



HAL
open science

On the applicability of currently existing digital tools in reconstructive textual editing

Philipp Roelli

► **To cite this version:**

Philipp Roelli. On the applicability of currently existing digital tools in reconstructive textual editing. 2022. hal-03718032v1

HAL Id: hal-03718032

<https://hal.science/hal-03718032v1>

Preprint submitted on 8 Jul 2022 (v1), last revised 25 Jan 2023 (v4)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

On the applicability of currently existing digital tools in reconstructive textual editing

Philipp Roelli

University of Zurich, Switzerland

roelli.sglp@yandex.com

Abstract

This article reviews some of the digital tools currently available for reconstructive textual editing. First the main idea of reconstructive textual editing is summarised, then its parts amenable to algorithmic description are compared to similar questions in evolutionary biology. Two Latin texts with a complicated transmission are then introduced to illustrate some available tools. The main focus is on stemma reconstruction. Some steps of the process can already be largely automated, especially collating texts, but it is found that tree-constructing software is of little help in the case of the *Liber Aurelii*, whereas it is somewhat more helpful for Plato of Tivoli's Latin translation of the *Centiloquium*. In a concluding part, the main problems for algorithmic approaches to the stemma are discussed: incomplete witnesses leading to only partly overlapping text samples, contamination of some witnesses, and rooting the tree.

keywords

Critical editing, Textual criticism, Stemmatology, Computer aids, Significant errors.

INTRODUCTION

I THE MAIN IDEA OF RECONSTRUCTIVE TEXTUAL EDITING

Reconstructive textual editing is an approach that aims to reconstruct a text known only in copies made from a lost original as closely as possible by using all available data (for a more detailed characterisation [Roelli, 2020: 3–4]). The latest common ancestor of all surviving copies is known as the 'archetype'. The first goal is to reconstruct this (often lost) archetype from the extant witnesses as far as possible, the second to examine it and to try to correct its errors using the available external information about the original, the author, and his time. In reality, the situation may become more complicated due to several factors: the author may have reworked the text and there may thus be more than one 'original' or some information may make its way into some witnesses from a pre-archetype witness now lost (a process known as 'extra-stemmatic contamination', cf. Paolo Trovato [Roelli, 2020: 123] and Marina Buzzoni [Roelli, 2020: 386]). In antiquity and the early middle ages there are often many centuries between the oldest surviving manuscript and the original; and worse: even between the archetype and the original. Depending on the text in question, there may be from

one to hundreds, occasionally even thousands of surviving copies of extant antique or medieval texts. Clearly, digital aids are promising to handle the data of especially abundant traditions. Texts surviving in a handful or less witnesses can just as well be dealt with manually. The mentioned goal of approaching the lost text (that is, as a first step, its archetype), is reached by determining the relationship between the surviving witnesses by evaluating and weighting their readings.

Fig. 1 depicts what is apparently the first printed *stemma codicum* in our field, dating from 1827 [Schlyter, 1827: appendix].¹ The stemma is the genealogical tree which explains the observed variation in the witnesses in the most economical way. Here, the archetype or original is situated at its top end without a label; indeed the concepts ‘archetype’ and ‘original’ were not yet differentiated by Schlyter. The lines represent the relation ‘was copied to’. The stemma provides information about the weight readings should be given in the process of editing the archetypal text depending on their witnesses’ position. If the depicted stemma is correct, we will have to give the readings of A more weight for the reconstructed text than the other witnesses. For argument’s sake, let us consider a locus with four readings that occur in the witnesses as indicated in the stemma. Let the readings be: reading 1 (H, K, M, No166), reading 2 (A, L), reading 3 (C), reading 4 (B, G, N). Without the stemma we might be tempted to use a majority criterion and adopt reading 1 or 4. If we have the stemma, we will conclude that readings 1, 3, and 4 are innovations and we will choose reading 2 for the archetypal text. In cases similar to this one, it becomes thus possible to choose the archetypal reading mechanically, that is in a certain sense more objectively than if a philologist chose it intuitively (which is what was usually done prior to the 19th century).

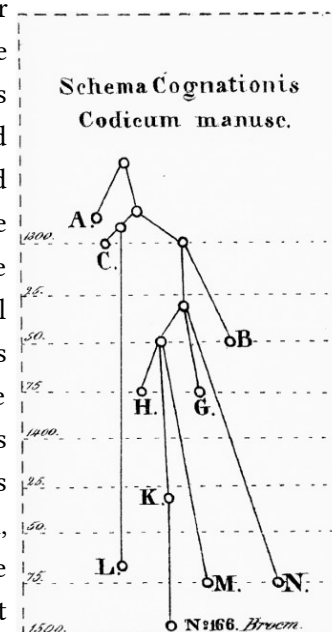


Figure 1: Carl Johann Schlyter’s stemma of Västgötalagen (1827).

The crucial question is now: how does one find the correct stemma or genealogical tree of a given set of witnesses? This question led scholarship to study variant readings and to the concept of significant errors or *Leitfehler* (this term was coined by Paul Maas, [Roelli 2020, p. 117 (Paolo Trovato)]). A significant error is an edit which can hardly or not at all be undone by a subsequent scribe without access to the original text form. Omissions of more than a word or two are good candidates; especially when they happen involuntarily in the form of eye-skips: a word or phrase occurs twice close-by in the original; the copyist copies up to the first occurrence, but then jumps by mistake to the second one, thereby omitting what stood in-between. Although this may happen to more than one scribe at the same locus, it is usually out of the question that the missing text can be reconstituted without access to a witness that still has it. A significant error should be directed: it must be clear which variant is the original one and which others are not. Only the secondary

1 Apart from a lack to consider contamination, this stemma seems to be largely correct: ‘Regarded as a schema that draws up the principal lines and disregards the contamination of the tradition it would actually seem to be almost accurate’ [Olrik Frederiksen, 2009, 129]. There are minor mistakes, so K and M are apparently independent of H [ibid., p. 135]. More on Schlyter’s method in [Froger, 1978] and [Holm, 1972].

readings form families, the original does not. In case of eye-skips the direction is usually clear,² although, unfortunately, the secondary reading (the omission) may happen more than once. The reconstruction of the stemma becomes much harder if some scribes used more than one source to compile a new manuscript, thus trying to improve their new copies. This phenomenon is called contamination [Roelli, 2020, ch. 4.4 (Tuomas Heikkilä)]; it will be discussed further below. The approach just outlined was first developed by German and French scholars in the 19th century. It was refined in many ways mostly by Italian scholars of the second half of the 20th century [Roelli 2020, esp. ch. 2.4 (Paolo Trovato)].

II PARALLELS WITH MOLECULAR BIOLOGY

The procedure just outlined can *in praxi* be summarised in these four steps:

- Transcription of available witnesses,
- Evaluating significant errors,
- Reconstructing the stemma,
- Editing the archetypal text according to it,
- Emending the archetypal text to approach the original as far as possible.

