

Some puzzles in Cao Lan

Kenneth J. GREGERSON & Jerold A. EDMONDSON

University of Texas at Arlington

There is a Latin phrase *tertium quid* for something that escapes easy division into two groups supposed to be exhaustive.¹ In this paper we outline why the Cao Lan language may be regarded as a *tertium quid* in regard to:

its *tonality*--Cao Lan appears to have merged its tones so that today as few as three may remain in live syllables and two in dead syllables, but this radical surgery has in part been compensated for by developing vowel quality and voice quality contrasts,

its *phonology*--Cao Lan demonstrates some phonological rules regarded as exclusively Northern Tai (NT) and some others that have been considered as Central Tai (CT) features,

its *vocabulary*--Cao Lan has a sizeable inventory of vocabulary that is common Tai, but also it has both Central Tai and Northern Tai forms as well.

The motivation for these variegated features we believe lies in the migratory history of the Cao Lan people. Cao Lan by their own account and from the study of experts is a group that has moved from its original homeland in NT territory in China to CT territory in Vietnam. In this paper we will present some evidence from our own fieldwork and computer-assisted analysis subsequent to it about the frailty of Cao Lan tone as well as the mottled phonological and lexical traits in Cao Lan and speculate about their origins. First of all though, we present some background about this group.

1. Introduction

The *Cao Lan-San Chây* people are one of the fifty-four officially recognized minority nationalities of Vietnam. The hyphenated designation as a part of their name is a recognition that this group speaks two languages and has two non-overlapping cultures. The Cao Lan, for one, use a language that belongs to the Tai Branch, whereas the San Chây make use of a form of Han Chinese. As a whole, the

¹The research reported on here was supported by a grant to the authors from the National Endowment for the Humanities and the National Science Foundation for a project called "Languages of the Vietnam-China Borderlands, 95-97."

Cao Lan-San Chây are reported to have a population of 114,000 according to the 1990 census, living concentrated in Tuyên Quang and Bắc Thái Provinces. They are also found in smaller numbers in: Yên Bái, Thái Nguyên, Vĩnh Phú, Bắc Giang, Lạng Sơn, and Quảng Ninh Provinces. Our Cao Lan informant, Mr. Hoàng Trường Vinh, used for his people the autonym *San Chây*, which surprised us in as much as San Chây is the official name not of his but of the Han-speaking group. From him we also learned that there is the unusual situation of two speech communities within one ethnicity, both using Chinese characters to record their folk customs and practices. And, as Mr. Vinh confirmed to us, at least some of the Cao Lan call themselves *San Chây*. MVNP 1978 reports that the Cao Lan and San Chây do not live in a classical diglossic situation of high language vs. low language, but as two groups with mostly different identities despite a small overlap today and a common link in the past. Looking at this history, scholars have concluded that at one time these people lived in China in territory along the border areas of Hunan, Guangdong, and Guangxi Provinces. They came to Vietnam circuitously arriving about 400 years ago and the unusual language use feature of the Cao Lan-San Chây appears to have resulted from events that occurred during this migration.

As we noted, the San Chây today speak mostly in a kind of Chinese, but some of their old people can also speak a Tai language. Similarly, seniors among the Cao Lan, such as Mr. Vinh, can speak and write in the Han language. Despite the language choice difference, the two have nevertheless common features of culture, language, and traditions. Some octogenarians among the Cao Lan suggest that they were originally San Chây; and in parallel fashion some older members of the San Chây community suggest they were once Cao Lan. They do share common folkloric traditions and common writing, as is portrayed in one song called “Sinh Ca” (Ca = song). According to the people themselves, they arrived in Vietnam as one group. The late French scholar, André Haudricourt 1973, explained that after they began their migrations, they stopped in Guangxi for a time and assumed a local form of the Zhuang (Tai) language. Later, on the way to Vietnam, some of them adopted the written language of Chinese and some began to use a spoken form of Chinese. To add further to the complication, it has also been suggested that the Cao Lan-San Chây have a connection to the Yao people as well. Indeed, there is a group of Yao in Fangcheng in Guangxi who call themselves [san³³ tɕai³³], Mao et al (1982:8). All these facts tell us that the two were in some sense one nationality with two partially overlapping speech communities whose original bilingualism has developed into separated mostly monolingualism through separation, as the majority of the San Chây live in Quảng Ninh and the Cao Lan live mostly in Tuyên Quang, Thái Nguyên, and Bắc Giang.

