountain Lalo are represented by one variety each, CE-YA and CW-QY; other Central varieties are not labeled as West Mountain, since those speakers do not use that loconym. Figure 2 shows a map of the locations of the 18 Lalo villages.

Both the comparative method and phonetic distance analysis require lexical data from the various varieties being compared. A 1,001-item wordlist was used, adapted from Pelkey 2008. Chinese glosses and photos were arranged in semantic categories, and a group of two to three bilingual speakers were asked to translate from Chinese into Lalo. A Lalo speaker from CE-YA conducted all recording sessions and identified any mis-translations. Selected participants were all fluent native speakers who had grown up in the village and had at least one parent from that village. Participation was voluntary, and participants were compensated for

their time. Participants were mostly male (36 out of a total 47), so further research on female speech in Lalo is needed.

| Group | Prefecture | County | Township | Village | Abbreviation |
|---|---|---|---|---|---|
| Central  (East Mt) | Dali | Weishan | Yongjian | Yong'an | CE-YA |
| Central (West Mt) | Dali | Weishan | Ma'anshan | Qingyun | CW-QY |
| | Dali | Weishan | Wuyin | Longjie | C-LJ |
| | Dali | Yangbi | Wachang | Wachang | C-WC |
| Central | Dali | Yongping | Shuixie | Leba | C-LB |
| | Dali | Nanjian | Xiaowandong | Chajiang | C-CJ |
| | Pu'er | Jingdong | Anding | Qingsheng | C-QS |
| | Dali | Dali | Shijiao qu | Diaocao | E-DC |
| Eastern | Dali | Dali | Fengyi | Houshan | E-HS |
| | Dali | Dali | Taiyi | Taoshu | E-TS |
| | Dali | Yangbi | Taiping | Dutian | W-DT |
| Western | Dali | Yangbi | Longtan | Shuizhuping | W-SZP |
| | Baoshan | Longyang | Wama | Shanglizhuo | W-SLZ |
| | Dali | Yongping | Changjie | Yilu | W-YL |
| Xuzhang | Baoshan | Longyang | Wafang | Xuzhang | XZ |
| Yangliu | Baoshan | Longyang | Yangliu | Yangliu | YL |
| Eka | Lincang | Shuangjiang | Heliu | Yijiacun | Eka |
| Mangdi | Lincang | Gengma | Hepai | Mangdi | MD |

*Table 1. Locations and abbreviations of Lalo data points.*



*Figure 2. Map of Lalo data points in western Yunnan.*

## 4. DIACHRONIC SUBGROUPING BASED ON SHARED INNOVATIONS

After the recordings were transcribed into a spreadsheet, the comparative method was used to reconstruct the Proto-Lalo (PLa) phonological system and over 900 lexical items. For the full reconstruction, see Yang (2010: 99-168). Reconstructions of Lalo's ancestor languages, as in Bradley's (1979) reconstruction of Proto-Ngwi and Matisoff's (2003) reconstruction of Proto-Ngwi-Burmese and Proto-Tibeto-Burman, were invaluable aids in tracking Lalo's diachronic development. The systematic reconstruction of PLa was foundational in determining if a feature shared between daughter varieties was a retention from PLa or a shared innovation of that particular Lalo subgroup. Criteria for member-ship in the Lalo cluster is given in Table 4 below.

The PLa syllable template is *(C)V(ŋ)T, in which one of 47 optional initials is followed by one of 9 obligatory rhymes and one of five tones, two of which have harsh phonation. The rhyme system includes eight monophthongs and one nasal final rhyme *aŋ. PLa is mostly monosyllabic with some disyllabic compound words, unlike its ancestor Proto-Ngwi, which is sesquisyllabic (i.e. weak syllable + strong syllable structure [Matisoff 1973]). Table 2 presents PLa's initial consonant inventory. PLa has palatalized labial and velar consonant clusters, and a labiovelar nasal cluster *ŋw. PLa distinguishes preglottalized and plain sonorants and has one preglottalized and plain fricative pair, *ʔv and *v. Preglottalized *ʔŋw is unattested, but theoretically possible given the symmetry between preglottalized and plain initials. PLa's preglottalized initials are the result of coalescence of the Proto-Ngwi *ʔə- prefix with the following sonorant. A voiceless glottal fricative *h was probably nasalized [h̃].

| *p | *pj | *t | *ts | *tʃ | *k | *kj | | |
| *pʰ | *pʰj | *tʰ | *tsʰ | *tʃʰ | *kʰ | *kʰj | | |
| *b | *bj | *d | *dz | *dʒ | *g | *gj | | |
| *f | | | *s | *ʃ | *x | | | *h |
| *v | | | *z | *ʒ | *ɣ | | | |
| *ʔv | | | | | | | | |
| *m | *mj | *n | | *ɲ | *ŋ | *ŋj | *ŋw | |
| *ʔm | *ʔmj | *ʔn | | *ʔɲ | *ʔŋ | *ʔŋj | *ʔŋw | |
| | | *l | | | | | | |
| | | *ʔl | | | | | | |

*Table 2. Proto-Lalo initial consonant inventory.*

Table 3 gives the PLa rhyme inventory. All open vowels were found in both modal and harsh phonation. Harsh (or tense) phonation, marked with an underscore on the vowel (e.g. a̱), is an aperiodic phonatory quality produced when the ventricular folds incur over the vocal folds, and the laryngeal sphincter is constricted (Edmondson & Esling 2006). PLa basically retains the Proto-Ngwi tonal categories. PLa has three levels of pitch height in non-harsh phonation (*1, high and modal, *2, low and breathy, *3 mid and modal), and two levels with harsh phonation (*H, mid; *L, low).

| *i | *y | *ɨ | *u |  |
|----|----|----|----|----|
| *e |    |    | *o |  |
| *ɛ |    |    | *a | *aŋ |

*Table 3. Proto-Lalo rhyme inventory.*

Subgrouping criteria are phonological, lexical, and morphological innovations that satisfy the following standards: (1) linguistic complexity of the innovation, (2) ecological distinctiveness of the innovation, and (3) sociohistorical plausibility (Toulmin 2009). Single sound changes are not considered strong enough evidence for subgrouping, since most types of individual sound changes can diffuse through contact (Pittayaporn 2009: 298). Rather, a set of innovations is required as evidence of shared history.
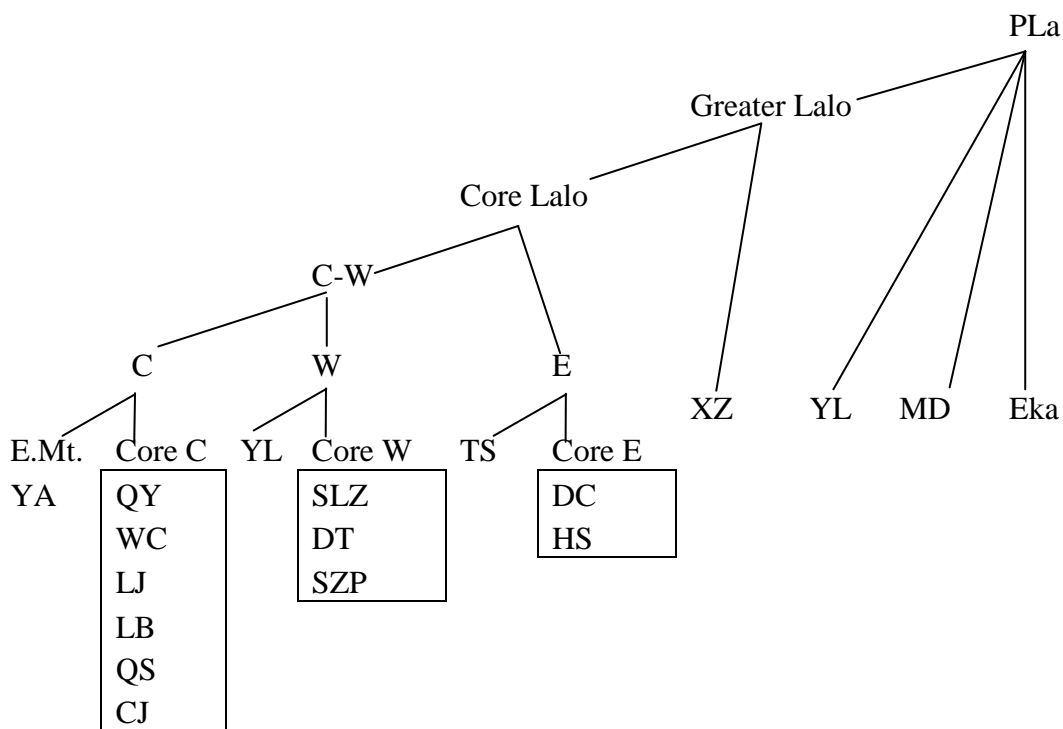