As many of these steps are nearly of an algorithmic nature, the idea to use computers to perform them arose early-on and quite naturally when computers became available. The first such approaches were already done in the late fifties [Roelli, 2020, ch. 5.1 (Armin Hoenen)]. For most of these steps there are today computer aids, but none of them can be fully automated as yet. The computerised approaches to these five steps differ widely in their advancement. Recent years have seen dramatic advances in the automatic reading of handwritten documents (i), for instance the open-source Kraken³ project in Paris has already reached levels of reading old handwritten material that would have seemed quite unthinkable a decade ago. The second step (ii) has been largely neglected, many digital scholars seem to hope that a big amount of non-significant errors will work just as well as a few rare significant ones. Below, we will see that this is in general not the case. We [Roelli & Bachmann, 2010]⁴ made a first attempt to automatically identify candidates of significant errors. The third step (iii) is very similar to what evolutionary biologists do when they construct genealogical trees from DNA sequences. Thus, it was possible to borrow much know-how from them. However, we will see that locating the archetype – finding the root of the tree – is a specific problem in our field for which the biologists’ tools are not helpful. The fourth step (iv), sometimes referred to as *Urtext* reconstruction, is still quite experimental; although given a stemma and a set of texts the task to find the most likely readings for the (lost) intermediary nodes would seem to be a relatively straight-forward task to program [Roelli, 2020, ch. 5.4.6 (Armin Hoenen)]. If the position of the archetype (in biological terminology: the ‘root’) is also known, the *Urtext*’s readings could then

2 Unless an addition (that may have stood *supra lineam*) moved into the copyist’s text and the anchor word was repeated; but the context usually makes clear whether this is the case or not.

3 <https://escripta.hypotheses.org/tag/kraken>.

4 The software (a perl script) is shared upon request.

be largely reconstructed probabilistically by the computer. A problem is that (as shall be seen below) the archetypal text is located between two nodes in case there are only two hyparchetypes, so the last step from these to the archetype remains unclear. Traditional editors face the same problem: it is largely up to their *iudicium* which of the two hyparchetypal readings they prefer where they disagree. The final step of emending the archetype's text has not been tackled by digital aids. In the remainder of this essay we shall mostly focus on the third step, the automatic finding of the stemma from transcribed data.

There is considerable similarity of our problem with that of molecular biologists. Where we deal with text strings, they deal with DNA or amino acid sequences, thus basically 'texts' consisting of four or twenty letters respectively. We both try to find the tree that explains the variation in the most economic way. The standard approach in biology is to measure the distances between each pair of taxa and to use the resulting distance matrix to find the best tree ([Roelli 2020, ch. 5.5 (Jean-Baptiste Guillaumin)] for a practical example). Alternatively, probabilistic Bayesian methods may omit the step of the distance matrix and optimise an initial tree directly in tree-space [Roelli, 2020, ch. 5.2 (Sara Manafzadeh and Yannick M. Staedler) and 5.3 (Teemu Roos)]. Either way, a point that has often not been paid its due attention is the metric one uses to measure the distance between items. In biology this point seems to be less crucial than in stemmatology. In textual criticism manual weighting has been attempted [Macé, De Vos, Geuten, 2012]. The mapping from text-space to matrix-space by means of a notion of distance can be formalised in this way:

$$\begin{array}{ccccc}
 & & m & & \textit{heuristic software} \\
 L^{(r \times n)} & \xrightarrow{\quad} & R^{(n \times n)} & \xrightarrow{\quad} & T_n \\
 \text{(witnesses' texts)} & & \text{(distance matrix)} & & \text{(tree graph)}
 \end{array}$$

In this formula, m is a metric (distance function) that calculates a 'distance' between any two texts $t_1, t_2 \in L^r$, $(t_1, t_2) \mapsto m(t_1, t_2) \in \mathbb{R}$.⁵ The idea is to make m correspond as closely as possible to our real-world notions of scribal accuracy or likelihood of copyists' mistakes happening. L stands for the texts' lexicon, r the number of readings (\leq total words of the longest witness). T_n is the set of unrooted binary trees with n leaves. The second step, from matrix to tree is done by means of a heuristic algorithm such as those found in the free Phylip software package used in molecular biology, for instance the one by Fitch and Margoliash.⁶ The most natural (naive) choice of the metric m is 'counting variants', i.e. $m(t_1, t_2)$ is the number of operations needed to transform t_1 into t_2 (known as 'editing distance'), where an operation is the deletion or insertion of a word or, better, a reading (that may consist of several words). The Unix `diff` command uses the longest common subsequence algorithm (LCS) repeatedly to detect differences in text-strings and can be used for such a task. Now the insertion or deletion of n consecutive words can be counted as a single operation, n distinct operations, or anything in between. Something in-between is likely the most fitting approach. But we still have not made a difference between a trivial edit, say of the synonyms *dominus* to *deus* or simple word transpositions, and a more significant error. To my knowledge there

5 For all n witnesses this produces a matrix $\mathbb{R}^{(n \times n)}$. As m is commutative, this matrix is symmetrical, which besides has only zeros in its diagonal (the distance between any text and itself being zero), thus only the subspace $\mathbb{R}^{n(n-1)/2}$ is relevant.

6 <https://evolution.gs.washington.edu/phylip/doc/main.html>.

is as yet only one attempt to automatically weight the significance of variants. Bachmann and I devised it in 2010 by assuming that if there is no contamination in a set of texts, then two significant errors' absence and presence will not be found in all four possible combinations in the witnesses. This approach uses the principle that significant errors usually happen only once in a specific textual tradition as their very idea is that they cannot be undone easily [Roelli & Bachmann, 2010, 317–318]. Variants that fulfil this criterion in combination with many other of the candidate variants can be weighted more strongly. This approach, thus, tries to assign a numerical value for the fitness of a variant as significant. Unfortunately, it cannot determine which of the variants is the original one as the presence and absence of words or readings is symmetrical. Therefore we use quotes and speak of the 'Leitfehler'-method. Clearly more work on this approach is required, but first results were promising.

III SOME EXISTING TOOLS WITH TWO COMPLICATED TEXTUAL TRADITIONS AS EXAMPLES

In order to work with concrete, real-life data I will now use data from two Latin texts that are currently under study to show how well some of the available software works. All software I used for what follows is freely available. I used it either online or on my rather low-power Linux laptop.

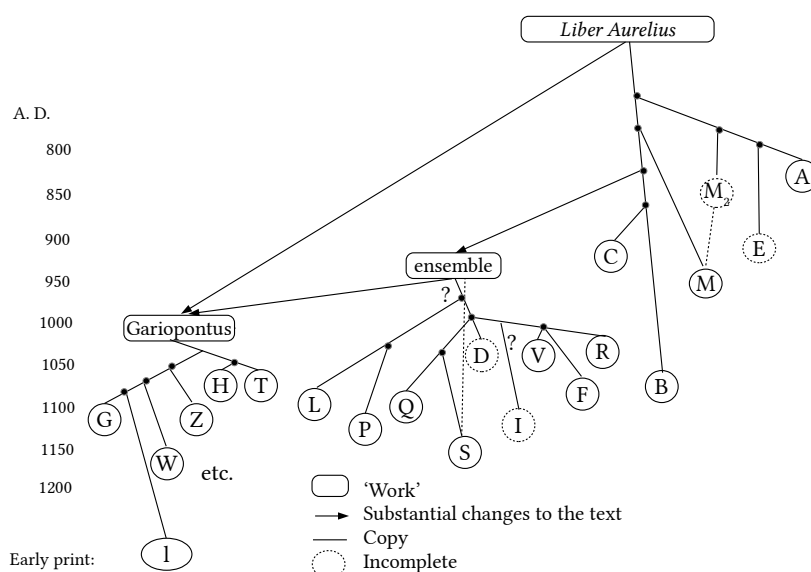


Figure 2: Manually determined stemma of the Liber Aurelii.

3.1 Liber Aurelii

The first one is a text I have recently edited critically for the first time using the traditional 'manual' approach [Roelli, 2021]. This anonymous text, erroneously known as *Liber Aurelii*, is a late antique Latin medical work on acute illnesses that goes back to lost Greek sources, mostly from the Methodic school and Soranus of Ephesus (fl. 2nd century AD). Its original title is unknown. It is extant in three recensions: the main text that had come down to the high middle ages in a seriously garbled text form and two reworkings by physicians of probably the 10th and 11th centuries.

