

Senri Ethnological Studies 98

Let's Talk about Trees:

Genetic Relationships of Languages and
Their Phylogenetic Representation

Edited by
KIKUSAWA Ritsuko
Lawrence A. REID

SES 98

Let's Talk about Trees:
Genetic Relationships of Languages and
Their Phylogenetic Representation

KIKUSAWA Ritsuko
Lawrence A. REID (eds.)

2018



National Museum of Ethnology

10-1 Senri Expo Park, Suita, Osaka 565-8511, Japan

CONTENTS

List of Figures	i
List of Maps	v
List of Plates	vi
List of Tables	vii
1. Introduction KIKUSAWA Ritsuko and Lawrence A. REID	1
2. Tree and Network in Systematics, Stemmatics, and Linguistics: Structural Model Selection in Phylogeny Reconstruction MINAKA Nobuhiro	9
3. Inferring Population Phylogeny from Genetic Data KIMURA Ryosuke	25
4. Jackknifing the Black Sheep: ASJP Classification Performance and Austronesian Soren WICHMANN and Taraka RAMA	39
5. Freeing the Comparative Method from the Tree Model: A Framework for Historical Glottometry Siva KALYAN and Alexandre FRANÇOIS	59
6. Modeling the Linguistic Situation in the Philippines Lawrence A. REID	91
7. Macrophyletic Trees of East Asian Languages Re-examined Weera OSTAPIRAT	107
8. The Family Tree Model and “Dead Dialects”: Eastern Middle Iranian Languages YOSHIDA Yutaka (translated by KIKUSAWA Ritsuko)	123
9. What the Tree Model Represents: Language Change, Time Depth, and Visual Representation KIKUSAWA Ritsuko	153

5. Freeing the Comparative Method from the Tree Model: A Framework for Historical Glottometry

Siva KALYAN

Australian National University

Alexandre FRANÇOIS

CNRS-LaCiTO, Australian National University

Abstract

Since the beginnings of historical linguistics, the family tree has been the most widely accepted model for representing historical relations between languages. While this sort of representation is easy to grasp, and allows for a simple, attractive account of the development of a language family, the assumptions made by the tree model are applicable in only a small number of cases: namely, when a speaker population undergoes successive splits followed by complete loss of contact. A tree structure is unsuited for dealing with dialect continua, and language families that develop out of dialect continua (“linkages”, as Ross 1988 calls them); in these situations, the scopes of innovations (their isoglosses) are not nested, but rather they constantly intersect, so that any proposed tree representation is met with abundant counterexamples. In this paper, we define “Historical Glottometry”, a new method capable of identifying and representing genealogical subgroups even when they intersect. We apply this glottometric method to a specific linkage, consisting of 17 Oceanic languages spoken in northern Vanuatu.

1. Introduction

The use of genealogical trees for the representation of language families is nearly as old as the discipline of historical linguistics itself; it was first prominently used by August Schleicher in 1853, six years before Darwin proposed a tree model in evolutionary biology (e.g., Minaka and Sugiyama 2012: 177). It has since been the dominant method of visualising historical relationships among languages, and for good reason: its simple structure allows any hypothetical representation of a language family to be interpreted unambiguously as a set of claims about the sequence of demographic and social events that actually occurred in the histories of the communities involved. These hypotheses can then potentially be falsified by new linguistic data or analysis, leading to a more valid representation. Since Schleicher’s time, there have been many other proposals for how to represent the historical relationships among languages, including Johannes Schmidt’s (1872) “Wave Model” (as illustrated e.g., in Schrader 1883: 99 and Anttila 1989: 305); Southworth’s (1964) “tree-envelopes” (which may well predate the similar-looking

“species trees” of phylogeography, e.g., Goodman et al. 1979, Maddison 1997); Hock’s (1991: 452) “truncated octopus’-like tree”; van Driem’s (2001) “fallen leaves”; and more recently NeighborNet (Hurlles et al. 2003; Bryant et al. 2005). However, to our knowledge, no alternative representation has yet combined precision and formalisation with direct interpretability in terms of historical events, to the same extent that the family tree model has.¹⁾

Yet there are important reasons to be dissatisfied with the family-tree model (as has long been noted; see e.g. Bloomfield 1933: 310–314). In particular, the family-tree model rests entirely on the assumption that the process of language diversification is one where language communities undergo successive splits—via migration or other forms of social disruption—with subsequent loss of contact. While this particular social scenario may have occurred occasionally (e.g., in the separation of Proto-Oceanic from the remainder of the Austronesian language family; see Pawley 1999), it can hardly be regarded as the general case.

The way language change arises is via a process of LANGUAGE-INTERNAL DIFFUSION (François 2014, 2017; cf. Labov 1963; Milroy and Milroy 1985; Croft 2000: 166–195; Enfield 2008), as speakers in a network imitate each other so as to jointly adopt an innovative speech habit. Once an innovation settles into a certain section of the social group, it becomes part of its linguistic heritage and can be transmitted to its descendants. This diffusion process is the underlying mechanism behind genetic relations among languages, as each subgroup is defined by the innovations its members have undergone together. Whereas contact-induced change takes place between separate languages, the process of *language-internal diffusion* that underlies language genealogy involves mutually intelligible speech varieties.

The tree model can represent such “genetic” (or better, to use Haspelmath’s (2004: 222) preferred term, “genealogical”) relations in just one type of case: when a language community has split into separate groups, each of which has later gone through its own innovations. It cannot properly handle the frequent case where adjacent speech communities remain in contact even after undergoing innovations that differentiate them from each other. In these cases, as long as the speech varieties remain mutually intelligible for some time, nothing prevents successive innovations from targeting overlapping portions of the network: e.g., one innovation may target dialects A-B-C, another one C-D-E, then B-C, then D-E-F, etc. In such cases of dialect chains or networks, the layering of partially overlapping innovations results in INTERSECTING genealogical subgroups—a situation which cannot be described by the tree model (Gray et al. 2010: 3229).

As is increasingly evident from the work of historical linguists, this sort of intersecting configuration typical of dialect continua is also the normal situation in most language families around the world (e.g., Geraghty 1983; Ross 1988; Toulmin 2009; Heggarty et al. 2010; Huehnergard and Rubin 2011): such families, characterized by an internal structure where genealogical subgroups intersect, are not compatible with a tree representation. Of course, one can force any set of data into a tree structure, but in most cases, this can only be done by selectively discarding some of the data, so as to retain

only those innovations which are compatible with a particular subgrouping hypothesis. Debates about which tree best represents the language family thus usually boil down to arguments (often pointless) over which parts of the data may be ignored.

In this paper, we start by elaborating on the arguments and claims made in the preceding paragraphs, by illustrating in greater detail how trees are used in historical linguistics, and discussing their advantages and disadvantages. We then propose a new method of representing genealogical relationships among languages, which we call Historical Glottometry. While ultimately inspired by the Wave Model which Schmidt (1872) proposed as an alternative to the family tree, our method also draws on the quantitative approach of dialectometry (Séguy 1973; Goebel 2006; Szmrecsányi 2011). We hope this method provides more realistic insights into language history than the tree model, while still combining precision and formalisation with historical interpretability. We conclude the paper by applying our method to a group of seventeen Oceanic languages spoken in Vanuatu.

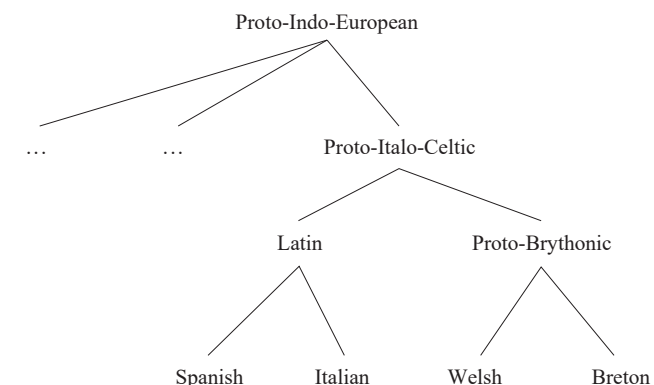
2. Subgrouping under the Tree Model

2.1 An Example from Indo-European

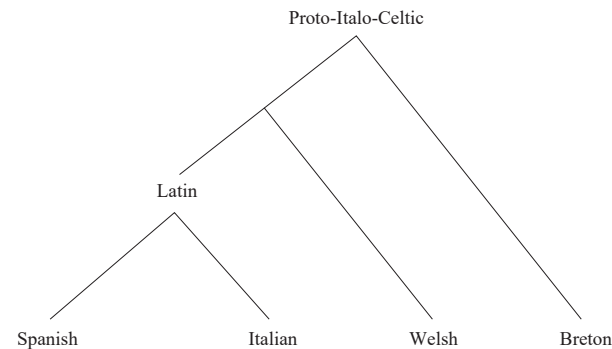
Consider the family tree shown in Figures 5-1, which represents a selection from the Indo-European language family. At the bottom are languages that are currently spoken; languages higher in the tree are ancestors of the languages that branch from them. Each nodal ancestor is called a proto-language, whose descendants together form a subgroup.

In some cases, ancestor languages have been preserved in writing; thus we have direct evidence that (some variety of) Latin is the common ancestor of Spanish and Italian. In other cases, the ancestors are hypothetical, and must be reconstructed by comparing their surviving descendants; thus it is merely a hypothesis that there was a unified Proto-Brythonic language from which Welsh and Breton descended, and the features of this proto-language are also hypothetical.