*Figure 3. Lalo phylogenetic tree.*

The proposed Lalo phylogenetic tree is seen in Figure 3 above. Eka, MD, and YL descend directly from Proto-Lalo and do not share any of the innovations that define Greater Lalo or subsequent subgroups. As noted in Section 1.1, the ancestors of Eka and MD speakers migrated out of the Lalo homeland at an early date (over 300 years ago for Eka, approximately 200 years ago for MD). While YL's migration history is unknown, its failure to share any post-Greater Lalo innovations suggests an early separation from other Lalo varieties. The Greater Lalo subgroup includes XZ and the Core Lalo varieties. Core Lalo includes the three dialect clusters E, C, and W. C and W in turn form a subgroup that excludes E. Within each Core Lalo cluster, some varieties (e.g. CE-YA) diverge from an inner core (Core E, Core W, Core C). Core Lalo varieties are all located in or

originate from southern Dali Prefecture, the probable area of origin for all Lalo, an area with the earliest historical accounts of a Lalo presence, the highest concentration of Lalo population, and the greatest linguistic diversity.

Distinct sets of innovations that define PLa and each Lalo subgroup are given in Table 4 below. Relative chronology is marked with ordinal numbers. All Lalo languages show evidence of the following innovations: Proto-Ngwi Tone *2 syllables shifted to *L when preglottalized obstruents were followed by the rhyme *-a; Proto-Ngwi *o and *u merged to *o after *b, and then *o became *wɨ after all labial stops. No other Ngwi language shows these innovations. For more details about these changes, see Yang (2010: 211-214). Eka, YL and MD descend directly from PLa and do not share in any Greater Lalo innovations. Instead, they show divergent developments that make them decidedly different from any other Lalo language. MD shows a set of rhyme innovations, Eka shows micro-splits in *1, *2, *H, and *L, and YL shows mergers of *1 and *3 to mid and *L and *H to high.

Greater Lalo (XZ plus Core Lalo) has two lexical innovations, for 'crow' and 'taro'. XZ merges *H, *L, and *3 to mid, a change that distinguishes it from Core Lalo. Core Lalo varieties share a morphologically innovative *-tsa³³ plural marker in the personal pronoun paradigm, e.g. ŋa⁵⁵ '1S' and ŋa³³tsa³³ '1PL'. C-W shows lenition of *g before low, back vowels, and a folk etymology of 'young' as 'soft in years'. C varieties shares a set of changes in rhymes, W varieties share a chain of tonal innovations, and E varieties shares a set of innovations in tones and rhymes. Within each Core Lalo cluster, later changes divide divergent varieties from the inner core. The changes listed below create complex synchronic correspondence sets that negatively impact cross-dialectal comprehension, to the extent that comprehension between the three clusters is negligible, as shown in Section 6.

| Group | Innovations |
|---|---|
| PLa | a) *2 > *L/*ʔ-obstruent + *a_ |
|  | 1) *o, *u > o/b_ |
|  | 2) *o > wɨ/labial stops_ |
| Eka | a) *1 > low-rising/default |
|  | b) splits in *1, *2, *H, *L to high level under various conditions |
| MD | a) *y > ɑ |
|  | b) *o > ɨ/velar initials_ |
|  | c) *e, *ɛ > ɛ |
| YL | a) *1 > mid/*[-voi]_, > low-rising/elsewhere |
|  | b) *L, *H > high (loss of harsh phonation) |
| Greater Lalo: Core plus XZ | a) 'crow' (n.) *a³nak$^H$ ('the black one') > *a¹ŋja$^H$bɛ$^H$ ('bird' + bɛ$^H$, possibly meaning 'untamed') |
|  | b) 'taro' *a¹tʃʰo̰$^H$ |
| XZ | *H, *L, *3 > mid (loss of harsh phonation) |
| Core Lalo: C-W-E | *-tsa³³ plural marker in personal pronoun paradigm |
| C-W | a) *g > ɣ/_*a, *aŋ, *o̰ |
|  | b) 'young in years' *tʰy²nu¹ > 'soft in years' *tʰy²nu² |
| C | a) *e > i |
|  | b) *ḛ > ɨ/[+high, -back]_ |
|  | c) *ɛ̰ > a̰/[+back]_ |
|  | d) *a > ɛ/ *C[+anterior, +strident]_$CV[-back] |
| Core C (excludes CE-YA) | a) *ɛ̰ > a̰/elsewhere |
| W | a) *L > high |
|  | b) *1/+ʔ_ > mid-high |
|  | c) *H > mid-high |
| Core W: (excludes W-YL) | a) *1 > low-rising/elsewhere |
|  | b) metathesis in 'grasshopper': *tʃɛ¹pu¹ > pɛ¹tʃu¹ |
|  | c) shared tone sandhi patterns: high > mid-high/_high |
| E | a) *H, *3 > mid (loss of harsh phonation in *H) |
|  | b) *y > ɨ/labials_ |
|  | c) *o > ɨ/labials_ |
| Core E: (excludes E-TS) | *y, *o > ɨ/elsewhere |

*Table 4: Innovations that define Lalo subgroups*

## 5. SYNCHRONIC GROUPING BASED ON PHONETIC DISTANCE

Phonetic distance between Lalo varieties is measured by applying a string edit distance algorithm called Levenshtein distance (LD), a method developed by dialectometrists at the University of Groningen (Heeringa 2004; Nerbonne 2009).

LD optimally aligns the phonetic segments of two cognates to compute the least cost of transforming one cognate into another in terms of substitutions, insertions and deletions. The LD algorithm was applied to all sets of Lalo cognates in the wordlists. Each set of cognates was compared, even if the cognate was shared by only two varieties. For example, Greater Lalo varieties reflect *a$^1$ŋja$^H$bɛ$^H$ for 'crow,' while non-Greater Lalo languages reflect *a$^3$nak$^H$, a retention from Proto-Ngwi. Reflexes of *a$^1$ŋja$^H$bɛ$^H$ formed one set of cognates, and reflexes of *a$^3$nak$^H$ formed a different set. In total, 955 sets of cognates were compared. The distance between all shared cognates was then averaged for each pair of varieties, resulting in a distance matrix.

Figure 4 below illustrates LD between the pronunciations for the word 'tiger' in CW-QY (/la$^{LE}$pa̱$^{LE}$/) and E-DC (/lɔ$^{LE}$pu$^{MF}$/). Tone is represented by onset and following contour (e.g. mid-falling as MF), a representation that Yang & Castro (2008) showed to have the strongest correlation with comprehension compared to other representations. To prevent longer words from having undue weight in the calculation of average distance, a normalization function is used in which the total cost is divided by the longest alignment, as in Gooskens & Heeringa 2004. In Figure 4, the alignment length is nine, as harsh phonation is counted as one element.