2. History of previous work

Cao Lan was analyzed first by Bonifacy. Later references are then in Haudricourt (1960) The most famous among the unpublished sources is the 1938 survey of the École française d'extrême orient (EFEO) of 1938. It is a compilation of about 500 items transcribed in Quốc Ngữ script by civil servants, members of the institute, and others gathered at various locations. Among these materials is the entry XVI.2 Mản Cao Lan of Yên Sơn Tuyên Quang, XI.7 Mản Cao Lan Phú Thọ, and XI. Mản Cao Lan Tieu Á. The first location is the same as the one being reported on here, moreover, the comparison below shows that the varieties of Cao Lan are the same.

Most of the post 1945 work has concentrated on the question of genetic affiliation, specially whether Cao Lan is a Northern Tai or a Central Tai language. Haudricourt 1960 suggested that Cao Lan was neither (Black or White) Thái nor Yao, and, moreover, to use the terminology of Li 1977 and Gedney it is neither Central Tai nor is it Northern Tai. Basically, it has phonological and lexical features of both. As far we know, there have not been other studies of the linguistic features of the language and very little information about it has found its way to print. In our modest paper on the Cao Lan language we will provide considerably more data than has heretofore been available on Cao Lan, examine some of its features with computer-assisted techniques, and highlight the features of the mixture. Our informant was Mr. Hoàng Trường Vinh of Tuyên Quang Province, Yên Sơn District, Kim Phú Village, and Giếng Tanh Hamlet. Mr. Vinh was born in 1926 and was 79 years of age at the time of the interview in 1995. He was a very good speaker of the Cao Lan language, using it on a daily basis with his family. He was able to read and write Chinese characters to render his language in written form.

We conclude from the study of Mr. Vinh's speech, very much in the sense of Haudricourt, that the special properties of Cao Lan suggest a profound and prolonged interfusion of Northern Tai and Central Tai influences with a significant adstratum of a kind of Chinese, presumably *Pinghua*, the first kind of Han spoken in Guangxi, cf. below. We did not find evidence of Yao influence in the language per se. The presence of Pinghua we believe to be crucially responsible for the special tonal and voice quality features of Cao Lao. In fact, we see many commonalities between Cao Lan and E, a mixed language reported on in Edmondson 1992.

3. Tonal system

The Cao Lan language has a very unusual prosodic system for a Tai language. In particular, we found that contemporary form showed a collapse of

many original categories into only a few tones. Indeed, the pitch alone seems to distinguish only three categories today, cf. below for the CECIL plots of the tones. Consider the following data:²

- A1: *tu*³¹ *mo*⁴² ‘pig’; *phon*⁴² ‘rain’; *tu*³¹ *ma*⁴² ‘dog’; *na*⁴² ‘thick’
 A2: *tu*³¹ *pa*⁴² ‘fish’; *tu*³¹ *ka*⁴² ‘crow’; *mai*⁴² ‘thread’; *tha*⁴² ‘eye’;
*tha:i*⁴² ‘to die’; *ka:i*⁴² ‘far’
 A3: *?bon*⁴² ‘sky, heavens’; *?dai*⁴² ‘good’; *?bau*⁴² ‘lightweight’;
*?ju*⁴² ‘medicine’
 A4: *tu*³¹ *va:i*³¹ ‘buffalo’; *ha*³¹ ‘cogongrass’; *na*³¹ ‘wet field’
 B1: *mo*⁴² ‘new’; *tu*³¹ *thu*⁴² ‘rabbit’; *pha*⁴² *mai*³¹ ‘palm hand’
 B2: *kai*⁴² ‘chicken’; *thau*⁴² ‘shuttle of loom’
 B3: *?ba*⁴² ‘shoulder’; *?o:n*⁴² ‘soft’; *zaj*⁴² ‘body’; *?da*⁴² ‘scold’
 B4: *ta*⁴² ‘river’; *pai*⁴² *?ba:u*⁴² ‘older brother’
 C1: *thai*³⁵ ‘intestine’; *heu*³⁵ ‘teeth’; *na*³⁵ ‘face’
 C2: *kon*³⁵ ‘hip’; *tai*³⁵ ‘below’; *hai*³⁵ ‘near’
 C3: *?ba:n*³⁵ ‘village’; *zaj*³⁵ ‘winnowing basket’
 C4: *ma*⁴² ‘horse’; *so*³¹ ‘early’; *kau*⁴² ‘owl’
 D1L: *ma:k*⁴⁵ ‘fruit’; *tha:p*⁴⁵ ‘carry on pole’
 D2L: *pɔt*⁴⁵ ‘lungs’
 D3L: *?be:k*⁴⁵ ‘carry on shoulder’; *zik*⁴⁵ ‘to peel’
 D4L: *lo:k*³⁴ ‘to skin’; *no:k*³⁴ ‘outside’; *lak*³⁴ ‘child’
 D1S: *phak*⁴⁵ ‘vegetable’; *lap*⁴⁵ ‘to pull’; *hap*⁴⁵ ‘to bite’
 D2S: *tap*⁴⁵ ‘liver’; *tu*³¹ *kop*⁴⁵ ‘frog’
 D3S: *?dip*⁴⁵ ‘alive’; *?uk*⁴⁵ ‘brain’
 D4S: *sak*³⁴ ‘to wash (clothes)’; *mot*³⁴ ‘ant’