Ancestral languages (whether attested or reconstructed) can themselves be compared, and their own ancestors hypothesized and reconstructed, in a recursive fashion. Thus,



Figures 5-1 A selection of Indo-European languages, organized as a tree



Figures 5-2 An incorrect tree of Italo-Celtic languages

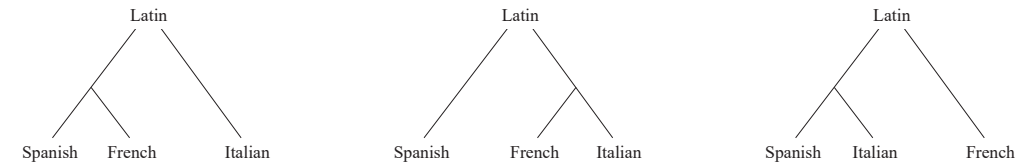
some linguists (e.g., Kortlandt 2007) believe that Latin and Proto-Brythonic ultimately descend from a language termed Proto-Italo-Celtic (PIC).² Repeatedly applying this process of comparison and reconstruction—called the Comparative Method—leads to proto-languages further and further back in time, ultimately ending in Proto-Indo-European (PIE).³

Granted that the uppermost node, Proto-Indo-European, is valid (since the Indo-European languages are indeed related to one another), on what basis are lower-level proto-languages (or equivalently, subgroups) posited? For example, why isn't Welsh grouped with Latin, separately from Breton, as in the fictitious Figures 5-2?

The reason is that this would imply that Latin and Welsh both exhibit certain changes (or *innovations*) from PIC (and hence, from PIE) that are not exhibited by Breton. But there are no notable innovations of this kind. Figures 5-2 would also imply that there are *no* innovations shared by Welsh and Breton which are not also shared by Latin (and all other members of the Italo-Celtic subgroup). This too is false: for example, the Brythonic languages changed **k^w* to *p*, and changed **s* to *h* at the beginnings of words (Schmidt 1993: 80–81); Latin, on the other hand, preserved these sounds intact. In sum, the representation in Figures 5-1 is more faithful to the empirical data we have from attested languages, than is Figures 5-2.

As we have just illustrated, in the Comparative Method, a subgroup is posited on the basis of *EXCLUSIVELY SHARED INNOVATIONS* among its members—a principle often attributed to Leskien (1876: xiii), but more accurately ascribed to Brugmann (1884: 231). In other words, a subgroup represents a hypothesis that all of its members share certain innovations that are not exhibited by any other language, and that any innovation that a member shares with a non-member is necessarily shared by *all* members. (This is similar to how, in phylogenetics, clades are interpreted as monophyletic groups defined by synapomorphies: see Skelton *et al.* 2002: 27–28.)

Let us now consider what happens when we add another language—French—to our tree. There is no question but that French is a descendant of Latin; hence it should ultimately be a daughter of the “Latin” node. However, there are multiple ways in which it could be put into a tree together with Spanish and Italian (Figures 5-3). Which of these



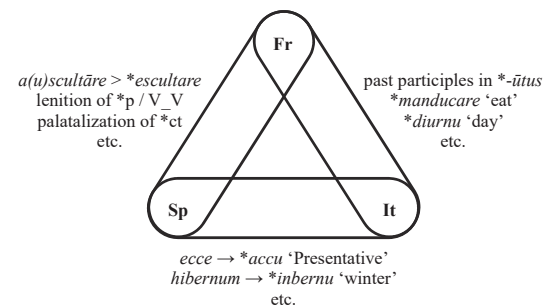
Figures 5-3 Three possible ways to represent the relations between Spanish, French and Italian

choices is correct?

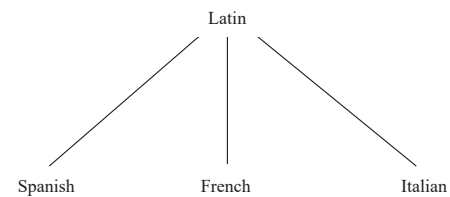
Choice 1, with Spanish and French forming a subgroup, seems justified by the innovations that are shared between these two languages, and not shared by Italian: for example, the irregular change of *a(u)scultāre* ‘listen’ to **escultāre* > Sp. *escuchar*, Fr. *écouter*, vs. It. *ascoltare* (Berger and Brasseur 2004:90); intervocalic lenition of **p*—e.g., *rīpa* ‘riverbank’ > Sp. *riba*, Fr. *rive*, vs. It. *ripa* (Posner 1996: 234); and the palatalisation of **ct* clusters—e.g., *factum* ‘done’ > Sp. *hecho*, Fr. *fait* vs. It. *fatto* (Hall 1950: 25). However, one can also find innovations shared by French and Italian but not by Spanish, which would argue in favour of choice 2: for example, the innovative weak past participle suffix **-ūtus* which affected many verbs—e.g. **sapūtus* ‘known’ > It. *saputo*, Fr. *su*, as opposed to Sp. *sabido* < **sapītus* (Alkire and Rosen 2010: 177); or numerous lexical innovations such as **diurnu* > It. *giorno*, Fr. *jour* ‘day’, replacing Lat. *diēs* (Sp. *día*), or **manducāre* ‘chew’ > It. *mangiare*, Fr. *manger* ‘eat’, replacing Lat. *comedere* (Sp. *comer*). Finally, one could cite evidence in favour of subgrouping Spanish and Italian together as opposed to French (as in choice 3), e.g., the irregular change of Lat. *ecce* to **accu* (Wüest 1994), as in the (feminine) distal demonstrative **accu-illa* > Sp. *aquella*, It. *quella*, where French preserves *ecce* (**ecce-illa* > Fr. *celle*); or the irregular insertion of */n/* in *hibernum* ‘winter’, yielding **inbernu* > Sp. *invierno*, It. *inverno*, vs. Fr. *hiver* (Alkire and Rosen 2010: 339). Many other examples of exclusively shared innovations⁴ could be found for each of the three language pairs. In all cases, the nature of the changes (especially phonological and morphological change, whether regular or irregular) is typical of the sort of evidence that is traditionally considered diagnostic of genealogical subgroups under the Comparative Method.

In this particular case, the data simultaneously support three intersecting subgroups (Figures 5-4): Spanish–French, French–Italian and Spanish–Italian. The tree model would force us to privilege one of these three groupings at the expense of the other two, but this would not do justice to the empirical evidence.⁵

One might be tempted to represent this problematic situation by resorting to the diagram in Figures 5-5, which does not necessarily commit us to any subgrouping hypothesis. This sort of diagram (cf. Ross 1997: 213) is sometimes used as an “agnostic” representation, which Pawley (1999) calls a “rake-like” structure, and van Driem (2001) likens to “fallen leaves”. (In phylogenetics this is known as “(soft) polytomy”:⁶ see Page and Holmes 2009: 13.) Yet it too is unsatisfactory, as it could be interpreted as claiming that there are *no* exclusively shared innovations between Spanish and French, between



Figures 5-4 Historical evidence supports three intersecting subgroups involving Spanish, French and Italian—a situation incompatible with the family tree model.



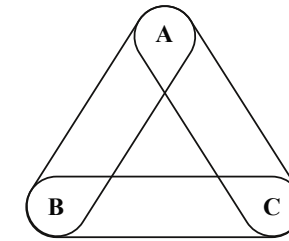
Figures 5-5 A rake (or “polytomy”)

French and Italian, or between Spanish and Italian, when—as we have seen—there is in fact solid, positive evidence for all of these groupings.

At best, a rake-like representation leaves us with the impression that science is simply incapable of unraveling the precise history of the language family. While this is sometimes the case due to lack of data, it is certainly not true in such a well-documented family as Romance. The history of individual changes across Romance dialects and languages is extremely well-known: if this family cannot be represented by a tree, then this cannot be due to a lack of data, but to the inherent flaws of the tree model itself—in particular, the axiom that genealogical subgroups defined by exclusively shared innovations are necessarily nested, and never intersect. This axiom results from an incorrect understanding of language change (cf. Bossong 2009; François 2014, 2017), namely that an innovation consistently results in total social isolation and lack of contact with communities that did not undergo the innovation—a wrong assumption in most of the world’s history. What we see, on the contrary, is that the spread of an innovation within part of a dialect network, insofar as it still allows mutual intelligibility with non-participating dialects, can perfectly well be followed by other innovations whose geographical scope may cross-cut its own (de Saussure 1995 [1916]: 273–278; Bloomfield 1933: 310–318), resulting in intersecting subgroups. We need a model of language relationships that is capable of accommodating such situations in a more accurate and faithful way than the tree model.

2.2 The Problem of Linkages

We can generalise our observations above by considering an abstract case, consisting of a family of three languages: A, B, and C. If A and B have some exclusively shared innovations, but neither B and C nor A and C do, then the situation is amenable to a tree representation (as in choice 1 in Figures 5-3 above). Historically, this represents a situation where the Proto-ABC speech community somehow split into two groups, one of which (the common ancestor of the modern A and B communities) underwent certain



Figures 5-6 When shared innovations intersect

linguistic innovations, separately from C; these innovations are said to have resulted in a hypothetical language “Proto-AB”. Later on, a similar split took place in the Proto-AB community, that resulted in the separate development of A and B.

But another situation is also possible, as we saw in the case of Romance languages. This is the situation where there are exclusively shared innovations not only between A and B, but also between B and C, and/or between A and C: that is, a situation in which shared innovations define intersecting groupings—see Figures 5-6 (and Figures 5-4 above).

This situation cannot be represented using the tree model, which assumes that a language can belong to one genealogical subgroup only. The only way to force the data into a tree—and posit, for example, a subgroup AB—would be to disregard the other two sets of innovations which contradict this grouping. Admittedly, such a procedure may be tenable in some cases. For example, C could have undergone some of the same innovations as A and B purely by chance, so that these are not really “shared innovations” in the relevant sense, but are rather “parallel innovations”. The trouble with this argument is that it is often extremely difficult to come up with positive evidence for it. In particular, if it is believed that C was still in contact with A and B at the time it underwent these innovations, it is unparsimonious to invoke independent, parallel development as an explanation: it is more probable that the changes they have in common reflect events of language-internal diffusion across dialects.