Especially the later one by the Salernitan author Gariopontus was very successful and has survived in more than 65 witnesses [Glaze, 2005]. Gariopontus still had a more comprehensible copy available, besides he also used the other recension. The author of this earlier recension had shortened the text drastically, often simply omitting hard to understand content. His text is less than half as long as the other two. As usually in medieval Latin, there is also quite a great deal of variability in spelling. The text is about the size of a typical antique book. We will use this text's chapter 12 (some 750 words in the main version, some 500 in the shortened one, here Gariopontus took his text mostly from the shortened recension) to provide data for what follows. For the edition I determined the stemma of the *Liber Aurelii* manually using the traditional methods of textual criticism (fig. 2; for a key to the sigla, see the edition: [Roelli 2021: lxiv–lxv]).

I had already transcribed all witnesses manually into a txt file in which each line contains the text of one witness; therefore I will skip the automated reading of the manuscript texts. While transcribing, one should note all orthographic and other peculiarities of the witnesses as this information can come in handy in later steps of editing. But for the input of tree-finding software, it is preferable to use standardised spelling as differences such as *hec* vs. *haec* are very unlikely to be relationship revealing, rather they are what could be termed 'orthographic noise'. Regex syntax in a text editor (e.g. the free Edit pad lite on Windows) can be used to quickly standardise the spelling. In our example, among other things, I removed all punctuation and all spelling that tends to be variable in medieval Latin. This can be done drastically and quickly by just removing for instance all h's altogether, or even changing all b's into u's. (The mixing up of v and b is a common feature of some of this text's witnesses.) The resulting txt-file may no longer look like Latin but the important information to determine the relationships between witnesses is preserved and much of the 'noise' is now removed. This step means that such differences are defined as zero distances: e.g. *humiliatio* = *umiliacio*. This step is fully automatic and takes no more than a few minutes. The exact list of what should be standardised and what not, of course, depends, on time, language, and textual tradition.

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | |
|----|-----------|-----------|-----------|-----------|-----------|-------------|---------------|------------|------------|------------|------------|------------|------------|------------|---------------|------------|-------------|-------------|-------------|-------------|
| 1 | A | B | C | E | F | G | H | L | M | M2 | P | Q | R | S | T | V | W | Z | I | S |
| 2 | | | | | | item | item | | | | | | | | | | | | | |
| 3 | de | de | de | de | de | de | de | de | | | de | de | de | de | de | de | de | | | de |
| 4 | sinance | | | sinance | sinancia | eadem | eadem | | | | | sinancia | sinancia | sinance | eadem | sinancia | sinancia | | | eodem mor |
| 5 | | sinancis | sinancis | | | sed allium | auctorem | sinancis | | | | sinancis | sinancis | | pasione | | | | | |
| 6 | sinance | sinance | sinance | sinance | sinancia | sinance | sinance | sinancia | sinance | sinance | sinancia | sinancia | sinancia | sinancia | sinance | sinancia | sinancia | sinancia | sinance | sinance |
| 7 | dicta est | dicta est | dicta est | dicta est | dicta est | dicta est | dicta est | dicta est | dicta est | dicta est | dicta est | dicta est | dicta est | dicta est | dicta est | dicta est | dicta est | dicta est | dicta est | dicta est |
| 8 | au eo | au eo | au eo | au eo | au eo | a | a | au eo | au eo | au eo | au eo | au eo | au eo | au eo | a | au eo | a | a | a | a |
| 9 | quod | quod | quod | quod | quid | | | quod | quod | quod | quod | quod | quod | quod | quod | quod | quod | quod | quod | quod |
| 10 | ueluti | ueluti | ueluti | ueluti | ueluti | | | ueluti | ueluti | ueluti | ueluti | ueluti | ueluti | ueluti | ueluti | ueluti | ueluti | ueluti | ueluti | ueluti |
| 11 | prefocac | profocac | prefocac | prefocac | prefocac | prefocac | prefocac | prefocac | prefocac | prefocac | prefocac | prefocac | prefocac | prefocac | prefocac | prefocac | prefocac | prefocac | prefocac | prefocac |
| 12 | quodam | | quodam | quodam | | quia sinanc | quia sinancis | quodam | quodam | | | | | | quia sinancis | | quia sinanc | quia sinanc | quia sinanc | quia sinanc |
| 13 | | quandam | | | quandam | | | quandam | quandam | quandam | quandam | quandam | quandam | quandam | | quandam | | | | |
| 14 | paciantur | paciantur | paciantur | paciantur | paciantur | | | paciantur | paciantur | paciantur | paciantur | paciantur | paciantur | paciantur | quandam | paciantur | | | | |
| 15 | | | | | | | | | | | | | | | | | | | | |
| 16 | qui laura | qui laura | qui laura | qui laura | qui laura | | | qui laura | qui laura | qui laura | qui laura | qui laura | qui laura | qui laura | qui laura | qui laura | qui laura | qui laura | qui laura | qui laura |
| 17 | ec enim | ec enim | ec enim | ec enim | ec enim | | | ec enim | ec enim | ec enim | ec enim | ec enim | ec enim | ec enim | ec enim | ec enim | ec enim | ec enim | ec enim | ec enim |
| 18 | grece | grece | grece | grece | grece | | | grece | grece | grece | grece | grece | grece | grece | grece | grece | grece | grece | grece | grece |
| 19 | | | | | | | | latine | latine | latine | latine | latine | latine | latine | latine | latine | latine | latine | latine | latine |
| 20 | | | | | | | | prefocacio | prefocacio | prefocacio | prefocacio | prefocacio | prefocacio | prefocacio | prefocacio | prefocacio | prefocacio | prefocacio | prefocacio | prefocacio |
| 21 | | | | | | | | latine | latine | latine | latine | latine | latine | latine | latine | latine | latine | latine | latine | latine |
| 22 | dicitur | dicitur | dicitur | dicitur | dicitur | | | dicitur | dicitur | dicitur | dicitur | dicitur | dicitur | dicitur | dicitur | dicitur | dicitur | dicitur | dicitur | dicitur |
| 23 | latine | latine | latine | latine | latine | | | latine | latine | latine | latine | latine | latine | latine | latine | latine | latine | latine | latine | latine |
| 24 | | | | | | | | | | | | | | | | | | | | |
| 25 | prefocare | profocare | prefocare | prefocare | prefocare | | | prefocare | prefocare | prefocare | prefocare | prefocare | prefocare | prefocare | prefocare | prefocare | prefocare | prefocare | prefocare | prefocare |
| 26 | est | et | et | est | | | | est | est | est | est | est | est | est | est | est | est | est | est | est |
| 27 | definicio | definicio | definicio | definicio | definicio | | | definicio | definicio | definicio | definicio | definicio | definicio | definicio | definicio | definicio | definicio | definicio | definicio | definicio |
| 28 | autem | ante | ante | autem | | | | autem | autem | autem | autem | autem | autem | autem | autem | autem | autem | autem | autem | autem |
| 29 | sinance | sinancem | sinance | sinance | sinance | | | sinance | sinance | sinance | sinance | sinance | sinance | sinance | sinance | sinance | sinance | sinance | sinance | sinance |

Figure 3: The beginning of the automatically generated alignment table of *Liber Aurelii*, chapter 12, automatically obtained by using CollateX.

The typical collation table that is used in traditional editing can now be generated automatically from a txt file within seconds using the software CollateX (<https://collatex.net/doc/#text-input>). As a command line tool, CollateX⁷ produces from a json-file (which is easily obtained from our txt-file) a csv-spreadsheet of a collation table as the one shown in fig. 3. In our case the standardisation of the spelling reduced the number of readings only slightly (from 883 to 860). This software functions well (although not perfectly in all instances) and can save a significant amount of time.