Figure 4 uses a simple phone representation, without comparisons based on features, so small and large phonetic differences are given the same weight. McMahon and McMahon (2005: 210-214) criticize the simple representation as blunt, but the procedure was externally validated through correlation with speakers' perceptions of phonetic distance, measured experimentally (Gooskens & Heeringa 2004). Heeringa et al. (2006) correlated feature-based and simple representations with speakers' perceived distance and found that the feature-based representation did not have a significantly higher correlation to speaker's perceived distance than the simple representation. Houtzagers et al. (2010) advocate the use of simple phone representation when the focus is on the aggregate distance between varieties, not on which differences are important. This type of LD also has a strong correlation with comprehension (Gooskens 2006; Yang & Castro 2008).

| Variety | 'tiger' | Operation | Cost |
|---------|---------|-----------|------|
| CW-QY | la$^{LE}$pa̱$^{LF}$ | | |
| | la$^{LE}$pa$^{LF}$ | delete ̱ (harsh phonation) | 1 |
| | la$^{LE}$pa$^{MF}$ | substitute M for L | 1 |
| | la$^{LE}$pu$^{MF}$ | substitute u for a | 1 |
| E-DC | lɔ$^{LE}$pu$^{MF}$ | substitute ɔ for a | 1 |
| | | total cost | 4 |
| | | normalized cost | 0.44 |

*Figure 4. Operations in calculating Levenshtein distance for 'tiger'.*

## 5.1 NeighborNet network analysis

The distance matrix generated by LD is then processed by the network-building program NeighborNet as well as multidimensional scaling (MDS). NeighborNet, first developed by Bryant and Moulton (2004) for use in evolutionary biology, is freely available in the SplitsTree 4 package (Huson & Bryant 2006). The use of NeighborNet here is phenetic, that is, based on varieties' overall phonetic similarity, rather than cladistic, i.e. based on historically significant shared innovations. NeighborNet phenograms (i.e. diagrams representing phenetic relations) present a synchronic snapshot of cross-varietal relations that includes all phonetic similarity, whether the similarity is a result of retentions, shared innovations significant for subgrouping, contact-induced change, or parallel developments. McMahon et al. (2007) and Maguire et al. (2010) use NeighborNet in a similar way to quantify the degree of difference between English dialects, with a focus on synchronic relationships. The phenogram given in Figure 7 below has much in common with the Lalo phylogenetic tree, but is not identical, due to inter-varietal contact and shared retentions.

As McMahon et al. (2007) note, one of the advantages NeighborNet brings is the ability to represent multiple trees in a single diagram. If there are similarities incompatible with one tree, NeighborNet still represents them through reticulated, or rectangular-like, lines that form a network. Ambiguities or mixed signals in the data are explicit, instead of being collapsed into a single line as they are with tree-building programs such as the neighbor-joining method (Saitou & Nei 1987). This advantage is illustrated through comparing Figure 5 and Figure 6 below. Figure 5 shows the collapsed tree built by the neighbor-joining algorithm, which groups C-WC and C-LB together in the upper left quadrant versus CW-QY and C-LJ in the lower right quadrant. Figure 6 shows the same basic grouping as Figure 5, but is further able to display an alternate grouping of C-WC and CW-QY versus C-LB and C-LJ. This alternate grouping is shown in Figure 6 through the shorter lines of the rectangle that push C-WC and CW-QY toward the upper right quadrant and push C-LB and C-LJ toward the lower left quadrant. In contrast with Figure 5, Figure 6 is able to show all sets of similarities between all pairs of varieties. NeighborNet will only produce a tree when the data actually fit a tree-like pattern; it is therefore an optimal way of representing dialect networks, which often have complex, partially conflicting isoglossic patterns.
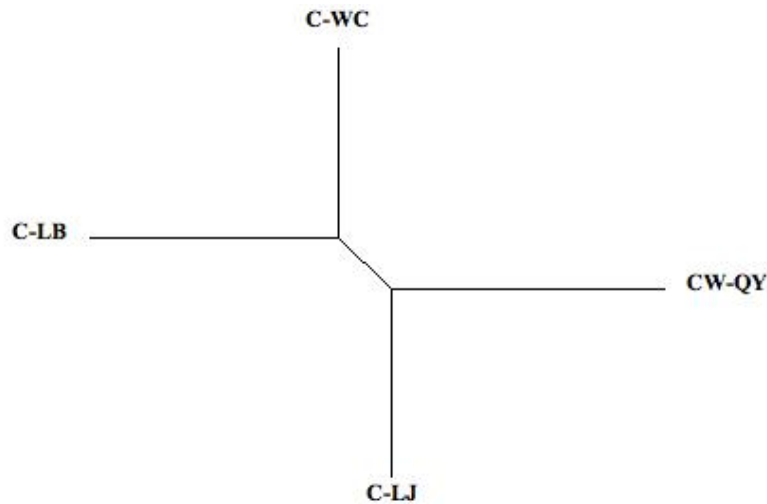
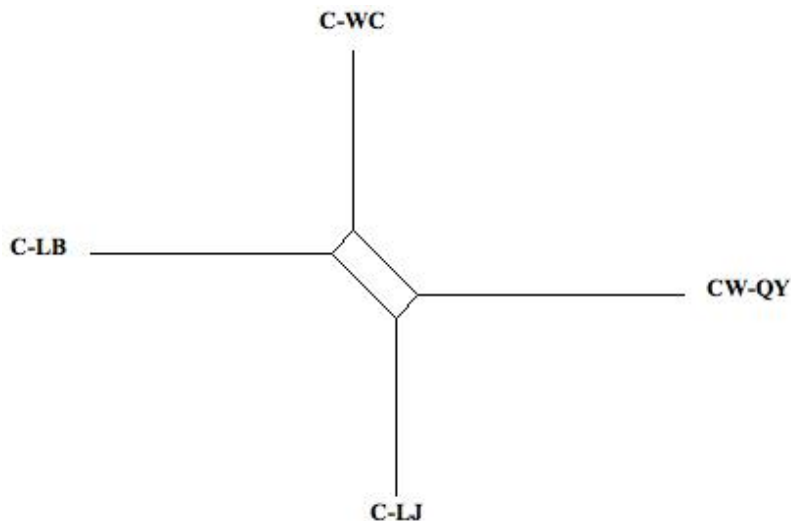*Figure 5. Neighbor-joining tree for four Central varieties.*



*Figure 6. NeighborNet reticulated rectangle for four Central varieties.*

Figure 7 below shows the NeighborNet phenogram using the equal angle method for all 18 Lalo varieties, plus Proto-Lalo. Fewer reticulated lines mean a clearer signal and thus a more clearly defined cluster. Compare, for example, the relatively narrow branches of the C or W clusters versus the many reticulated lines that pull the E varieties away from each other and towards other varieties. The lines' relative lengths depict relative difference: the longer the line, the more different the variety is from other varieties. For example, YL's relatively long line marks it as very different from all other Lalo languages, while the shorter lines connecting C varieties' depict a relatively small degree of difference within the C cluster.

The three Core Lalo dialect clusters C, W, and E are clearly identifiable in Figure 7. C varieties cluster together on the left side of the phenogram, with fewer reticulations and shorter lines. The C cluster is also closest to Proto-Lalo due its conservative nature, e.g. its retention of the Proto-Lalo tonal system. W