Another situation in which it may be reasonable to disregard the B–C and A–C innovations is when there is good reason to believe that these all occurred historically *after* the A–B innovations, and at a point in time when C had already become mutually unintelligible with A and B (i.e. had become a separate language). In this case, many historical linguists would label the B–C and A–C innovations as effects of “language contact”, and would disregard them for the purpose of representing genealogical relationships. This sort of reasoning only works under the assumption that it is possible to draw a principled line between diffusion across language boundaries (“contact”) and diffusion within them (“internal change”). This seems unlikely, given that the concept of a “language boundary” (i.e. whether two speech varieties are separate languages or simply dialects of the same language) is itself a gradient notion. However, the argument of contact is usually proposed in good faith, and may be accepted in some obvious cases,

namely when the genealogical distance between the speech varieties involved was already much too great at the time of contact for mutual intelligibility—e.g. lexical borrowings from Old Norse into Old English, or from Polynesian languages into other Oceanic languages (Biggs 1965).

In sum, given a set of changes with overlapping distributions, there are occasionally *bona fide* reasons for arguing that some of them are *not* genealogical in nature, and thus should be discarded for the purpose of subgrouping. In general, though, there is often no legitimate basis for deciding which changes may be so discarded. Sometimes, this is merely due to lack of evidence (historical or linguistic) about which set of changes predates the other. But in many cases, the problem is simply that the tree model fails to capture the fact that innovations do spread in entangled patterns across sets of mutually intelligible dialects, resulting in intersecting genealogical subgroups. This is what happens in dialect chains and networks, as well as in full-fledged language families that evolve out of dialect networks—which Ross (1988: 8; 1997: 213) calls *linkages*. The relationships among Spanish, French and Italian—or among other Romance languages, for that matter (with the possible exception of Romanian)—are typical of a linkage. Crucially, linkages are common throughout the world: similar configurations have been described, under various names, for Sinitic (Hashimoto 1992; Chappell 2001), Semitic (Huehnergard and Rubin 2011), Indo-Aryan (Toulmin 2009), Athabaskan (Krauss and Golla 1981; Holton 2011), Oceanic (Geraghty 1983; Ross 1988), and many other language families. In Section 4, we will be presenting a detailed example from a section of the Oceanic linkage.

When dealing with linkages, decisions about which innovation-defined groupings should be ignored for the purpose of representing genealogical relationships tend to be *ad hoc*, and debates rage with no sign of resolution. In our view, such problems are mere artefacts of the assumptions present in the tree model, and lack any legitimate basis as far as language change is concerned. In fact, there is no justification to the assumption that dialects and languages evolve primarily by splitting in a tree-like fashion: the more is known about language change, the more it becomes obvious that this model is a poor approximation of reality, resting as it does on a misleading metaphor.

In the remainder of this paper, we advance a more flexible approach: Historical Glottometry. It elaborates on the principles of the Comparative Method, yet attempts to liberate it from the misleading influence of the family-tree model, by proposing a representation that reflects historical reality more faithfully.

3. Defining Historical Glottometry

3.1 Intersecting Subgroups

Insofar as our method is meant to represent (past or present) dialect networks, it is useful to start by looking at how dialects are visualised by dialectologists. A key concept in dialectology is that of the *isogloss* (Chambers and Trudgill 1998: 89). Given a linguistic property that is distributed across a dialect network in a certain way, an isogloss is a line delimiting the dialects that share that property. Isoglosses can be represented on

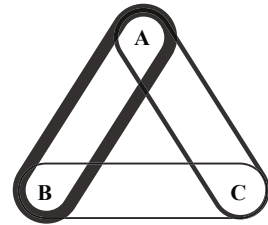
geographically realistic maps, or on more abstract figures. The lines in Figures 5-4 above are examples of isoglosses, showing the distribution of certain linguistic properties in (part of) the Romance family.

In principle, an isogloss may involve any property that is shared among languages, regardless of its historical origin. And indeed, because dialectology traditionally examines modern speech varieties from a purely synchronic perspective, isogloss maps often fail to distinguish between those similarities that result from shared innovations (*synapomorphies*) and those that are simply shared retentions from a common ancestor (*symplesiomorphies*), or even parallel innovations (*homoplasies*) and accidental similarities.⁷ From the perspective of historical linguistics, however, it is indispensable to restrict our observations to *shared innovations*: this is a pillar of the Comparative Method, also known as Leskien's (or Brugmann's) principle. The methodology we propose can be described as a *dialectological approach to language history*; it combines the precise descriptive tools of dialectology and dialectometry (Goebel 2006; Nerbonne 2010; Szmrecsányi 2011) with the powerful concepts of the Comparative Method—notably the stress on shared innovations.

One problem with isogloss maps (and admittedly the main reason why they have not been adopted more widely outside of dialectology) is that they become visually messy very quickly as more and more intersecting isoglosses are added; furthermore, they do not lend themselves to straightforward storytelling as much as a tree diagram would. The former issue, at least, can be addressed if we choose to use isoglosses to represent not individual innovations, but rather language groupings defined by one or more exclusively shared innovations—in other words, subgroups. A genealogical subgroup is a grouping of dialects or languages identified by a bundle of (innovation-defined) isoglosses; note, in passing, that nothing in this definition entails that subgroups should be discrete or nested as they are in a tree: they can perfectly intersect (François 2014: 170).

The thickness of the isogloss line can then be used to represent the strength of the evidence for each language grouping. For example, Figures 5-7 translates visually the fact that, while the three subgroups AB, AC and BC are all empirically supported, BC is the weakest pairing, and AB the strongest.

With such a configuration of the data, historical linguists who take the tree model for granted might be tempted to favour AB as the only valid subgroup, and dismiss the evidence for the two other subgroups altogether, under the assumption that these “weaker” groupings must be mere illusions—whether their similarities be due to “contact”, or to “parallel innovation”, etc. However, unless there is indeed a principled way of ruling out these isoglosses, it is wiser to keep them in the picture: the idea is that those innovations that are shared between A and C, or B and C, reflect historical events of shared linguistic development just as much as do those between A and B. It is just that the social relations between communities A and B, over the entire course of the history of the ABC family, have been stronger, more frequent or more sustained than those between other pairs of communities. Historical Glottometry can be used precisely as a means to explore and evaluate the strengths of historical connections between social groups, based on the linguistic traces they left in modern languages.



Figures 5-7 A representation of intersecting subgroups with relative weighting

In sum, linguistic linkages make it necessary to accept the idea of a language family in which genealogical subgroups have different strengths, and can cross-cut. Rather than a simplistic binary answer (*X forms vs. does not form a subgroup with Y*), subgrouping studies should allow for the possibility of *stronger vs. weaker subgroups*. Just as a village A may have more frequent mutual interaction with another village B than with C, likewise languages A and B can be said to form a stronger subgroup together (i.e., be “more subgroupy”) than languages A and C. Ideally, such claims could even be quantified—as in “A subgroups *n* times as strongly with B as it does with C”.⁸⁾

The crucial question is now: how can we define, and calculate, the “strength” of a subgroup? This is the object of the next subsection.

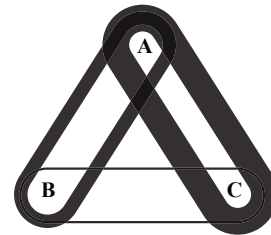
3.2 The Cohesiveness of Subgroups

The most obvious way to represent the strength of a subgroup using isoglosses would be to simply make the thickness of isoglosses directly proportional to the number of innovations defining the respective groupings. For example, suppose that in the above example of languages A, B and C, there were 12 innovations exclusively shared between A and B, 4 between A and C, and 2 between B and C: then our diagram would look exactly as in Figures 5-7 (where 1 shared innovation = 1 pixel).

However, suppose that instead of 4 exclusively shared innovations between A and C, there were 24. Our diagram would then be as in Figures 5-8.

Insofar as the thickness of their lines is exactly proportional to the number of exclusively shared innovations between each pair of languages, Figures 5-7 and 5-8 are accurate, fully-detailed representations of their respective data. However, they fail to represent an important fact: that the strength of the AB grouping in the first situation is *greater*, relative to the other isoglosses, than the strength of the same grouping in the second situation—*despite* the fact that the same number of defining innovations ($n=12$) is involved in both cases.

Interestingly, Pawley (2009: 13), discussing the factors that provide evidence for a particular subgrouping hypothesis, notes that “The weight of this evidence depends on the number and quality of the innovations concerned *and on the number and quality of innovations that have conflicting distributions*” (our emphasis). We thus need to quantify



Figures 5-8 Intersecting isoglosses, with more support for AC than for AB.

the strengths of groupings in a way that takes into account not only the absolute number of innovations that *support* the grouping, but also the number that *conflict* with it. An isogloss *x* is said to “conflict” with a subgroup *y* if they cross-cut each other—i.e. if and only if *x* contains some but not all members of *y*, and also contains members outside *y* (mathematically speaking: $x \cap y$, $x \setminus y$ and $y \setminus x$ are all nonempty). In our case, even though the AB grouping is supported by 12 innovations in both cases, it is *more strongly* supported in the first case (where the 12 innovations of AB conflict with only 4 isoglosses for AC plus 2 for BC) than in the second (where the number of conflicting isoglosses is $24 + 2$).

Drawing on the use of “relative identity weight”—also known as the “Jaccard coefficient”—in dialectometry (e.g. Goebel 2006: 412), we propose to define the “cohesiveness” of a subgroup as the *proportion* of supporting evidence with respect to the entire set of relevant evidence. Thus, for each given subgroup *G*, let *p* be the number of supporting innovations, and *q* the number of conflicting innovations. The total amount of evidence that is relevant for assessing the cohesiveness of *G* is $(p+q)$.⁹⁾ Now, if we call k_G the cohesiveness value of *G*, we have:

$$k_G = \frac{\text{number of supporting innovations}}{\text{total number of relevant innovations}} = \frac{p}{(p + q)}.$$

In the situation depicted in Figures 5-7, the cohesiveness of AB would be calculated as:

$$k_{AB} = \frac{12}{12+(4+2)} = \frac{12}{18} = \frac{2}{3} \approx 67\%.$$

This result can be translated into plain language by saying that, out of all the innovations that affected the subgroup AB (i.e. either encompassed the subgroup as a whole, or affected one of its members together with an external member), exactly two thirds confirmed the cohesion of AB as a subgroup, while one third contradicted it. More simply, A and B *evolved together* two-thirds of the time, and *evolved apart* one-third of the time.