Another useful piece of software is Juxta (<http://www.juxtasoftware.org>): it allows to generate ‘critical’ apparatuses based on a set of witnesses. It compares the witnesses but does not evaluate them, so one has to provide a manual critical text in order to be able to get a critical apparatus from the software. Such lists can be useful to find readings that are shared by a given group of witnesses. The output is an html-file that can be easily searched in any web-browser. The latter is shown for a part of our text in fig. 4.

```

Critical Apparatus

1 ^item G, H
1 de sinance~ M, M2, Z
1 sinance ] sinancis B, C
1 sinance sinance ] sinancia sinancia F, L, P, Q, R, V
1 sinance ] eadem sed alium auctorem G
1 sinance ] eadem H
1 sinance ] eodem moruo I
1 sinance ] eadem passione T
1 sinance sinance ] sinancia sinancis W
1 sinance ] sinancia S
1 au eo quod ueluti prefocacionem quandam paciuntur qui laurant ec enim ] a prefocacione quia sinancis G, I, T, W, Z
1 au eo quod ueluti prefocacionem quandam paciuntur qui laurant ec enim grece dicitur latine prefocare et definicio autem s
impetu circa inguine nam oc tonsillarum impetu differ t quod ea acuta passio est itemque et illo quod non profocant tonsille eis
inspiciamus ostendendum est ] a prefocacione quia sinancis grece prefocacio latine dicitur sinance est difficilis translacio spirit
H
1 quod ] quid F
1 prefocacionem ] profocacionem B
1 prefocacionem ] prefocacione V
1 quandam ] quodam A, E, M, M2
1 quandam ] quadam C
1 ^paciuntur S
1 paciuntur ] paciantur Q

```

Figure 4: Automatically generated apparatus by Juxta.

An entire set of tools is offered by Tara Andrews’ Stemmaweb (<https://stemmaweb.net>). Besides an aligned text, one can also visualise a flow diagram with readings as nodes and witnesses as edges (fig. 5). One can also create genealogical trees using several approaches. Fig. 6 shows such a tree for our example text using the intuitive but simplistic NeighborJoining approach. It does not weight readings philologically, but the result is nonetheless not too bad: The three recensions are kept apart, but B and C should form a group and stand below M; V should group with F and R. A negative point of this approach is that all edges are assigned the same length. More sophisticated approaches provide branch length which enables the user to gauge to what extent texts differ and to decide which bifurcations should actually be understood as multifurcations.

⁷ A simple java -jar collatex-tools-1.7.1.jar -f csv DATA.json > OUTPUT.csv. For large samples it may be necessary to split the file into smaller ones.

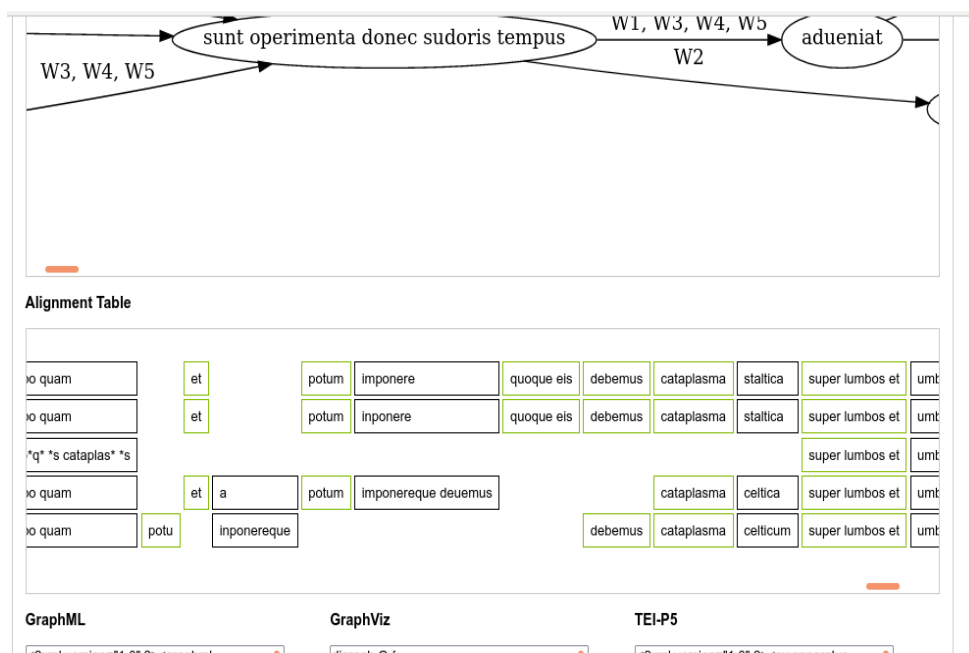


Figure 5: StemmaWeb.

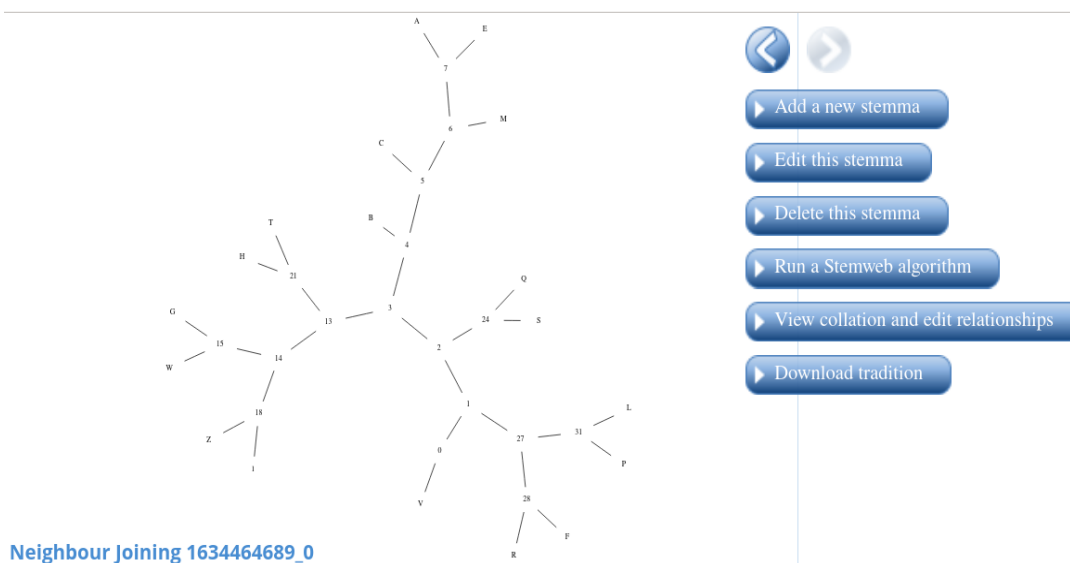


Figure 6: Tree from StemmaWeb based on NeighbourJoining.

The Roelli-Bachmann script generates the tree shown as fig. 7 with automatically determined ‘Leitfehler’; it does include variable branch-length. The clear pairs of very similar manuscripts are mostly correctly found: AE, LP, FR, GI, HT, BC. Three lines are not correct, they are depicted as manually added red arrows in the figure. The most serious mistake (in both plots) is that the abbreviated recension (on top in fig. 7) should belong close to B and C. This artefact is a common problem of most similar software: it moves shorter texts towards the centre of the plot.

