

Intertextual Distance and *Grand Siècle* theatre

CHARLES BERNET [English translation by Luke Arnason]

1. Corneille and Molière (the affair)

In the first years of this century, a question has animated the community of researchers in quantitative lexicology: do statistical methods validate the hypothesis suggesting that Pierre Corneille is the author of Molière's comedies? The authors of an article entitled "Inter-Textual Distance and Authorship Attribution. Corneille and Molière" (Labbé and Labbé : 2001) argue that they do, based on a calculation of "intertextual distance" in the vocabulary of classical theatre.

This is an idea that had some supporters in the 20th century. Pierre Louÿs first formulated it in 1919¹ after having perceived certain "cornelian" characteristics in passages of [Molière's] *Amphitryon*, giving rise to a first Corneille-Molière polemic. The same hypothesis was then revisited, by Henry Poulaille (1951 and 1957), then by Hippolyte Wouters and Cécile de Ville de Goyet (1990), without convincing the community.

Despite the controversy that followed the publication of the 2001 article, we observe in recent years an increased interest for this hypothesis, now defended by arguments based on scientific methods. The large number of publications² dealing with this question is a testament to this renewal of interest in the subject.

2. Application of Cyril and Dominique Labbé's method to a broader corpus of works

In this article, we will not dwell on the questions that deal with literary history and criticism. It seems to me that the essential has already been said, notably by Georges Forestier (2003).

What remains is to pursue the expertise of Cyril and Dominique Labbé's study on the points that merit further attention today. This contribution will consist of filling in a shadowy area, testing their method on other plays from the same chronological period, so as to widen the field of comparisons within a homogeneous group and, secondly, to illustrate the calculation of intertextual distances in the lexical sub-group of rhyming words.

Since Dominique Labbé was careful to make his data and programs available to researchers (in JADT-2004), it is possible to reproduce the procedure step-by-step and to continue the experiment on other texts. The initial corpus included the plays of Corneille, Molière and Racine, to which we have added a selection of six other plays. Three Comedies: *L'Intrigue des filous* by Claude de L'Estoile (1648) as well as *Le Distrain* (1698) and *Le Légataire universel* (1708) by Jean-François Regnard, and three tragedies: *Astrate, roi de Tyr* by Philippe Quinault (1665), *Ariane* by Thomas, younger brother of Pierre Corneille (1672) and *Électre* by Longepierre (early 18th century). These texts belong to the genres already represented, they cover a chronological period stretching from the middle of the 17th century to the first years of the 18th century and diversify the group of authors taken into account.

The principal difficulty of the procedure resides in the preparation of the data, which is quite involved. It is essential, before any comparison can be made, to unify the spelling of the texts that do not follow the same editorial principles. The texts analysed by Cyril and Dominique Labbé follow "modernised" spelling from the 19th century. Some of the texts that we have added are scholarly editions maintaining the spelling of the period, while others are aimed at a broader public

1 "I am certain about Molière, every day brings me new evidence and never the slightest difficulty" (Letter to Frédéric Lachèvre, November 27, 1919 in Louÿs : 1938, p. 88). Pierre Louÿs was certain about something else, in relation to the *Histoire comique de Francion*: "Dear Sir, you shall know my secret before I publish it: it is Francion. Pierre Corneille wrote Francion in class, under the Jesuits, during his second year of college, at the age of sixteen years and in the space of three months" (Letter to Frédéric Lachèvre, November 20, 1919, in Louÿs : 1938, p. 86).

2 See in particular Boissier (2004), Goujon and Lefrère (2006) and Ferrand (2008).

and use the spelling of today.

In the protocol of preparation of data used by Dominique Labbé, the corpus is lemmatised before any statistical analysis is carried out. This procedure consists of attaching, in coded form, a categorical tag to each word of the text, so as to group different occurrences of a same lexeme and to distinguish between homographs. Our interventions consisted essentially of prescribing a lemma of “modern French” to archaic spellings and occasionally correcting erroneous tags.