In the situation depicted in Figures 5-8, the cohesiveness of AB would be:

$$k_{AB} = \frac{12}{12+(24+2)} = \frac{12}{38} \approx 32\%.$$

That is, in Figures 5-8, AB as a subgroup is confirmed 32% of the time, and contradicted 68% of the time.

These rates of 67% and 32% should be compared with the theoretical cohesiveness value which all subgroups are supposed to have in a “well-behaved” family tree, namely 100%. In an ideal tree, any group of languages defined by even a single shared innovation is supposed to *always* behave like a subgroup: that is, 100% of the innovations that affect it should confirm its cohesion, and there should be no genealogical innovation involving some (but not all) of its members together with some non-members. As we will see below with real data, this extreme figure of 100% is a convenient fiction

that is virtually never met with among real-life languages—at least not in the situation of a linkage. Rates of cohesiveness in most subgroups typically fall far short of this “ideal” (in our data, most cohesiveness values are between 10% and 30%). This does not mean that we are not dealing with genealogical subgroups at all; but rather, that this very notion must be redefined so as to accommodate the heterodox notion of the *strength* of a subgroup.

3.3 Subgroupiness

Given this measure of cohesiveness, we could use these values to determine the thicknesses of our isogloss lines. However, cohesiveness alone is not sufficient to provide an accurate representation of each subgroup’s strength: as we will see now, it is necessary to also consider the absolute number of exclusively shared innovations.

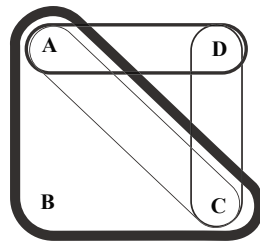
Consider now a family of four languages, A, B, C and D, where there are 12 innovations shared by ABC; 4 by AD; 2 by CD; and 1 by AC, as in Figures 5-9.

Note here that the number of innovations shared by AC (n=1) is irrelevant to the calculation of the cohesiveness of ABC, since this innovation neither confirms the subgroup nor contradicts it (see fn. 9). In order to assess the cohesiveness of ABC, what we need is to observe the number of innovations that confirm the subgroup (n=12) and those that clearly conflict with it—i.e. the innovations of AD (n=4) plus those of CD (n=2). The cohesiveness of ABC is thus:

$$k_{ABC} = \frac{12}{12+(4+2)} = \frac{12}{18} = \frac{2}{3} \approx 67\%.$$

Let us now calculate the cohesiveness of AC. This grouping is confirmed not only by the innovations that are exclusively shared by A and C (n=1), but also by those which they share non-exclusively, since these too show that languages A and C tend to undergo the same linguistic changes together. This includes, in Figures 5-9, the 12 innovations shared by ABC. As a result, the cohesiveness of the grouping AC should be as follows:

$$k_{AC} = \frac{12+1}{(12+1)+(4+2)} = \frac{13}{19} \approx 68\%.$$



Figures 5-9 A family of four languages.

In sum, the cohesiveness of AC is even *greater* than that of ABC. Yet we would not want to say that AC is a “stronger” subgroup than ABC, because the latter has a far greater number of *exclusively* shared innovations.

Our proposed solution to this problem is to use the absolute number of *exclusively shared innovations* as the main point of reference, and *weight* it using the subgroup’s cohesiveness rate (k). For each given subgroup G , let ε (‘epsilon’) be its number of exclusively shared innovations; p its number of supporting innovations (i.e. shared innovations, whether exclusively or not), and q the number of conflicting innovations. We already saw that the *cohesiveness* rate is $k = \frac{p}{(p+q)}$. We now propose to define the *subgroupiness* of a language cluster (call it ‘sigma’, ζ) as the product of the cohesiveness rate (k) with the number of exclusively shared innovations (ε):

$$\zeta = \varepsilon \times k = \varepsilon \times \frac{p}{(p+q)}.$$

For example, if we come back to the comparison of Figures 5-7 and 5-8, we can now weight the absolute number of innovations exclusively shared by A and B (ε_{AB}) using AB’s cohesiveness rate k_{AB} (given above), and thus calculate its subgroupiness ζ_{AB} .

$$\text{In Figures 5-7: } \zeta_{AB} = 12 \times \frac{12}{18} = 8.$$

$$\text{In Figures 5-8: } \zeta_{AB} = 12 \times \frac{12}{38} \approx 3.79.$$

These numbers constitute exact measurements of the extent to which AB is a more strongly-supported subgroup in the first case than in the second case. (In other words, we can now say, “AB is more than twice as strongly supported—or more simply, *more than twice as subgroupy*—in Figures 5-7 than in Figures 5-8”.) As for Figure 5-9, we find that

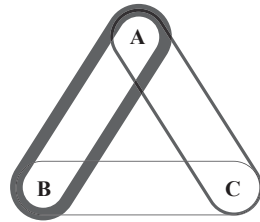
$$\zeta_{ABC} = 12 \times \frac{12}{18} = 8 \text{ and } \zeta_{AC} = 1 \times \frac{13}{19} = \frac{13}{19} \approx 0.68;$$

in other words, ABC is more than eleven-and-a-half times as subgroupy as AC. These results are consistent with the intuition that the subgroup ABC is more strongly supported than AC. In conclusion, subgroupiness constitutes the best criterion we have found for assessing the relative strengths of the genealogical subgroups in a language family.

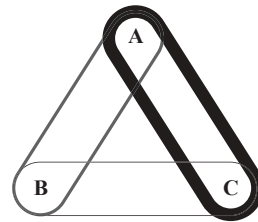
3.4 A Visual Representation

In terms of visual representation, we propose to draw lines around subgroups, and make their thickness proportional to their calculated subgroupiness ζ . Figures 5-7’–5-9’ show our proposed representations of the situations depicted in Figures 5-7–5-9, respectively. We call these kinds of figures *historical glottometric diagrams* (or ‘glottometric diagrams’ for short).

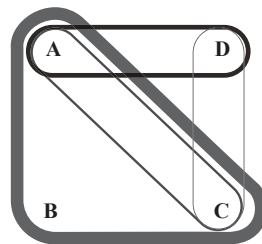
The examples given in this section were abstract, and simple in the sense that they involved small numbers of languages and of innovations. But the same tools can be



Figures 5-7' Illustration of subgroupiness-based isogloss thickness for the situation depicted in Figure 7. Subgroupiness rates: $\zeta_{AB}=8$; $\zeta_{AC}=0.89$; $\zeta_{BC}=0.22$.



Figures 5-8' Illustration of subgroupiness-based isogloss thickness for the situation depicted in Figure 8. Subgroupiness rates: $\zeta_{AB}=3.79$; $\zeta_{AC}=15.16$; $\zeta_{BC}=0.11$.



Figures 5-9' Illustration of subgroupiness-based isogloss thickness for the situation depicted in Figure 9. Subgroupiness rates: $\zeta_{ABC}=8$; $\zeta_{AD}=0.84$; $\zeta_{AC}=0.68$; $\zeta_{CD}=0.21$.

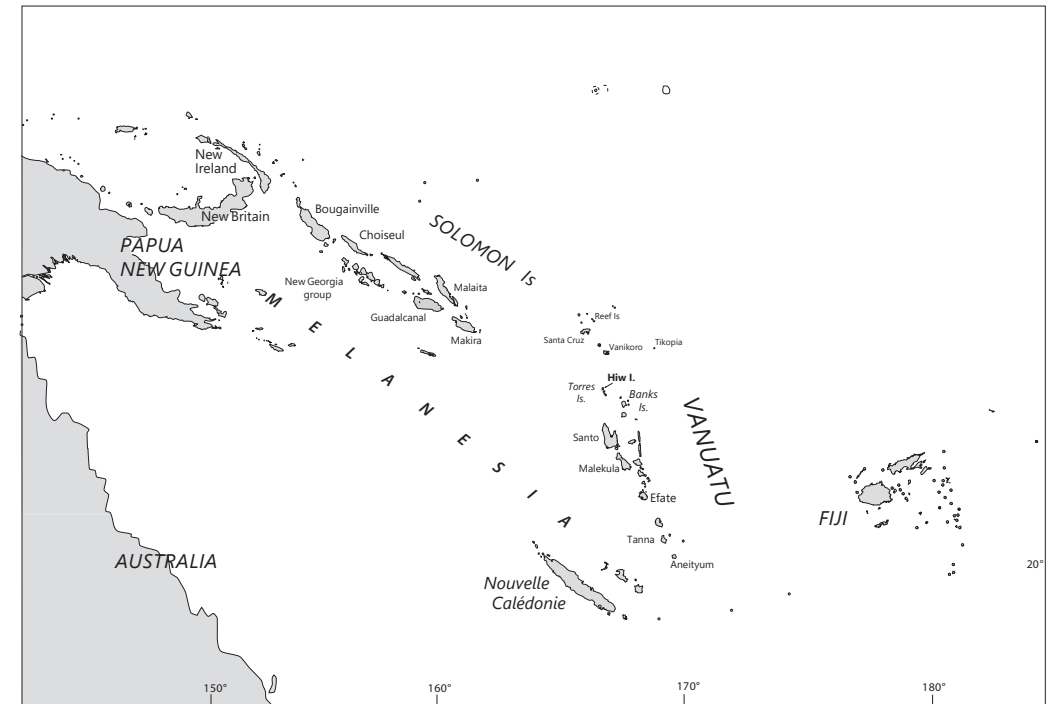
profitably applied to a much richer set of data. The next section will show precisely how Historical Glottometry can be applied to a real dataset involving 17 languages, and a total of 474 innovations.