varieties also form a clearly defined cluster. All E Lalo varieties are found in one area of the phenogram, with reticulations connecting them to CE-YA and to XZ. In contrast to the tight bundle of C Lalo, the E cluster is much looser, but still identifiable. NeighborNet's identification of the Core Lalo clusters corroborates the lower-level diachronic subgrouping in Section 4. This corroboration is expected, as shared history is a source for synchronic similarity.
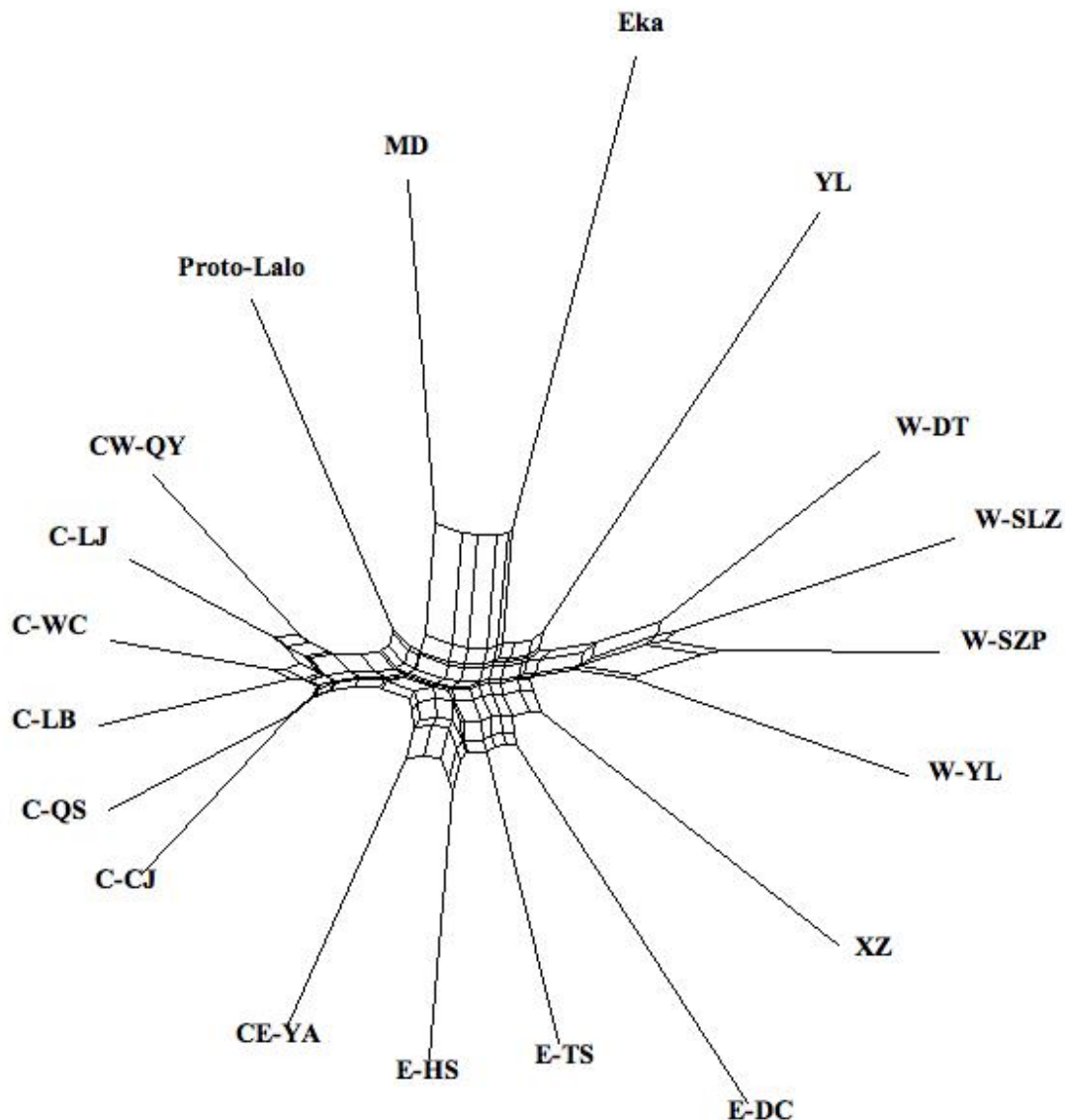


*Figure 7. NeighborNet phenogram of Lalo varieties, based on phonetic distance.*

However, Figure 7 shows two main discrepancies with the phylogenetic tree: (1) some varieties are clustered together because of contact-induced change or shared retentions, not shared innovations, and (2) NeighborNet does not show any upper-level groupings like Core Lalo or Greater Lalo. CE-YA, a C variety located in northeast Weishan County, just to the south of the E varieties, clusters with E varieties. The connection to E is a result of the close contact between CE-YA and E speakers, including frequent intermarriage. CE-YA shares the bulk of its innovations in initials and finals with other C varieties, but also shares the

Tone *1 split with the E cluster. The Tone *1 split affects a large number of words in the lexicon, and thus has a relatively large effect on the LD between CE-YA and other E varieties.

MD and Eka appear to form a cluster through a long, reticulated branch, but are then separated from each other by long individual lines. The apparent MD-Eka cluster is a result of shared retentions and possible areal influence, not innovations, and they are not grouped together on historical grounds. Synchronically, though, both Eka and MD are located in the same area of southern Lincang Prefecture and thus share a similar linguistic ecology. Both Eka and MD are in contact with the languages of this region, such as Dai and Black Lahu. From a language planning perspective, Eka and MD's geographical proximity and shared areal influence make them candidates for partnership in community-based language planning efforts.

The NeighborNet phenogram does not show any of the higher-level groups such as C-W, Core Lalo, or Greater Lalo. There are two reasons for this discrepancy. First, the shared innovations that characterize the upper-level subgroups do not have a large effect on the lexicon; they are only found in a small subset of words, and therefore their contribution to phonetic distance is small. Second, later changes in the tonal systems that characterize the lower-level clusters affect a large set of words and therefore have a large impact on the phonetic distance. These subsequent changes at the dialect cluster level have lessened the subgroups' phonetic similarity; for example, even though C and W are historically linked, their modern forms are now very different.

## 5.2 Multidimensional scaling

MDS, like NeighborNet, displays the relationships between Lalo varieties without forcing them into a tree. MDS does not explicitly cluster the varieties together at all, but instead presents them as they interrelate to all other varieties. Varieties may form a visual cluster by all being located near each other, as is the case with Central varieties in Figure 8 below. Two-dimensional, Euclidean space is used for representing Lalo varieties. Kleiweg (2004), in his tutorial on the RuG-L04 software, recommends Kruskal's (1964) method as often giving the best results. Kruskal's method finds a configuration of varieties that best matches the rank order of phonetic distances between Lalo varieties, so that the most distant varieties in the distance matrix are the most distant in the MDS space, and the most similar varieties are the closest (UNESCO 2008).

Figure 8 shows the results of using Kruskal's method on the phonetic distance matrix. In general, the results are consistent with NeighborNet's network diagrams given in Section 5.1. MDS, in contrast to NeighborNet, identifies the upper-level grouping of Greater Lalo. Greater Lalo varieties (Core Lalo plus XZ) form a cluster in the middle of the diagram, with peripheral varieties YL, Eka, and MD on the outside, far apart from each other and from everyone else. The MDS distance echoes the geographic distance these Lalo varieties have with Core Lalo, with Eka and MD to the far south and YL to the far west of Core Lalo

distribution. C and W clusters are clearly defined, while the E cluster is less defined, similar to the results given in Section 5.1. CE-YA is part of the C cluster, but still closer to E varieties than any other C variety. E-DC, although still relatively close to other E varieties, appears on the edge of the Greater Lalo cluster.
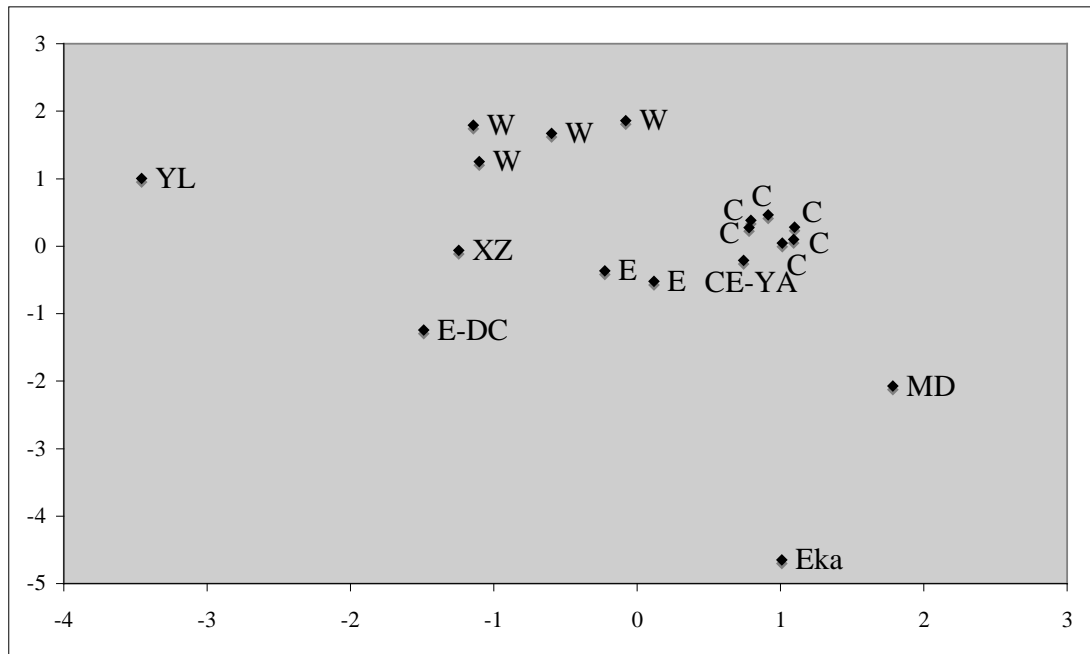


*Figure 8. Multidimensional scaling of Lalo varieties, based on phonetic distance.*

In summary, NeighborNet and MDS diagrams based on phonetic distance agree with the diachronic subgrouping at a relatively shallow time depth, i.e. at the level of dialect clusters and individual varieties. This supports the classification of the three dialect clusters and four peripheral languages. MDS shows Greater Lalo varieties closer together, but NeighborNet does not show any of the higher-level subgroups identified in the diachronic subgrouping. For deeper time depth relationships in the Lalo language cluster, and for a more precise distinction between contact and true phylogenetic relationships, the comparative method must be called on. However, in terms of language planning, identifying more recent relationships is an important step, and the dialectometric tools are a useful complement to the comparative method in this function.