Since our objective was simply to expand on a previous experiment, we will not dwell on the protocol of tagging set out by Dominique Labbé. We would make the single observation, however, that the demarcation between the verbal and adjectival past participles follows a rather summary mechanic that inflates the class of adjectives and increases the number of lemmata. The grammatical taxonomies are never definitive or perfect, and no norm can fully satisfy linguists, but it must be said, to Dominique Labbé's credit, that he did take the trouble of lemmatising all of the texts studied, thus avoiding founding his research on spellings in their raw state. The criticisms levelled at him on that subject are unjustified.

The method of calculation of the standardised measure of the actual distance between texts has been laid out in the article mentioned above, which we invite the reader to consult. The “intertextual distance” between two texts is formulated by an index whose value is situated between 0 and 1. For Cyril and Dominique Labbé, this index does not simply provide an evaluation of the distance that separates the vocabulary of the texts studied, it also makes it possible to identify the author of a text.

The critical levels of the index were established empirically, the calculation having “been applied to several thousand texts of all origins” (Labbé : 2003, p. 14). Dominique Labbé affirms (pp. 14 and 15) that:

A distance equal or inferior to 0.20 indicates with certainty that the texts are by the same author. Even when a writer is “pastiche” another, the distance between the pastiche and the originals is always superior to this level.

Between 0.20 and 0.25, it is virtually certain that the author is the same. [...] In the case of literary works that belong to two different authors, a distance equal or inferior to 0.25 indicates either a close collaboration or the certain plagiarism of the first by the second (when it is possible to know with certainty who the first is!).

Above 0.25, we enter into a grey area where two hypotheses are possible: a same author dealing with different themes or two contemporary authors a same theme in their own style... Such that the higher the value rises above 0.25, the more difficult it becomes to attribute an anonymous text to the considered author, without excluding the possibility altogether.

Above 0.40, the authors are certainly different or, for a single author, the texts are of radically different genres, for example, oral and written.

3. Values of the index in light of the scheme proposed by Dominique Labbé

Table 1 presents two distinct sub-groups, comedies and tragedies. As Jean-Marie Viprey (2003) has shown, the index is sensitive to differences in length. Since the interpretation of distance runs the risk of being biased when differences in length are too dramatic, the shortest plays were excluded.

The left part of the table gives detailed results for Corneille's *Le Menteur* and *La Suite du Menteur* — let us recall that it was the relative proximity of these two comedies by Corneille to several plays by Moliere that led Cyril and Dominique Labbé to endorse Pierre Louÿs' idea — and for *L'Intrigue des filous* as well as *Le Distrait* and *Le Légataire universel* by Regnard.

We observe that the values of the index are, for the three comedies added to the corpus, inferior to 0.3; between these and the two *Menteur* plays by Corneille, the index is below crucial 0.25 mark; and it is inferior or equal to this same value 5 times between the *Intrigue des filous* and the plays by Moliere, and 8 times between each of the plays by Regnard and those by Moliere.

For these three plays, as for the *Menteur* plays, the lowest values are obtained in the comparisons with Moliere's verse plays (see especially *L'Étourdi*, *Le Dépit amoureux*, *Sganarelle*, *L'École des maris*, *Les Fâcheux*, *L'École des femmes*, *Tartuffe*, *Le Misanthrope*, *Les Femmes*

savantes) and, among those in prose, with *L'Avare*. In all cases, the values of the index of the comedies by Regnard are extremely close to those of the *Menteur* plays, despite a chronological gap of over a half-century. Between the plays by Moliere and those of Regnard, there are similarities comparable to those observed between Corneille and Moliere.

The right side of the table presents the results concerning the sub-group that includes tragedies, some tragicomedies and Corneille's "heroic comedies". The extract gives the values of the index in all comparisons concerning Quinault's *Astrate*, Thomas Corneille's *Ariane* — these two tragedies being contemporary of Pierre Corneille's last tragedies and Racine's secular tragedies — and Longepierre's *Électre*, a much later play.

Here we remark greater contrast. For Longepierre's play, the values of the index are distinctly higher — they often pass 0.30. On the other hand, in the two other cases, these values are considerably inferior to the 0.25 mark, and on several occasions extremely close to 0.20. In the comparisons with Pierre Corneille, the index remains under the 0.25 mark 12 times out of 24 for *Ariane* and 19 times out of 24 for *Astrate*.