4. A Case Study from North Vanuatu

4.1 The Languages

We can now illustrate Historical Glottometry using actual data from the languages of Vanuatu, an archipelago in the south Pacific (see Map 5-1).

Vanuatu is home to 138 indigenous languages, all members of the Oceanic branch of the Austronesian language family (François et al. 2015). The evidence for Oceanic being a (classical, nearly 100% cohesive) subgroup of Austronesian is massive (Pawley and Ross 1995; Ross 1988). It is widely accepted that there was at some point a more or



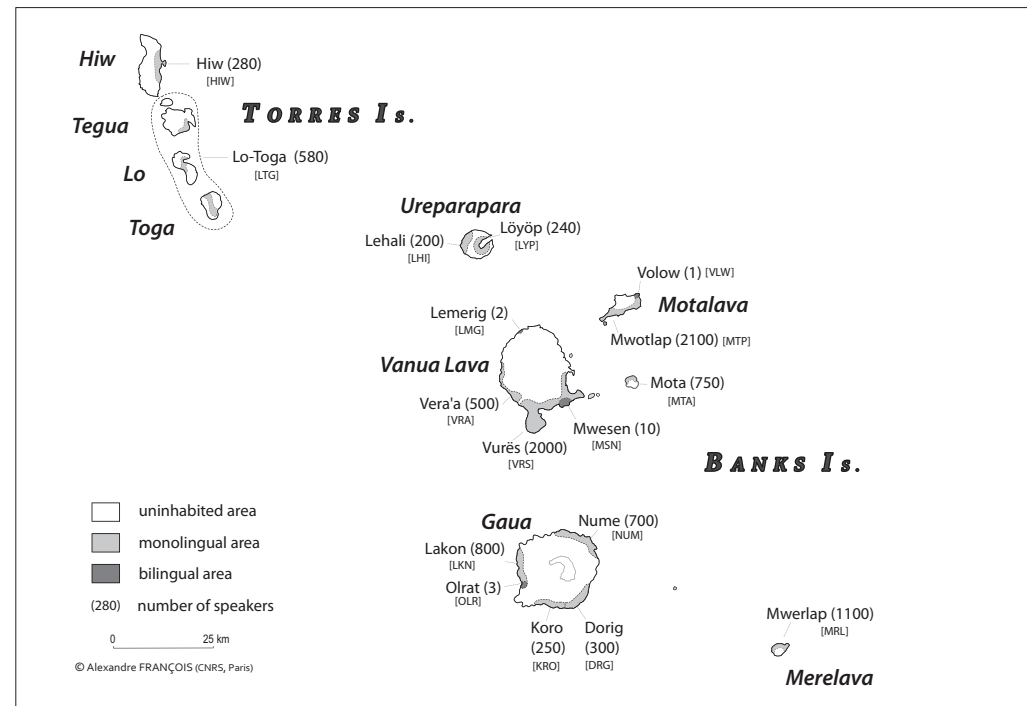
Map 5-1 The archipelago of Vanuatu, in the South Pacific [François (2011a: 181)]

less homogeneous Proto-Oceanic language spoken throughout most of the south Pacific (Pawley and Green 1984; Pawley 2008, 2010), which gradually fragmented into dialects and then independent languages—following a scenario quite similar to the history of the Romance languages. Over the decades, there have been many attempts to fit the modern-day languages of Vanuatu into a tree model. Clark (2009: 49) lists as many as nine conflicting subgrouping hypotheses, none of which has reached consensus. This tends to confirm our hypothesis that the genealogical relations among Vanuatu languages cannot be rendered by a tree: they constitute a *linkage*, i.e. a group of modern languages which emerged through the *in-situ* diversification of an earlier dialect network (Tryon 1996; François 2011a; 2011b; 2016; François et al. 2015).

We will be focusing on the two northernmost island groups of the Vanuatu archipelago, the Torres and Banks Islands. The second author (François) has been conducting fieldwork there since 1997, and has collected extensive data on the 17 languages still spoken in this small area, many of which are endangered (see François 2012). The names of these languages are given on Map 5-2, together with three-letter abbreviations and numbers of speakers.

4.2 Intersecting Isoglosses in North Vanuatu

The communalects (to use a term that is neutral between “language” and “dialect”) of

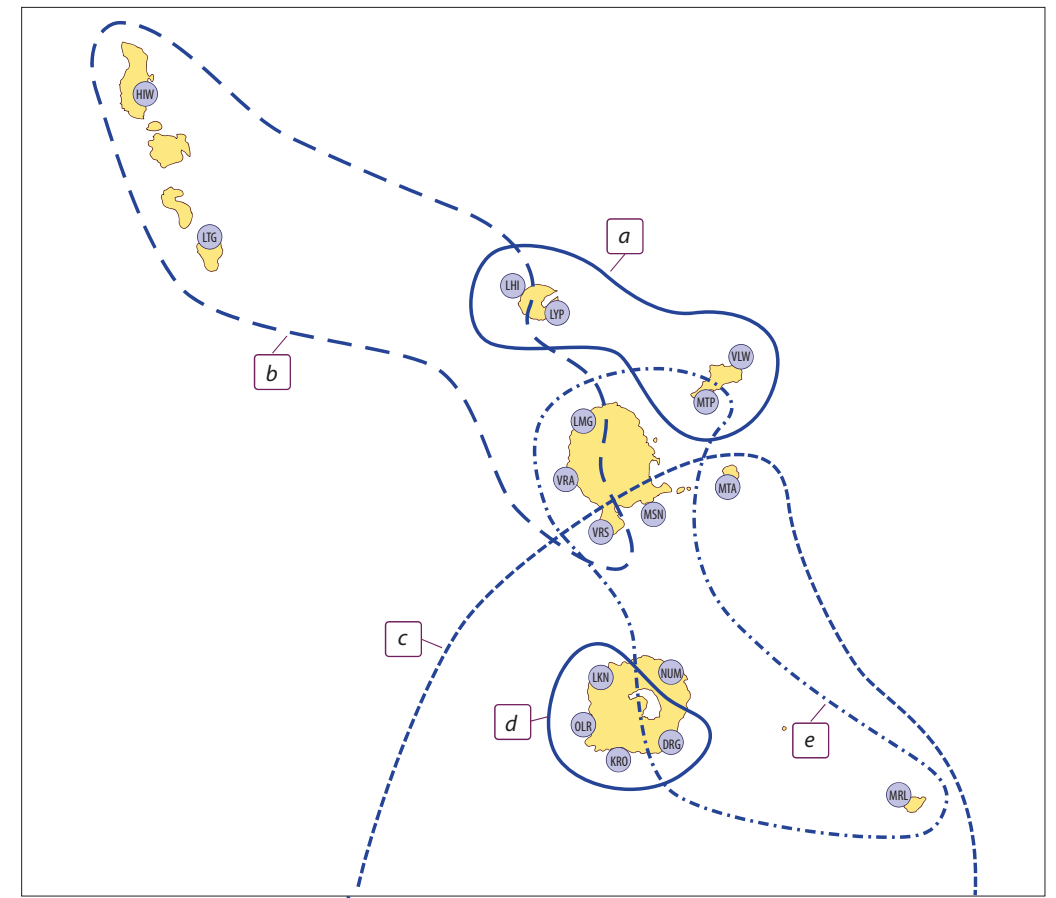


Map 5-2 The 17 languages of the Torres and Banks Islands, in northern Vanuatu [François (2011a: 181)]

northern Vanuatu have now lost mutual intelligibility, and constitute distinct languages. The Comparative Method makes it possible to unravel the various linguistic changes that took place since the time of earlier linguistic unity, and brought about the present linguistic diversity (François 2005; 2011a; 2011b; 2016). Even though some changes affected a single communalect in isolation, the most typical case was for a given innovation to emerge in some location, and diffuse via social interaction from one dialect to its neighbors, until it settled down to a certain portion of the dialect network. Some isoglosses encompassed the entire area, while others only targeted a set of four or five villages. And of course, in a manner similar to Romance dialects, what we see is that the isoglosses defined by the various innovations cross-cut each other.

The innovations under discussion here are of various kinds (François 2011a: 192–211). They include regular phonological change; irregular sound change (which affects one or a few words rather than applying across the lexicon); morphological change; syntactic change; and lexical replacement. Map 5-3 shows a selection of isoglosses for the following five innovations:

- a) Regular sound change: $*r > /j/$
- b) Irregular sound change: $*malate \rightarrow *malete$ ‘broken’
- c) Irregular sound change: $*\text{ʔaŋa}Ri \rightarrow *ʔaŋai$ ‘almond’



Map 5-3 Five isoglosses in the Torres–Banks Islands

- d) Morphological change: metathesis in trial pronouns
(Plural+three \rightarrow three+Plural)
- e) Morphological change: $*toya$ ‘stay’ \rightarrow Prohibitive

Map 5-3 makes it clear that isoglosses in the Torres and Banks languages—like those in the Romance family—constantly intersect.¹⁰ There is no way in which the genealogical relations among these languages could be represented by a tree. François (2004) was an attempt to do precisely this; while a tentative tree was indeed proposed, the number of issues raised (conflicting evidence, intersecting isoglosses, the need to constantly resort to *ad hoc* hypotheses to preserve the tree structure) was a preliminary sign of the inadequacy of the cladistic approach in this part of the world.

What we need here, then, is a Historical Glottometry approach, which will tell us, amongst the 131,070 ($=2^{17-2}$) potential groupings involving these languages, which ones actually exist, and which ones constitute the strongest subgroups. That these subgroups will probably intersect is to be expected, and is no longer a problem: as we have argued,

there is good reason to believe that this is the default situation in most language families. What we need is simply to go beyond the observation of individual isoglosses as in Map 5-3, and to be able to base our calculations on a rich database.