We have prepared plots using WINCECIL and a compositing program. In Figures 1-5 we show the pitch trajectories for representative vocabulary illustrating A1 vs. A4 in Figure 1; B3 vs. B4 in Figure 2; C2 vs. C4 in Figure 3, D2S, D4S, D3L vs. D4L in Figure 4, and finally in Figure 5 we compare A1, B4, and C4 tone examples, *ma*⁴² ‘dog’; *ta*⁴² ‘river’; and *ma*⁴² ‘horse’.

²CECIL and WINCECIL are hardware and software speech analysis systems (JAARS International, Inc., Waxhaw, NC) that allow reliable extraction of pitch in a field setting.

The following list uses the system of Gedney to portray the potential sources for tones. Gedney assumes five original tones called *A*, *B*, *C*, *DL*, and *DS* and four consonant types for the initials, called 1=voiceless friction (aspiration and voiceless sonorants and fricatives, 2=plain voiceless stops, 3=“preglottalized” voiced stops and glottal stop initials, and 4=original fully voiced stops. We use then *A1* to categorize vocabulary with voiceless friction initials in proto tone A, etc.

Cao Lan Tone A

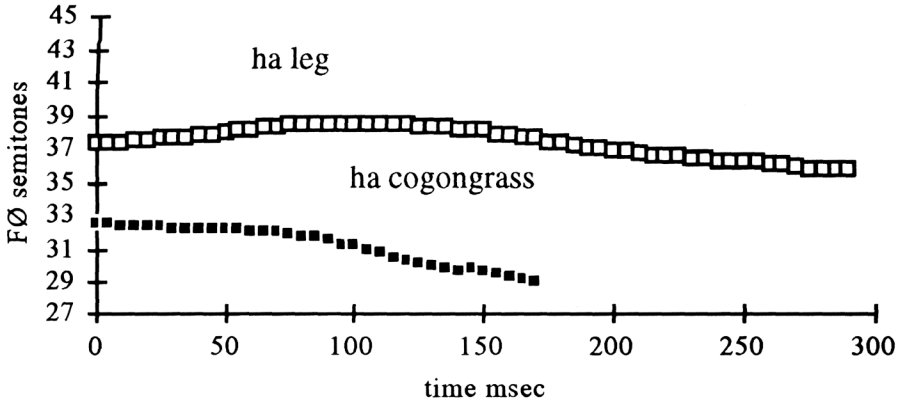


Figure 1. Cao Lan reflexes of Proto-Tai A tone

Cao Lan Tone B

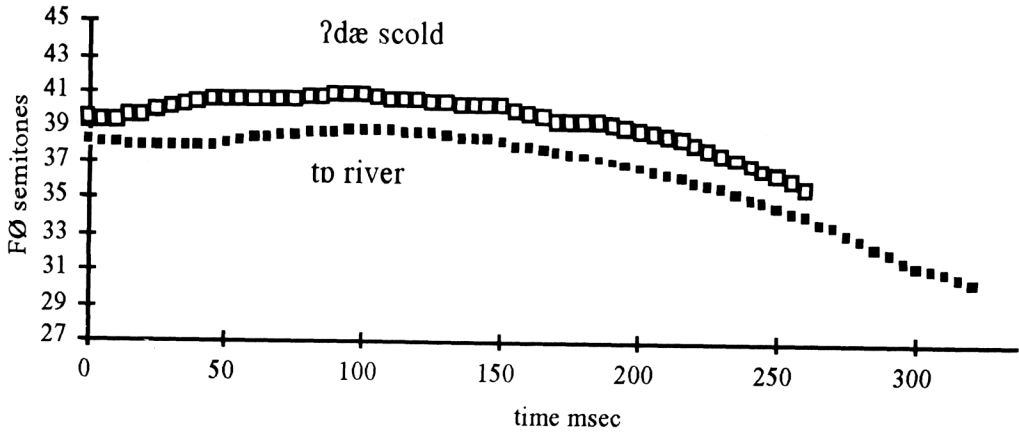


Figure 2. Cao Lan reflexes of Proto-Tai B tone

Cao Lan Tone C

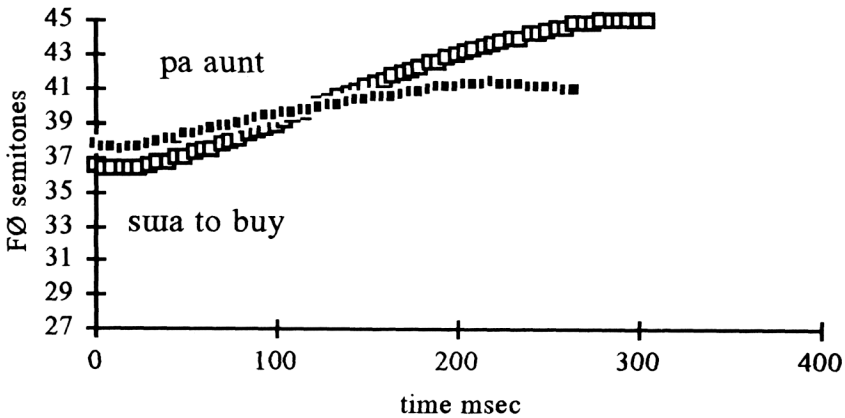


Figure 3. Cao Lan reflexes of Proto-Tai C tone

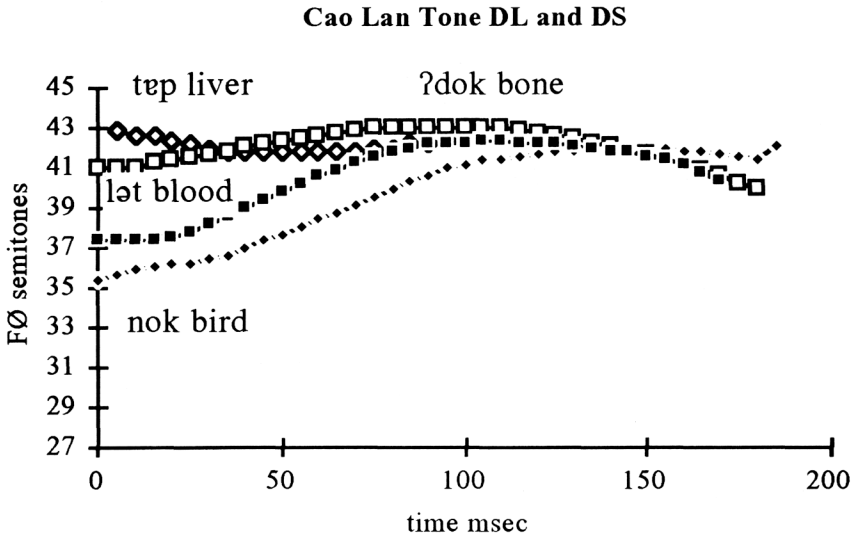


Figure 4. Cao Lan reflexes of Proto-Tai D tone

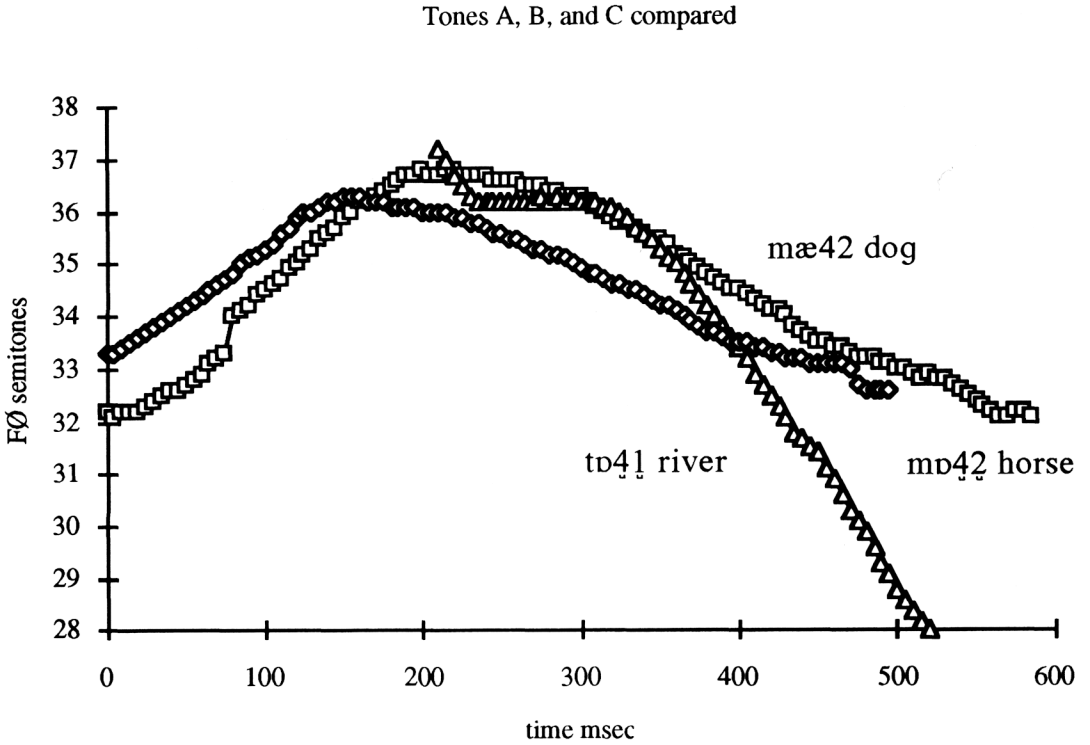
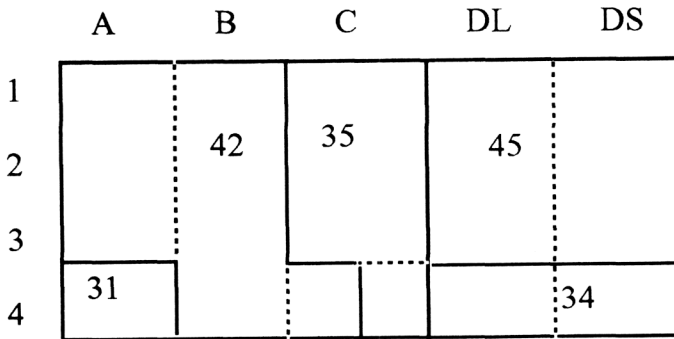


Figure 5. Comparison of *mæ* ‘dog’ and *mɔ* ‘horse’

We can display the results of this comparison with a Gedney diagram showing the pattern of tone splitting, as is shown in Figure 6;



It is to be noted that some lexical items from C4 have become like C123, e. g. *məi*³⁵ ‘tree’ whereas others resemble the pitch trajectories of the categories A123 and B1234. While the pitch may be similar, that is not to say that the words are homophonous. Indeed, Mr. Vinh repeatedly pointed our errors in pronouncing the items *mɑ*⁴² ‘dog’ vs. *mɑ*⁴² ‘horse’. We finally realized that in his speech an expected tonal contrast would better be called one of voice quality and vowel quality transcribable as *mæ*⁴² ‘dog’ vs. *mɒ*⁴² ‘horse’. The vowel in dog is higher and further front than that in horse. In below we display a plot of vowel formants determined by examining spectral data (F2-F1 vs. F1), which clearly demonstrates for several locations in each syllable the quality of the two vowels involved.