## 6. COMPREHENSION TESTING

Recorded Text Testing (RTT), first developed by Casad (1974) and adapted by Kluge (2007), measures listeners' comprehension of a particular text, which is then used as an estimate of the overall degree of intelligibility of that variety. An RTT is a short narrative recorded in variety A and played to a listener from variety B. The variety B listener then listens to the story in short sections and retells the content of each section. The comprehension score is measured by

checking the retold content against a baseline of core elements previously identified by a pilot test panel of native listeners.

To conduct the pilot test, a speaker native to village A relates a short narrative, typically one to three minutes long, which is recorded and then translated into Chinese. A panel of eight native speakers is invited to listen to the text. Panel participants grew up in the village, spoke the language fluently, and had parents native to the locale. Although fluency in Chinese was not a requirement for participation, the rate of bilingualism among Lalo is such that all participants were also fluent in Chinese. Thus, participants answered in Chinese without the use of a translator. Native speakers first listened to the whole text and then listened again section-by-section, with a pause after each section in which the listener retold the section's content in Chinese. Elements that all native listeners retold formed the baseline for scoring responses for participants from other varieties. By first pilot testing the text, comprehension levels seen in cross-dialectal listeners' scores were comparable to a native speaker's comprehension of the same text.

When administering an RTT in dialect B, the same selection procedures were used as for the pilot test participants, with the added stipulation that the participant must not have lived in the area where the tested variety was spoken for more than one month. Although this could not eliminate casual contact at markets or festivals, it still screened out participants who had acquired comprehension through substantial contact. Eight to ten participants were selected in each village, with roughly half the participants male and the other half female. If comprehension scores were close to zero (below 10%) for more than two participants, the test was discontinued.

After screening, each participant followed the same testing procedure as the panel of native listeners. Played sections were never longer than 10 seconds, ensuring that the RTT was not testing memory but rather comprehension. A listener's comprehension score was the number of core elements mentioned divided by the total number of core elements identified by the panel of native listeners.

Not all RTTs were tested at each village, due to time constraints and the fatiguing nature of the RTT process for the participants. CW-QY was tested at 16 out of 18 locations, because of its candidacy as a reference dialect for orthography development. C-LJ was tested in seven locations, E-HS and E-DC in five, W-DT in three, CE-YA and C-CJ in two, and C-WC in one, with a total of 122 participants.

In general, RTT results match expectations based on the phylogenetic tree and dialectometric groupings. Cross-cluster comprehension is low, unless a variety has significant cross-cluster contact, while intra-cluster comprehension is usually intermediate or high. Figure 9 shows the mean RTT result for each village when responding to the CW-QY text. Peripheral varieties all showed low comprehension of the CW-QY text. E listeners showed slightly higher comprehension than peripheral varieties, but still at or below 40%. W varieties range widely on their comprehension, from 10% in W-DT to 62% in W-YL. W-

YL is on the border between C and W varieties and is in contact with C varieties such as C-LB. The unexpectedly high score for W-YL may reflect acquired rather than inherent comprehension, i.e. may be due to contact rather than structural similarity. The standard deviation of W-YL's scores was 17, with scores ranging from 30% to 77%. A large standard deviation (over 12) suggests comprehension acquired through contact (Blair 1990: 25). Within the C cluster, C-WC and C-LJ score the highest, which reflects their geographic and genetic proximity to CW-QY. Other C varieties also score relatively high (above 60%), especially when compared with listeners from the E and W clusters.
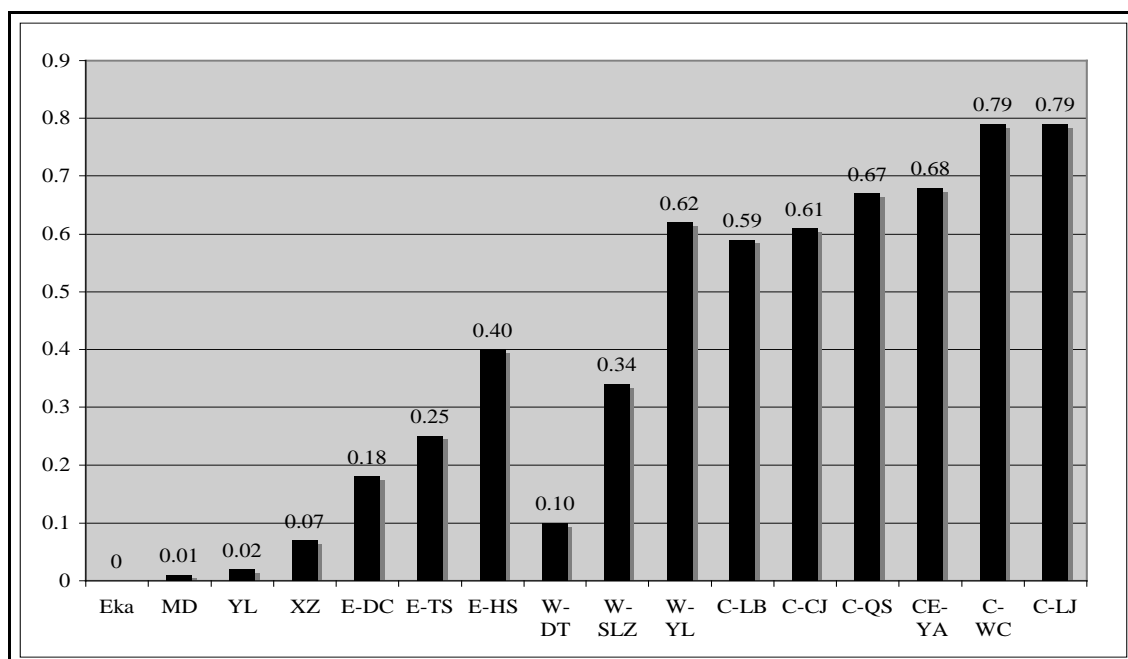