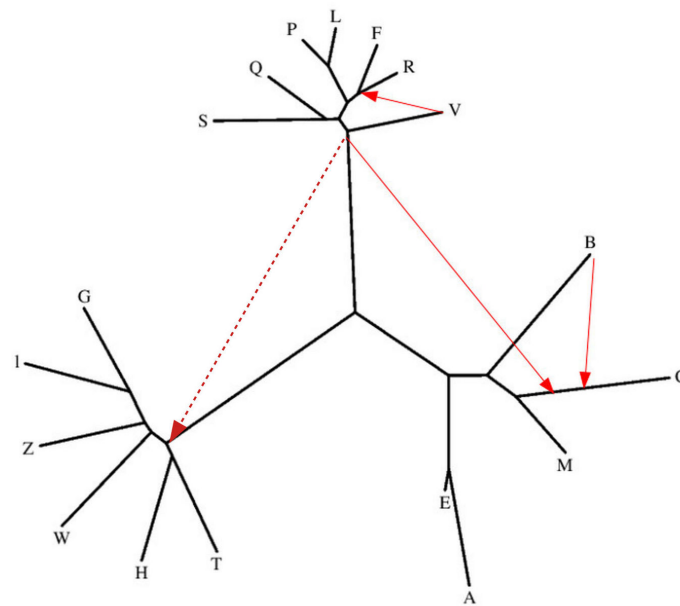


Figure 7: Tree by Roelli–Bachmann software. Red arrows show where the plot is mistaken. The dotted line represents contamination. Minor cases of contamination (especially in M and S are not shown).

The weighting of variants changes the result significantly. This is less the case for the kinship of leaves in the tree, especially when they are rather distinctive, but the upper branches in the tree are prone to shift easily when the weights are changed. Unfortunately, these are the parts of the evaluation that are also the most difficult for a traditional human textual critic. The result of fig. 7 could be improved by using hand-picked significant errors, but this, of course, will easily become a case of finding again what one had previously hidden behind the bush (as Nietzsche used to say).⁸ The ‘*Leitfehler*’ script functions better than simple NeighbourJoining⁹ but would still need significant improvement. Some improvement would be achieved by using entire readings or possibly n-grams in order to be able to include series of common words that differ in their combination as potentially significant. Even then, in this case, the improved result would not be very useful in practical life as an editor because it is of no help for the two most serious problems, rooting and contamination, as will be discussed below.

3.2 Centiloquium

The second text I used is the Latin version of Pseudo-Ptolemy’s *Centiloquium* by Plato of Tivoli, currently being edited by my colleague Emanuele Rovati, who kindly allowed me to use his transcription data and some of his new insights into the text’s transmission. The *Centiloquium*¹⁰ is a collection of one hundred aphorisms on astrology that was considered to be written by the famous

8 The following list of hand-picked errors: *definitio, tumore, cognoscere, proiciunt, afixia, introrsum, irruit, uisceribus, retentione, aperiant, molitam, catarticum, efficitur, sursum* improves the result somewhat (plot not shown).

9 RHM’s result on StemmaWeb is completely wrong. Does it fail here because of different string lengths or is there a problem in its implementation?

antique astronomer Claudius Ptolemaeus. In Arabic it is often transmitted together with the commentary of the tenth-century astrologer Abū Ja‘far Aḥmad ibn Yūsuf ibn Ibrāhīm ibn al-Dāya (on whom: [Lemay, 1978]). The extant Greek version of the text does not seem to be the original version, but rather a Byzantine translation from the Arabic. This short text (of some 16’000 words) was translated several times from Arabic into Latin, sometimes with, sometimes without Abū Ja‘far’s commentary, and remained influential into early modern times. The text we use here was translated into Latin by Plato of Tivoli in 1136 and includes the commentary. There are 101 known manuscripts and three early prints.¹¹ In total, some 80 of the manuscripts contain a substantial part or the entire text. The sample for our plots uses twelve of the one hundred aphorisms and their commentaries (the numbers 9–10, 30, 51–53, 60–62 and 98–100), and spans some 3’000 words per witness. The witnesses’ text is much more homogeneous than the one of the *Liber Aurelii*. The most common edits are in orthography, word transpositions, changes of mood, tense, number, or person in verbs, and the substitution of pronouns (e.g. *eius* vs. *illius* vs. *istius*) without changing the meaning.

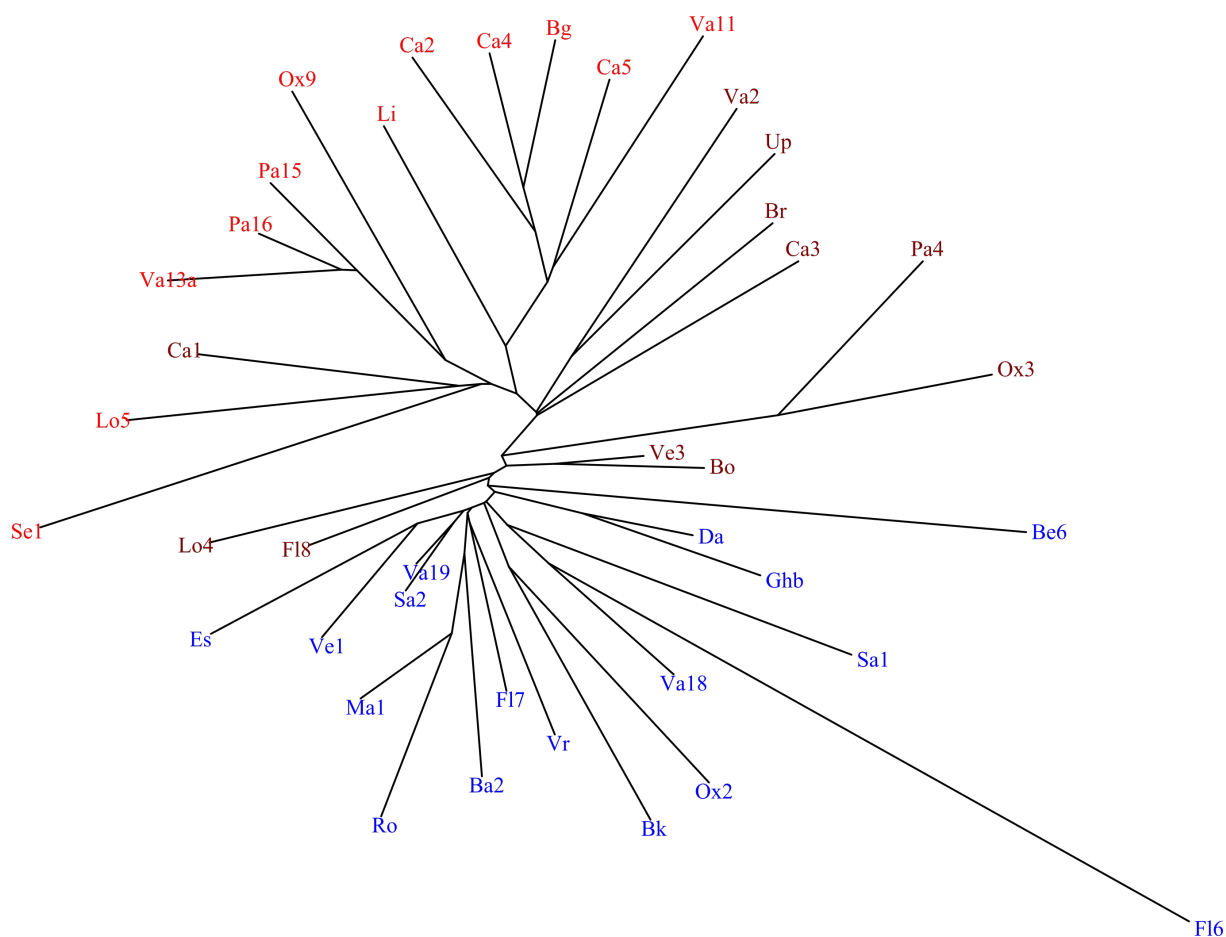


Figure 8: Centiloquium plot, the α group is depicted in red (manuscripts from the threefold version in light red), β in blue.