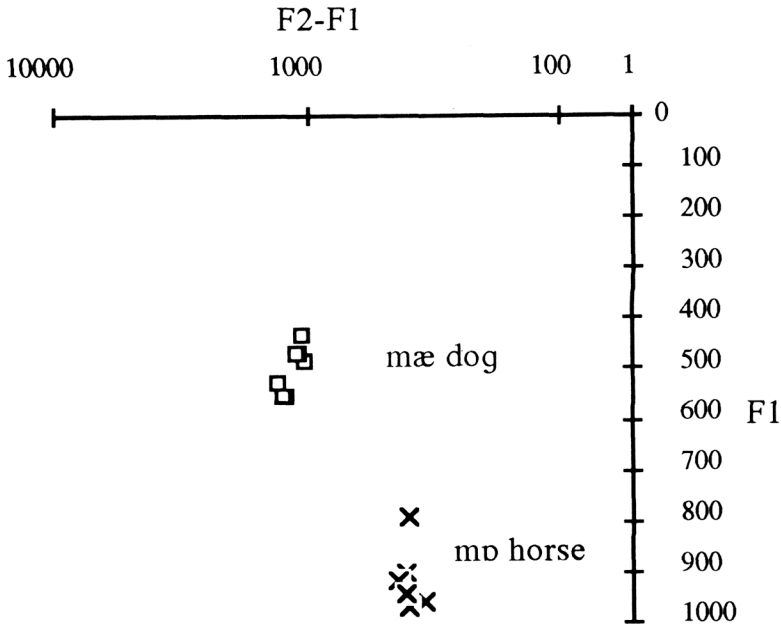


Figure 7. Comparison of the vowel formants of *mæ* ‘dog’ vs. *mɒ* ‘horse’

Once alerted, we examined the voice quality (tense vs. lax) and vowel quality (high-front vs. low-back) to see if there is systematic use of these phonological contrasts to enhance or replace a weakened tonal system. We have concluded at this point that though there is a tendency for high tone items to have tense voice and higher/fronter vowels, there are also numerous exceptions to this tendency. It cannot be said that Cao Lan has developed a register system of the manner found in some Mon-Khmer languages.

4. The phonological and lexical puzzles

Beyond the strange tone features, the Cao Lan language has also developed unexpected segmental features in its phonology and vocabulary. Cao Lan—like Nùng An of Cao Bằng and Ts'un -wa of Quảng Ninh Province, as Strecker reports—possesses characteristics that are geographically out of place; they seem half like the languages of the Northern Tai (NT) Subbranch of Tai and half like those of the Central Tai (CT) Subbranch. Specifically, Strecker wants to argue that Cao Lan belongs to the Northern Tai Branch of Tai or Ddoi, because in Cao Lan, like NT and unlike CT languages **dl* clusters, as in **dluon^A* -> lom³¹ 'wind' result in the same initial as **l* or **r*, as in **ruan^A* -> laan³¹ 'house'. Strecker regards the development of **thr-* and **tr-* to *th-* (according to Li Fang Kuei's 1977 reconstruction), as in **tra^A* -> tha⁴² 'eye' and **thrau^A* -> thau⁴² 'louse' to be "areal" properties and not genetically inherited properties. However, from our data the picture of genetic affiliation of Cao Lan seems much murkier. Below we present fuller data than Strecker had at his disposal on the developments of the proto **r-* and **l-* as well as those of the proto clustered *r*: **dl-*, **thr-*, **nl-*, **kr-*, **khr-*, **gr-*, **pr-*, **pl-*, as well as **f-*, **fr-/vr-/vl-*, **hr-*, and **hw-*. We have provided in many cases corresponding data from a clear example of a NT language in Vietnam, namely, Giáy of Lào Cai Province and a clear example of a CT language of Vietnam, namely, Tày Cao Bằng.³

The development of **dl-*

Gloss	Cao Lan	Giáy Lào Cai (NT)	Tày Cao Bằng (CT)
'wind'	lom ²	đum ²	lom ²
'fingernail'	lip ⁸	zit ⁸ fuŋ ²	lep ⁸

³Tones are marked following the Chinese system for recording tones in Tai languages as: A (high)=1; A (low)=2; B (high)=5; B (low)=6; C (high)=3; C (low)=4; DL (high)=9; DL (low)=10; DS (high)=7; and DS (low)=8.

The development of **dr-*

Gloss	Cao Lan	Giấy Lào Cai (NT)	Tày Cao Bằng (CT)
'boat'	lu ²	đua ²	lũa ²
'day a. tomorrow'	lɔi ²	đu ²	lũa ²

The development of **r-*

'house'	laan ²	đaan ²	ruən ²
'rice husk'	lam ²	đam ²	ram ²
'dry field'	lai ⁶	đi ⁶	rei ²

The development of **hr-*

'to bark'	lau ⁵	đau ⁵	hau ⁵
'mushroom'	let ⁷	đat ⁷	---

The development of **hw-*

'comb'	loi ¹	đoi ¹	wi ¹
--------	------------------	------------------	-----------------

The development of **nl/r-*

'water'	nom ⁴	đəm ⁴	nəm ⁴
'outside'	nɔ:k ¹⁰	đɔ:k ¹⁰	nɔ:k ¹⁰
'bird'	nok ⁸	đok ⁸	nok ⁸