4.3 Identifying Innovations

4.3.1 Applying the Comparative Method

Our dataset consists of a table of 474 separate innovations which A. François identified in these 17 languages. For each linguistic feature considered, systematic comparison was conducted among languages of the sample as well as with other Oceanic languages, following principles of the Comparative Method, to establish the ancestral state of each property in the languages' shared ancestor (Proto-Oceanic, or a close variant thereof) as well as the direction of change.

Some cases make it relatively easy to determine what the innovation was. For example, consider the words for 'almond': whereas the eight languages to the north reflect the proto-form *ʔaŋaRi (e.g., Vera'a *ŋar*), the languages further south reflect a form *ʔaŋai (e.g., Vurës *ŋe*). The latter protoform shows the irregular loss of *R, a frequent yet lexically-specific sound change in the area (François 2011b). It is clearly an innovative form, whose distribution in the Banks Islands is represented by isogloss (c) in Map 5-3.

In other cases, identifying the innovation requires more reflection. For example, most of the northern Vanuatu languages have an adjective meaning 'broken', with forms that are cognate with each other:

(1) 'broken': HIW *mjit*; LTG *məlit*; LHI *melet*; LYP *malat*; VLW *malat*; MTP *malat*; LMG *mele?*; VRA *muli?*; VRS *mlit*; MSN *malat*; MTA *malate*; NUM *malat*; DRG *mlat*; MRL *melet*.

One can show that these modern forms go back to two distinct proto-forms: *malate and *malete. This conclusion is based on our knowledge of regular sound changes in this area, established using the Comparative Method (François 2005). This allows us to discern even those cases where two cross-linguistic homophones derive from different etyma: for example, while Lehali /melet/ necessarily reflects *malete, the same surface form /melet/ in Mwerlap is a regular reflex of *malate, because a stressed /a/ followed by an unstressed /e/ in the next syllable regularly underwent umlaut in this language (*aCe > /εC/). Knowledge of each language's phonological history likewise enables us to link each modern form in (1) to one, and only one, of the two proto-forms—either *malate or *malete. The next, crucial step consists in determining which of these two is conservative, and which one is innovative. External evidence is indispensable here, and shows that other Oceanic languages outside the Torres–Banks area point to the form with /a/: e.g., Araki /ŋalare/ 'broken' < *malate (François 2002: 270). In sum, the innovation we are concerned with here is a lexically-specific, irregular sound change whereby *malate became *malete, and not the other way around. The languages that participated in this particular innovation are: Hiw, Lo-Toga, Lehali, Lemerig, Vera'a and Vurës. This innovation is represented with isogloss (b) in Map 5-3.

4.3.2 Creating the Dataset

The sort of reasoning illustrated above, which follows a rigorous application of the Comparative Method, was used to identify all 474 innovations. The distribution of innovations into various types is presented in Table 5-1.

Among these types of changes, we consider irregular sound change and morphological change to be the most diagnostic of historical relatedness (following Greenberg 1957: 51; Ross 1988: 12), because they are least likely to be independently innovated. Lexical material is often excluded from subgrouping studies under the assumption that it is easily borrowable; to avoid this (perceived) problem, we have included here only those lexical replacements which can be shown to predate events of (regular or irregular) sound change.¹¹⁾

Figures 5-10 shows what the final database looks like. The 17 languages are sorted from north-west to south-east; each row corresponds to one innovation, and indicates whether there is positive evidence that a language participated (1) or did not participate (0) in that innovation. An empty box (–) was used when the data is inconclusive, non-applicable, or simply lacking. Altogether, the database contains 2728 positive ('1'), 5040 negative ('0') and 290 agnostic ('–') data points.

Note that each pattern of 1s and 0s corresponds to a diffusion area, and would be represented with an isogloss. We will now illustrate the application of Historical Glottometry to this database, following the methods explained in the previous section.

4.4 The Results

4.4.1 Numerical Results

The first thing we can do with this dataset is to measure cohesiveness for clusters of two languages. This measure of "pairwise cohesiveness",¹²⁾ applied to all pairs of languages (17²=289), yields the results in Table 5-2.

The figures of 100% in the diagonal simply say, as it were, that a language always subgroups perfectly with itself; these can thus be disregarded. More instructive is the observation that the cohesiveness *k* of language pairs tends to vary a lot, but with the highest figure being only 92%. The colored (yellow and orange) cells indicate rates of 50% and above, i.e. pairs with relatively high cohesiveness.

To illustrate the proper interpretation of the table, the figure of 92%, between Volow

Table 5-1 Typology of innovations represented in our North Vanuatu database

NATURE OF CHANGE	NUMBER	PROPORTION
Regular sound change	21	4%
Irregular sound change	116	25%
Morphological change	91	19%
Syntactic change	10	2%
Lexical replacement	236	50%
<i>Total</i>	474	100%

and Mwotlap, indicates that when either of these languages underwent a change (together with some other language), it shared this change with the other member of the pair 92% of the time. Table 5-2 thus shows that languages share innovations with their immediate neighbours a lot of the time—yet they do so at varying rates.

These figures, incidentally, are a valuable result in themselves, as they provide an empirical measurement of how much two languages have evolved together throughout their history. For example, the fact that Lo-Toga (#2) and Lehali (#3) shared only 41% of their innovations together points to a rather strong social divide between the Torres islands on the one hand, and the Banks islands on the other hand: clearly, the Lo-Toga community has had much less social interaction with Lehali ($k=41\%$) than with Hiw ($k=83\%$). Likewise, it is instructive to observe that, even though the language Vurës is spoken only a couple of hours' walk geographically from Vera'a (see Map 2), the two languages share together no more than 58% of their innovations; the historical links were much stronger, on the one hand, between Vera'a and Lemerig ($k=75\%$), and on the other hand, between Vurës and Mwesen ($k=85\%$). Interestingly, these figures closely match the intuitive feel one gets when learning and comparing the languages of Vanua Lava, as well as the islanders' own impressions; except that the figures have the advantage of being precise, and directly comparable with one another.

In order to deserve the status of genealogical subgroup, a cluster of languages needs to be “attested” historically, i.e. have at least one exclusively shared innovation ($\epsilon \geq 1$). A subgroup uniting Volow and Löyöp, for example, would have high cohesiveness (73%) if it existed; but because no innovation happens to be shared exclusively by these two languages, they cannot count together as a subgroup. Pairings that are not supported by at least one isogloss appear here in orange. Conversely, the yellow cells in Table 5-2 correspond to those higher-cohesiveness pairings ($k \geq 50\%$) which are actually attested as subgroups: e.g. Hiw-Lo-Toga with 83%, Lehali-Löyöp with 71%, etc.

We applied the same formula to calculate the cohesiveness (k) of all attested clusters of North Vanuatu, of any size. In total, the number of unique innovation-defined subgroups was 143. This figure includes the 15 pairs of languages shown in yellow in Table 5-2 above, but also clusters of various sizes, up to 15 members. The results, which cannot all be presented here for lack of space, were useful for the next stage: the calculation of subgroupiness values (ζ).

4.4.2 A Glottometric Diagram

We calculated the subgroupiness of all 143 attested language clusters, by applying the principles presented in §3 above. The 15 subgroups with the highest subgroupiness values are listed in Table 5-3.

In terms of visual representation, the abundance of subgroups of varying strengths made it necessary to represent only the strongest ones—we chose to show only those whose subgroupiness value is greater than or equal to 1 ($\zeta \geq 1$). This includes the 15 subgroups listed in Table 5-3, plus 17 others. We then represented each subgroup's strength by having line thickness proportional to its subgroupiness. In addition, the degree of redness (brightness value of the contour line) was made proportional to its

Table 5-3 The 15 strongest subgroups in the Torres–Banks linkage

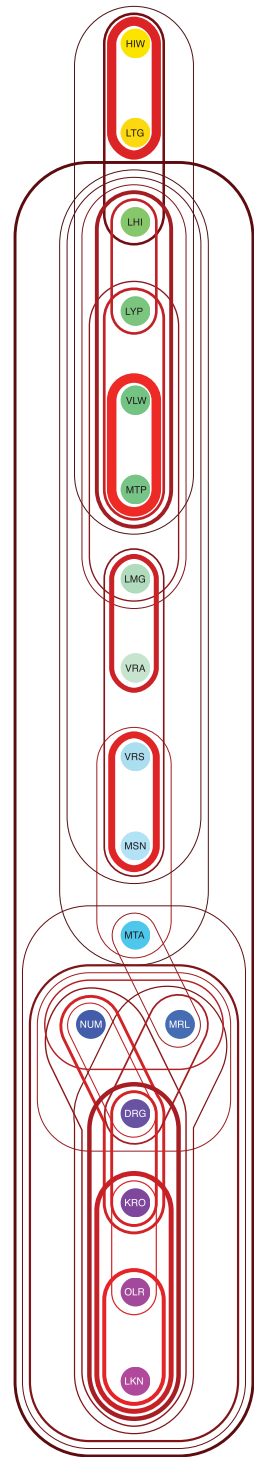
subgroups	subgroupiness
Volow–Mwotlap	12.82
Hiw–LoToga	12.45
Vurës–Mwesen	9.34
Lemerig–Vera'a	6.78
Koro–Olrät–Lakon	6.63
Dorig–Koro–Olrät–Lakon	6.01
Olrät–Lakon	5.34
Lehali–Löyöp–Mwotlap–Volow	5.22
15 Banks languages (LHI → HLKN)	3.92
Dorig–Koro	3.90
Löyöp–Volow–Mwotlap	3.64
Lehali–Löyöp	3.53
Hiw–LoToga–Lehali	3.43
southern Banks (Mwerlap + Gaua)	2.99
Dorig–Mwerlap	2.37

cohesiveness, with more cohesive subgroups appearing more intensely red. The final result was a comprehensive *glottometric diagram* of the whole region (Figures 5-11).