10 Further data about the text [Juste, 2022]. Edition of the Arabic text [Martorello & Beza, 2013], and the *Ptolemaeus Arabus et Latinus* project: <https://ptolemaeus.badw.de>.

11 The first of these (the editio princeps): Venice, Erhardus Ratdolt, 1484. Ratdolt was a competent editor. The three prints are related to one another and to manuscript Va19.

There are two major groups: α is the text translated by Plato of Tivoli, β a reworked version by Gerard of Cremona, as Rovati has been able to show.¹² The latter is more stable than the former. Gerard checked the Arabic original and changed some 400 passages in Plato's translation. Later copyists often conflated the two versions. The majority of witnesses of Plato's unaltered translation belong to what is known as the threefold version which quotes each aphorism in three translations (Plato, Adelard of Bath, and the anonymous 'Mundanorum 1' translation). Obviously only the text of the first of these translations is used here. An example can illustrate the typical differences between α and β : Plato wrote *Obtinebit, inquit, locum patris et erit 10 annis fere in regno, sed erit sicut ille cui iubetur*. Gerard compared the Arabic original and added, apparently in the margin: *in alio* (i.e. in another Arabic manuscript): *sub potestate vel regimine alterius*. Most β manuscripts now read exclusively: *Obtinebit, inquit, regnum patris et erit 10 annis fere in regno, sed erit sub potestate alterius*.

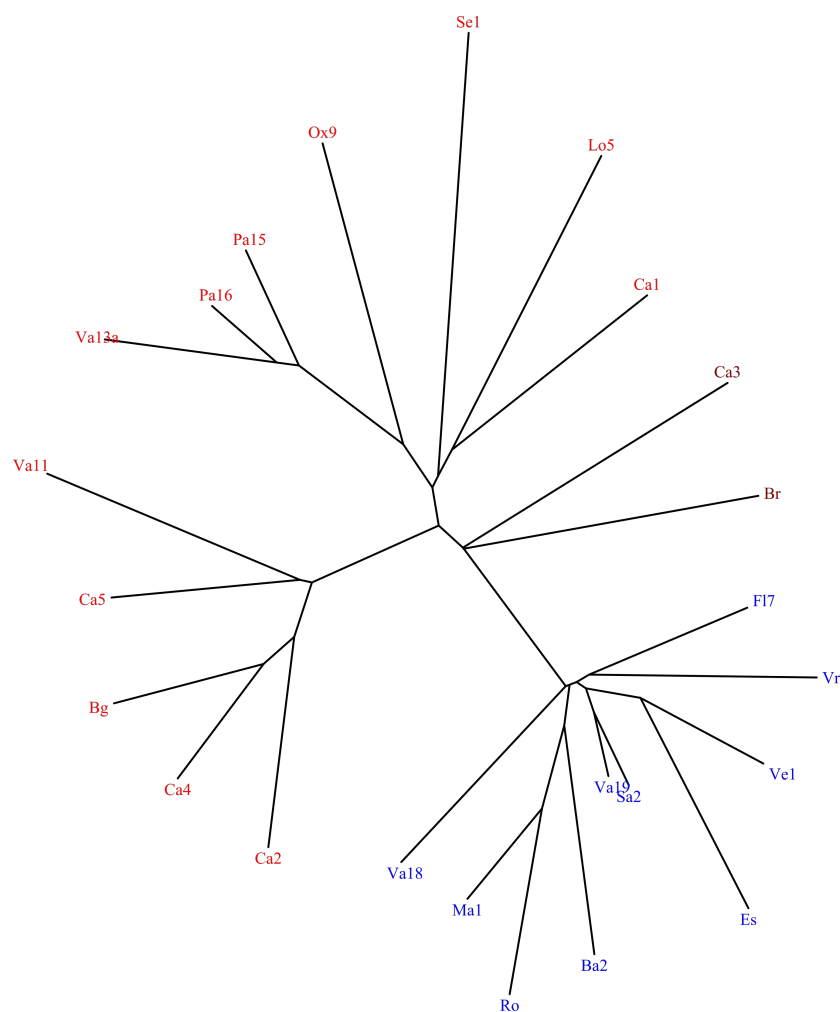


Figure 9: A second plot without the most obvious cases of contamination. Colours as in fig. 8.

I will focus on the automated trees I plotted based on Rovati's transcriptions, thus I will skip the process leading to them as it is identical to what was described above for the *Liber Aurelii*. Suffice it

¹² The data in this paragraph comes from Emanuele Rovati, to be published in detail in his critical edition (around 2025).

to say that CollateX finds 3'190 readings for this sample. A first plot of forty manuscripts (mostly the oldest ones) produced the plot in fig. 8. Many of the witnesses can be identified as contaminated at first sight as they contain glosses with readings from other branches or pairs of both the α and the β reading within the main text joined with *uel* or similar. Removing the most obviously contaminated manuscripts from the sample yielded clearer groups (fig. 9). Among β Va18 indeed seems to be the closest to the hyparchetype, although it is not the best witness as it has quite a lot of *Eigenfehler*. This led to its long branch. In contrast, Va19's text has few *Eigenfehler*, also as depicted. The other groups in the plot are clustered correctly: Ma1–Ro–Ba2, Va19–Es–Ve1. Fl7 is a reworking hence naturally removed from the rest, it seems that Vr should be closer to Ma1–Ro–Ba2. Only Ma1–Ro–Vr share one of the few eye-skips in the sample. The position of Sa2 is as yet unclear, it reads mostly with Va19–Ve1 and sometimes with Va18–Ox4 or even Vr–Fl7. However the relationship between these groups does not seem to be as depicted: they do not all stem from one ancestor as the plot might insinuate.

The most original manuscripts of α are Br and Ca3. Contamination is common in α , especially in the non-threefold witnesses, but many details are not yet clear to Rovati. The second sample produced the following twenty best scoring '*Leitfehler*' (including their relative score):

| | |
|-----------------------|--------------------------|
| removebis -- 100% | proximos -- 81% |
| auctor -- 85% | contracta -- 81% |
| accepta -- 85% | fortassis -- 81% |
| <u>fecit</u> -- 82% | penitus -- 81% |
| divise -- 82% | pervenerunt -- 79% |
| <u>diutius</u> -- 82% | inquisivi -- 79% |
| libros -- 82% | <u>nequit</u> -- 76% |
| perficitur -- 82% | <u>relationis</u> -- 76% |
| sumitas -- 82% | potavit -- 75% |
| sibique -- 81% | <u>perpendi</u> -- 74% |

The underlined ones are among the best significant errors Rovati has found manually. But the two highest scoring ones in the list are not good significant errors, e.g. *removebis* stands against *removebit*, *removeberit*, forms that can easily be changed. Some other promising significant errors, like *rei pro qua fit electio*, could not be found by the algorithm as each word occurs also elsewhere. This would be different if readings were used instead of words.

Despite the greater number of witnesses, the situation of the *Centiloquium* transmission is less complex and more amenable to computer-aided study than for the *Liber Aurelii*. Despite the widespread contamination, the two main groups were found correctly (especially in the second plot), the known sub-groups in the β version were also correct. Still, the most contaminated witnesses had to be removed manually and the root cannot be detected in the plots. Rovati's forthcoming edition intends to print one of the two versions with an extra apparatus for the other's changes.