*Figure 9. Mean RTT scores of CW-QY text by village.*

Figure 10 below shows the mean RTT result for E-HS (marked in black) and E-DC (marked in gray) texts. In general, non-E varieties show low comprehension. The intra-cluster score of E-DC listening to E-HS is high at 88%. CE-YA's high score of 70% when listening to E-DC is probably due to the close social contact between the two villages, which has also influenced CE-YA's phonological development (e.g. CE-YA shares the Tone *1 split with E-DC, as noted in Section 5.1).
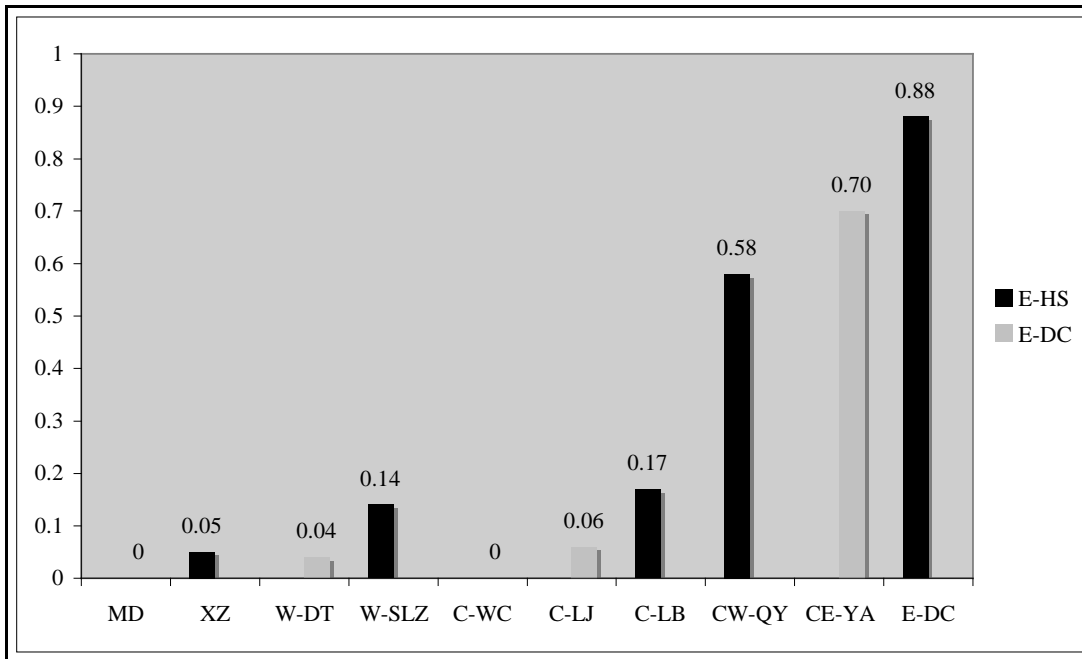
*Figure 10. Mean RTT scores of E-HS and E-DC texts by village.*

Finally, Figure 11 shows the mean RTT scores for W-DT (marked in black) and W-YL (gray) texts. These texts were collected towards the end of fieldwork and so were tested in the least number of locations. Still, a few observations can be made. Not surprisingly, C-LB shows low comprehension of W-DT, while W-YL shows high comprehension. Unexpectedly, W-SLZ shows low comprehension when listening to the W-YL text, but this may reflect the peripheral nature of W-YL's membership in the W cluster. XZ's scores are low for the W-YL text, but rather high for W-DT. Both XZ and W-DT share a chain shift in which $*a > o$ and $*aŋ > a$, so this may have aided XZ's comprehension.
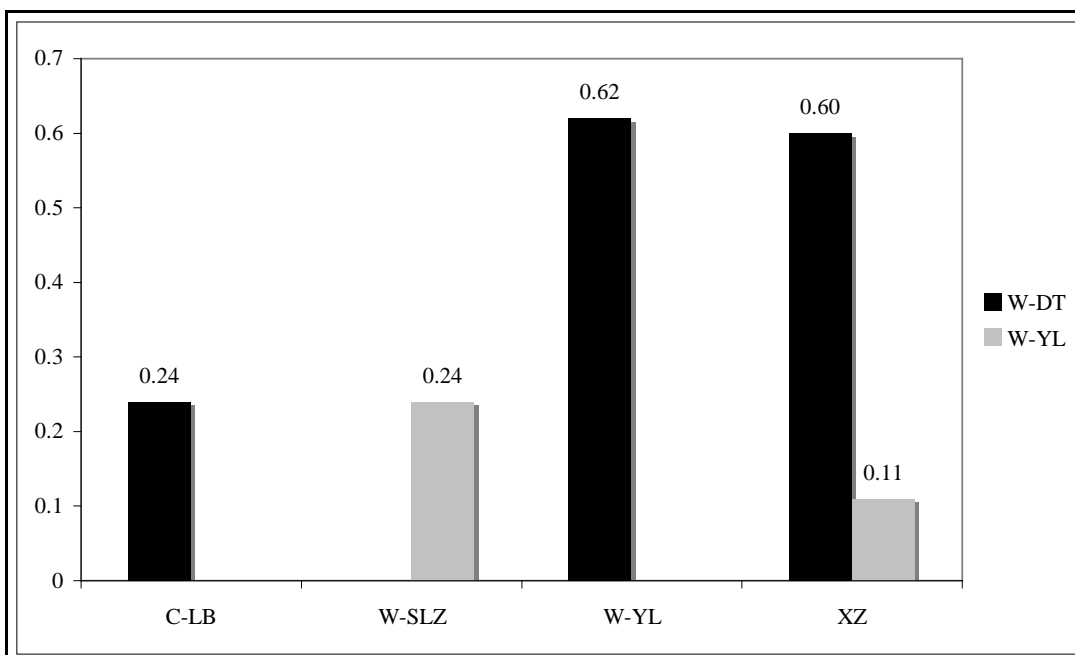


*Figure 11. Mean comprehension scores of W-DT and W-YL texts by village.*

In summary, comprehension test results support the classification of seven Lalo languages: C, E, W, Xuzhang, Yangliu, Eka, and Mangdi. Intra-cluster comprehension tends to be high, as seen in C Lalo varieties' scores when listening to CW-QY, and in E-DC's high score on the E-HS text. Cross-Core cluster comprehension is low, unless there is significant contact, as in the case of CE-YA listening to E-DC. Peripheral-Core comprehension is also low, as seen in the peripheral languages' negligible intelligibility of CW-QY and E-DC.

# 7. CORRELATION BETWEEN PHONETIC DISTANCE AND COMPREHENSION

The strong, significant correlation between RTT results and LD gives further validation for LD as an approximation of comprehension. Table 5 shows the correlation between LD and comprehension for the CW-QY text. Other RTT results are not included in the correlation, as the content of each RTT text is different, and results are therefore not comparable. N is the number of observations, R is Pearson's correlation coefficient, R squared is the proportion of variance explained by the model, and P is the level of significance. The closer R is to 1 or -1, the stronger the correlation between the two variables. The negative R value in Table 5 indicates a negative relationship between LD and comprehension, i.e. the greater the phonetic distance, the lower the comprehension score. R squared indicates how well a regression line approximates real data points, i.e. how much variance in comprehension can be explained by LD. The closer to 1 R squared is, the better LD is at predicting comprehension. The P value is the probability of finding the current R if the real R were in fact zero. P values less than 0.05 indicate that the current R is unlikely to have occurred by chance.

| N | R | R squared | P |
|---|---|---|---|
| *16* | *-0.88* | *0.77* | *0.000007* |

*Table 5. Correlation between LD and comprehension.*

The correlation given in Table 5 is strong and statistically significant. R squared indicates that LD may be able to explain a large proportion of the variance in comprehension scores. Note, however, that there is a degree of interdependence between LD and comprehension test results, since the representation of tone in LD was adopted according to findings in Yang & Castro 2008, which determined the optimal representation of tone through highest correlation to comprehension test results. The representation of tone, though, is a necessary modification of standard LD methodology when dealing with tonal languages. There are, of course, a multitude of other factors not included in LD that affect comprehension, such as differences in discourse patterns, lexicon, rhythm, prosody, syntax, as well as participants' language attitudes and reaction to the testing procedure. Further studies using multiple regression analysis are

needed to determine the weighting of phonetic distance among these other factors.

Figure 12 shows the scatter plot of RTT scores versus LD for CW-QY RTT results, with RTT scores as the Y variable and LD as the X variable. The regression line fits well, with only W-YL appearing as an outlier, due to contact with C varieties. C Lalo varieties cluster in the top left with high comprehension and low LD, E and W show lower comprehension and higher LD, and peripheral languages show almost no comprehension and high LD.