This result would warrant more commentary than is possible in this paper;¹³ we will stick to the essentials. First of all, the subgroupiness values, as well as the map derived from them, confirm the statement in §4.2, that the languages of northern Vanuatu form a *linkage* in which isoglosses, and hence subgroups, constantly intersect. Lehali (LHI), for example, subgroups both with the two Torres languages to its north ($\zeta=3.43$) and with the other Banks languages to its south ($\zeta=3.92$). Similarly, Mota (MTA) forms the bridge, as it were, between a northern Banks subgroup (running from Lehali to Mota, $\zeta=1.03$) and a distinct southern Banks subgroup (running from Mota to Lakon, $\zeta=1.30$). No family tree could ever account for this situation.

It is worthy of notice that the glottometric approach can *also* detect and represent those situations which are “tree-like”: for example, Volow and Mwotlap form a subgroup clearly separated from Löyöp; Vurës and Mwesen also clearly belong together. But evidently, these tree-like patches are a rarity in a language network which is strongly non-tree-like.

Another important result is the observation that Torres–Banks languages generally pattern in a geographically coherent way: all languages adjacent on the glottometric diagram are also adjacent geographically (though not vice versa; see below). This is even true for the non-linear part of the map, involving the four languages Mota–Nume–Dorig–Mwerlap: all the language pairs attested there (MTA–MRL, NUM–MRL, NUM–DRG, MRL–DRG) correspond to adjacent languages on Map 5-2. It is impossible to capture



Figures 5-11 A glottometric diagram of the Torres-Banks languages

such tight geographical organisation using a tree: any binary tree of 17 languages will allow 65,536 ($=2^{16}$) possible linear orderings of languages.

Expected though it may be, this consistency between language history and geography is a valuable result: for it shows that the languages' anchoring in space must have remained stable over the three millennia of their historical development, with limited inter-island migration (François 2011b: 181). Applying Glottometry to historically more turbulent families would make it possible to detect the genealogical relations that hold between languages *in spite* of their geographic locations, as accurately as the Comparative Method on which this method is based.

And indeed, a finer grain of observation reveals certain non-trivial patterns in our data that do more than just index geography. For example, even though Volow's location is closer to Mota than to Löyöp (Map 5-2), the position of the three languages in the diagram shows that Volow and Mota are genealogically quite remote ($k=36\%$). Evidently, the ancient societies of Motalava and Mota islands had very few direct social interactions with each other, and much more with the other islands (Ureparapara, Vanua Lava) located to their west. Such a result illustrates the potential of the method to reconstruct the shape of past social networks.

5. Conclusion

In conclusion, our newly proposed method of Historical Glottometry allows us to escape the false dichotomies of the tree model by allowing us to posit intersecting subgroups, and to quantify the *strength* of the genealogical evidence in favour of each language cluster.

If we were to use a tree to represent our data, we would certainly be able to capture certain salient organising features, e.g. the split between the two Torres languages (Hiw and Lo-Toga) and

all the languages to the south. But a tree would only be able to provide a very distorted picture of the social history of the region—as an orderly sequence of migrations with loss of contact—while the story told by the data (made visible to us by the glottometric diagram) is a much richer and more varied narrative of social interaction in which languages converge as much as they diverge. Far from the approximations imposed by the assumptions of the tree model, we hope to have shown the way towards a more accurate and realistic representation, which stays true to the most valuable insights of the Comparative Method.

Notes

- 1) The authors are grateful to Malcolm Ross, Mark Donohue and Martine Mazaudon for their comments on an earlier draft of this paper. They would also like to thank the participants of the symposium *Let's Talk about Trees* for their valuable questions and feedback. This work is part of the ANR “Labex” program *Empirical Foundations of Linguistics* – and of its axis *Typology and dynamics of linguistic systems*. Unless otherwise stated, all tables, figures, and maps were compiled by the authors.
- 2) Brythonic is a branch of Celtic, which in turn is a branch of Italo-Celtic; likewise, Latin is a member of the Italic branch of Italo-Celtic. The fact that the existence of Proto-Italo-Celtic is controversial is irrelevant to the present demonstration—what is important is that Latin and the Brythonic languages do in fact have a common ancestor (even if that ancestor turns out to be nothing other than Proto-Indo-European itself).
- 3) On general principles of the comparative method, see Hock (1991), Campbell (2004), Crowley and Bowerman (2010), among many others.
- 4) Obviously, the term “exclusively” must be understood within the restricted set of three languages taken here for the sake of discussion. Some of the innovations shared by French and Spanish are also shared with Catalan, Portuguese, etc., but this is not relevant for the present demonstration. (Interestingly, Catalan seems to exhibit most of the innovations mentioned.)
- 5) This is what Hall (1950) does: his assumption that languages must evolve following a cladistic model has him force the data into a tree structure. His “Western Romance” node, by grouping French and Spanish together, arbitrarily favours only one of the three groupings outlined here, and deliberately ignores any conflicting evidence.
- 6) We are grateful to Nobuhiro Minaka (p.c.) for pointing this out. In phylogenetic terms, a rake is ambiguous between “soft polytomy” (where the rake structure reflects lack of data or lack of certainty) and “hard polytomy” (which involves a claim that the actual structure of the data is inherently rake-like).
- 7) Important exceptions include the “dialect map of the Indo-European languages” in Anttila (1989: 305), which is extremely similar in spirit to the representation we will be proposing below, as well as the diagrams in Southworth (1964), which are less so. We are grateful to Malcolm Ross for having brought these works to our attention.
- 8) A further extension of our model, which we will not have room to develop in this study, could be to provide both *quantification* and *qualification* to genealogical relations. Thus, one could

imagine statements along the lines of “A subgroups with B twice as strongly as it does with C as far as regular sound change is concerned; but it does so 1.6 times more with C than with B with respect to verbal morphology, 3 times with respect to lexical replacement in basic vocabulary”, etc.

- 9) Those innovations that are entirely nested within a subgroup (e.g., those that affected only the language B within AB, and no language outside AB) are irrelevant to the cohesiveness of that subgroup, and therefore do not take part in the calculations.
- 10) Note that one innovation, namely (c), involves not only a subset of the Banks languages, but also languages further south in Vanuatu (François 2011b: 157). The metathesis of pronouns (innovation d) is described in François (2016: 51).
- 11) This is the same reasoning that validates **manducāre* ‘eat’ as a legitimate example of an early lexical innovation shared by French and Italian (§2.1), because it reflects regular sound changes diagnostic of inherited vocabulary (compare French *manger* /mãʒe/ < **manducāre* with *venger* /vãʒe/ ‘avenge’ < **vindicāre*). By contrast, a recent Italian loanword such as *caporal* (‘corporal’), which does not exhibit any such sound changes, would not normally qualify as diagnostic evidence for subgrouping.
- 12) This is similar to the concept of “Relative Identity Weight” in the Salzburg school of dialectometry (Goebel 2006: 412); it is also sometimes known as the “Jaccard coefficient”.
- 13) The colors of the dots representing the languages are also significant; they are obtained by transforming the pairwise cohesiveness values into distances, performing multidimensional scaling in three dimensions, and assigning the three axes to red, green and blue. This is a procedure commonly used in dialectometry (e.g. Heeringa 2004: 161).

References

- Alkire, T. and C. Rosen
2010 *Romance Languages: A Historical Introduction*. Cambridge: Cambridge University Press.
- Anttila, R.
1989 *Historical and Comparative Linguistics* 2nd ed. Amsterdam: John Benjamins.
- Berger, R. and A. Brasseur
2004 *Les Séquences de Sainte Eulalie*. Geneva: Droz.
- Biggs, B.
1965 Direct and Indirect Inheritance in Rotuman. *Lingua* 14: 383–415.
- Bloomfield, L.
1933 *Language*. New York: Henry Holt.
- Bossong, G.
2009 Divergence, Convergence, Contact: Challenges for the Genealogical Classification of Languages. In K. Braunmüller and J. House (eds.) *Convergence and Divergence in Language Contact Situations*, pp.13–40. Amsterdam: John Benjamins.
- Brugmann, K.
1884 Zur Frage nach den Verwandtschaftsverhältnissen der indogermanischen Sprachen.

Internationale Zeitschrift für allgemeine Sprachwissenschaft 1: 226–256.

- Bryant, D., F. Filimon, and R. D. Gray
2005 Untangling Our Past: Languages, Trees, Splits and Networks. In R. Mace, C. Holden, and S. Shennan (eds.) *The Evolution of Cultural Diversity: Phylogenetic Approaches*, pp.69–85. London: UCL Press.
- Campbell, L.
2004 *Historical Linguistics: An Introduction*. Cambridge: MIT Press.
- Chambers, J. K. and P. Trudgill
1998 *Dialectology* (Cambridge Textbooks in Linguistics.) Cambridge: Cambridge University Press.
- Chappell, H.
2001 Language Contact and Areal Diffusion in Sinitic Languages. In A. Aikhenvald and R. M. W. Dixon (eds.) *Areal Diffusion and Genetic Inheritance: Problems in Comparative Linguistics*, pp.328–357. Oxford: Oxford University Press.
- Clark, R.
2009. **Leo Tuai: A Comparative Lexical Study of North and Central Vanuatu Languages* (Pacific Linguistics 603). Canberra: Australian National University.
- Croft, W.
2000 *Explaining Language Change: An Evolutionary Approach*. London: Pearson Education.
- Crowley, T. and C. Bower (eds.)
2010 *An Introduction to Historical Linguistics*, 4th ed. Oxford: Oxford University Press.
- van Driem, G.
2001 *Languages of the Himalayas: An Ethnolinguistic Handbook of the Greater Himalayan Region*.
- Enfield, N.
2008 Transmission Biases in Linguistic Epidemiology. *Journal of Language Contact* 2(1): 299–310.
- François, A.
2002 *Araki: A Disappearing Language of Vanuatu* (Pacific Linguistics 522). Canberra: Australian National University.
2004 Subgrouping Hypotheses in North Vanuatu. Paper presented at the Sixth International Conference on Oceanic Linguistics (COOL6). University of the South Pacific, Port Vila, Vanuatu. July 8, 2004.
2005 Unraveling the History of the Vowels of Seventeen Northern Vanuatu Languages. *Oceanic Linguistics* 44(2): 443–504.
2011a Social Ecology and Language History in the Northern Vanuatu Linkage: A Tale of Divergence and Convergence. *Journal of Historical Linguistics* 1(2): 175–246.
2011b Where *R They All? The Geography and History of *R-Loss in Southern Oceanic Languages. *Oceanic Linguistics* 50(1): 140–197.
2012 The Dynamics of Linguistic Diversity: Egalitarian Multilingualism and Power Imbalance among Northern Vanuatu Languages. *International Journal of the Sociology of Language* 214: 85–110.
2014 Trees, Waves and Linkages: Models of Language Diversification. In C. Bower and B.