IV LIMITATIONS OF DIGITAL TOOLS TODAY

Often, basically every witness has its own individual problems. They may range from physically missing parts to cancelled or otherwise lost text, missing chapters, unreadable characters or abbreviations. One is tempted to say '*chaque manuscript a son histoire*'. Fig. 10 shows a page of

manuscript E from the Aurelius tradition. It had some pages cut off and lacks about a third of each line of text. Another manuscript (M, see ill. 1 in [Roelli, 2021: 173]) was vandalised by a medieval physician who wrote the text from another branch over the original text where he thought it was better or more complete. In both cases there are many passages that are irretrievably lost. For the trees I plotted, I pragmatically inserted the readings of the closest relative where there were lacunae. For E this was A, for M it was C. If I had just left the missing words out, the resulting tree would have become much worse, as shorter witnesses tend to move to the centre of the tree. Such problems make a fully automated computerised treatment of the data at present illusory. Apart from such material problems, currently available stemma generating software struggles most with two problems: rooting and contamination (especially between various recensions).

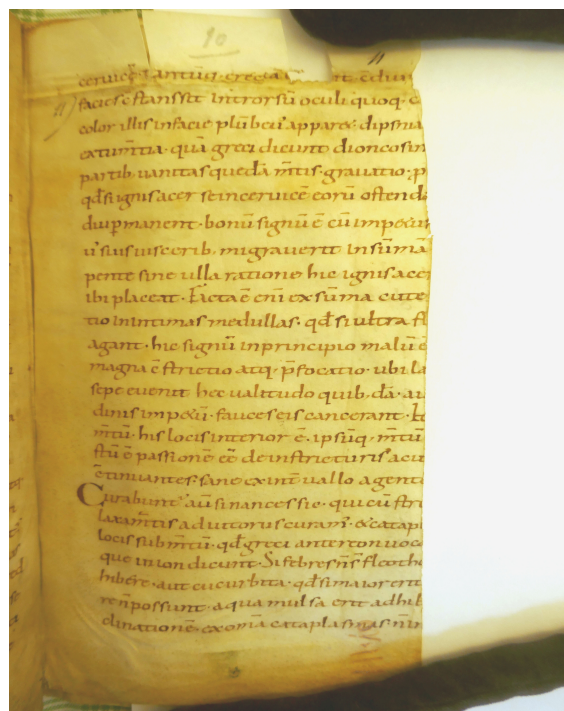


Figure 10: Cut manuscript page, Einsiedeln, Stiftsbibliothek 363, fol. 11r, of the *Liber Aurelii* text. Photograph by Philipp Roelli, published by permission.

4.1 Contamination

Contamination distorts the plots. If we add manuscript M's (*Liber Aurelii*) second contaminating hand to our plot,¹³ this 'M²' correctly groups with A and E, but the group MBC loses cohesion and the entire plot becomes distorted (fig. 11). Even at the very other end of the plot, G and I no longer form a group as they should. The basic idea of the Roelli-Bachmann algorithm for weighting variants ceases to function with clearly contaminated witnesses in the sample. But other software also suffers from this problem as it is constrained to find a tree and not a network, which would link the contaminated witness to its two (or more) sources. In a tree the contaminated witness is thus fit into a tree somewhere between its two ancestral groups. Mathematically, the correct solution would

¹³ For the sample text M² I used M's original text unless it was changed by the contaminating scribe.

be to plot not trees but networks that allow several incoming lines to nodes. There is software, such as SplitsTree, that can produce such networks (an example by Jean-Baptiste Guillaumin [Roelli, 2020: 351]). But first experiments along this line are not yet very promising. The degrees of freedom for the algorithm are too great, one would have to penalise extra edges strongly, otherwise the algorithm will plot them for any random identical spelling between unrelated witnesses. A pragmatic solution is to identify contaminated witnesses and add them later manually to the automatically plotted tree of the non-contaminated ones. However, it is not always easy to find out that a manuscript is contaminated and in some cases, like the *Centiloquium*, most witnesses are contaminated.

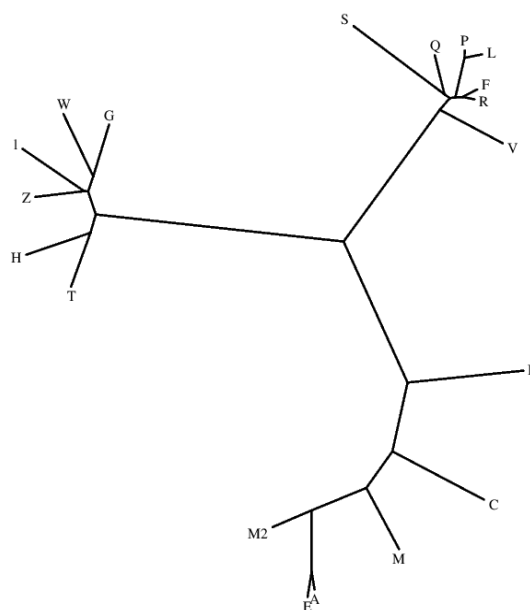


Figure 11: As fig. 7 but including the contaminated text-layer M^2 .

4.2 Rooting

The second great problem for tree-finding software is rooting. Whereas biologists can usually use an outgroup to root their trees, for instance a rodent for a genealogical tree of primates, this is in most cases not an option in our field as texts are created (more or less) *ex nihilo* by their authors.¹⁴ Procedures using textual distances can by definition not serve for identifying the root: distance is commutative, the information of the direction of change is not contained in the distance data. One would have to add a mathematical notion of direction: In order to decide where the root is in a tree one needs to know for at least some variants which one is primary and which ones are secondary. It is easy to see that the naive first guess that the root will be in the centre of the tree plot is in general wrong. Fig. 12 shows a constructed example with six loci symbolised by upper and lower-case letters of the alphabet. The obvious idea to start grouping close relatives together again and again in order to bring direction into the stemma, does not work in general either (even if we take for granted that the archetype is not among the witnesses). In the example, one would easily guess that the two main

¹⁴ In some cases early translations can serve as outgroups, cf. Caroline Macé [Roelli, 2021, ch. 3.2].

groups of readings ABC and bcd are significant variants shared by about half the tradition each and one will be tempted to posit the archetype where they intersect. This would be between α and γ (black line on the right hand side), which is not at all where the archetype is found in truth. We note in passing that in case of a stemma with a bifurcation at the top (a common case, as [Bédier, 1928] showed), the root or archetype is not found at a junction of lines but along one of the lines in the tree, like in this example (red arrow). As long as the information of directionality is not in the computational data, the root cannot be found by any approach without further data. An idea for future development could be to develop an algorithm that can detect eye-skips and use them to add an element of directionality for parts of the tree by assuming that the texts with the eye-skip are the derived text-state.

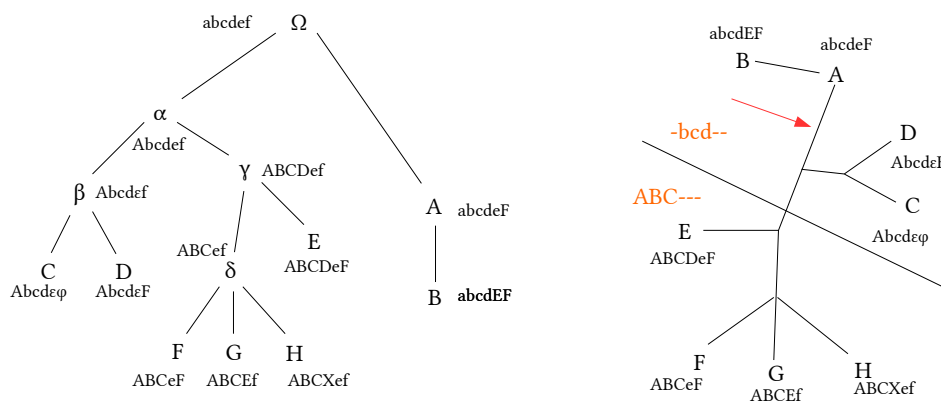


Figure 12: Constructed schematic example: stemma vs. unoriented tree.