The strong, significant correlation between comprehension and LD has already been noted for Scandinavian languages (Gooskens 2006) and East Asian tone languages such as Nisu (Yang 2009), Bai and HSH Zhuang (Yang & Castro 2008). This result is consistent with Gooskens' (2006) findings for correlation between LD and comprehension of Scandinavian languages: r=-0.82, p<0.01. These correlations suggest that LD performs consistently as a dialectometric tool in both Indo-European and East Asian tone languages.
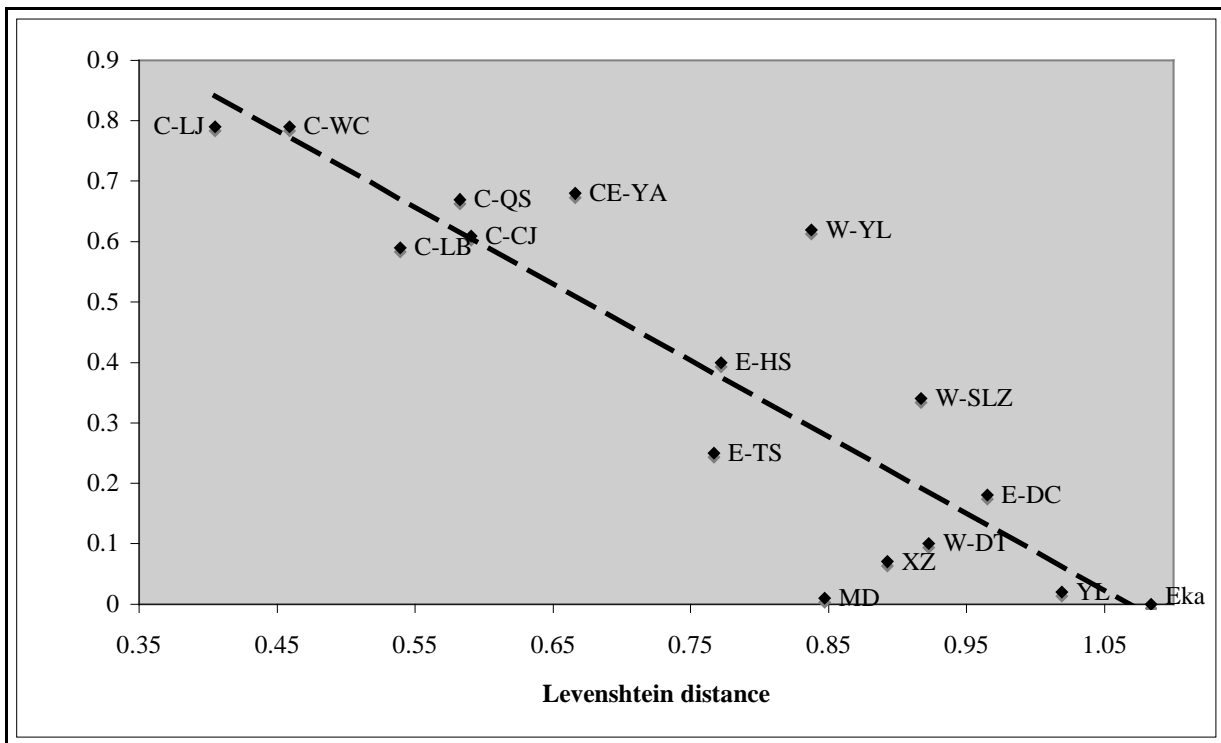


*Figure 12. Scatter plot of comprehension of CW-QY text versus LD.*

## 8. CONCLUSION

Findings from diachronic subgrouping, phonetic distance and comprehension agree in their identification of C, W, and E clusters and the four geographically peripheral languages. The phylogenetic tree, NeighborNet network analysis, and MDS all show the distinct groupings of C, W, and E, as well as the considerable divergence of the peripheral languages with all other Lalo languages. Comprehension test results indicate that the perceptual impact of the differences

between each language is a significant impediment to intelligibility, thus underscoring the distinctions between them. Additionally, the statistical correlation between comprehension and phonetic distance is strong and significant. The triangulation between the three different methodologies gives support to the classification of the seven Lalo languages presented in this paper.

Though each has its own limitations, the three methodologies complement each other by answering different questions. Historical analysis identifies important differences between varieties and postulates genetic subgroups; dialectometric analysis measures cumulative degrees of difference; comprehension testing measures the perceptual impact of those differences on comprehension. Diachronic subgrouping provides insights into Lalo history, but its qualitative, detail-focused nature invites an aggregate measure as a cross-check. Synchronic, dialectometric analysis clearly distinguishes groups at a shallower time depth, giving support to the diachronic subgrouping. However, it fails to identify higher-level groupings such as Core or Greater Lalo and conflates contact-induced change with shared innovations, limitations that call for the comparative method as a cross-check. Comprehension testing, while unable to identify specific differences between varieties, gauges the perceptual effect of those differences.

The discrepancies between synchronic and diachronic groupings have fruitful implications for language planning. Critically for endangered languages like the Lalo cluster, historical linguistics and language development efforts must be considered in tandem. NeighborNet network analysis strikingly displays the impact of contact between certain varieties, particularly in the placement of CE-YA in the E Lalo cluster. The comparative method can identify contact-induced changes in CE-YA, such as the development of a contrastive low-rising tone identical to E Lalo's. However, NeighborNet network analysis, as opposed to a phylogenetic tree, is able to show the degree of influence of the contact on the overall phonetic structure of CE-YA. From a diachronic perspective, the NeighborNet phenogram incorrectly portrays CE-YA as belonging to the E Lalo cluster. But from a synchronic perspective, CE-YA has become a hybrid Central/Eastern variety, and this result is accurately reflected in both NeighborNet's network diagram and in the high levels of comprehension that CE-YA speakers show when listening to an E text. Therefore, language planners should consider including CE-YA and E speakers in the same community based language planning efforts. For CE-YA speakers, non-print media in Lalo, such as a video on AIDS prevention, may be more effective in an E variety than, for example, in a C variety such as CW-QY.

The classification presented here has implications for our understanding of Lalo history. The center of Lalo's rich linguistic diversity is located in southern Dali Prefecture, a finding that supports the Lalo's traditional belief that this area has been their homeland for millennia. These findings support the claim of many Lalo to a historical link with the Meng clan, who became the leaders of the Nanzhao Kingdom (737-902 A.D.) (Fan 1961). Historical records indicate that the Meng clan first rose to prominence in southern Weishan and northern Nanjian

counties (Backus 1981), right in the heart of the Central Lalo area. Any link to the Nanzhao Kingdom has a large potential impact on the tourism industry in the area; indeed, exploitation of such a link can already be seen in the 2009 renaming of Weishan's county seat as "Nanzhao". While the linguistic evidence cannot directly confirm the Lalo people as descendants of the Meng clan, the diversity found in the Core Lalo area affirms the Lalo's origin in the land where the Meng clan first came to power.

The dialectometric tools used in this study (i.e. Levenshtein distance as input for NeighborNet and MDS) have only recently been developed and as yet have not been applied to many indigenous languages in East Asia. This study provides external validation for these tools through the strong, significant correlation between Levenshtein distance and comprehension, and the convergence between NeighborNet and MDS's clustering and diachronic subgrouping at a shallow time depth. When there are differences between historical and dialectometric analysis, the dialectometric results provide additional avenues of inquiry for language planning efforts. Therefore, these methods are appropriate for use with the under-documented and, in many cases, endangered languages of East Asia, which have an urgent need for further documentation and language maintenance work.

## REFERENCES

Backus, Charles. 1981. *The Nan-chao kingdom and T'ang China's southwestern frontier*. Cambridge; New York: Cambridge University Press.

Ben Hamed, Mahé. 2005. Neighbour-nets portray the Chinese dialect continuum and the linguistic legacy of China's demic history. *Proceedings of the Royal Society Biological Sciences* 272. 1015-1022. DOI: 10.1098/rspb.2004.3015.

Ben Hamed, Mahé & Feng Wang. 2006. Stuck in the forest: Trees, networks and Chinese dialects. *Diachronica* 23(1). 29-60.

Björverud, Susanna. 1998. *A grammar of Lalo*. Lund: Lund University PhD dissertation.

Blair, Frank. 1990. *Survey on a shoestring (SIL and UTA Publications in Linguistics 96)*. Dallas: SIL.

Bradley, David. 1979. *Proto-Loloish (Scandinavian Institute of Asian Studies Monograph No. 39)*. London: Curzon Press.

Bradley, David. 1979a. *Lahu dialects*. Canberra: Australian National University Press.