- Evans (eds.) *The Routledge Handbook of Historical Linguistics*, pp.161–189. New York: Routledge.
- 2016 The Historical Morphology of Personal Pronouns in Northern Vanuatu. In Konstantin Pozdnyakov (ed.) *Reconstruction et classification généalogique: Tendances actuelles. Faits de Langues*, pp.25–60. Bern: Peter Lang.
- 2017 Méthode comparative et chaînages linguistiques: Pour un modèle diffusionniste en généalogie des langues. In J. L. Léonard (ed.) *Diffusion: implantation, affinités, convergence* (Mémoires de la Société de Linguistique de Paris, XXIV), pp.43–82. Louvain: Peeters.
- François, A., S. Lacrampe, M. Franjeh, and S. Schnell (eds.)
- 2015 *The Languages of Vanuatu: Unity and Diversity* (Studies in the Languages of Island Melanesia). Canberra: Asia Pacific Linguistics Open Access.
- Geraghty, P. A.
- 1983 *The History of the Fijian Languages* (Oceanic Linguistics Special Publication 19). Honolulu: University of Hawaii Press.
- Goebel, H.
- 2006 Recent Advances in Salzburg Dialectometry. *Literary and Linguistic Computing* 21(4): 411–435.
- Goodman, M., J. Czelusniak, G. W. Moore, A. E. Romero-Herrera, and G. Matsuda
- 1979 Fitting the Gene Lineage into Its Species Lineage: A Parsimony Strategy Illustrated by Cladograms Constructed from Globin Sequences. *Systematic Biology* 28(2): 132–163.
- Gray, R. D., D. Bryant, and S. J. Greenhill
- 2010 On the Shape and Fabric of Human History. *Philosophical Transactions of the Royal Society London B* 365: 3923–3933.
- Greenberg, J. H.
- 1957 *Essays in Linguistics*. Chicago: University of Chicago Press.
- Hall, R. A. Jr.
1950. The Reconstruction of Proto-Romance. *Language* 26(1): 6–27.
- Hashimoto, M. J.
- 1992 Hakka in Wellentheorie Perspective. *Journal of Chinese Linguistics* 20: 1–49.
- Haspelmath, M.
- 2004 How Hopeless is Genealogical Linguistics, and How Advanced is Areal Linguistics? *Studies in Language* 28(1): 209–223.
- Heeringa, W. J.
- 2004 Measuring Dialect Pronunciation Differences using Levenshtein Distance. Ph.D. thesis, Rijksuniversiteit Groningen.
- Heggarty, P., W. Maguire, and A. McMahon
- 2010 Splits or Waves? Trees or Webs? How Divergence Measures and Network Analysis Can Unravel Language Histories. *Philosophical Transactions of the Royal Society B: Biological Sciences* 365(1559): 3829–3843.
- Hock, H. H.
- 1991 *Principles of Historical Linguistics*, 2nd ed. Berlin: Mouton de Gruyter.

- Holton, G.
- 2011 A Geo-linguistic Approach to Understanding Relationships Within the Athabaskan Family. Paper presented at the International Workshop of Language in Space: Geographic Perspectives on Language Diversity and Diachrony. Boulder, Colorado. July 23, 2011.
- Huehnergard, J. and A. Rubin
- 2011 Phyla and Waves: Models of Classification of the Semitic Languages. In S. Weninger (ed.) *The Semitic Languages: An International Handbook* (Handbücher zur Sprach- und Kommunikationswissenschaft 36), pp.259–278. Berlin: de Gruyter Mouton.
- Hurles, M. E., E. Matisoo-Smith, R. D. Gray, and D. Penny
- 2003 Untangling Oceanic Settlement: The Edge of the Knowable. *Trends in Ecology and Evolution* 18(10): 531–540.
- Kortlandt, F.
- 2007 *Italo-Celtic Origins and Prehistoric Development of the Irish Language*. Amsterdam: Rodopi.
- Krauss, M. E. and V. Golla
- 1981 Northern Athapaskan languages. In J. Helm (ed.) *Handbook of North American Indians*, vol. 6: *Subarctic*, pp.67–85. Washington D.C.: Smithsonian Institution Scholarly Press.
- Labov, W.
- 1963 The Social Motivation of Sound Change. *Word* 19(3): 273–309.
- Leskien, A.
- 1876 *Die Declination im Slavisch-Litauischen und Germanischen*. Leipzig: Hirzel.
- Maddison, W. P.
- 1997 Gene Trees in Species Trees. *Systematic Biology* 46(3): 523–536.
- Milroy, J. and L. Milroy
- 1985 Linguistic Change: Social Network and Speaker Innovation. *Journal of Linguistics* 21(2): 339–384.
- Minaka, N. and K. Sugiyama
- 2012 *Keitouju Mandara (Phylogeny Mandala: Chain, Tree, and Network)*. Tokyo: NTT Publishing. (In Japanese)
- Nerbonne, J.
- 2010 Measuring the Diffusion of Linguistic Change. *Philosophical Transactions of the Royal Society B: Biological Sciences* 365(1559): 3821–3828.
- Page, R. D. M. and E. C. Holmes (eds.)
- 2009 *Molecular Evolution: A Phylogenetic Approach*. Oxford: Blackwell.
- Pawley, A.
- 1999 Chasing Rainbows: Implications of the Rapid Dispersal of Austronesian Languages for Subgrouping and Reconstruction. In E. Zeitoun and P. J.-K. Li (eds.) *Selected Papers from the Eighth International Conference on Austronesian Linguistics* (Symposium Series of the Institute of Linguistics), pp.95–138. Taipei: Institute of Linguistics, Academia Sinica.
- 2008 Where and When was Proto-Oceanic Spoken? Linguistic and Archaeological Evidence.

- In Y. A. Lander and A. K. Ogloblin (eds.) *Language and Text in the Austronesian World: Studies in Honour of Ülo Sirk*, pp.47–71. München: Lincom Europa.
- 2009 Polynesian Paradoxes: Subgroups, Wave Models and the Dialect Geography of Proto Polynesian. Paper presented at the Eleventh International Conference on Austronesian Linguistics, Aussois, France, June 22, 2009.
- 2010 Prehistoric Migration and Colonisation Processes in Oceania: A View from Historical Linguistics and Archaeology. In J. Lucassen (ed.) *Migration History in World History: Multidisciplinary Approaches* (Studies in Global Social History 3), pp.77–112. Leiden: Brill.
- Pawley, A. and R. C. Green
1984 The Proto-Oceanic Language Community. *Journal of Pacific History* 19(3): 123–146.
- Pawley, A. and M. Ross
1995 The Prehistory of Oceanic Languages: A Current View. In P. S. Bellwood, J. J. Fox, and D. Tryon (eds.) *The Austronesians: Historical and Comparative Perspectives*, pp.39–80. Canberra: Australian National University.
- Posner, R.
1996 *The Romance Languages*. Cambridge: Cambridge University Press.
- Ross, M.
1988 *Proto-Oceanic and the Austronesian Languages of Western Melanesia*. Canberra: Pacific Linguistics.
1997 Social Networks and Kinds of Speech-community Event. In R. Blench and M. Spriggs (eds.) *Archaeology and Language 1: Theoretical and Methodological Orientations*, pp.209–261. London: Routledge.
- de Saussure, F.
1995 [1916] *Cours de linguistique générale*. Paris: Éditions Payot & Rivages.
- Schmidt, J.
1872 *Die Verwandtschaftsverhältnisse der indogermanischen Sprachen*. Weimar: Hermann Böhlau.
- Schmidt, K.-H.
1993 Insular Celtic: P and Q Celtic. In M. J. Ball and J. Fife (eds.) *The Celtic Languages*, pp.64–99. New York: Routledge.
- Schrader, O.
1883 *Sprachvergleichung und Urgeschichte: linguistisch-historische Beiträge zur Erforschung des indogermanischen Altertums*. Jena: Hermann Costenoble.
- Séguy, J.
1973 La dialectométrie dans l'Atlas linguistique de la Gascogne. *Revue de linguistique romane* 145–146: 1–24.
- Skelton, P., A. Smith, and N. Monks
2002 *Cladistics: A Practical Primer on CD-ROM*. Cambridge: Cambridge University Press.
- Southworth, F. C.
1964 Family-tree Diagrams. *Language* 40(4): 557–565.
- Szmrecsányi, B.
2011 Corpus-based Dialectometry: A Methodological Sketch. *Corpora* 6(1): 45–76.

- Toulmin, M.
2009 *From Linguistic to Sociolinguistic Reconstruction: The Kamta Historical Subgroup of Indo-Aryan* (Pacific Linguistics 604). Canberra: Australian National University.
- Tryon, D.
1996 Dialect Chaining and the Use of Geographical Space. In J. Bonnemaïson, K. Huffman, C. Kaufmann, and D. Tryon (eds.) *Arts of Vanuatu*, pp.170–181. Bathurst: Crawford House Publishing.
- Wüest, J.
1994 La restructuration du système des démonstratifs en protoroman. In J. Cerquiglini-Toulet and O. Collet (eds.) *Mélanges de philologie et de littérature médiévales offerts à Michel Burger*, pp.41–49. Genève: Droz.