VI A PRAGMATIC APPROACH FOR THE TIME BEING

The length of the text and the number of manuscripts are the main constraints to determine the time (and thus money) it takes to edit it if we take a skilled editor for granted. Thus it may be that a long text survives in a hundred witnesses but the time and money are only sufficient to use half a dozen manuscripts for the entire text. Similarly the entire variance of a tradition (like the *Centiloquium*) may not be worth recording as it is largely derived and repetitive. In such cases it is crucial to choose the right sample of witnesses. A relatively small text sample – or better two from different parts of the texts – are transcribed, potentially by a team of assistants. Software can then be used to plot trees; together with the other shown pieces of software true significant errors can be identified and with them *loci critici*. Analysing the readings of only these will help to draw a stemma for the sample texts. A large percentage of witnesses will often quickly be found to be derivative and of little importance for the reconstruction of the archetypal text (although they may well be interesting for other research questions). The remaining witnesses can be studied more in depth and a further sub-sample of them can be chosen for the edition. The archetype text is then reconstructed as far as possible and examined, *emendatio* may have to be applied to passages that can be shown to have been wrong in the archetype. As this is quite speculative, it should be done with circumspection and

clearly stated for each passage in the edition so that the reader can decide for himself whether he agrees with the emendation.

For the *Liber Aurelii* edition I decided that I could just as well do the work manually except for the Gariopontus recension where I pragmatically chose five of the oldest manuscripts and the early print from a sample of about a dozen. For a critical edition of the entire Gariopontus one will have to study the tradition in more depth but for the *Liber Aurelii* passages (which make up some 12% of Gariopontus' compilation) this seemed fine as the variation in the Gariopontus witnesses was relatively minor and the readings of the full *Liber Aurelii* could be compared. For the *Centiloquium* (having more manuscripts but less variance) software may be of greater use. The approach sketched above can help save time if there are many witnesses of similar length and many of them not contaminated – although there is still a lot of crucial manual work and critical thinking left for the editor. At least the alignment can be done fully automatically and there is no need to consider readings while transcribing. Indeed, several people can easily transcribe witnesses at the same time. Nonetheless transcription remains the bottle-neck of labour; OCR software for hand-written documents may bring help at this point in the not-too-distant future. Contamination must be found manually¹⁵ and the tree must be manually rooted.

VII POSSIBLE FUTURE DEVELOPMENTS

OCR software today already uses dictionaries with the help of which *dominus* gets weighted far more strongly than similar looking but nonsensical *dorunus* if the text is known to be in Latin. This approach could be strongly improved once software can construct syntactic trees automatically and then weight their probabilities. For instance the software would then expect a verb in the subjunctive in an *ut* clause. An automated analysis of vocabulary and style is a further possibility to improve results: if an author always says *scapulae* never *palae* for 'shoulder', an algorithm should consider the first word as the likely primary reading. In case other works of an author are known, their vocabulary and their syntactical trees could additionally be used to weight possible readings or trees more accurately for the author in question. These are basically the things an experienced philologist does intuitively when editing a text. Thus, the goal for software should be to mimic the behaviour of an experienced editor. Computerised methods should thus strive to approach the gold standard of our reconstructive textual editing – Neo-Lachmannian critical text editing [Trovato, 2017] – one should avoid using software from other fields trusting that the results will be fine even though the software was tailored for other problems involving 'descent with modification'.

References

- Bédier, Joseph (1928). La Tradition manuscrite du Lai de l'ombre: Réflexions sur l'art d'éditer les anciens textes. *Romania* 54: 161–196, 321–356. [reprint Paris: Champion, 1970. gallica.bnf.fr/ark:/12148/bpt6k8980]

15 There are special approaches that try to cope with highly contaminated transmissions, such as versions of the Bible: [Mink, 2011].

- Froger, Jacques (1978). Westgöta-Lagen, Edited by H. S. Collin and C. J. Schlyter. Facsimile edition with an addendum by Otto von Friesen, Our Oldest Manuscript in Old Swedish. Edited by Gösta Holm. *Scriptorium* 32/1: 183–188.
- Glaze, Eliza (2005). Galen refashioned: Gariopontus in the Later Middle Ages and Renaissance. In: Elizabeth Lane Furdell (ed.). *Textual healing: Essays on Medieval and Early Modern medicine*. Brill (Leiden): 53–75.
- Holm, Gösta (1972). Carl Johan Schlyter and Textual Criticism. *Saga och Sed* 1972: 49–80.
- Juste, David (2022). Pseudo-Ptolemy, Centiloquium (tr. Plato of Tivoli) (update: 10.03.2022). *Ptolemaeus Arabus et Latinus. Works*. Online: <http://ptolemaeus.badw.de/work/41>.
- Lemay, Richard Joseph (1978). Origin and Success of the Kitāb Thamara of Abū Ja‘far Aḥmad ibn Yūsuf ibn Ibrāhīm from the Tenth to the Seventeenth Century in the World of Islam and the Latin West’. In: *Proceedings of the First International Symposium for the History of Arabic Science* (Aleppo, April 5-12, 1976). Aleppo University Press (Aleppo): II, 91–107.
- Macé, Caroline, Ilse De Vos, and Koen Geuten (2012). Comparing Stemmatological and Phylogenetic Methods to Understand the Transmission History of the ‘Florilegium Coislinianum.’ In: Alessandra Bucossi and Erika Kihlman. *Ars Edendi Lecture Series*. Stockholm University Library (Stockholm): 2:107–129.
<http://www.diva-portal.org/smash/get/diva2:551286/FULLTEXT01.pdf>
- Martorello, Franco, and Giuseppe Bezza (eds.) (2013). Aḥmad ibn Yūsuf ibn al-Dāya. *Commento al Centiloquio Tolemaico*, Mimesis (Milano).
- Mink, Gerd (2011). Contamination, Coherence, and Coincidence in Textual Transmission. In: Klaus Wachtel and Michael W. Holmes. *The Textual History of the Greek New Testament. Changing Views in Contemporary Research*. SBL (Atlanta): 141–216.
- Olrik Frederiksen, Britta (2009). Stemmaet fra 1827 over Västgötalagen – en videnskabshistorisk bedrift og dens mulige forudsætninger. *Arkiv för nordisk filologi* 124: 129–150.
- Roelli, Philipp, and Dieter Bachmann (2010). Towards Generating a Stemma of Complicated Manuscript Traditions: Petrus Alfonsi’s Dialogus. *Revue d’histoire des textes* n.s. 5: 307–321.
- Roelli, Philipp (ed.) (2020). *Handbook of Stemmatology. History, Methodology, Digital Approaches*. De Gruyter (Berlin).
<https://doi.org/10.1515/9783110684384>
- Roelli, Philipp (2021). Liber Aurelii ‘On Acute Diseases’, critical edition. *Beihefte zum Mittellateinischen Jahrbuch* 21. Hiersemann (Stuttgart). <https://doi.org/10.36191/9783777222035>
- Schlyter, Carl Johan, and Hans Samuel Collin (eds) (1827). *Westgöta-Lagen*. Häggström (Stockholm).
- Trovato, Paolo (2017). *Everything You Always Wanted to Know about Lachmann’s Method: A Non-Standard Handbook of Genealogical Textual Criticism in the Age of Post-Structuralism, Cladistics, and Copy-Text*. Libreriauniversitaria.it (Padova).