Bradley, David. 1997. Onomastic, orthographic, dialectal, and dialectical borders: The Lisu and the Lahu. *Asia Pacific Viewpoint* 38.2. 107-117.

Bradley, David. 2002. The subgrouping of Tibeto-Burman. In Christopher Beckwith & Henk Blezer (eds.), *Medieval Tibeto-Burman languages*, vol. 2, 73-112. Leiden: Brill.

Bryant, David & Vincent Moulton. 2004. NeighborNet: An agglomerative algorithm for the construction of planar phylogenetic networks. *Molecular Biology and Evolution* 21. 255-265.

Casad, Eugene H. 1974. *Dialect intelligibility testing*. Dallas: SIL.

Chen Shilin, Bian Shiming & Li Xiuqing. 1985. *Yiyu jianzhi [Outline of the Yi language] (Zhongguo Shaoshu Minzu Yuyan Jianzhi Congshu [Series of Outlines of China's Minority Nationality Languages])*. Beijing: Minzu Chubanshe.

Edmondson, Jerold & John H. Esling. 2006. The valves of the throat and their functioning in tone, vocal register and stress: Laryngoscopic case studies. *Phonology* 23. 157-191.

Fan Chuo. 1961. *The Man shu: Book of the southern barbarians*. Translated by Gordon H. Luce. Ithaca: Cornell University.

Goebl, Hans. 1982. *Dialektometrie: Prinzipien und methoden des einsatzes der numerischen taxonomie im bereich der dialektgeographie [Dialectometry: principles and methods of the use of numerical taxonomy in dialect geography]*. Wien: Österreichische Akademie der Wissenschaften.

Goebl, Hans. 2006. Recent advances in Salzburg dialectometry. *Literary and Linguistic Computing* 21(4). 411-435. http://llc.oxfordjournals.org/cgi/content/abstract/21/4/411 (accessed 24 May 2007).

Gooskens, Charlotte. 2006. Linguistic and extra-linguistic predictors of inter-Scandinavian intelligibility. In Jeroen van de Weijer & Bettelou Los (eds.), *Linguistics in the Netherlands 2006*, 101-113. Amsterdam: John Benjamins.

Gooskens, Charlotte & Wilbert Heeringa. 2004. Perceptive evaluation of Levenshtein dialect distance measurements using Norwegian dialect data. *Language Variation and Change* 16(03). 189-207.

Heeringa, Wilbert. 2004. *Measuring pronunciation differences with Levenshtein distance*. Groningen: University of Groningen PhD dissertation.

Heeringa, Wilbert, Peter Kleiweg, Charlotte Gooskens & John Nerbonne. 2006. Evaluation of string distance algorithms for dialectology. In John Nerbonne & E. Hinrichs (eds.), *Linguistic Distances Workshop at the joint conference of International Committee on Computational Linguistics and the Association for Computational Linguistics*, 51-62. Sydney: Australia Association for Computational Linguistics.

Houtzagers, John, John Nerbonne & Jelena Prokić. 2010. Quantitative and traditional classifications of Bulgarian dialects compared. *Scando-Slavica* 56(2). 163-188.

Huson, Daniel H. & David Bryant. 2006. Application of phylogenetic networks in evolutionary studies. *Molecular Biology and Evolution* 23(2). 254-267.

Kleiweg, Peter. 2004. RuG/L04, software for dialectometrics and cartography. http://www.let.rug.nl/~kleiweg/indexs.html. (accessed 24 May, 2008).

Kluge, Angela. 2007. RTT retelling method: An alternative approach to intelligibility testing. *SIL Electronic Working Papers* 2007(006).

http://www.sil.org/silewp/abstract.asp?ref=2007-006 (accessed 15 Jan 2008).

Kruskal, Joseph. 1964. Nonmetric multidimensional scaling: A numerical method. *Psychometrika* 29(2). 115-129.

Maguire, Warren, April McMahon, Paul Heggarty & Dan Dediu. 2010. The past, present and future of English dialects: Quantifying convergence, divergence, and dynamic equilibrium. *Language Variation and Change* 22(1). 69-104.

Matisoff, James A. 1973. Tonogenesis in southeast Asia. In Larry Hyman (ed.), *Consonant Type and Tone*, 71-95. Los Angeles: University of Southern California.

Matisoff, James A. 1973 [1982]. *The grammar of Lahu (2nd ed., University of California Publications in Linguistics, No. 75)*. Berkeley: University of California Press.

Matisoff, James A. 1988. *The dictionary of Lahu*. Berkeley: University of California Press.

Matisoff, James A. 2003. *Handbook of Proto-Tibeto-Burman: System and philosophy of Sino-Tibetan reconstruction*. vol. 135 *(UC Publications in Linguistics)*. Berkeley: University of California Press.

McMahon, April, Paul Heggarty, Robert McMahon & Warren Maguire. 2007. The sound patterns of Englishes: representing phonetic similarity. *English Language and Linguistics* 11(1). 113-142.

McMahon, April & Robert McMahon. 2005. *Language classification by numbers*. New York: Oxford University Press.

Nerbonne, John. 2009. Data-driven dialectology. *Language and Linguistics Compass* 3(1). 175-198. http://www.let.rug.nl/nerbonne/teach/ling-tech/literature/Nerbonne-Compass-2008.pdf (accessed 28 Jan 2009).

Nerbonne, John & William Kretzschmar, Jr. 2006. Progress in Dialectometry: Toward Explanation. *Literary and Linguistic Computing* 21(4). 387-397. http://llc.oxfordjournals.org/cgi/content/abstract/21/4/387

Pelkey, Jamin R. 2008. *The Phula languages in synchronic and diachronic perspective*. Melbourne: La Trobe University PhD dissertation.

Pelkey, Jamin R. 2011. *Dialectology as Dialectic: Interpreting Phula Variation (Trends in Linguistics: Studies and Monographs)*. Berlin: Mouton de Gruyter.

Pittayaporn, Pittayawat. 2009. *The phonology of Proto-Tai*. New York: Cornell University PhD dissertation.

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

Séguy, Jean. 1971. La relation entre la distance spatiale et la distance lexicale [The relationship between spatial distance and lexical distance]. *Revue de Linguistique Romane* 35(138). 335-357.

Stanford, James N. in prep. Dialectometry, rice paddies, and clans in rural China.

Tang Chaoju. 2009. *Mutual intelligibility of Chinese dialects: An experimental approach*. Leiden: Leiden University PhD dissertation.

Toulmin, Matthew. 2009. *From linguistic to sociolinguistic reconstruction: The Kamta historical subgroup of Indo-Aryan (Studies in Language Change)*. Canberra: Pacific Linguistics.

UNESCO. 2008. Non-metric multidimensional scaling. *WinIDAMS 1.3 Reference Manual*, Chapter 8.1. Paris: UNESCO. http://www.unesco.org/webworld/idams/advguide/Chapt8_1.htm (accessed Feb 17, 2010).

Wang Chengyou. 2003. *Yi yu fangyan bijiao [Comparative study of Yi dialects]*. Chengdu: Sichuan Minzu Chubanshe.

Yang, Cathryn. 2009. Nisu dialect geography. *SIL Electronic Survey Reports* 2009(007). http://www.sil.org/silesr/abstract.asp?ref=2009-007 (accessed 30 Apr 2009).

Yang, Cathryn. 2010. *Lalo regional varieties: Phylogeny, dialectometry, and sociolinguistics*. Melbourne: La Trobe University PhD dissertation. http://arrow.latrobe.edu.au:8080/vital/access/HandleResolver/1959.9/1530 15.

Yang, Cathryn & Andy Castro. 2008. Representing tone in Levenshtein distance. *International Journal of Humanities and Computing* 2(1-2). 205-219. http://www.euppublishing.com/doi/abs/10.3366/E1753854809000391 (accessed 30 Oct 2009).

Zhu Wenxu. 2005. *Yiyu fangyan xue [Yi dialect studies]*. Beijing: Zhongyang Minzu Daxue Chubanshe.

*Author's address:*

8584 N Private Rd 600 W
Brazil, IN 47834
USA

cathryn_yang@sil.org