The development of **thr-*

'carry on a pole'	thaap ⁹	đaap ⁹	thaap ⁹
'tail'	thuŋ ¹	đuŋ ¹	thaŋ ¹
'stone'	thən ¹	đin ¹	thin ¹
'loom'	thok ⁷		thuk ⁷
'louse'	thau ¹	đau ¹	thau ¹
'to cook'	thoŋ ¹	đuŋ ¹	thoŋ ¹
'carry hanging'	thiu ³	điu ³	thiu ³
'hailstone'	thet ⁷	---	thep ⁷

As one can see from this data, the NT language Giấy demonstrates the interdental fricative [ð-] as a reflex of all of these proto-initials, whereas Cao Lan has **dl-*, **dr-*, **r-*, **hr-*, **hw-* → *l-*; but **nl/r-* → *n-*; and **thr-* becomes *th-*. Finally, Tày Cao Bằng has the richest set of reflexes: **dl-* and **dr-* → *l-*; **r-* remains **r-*; **nl/r-* → *n-*; and **thr-* becomes *th-*. So, it appears to us that Cao Lan occupies a

middle ground between maximum merger--in the case of Giáy--and maximum preservation--in the case of Tày Cao Bằng. Or it seems about half way between a NT and a CT language and it would be somewhat arbitrary to see some developments as “areal” and some inherited.

In regard to aspiration, the loyalties of Cao Lan are again split. As Li 1977 has pointed out, NT languages seem to have lost aspiration in stops at an early time without engendering any changes in the tonal system. On this test too, Cao Lan shows itself to be neither fish nor fowl, aspirated or non-aspirated. Consider:

‘near’	khai ³	tɕau ³	səu ³
‘vegetable’	phak ⁷	pjik ⁷	phjak ⁷

But the number of unaspirated or NT-like forms predominates in the lexicon so far studied:

‘kill’	ka ³	ka ³	kha ³
‘walk’	pəi ³	pjai ³	thai ³

On the scorecard then, the behavior of aspiration in Cao Lan basically argues for a NT heritage. However, there are equally well rules of CT that seem much in evidence in Cao Lan as well. For example, there is a rule that **f-* -> *ph-*. This rule operates widely in CT and SW Tai, e.g. in Shan, but it is virtually the norm in the CT languages of Vietnam.

The development of **f-* -> *ph-*

‘cloud’	phu ³	vuu ³	pha ⁴
‘cotton’	pha:i ⁵	vai ⁵	pha:i ³
‘rain’	phon ¹	pun ¹	phən ¹
‘hand palm’	pha ⁵	---	---

Another feature of NT vis-à-vis CT is that a number of common items are found in the lower tone set, whereas CT and SWT have these in the high set. On this litmus Cao Lan sides with the North.

‘excrement’	həi ⁴	ʔe ⁴	khi ³
‘rice’	həu ⁴	hau ⁴	khəu ³

There are also some differences of tone or segmental elements that seem to be restricted to individual words. First consider the items in which Cao Lan resembles the CT languages (as in Tày Cao Bằng in the third column):

'ladder'	ʔdɔi ¹	lai ¹	ʔduɛi ¹
'this'	nəi ³	ni ⁴	nəi ³
'a fly'	mɛŋ ² ʔiŋ ²	nɛŋ ²	mɛɛŋ ² wɛn ¹
'to sleep'	nə:n ²	nin ²	no:n ²

The following items, on the other hand, show similarities between Cao Lan and the NT language Giáy.

'rope'	saap ⁹	sak ⁹	zuək ⁹
'come back'	ma ¹	ma ¹	ma ²
'snake'	ŋu ²	ŋu ²	ŋu ²
'widow'	mai ⁵	---	mai ³
'meat'	no ⁶	no ⁶	nua ⁴
'snow'	nai ¹		mui ¹

Finally there are many items, reported in Zhang/Wei 1997 in which NT languages and CT languages have a lexical difference. Consider first these items:

'cloth'	pha:i ⁵	paŋ ²	pha:i ³
'sheep'	ʔbe ³	ji:ŋ ²	ʔbe ³
'goose'	pə:n ⁶	ha:n ⁵	pu:n ⁶
'slippery'	mɛ:k ⁸	mja:k ⁸	ma:k ⁸
'yellow'	lu:ŋ ¹	he:n ³	lu:ŋ ¹
'mud'	tom ¹	na:m ⁶	tom ¹
'bowl'	an ³ toi ³	wa:n ³	tui ³

But then again there are the items:

'sun'	thaŋ ¹ ŋin ²	tɕan ³³ wan ²	tha ¹ wɛn ²
'clothing'	pəu ⁶	pu ⁶	ʔua ³
'sky'	ʔbon ¹	ʔbu:n ¹	fa ⁴
'horn'	kɔk ⁷	kok ⁷	kau ¹
'yesterday'	ŋin ² lu:n ²	lu:n ²	wa ²
'tiger'	kok ⁹	kuk ⁹	ʔua ¹
'flower'	va ²	ʔdok ⁹ va ²	ʔbjok ⁷

5. Conclusion

In this paper we have reported on a very strange-looking tonal system and a phonology and lexicon with—on balance—equal indications of a NT and CT heritage. Admittedly, our ignorance about matters Cao Lan is still profound. Nevertheless, it seems to us that answers to the puzzling features of this language are to be sought in the fusion of three languages that must have existed in a speech community that was multilingual for some time. This would be a situation similar to that Edmondson 1992 reported for the language E of Sanjiang and Luo Cheng Counties, Guangxi. In that language we also found that the organization of tone contrasts was one of the casualties when fusion creolization occurs. The new system that results shows discontinuities in development without total language breakdown as is seen in fission creoles of island nations. Moreover, E also has a similar sort of voice quality development, presumably imported from Tuguai Hua, a local kind of Pinghua Chinese spoken in Guangxi Province.⁴ The resemblance of the voice and vowel quality in Cao Lan and E is uncanny. So, we would venture to speculate that Pinghua figures in the history of Cao Lan. We suppose, the Cao Lan people were once a NT people moving southward who stopped and interfused with a CT speaking group at a time early enough for some late CT rules still to operate. This contact must have been more direct than that which has occurred in most languages of the area, involving living together and intermarriage. That is why it seems to us procrustean to force Cao Lan onto either a NT or CT category. Cao Lan is what it has become, namely mottled and mixed at a very deep level. What's wrong with being a tertium quid anyway?

⁴Liang Min, in a paper read at the 30th ICSTLL in Beijing, suggested that *Pinghua* be recognized as a new type of the Han language. *Pinghua* is the kind of Chinese spoken by descendents of the garrison troops on roads and waterways in Guangxi beginning from Qin times. It is now not accorded the status of Yue, Kejia, Min, etc among Chinese dialectologists. Moreover, it is easily mistaken for Yue (Cantonese), because of its rich tonal system (8 or 9 tones) and codas (-p, -t, -k are still present). However, *Pinghua*, on close examination, is distinct from Cantonese and, moreover, antedates the arrival of Cantonese by many centuries as Liang Min showed. It is also the donor language for the loan words in the minority languages of Guangxi Province.

REFERENCES

- Bonifacy, August. 1907. "Étude sur les Cao Lan." *T'oung Pao ou Archives concernant l'histoire, les langues, la géographie et l'ethnographie de l'Asie orientale*. Serie ii. Vol. viii. Leiden: E. J. Brill, 429-438.
- Edmondson, Jerold A. 1992. "Fusion and diffusion in E, Guangxi Province, China." In T. Dutton et al eds., *The Language Game: Papers in Memory of Donald C. Laycock* (Pacific Linguistic Series, C-110). Department of Linguistics, Australian National University, Canberra, 131-140.
- Haudricourt, André. 1960. "Note sur les dialectes de la région de Moncay." *BEFEO* 50:167-177.
- Haudricourt, André. 1973. "Một số nhận xét về lý luận và thực tiễn nhân một chuyến đi tham các dân tộc thiểu số ở Tây Bắc Việt Bắc." *Ngôn Ngữ* 3.5. [Some remarks about theory and practice on the occasion of a trip to visit the ethnic groups of northwest and northeast Vietnam.]
- Li Fang Kuei. 1977. *A Handbook of Comparative Tai*. Honolulu: The University Press of Hawaii.
- MVNP, Vietnam Social Sciences Committee (ed.) 1978. *Các dân tộc ít người Việt Nam (các tỉnh phía bắc)*. [Minorities of Vietnam, the northern provinces.] Hanoi: Social Sciences Publishing House.
- Strecker, David. 1985. "The classification of the Caolan languages." In S. Ratanakul et al eds., *Southeast Asian Linguistic Studies Presented to André-G. Haudricourt*. Bangkok: Mahidol University.
- Zhang Yuansheng and Wei Xingyun. 1997. "Regional variants and vernaculars in Zhuang." In Edmondson and Solnit eds., *Comparative Kadai: The Tai Branch*. Dallas: Summer Institute of Linguistics.