It is apparent from these observations that the lexical proximity, altogether relative, between Corneille and Moliere is nothing exceptional. Applied to the dramatic texts examined here, the scheme proposed by Dominique Labbé would lead us to unfounded speculations and unrealistic hypotheses.

The index of intertextual distance cannot be considered as a valid tool in the matter of the attribution of texts³. The parameters subject to questioning are numerous: the subject, the themes, the genre, the author, the period, the length of the text, its individual dynamic, etc., and it is impossible to isolate these in order to measure precisely the influence of each factor. Furthermore, it is highly probable that the personal mark of the author is of insufficient quantitative weight to cause the index of intertextual distance to vary appreciably. By contrast, searching for markers, even isolated or sporadic ones, specific to each author constitutes a more fruitful approach⁴.

4. Tree diagrams of intertextual distance

Nevertheless, we should not throw out the baby with the bath water. Dominique Labbé's method, which Jean-Marie Viprey presented in a polemical fashion as being "crude and outdated" ("un appareillage fruste et suranné"), is nevertheless susceptible to giving statistically significant results when it is used within the limits of its validity, and more precisely when the plays examined are not too different in length.

The tree diagram of intertextual distance established in the corpus of tragedies and other plays represents, through the grouping of paths, the particularities clearly shown in other studies on quantitative lexicology (in particular Muller: 1967 and Bernet: 1983).

Inside this group, the extremities of the tree (figure 1) are occupied by Corneille (upper part of the table) and Racine (lower part of the table). The works of Corneille are distributed over several groups: one relatively dense group in the upper right corner containing the first plays among which we recognise the classical masterpieces, another, more spare, in the upper left, including his later plays. Under these two well-defined groups, the plays written between *Le Menteur* and *Pertharite* dominate, a period at the end of which Corneille temporarily interrupted his career as playwright. Here we find the plays through which Corneille tried to renew himself, such as *Andromède* or *La Toison d'or* (which both include sung parts), two christian plays including *Polyeucte* (which belongs to the classical period) and *Théodore* which was, like *Pertharite*, a failure.

In the lower part of the tree, the tragedies by Racine, less numerous than Corneille's, occupy

³ Étienne Brunet (2004) wrote : "[...] in the absence of an adequate number of trials, we refuse the idea of a fixed scale, of an arbitrary scheme, attached to a single, global and undifferentiated measure, applied to a sole – lexical – aspect of language no less."

⁴ Thus, when we compare Corneille and Moliere, the frequency of appellatives or of pejorative qualifications (*coquin*, *faquin*, *fripouille*, *jocrisse*, *pécore*, *pendard*, etc.) or, a less obvious but more symptomatic case, the use of the word *parfois*. See Bernet (2004), pp. 153-154.

a greater amount of space; as it is well known, the lexical diversity is more marked in Racine. His first two plays, far from one another, are situated on intersections to the middle part of the diagram. The secular masterpieces are grouped together in quite a compact fashion, from *Andromaque* to *Mithridate*, plays whose vocabulary is quite limited. A group of more loosely connected pieces includes the last two Greek tragedies and his biblical tragedies, that is to say the more lexically rich plays. We note in this group the presence of Longepierre's *Électre* which is, of the whole group, the secular tragedy farthest from Corneille's plays. Thomas Corneille's *Ariane*, a mythological play reputed to be “racinian”, and Quinault's *Astrate*, a “gallant” tragedy (like Racine's *Alexandre le Grand*), relatively close to one another, are situated in the median part of the tree, midway between the two great classics. It seems reasonable that by fleshing out the corpus with other plays by these two authors, a third pole could emerge from this median space.

5. Intertextual distance and rhyming words

It is also possible to apply the calculation of intertextual distance to lexical data other than overall vocabulary. This is what we have done with rhyming words in a group of 16 comedies by Corneille and Moliere.

The interest of such research resides in the fact that the selection of rhyming words relates to an investment of a different nature for the author than the ordinary choice of words. One can thus expect results representative of this situation.

Since it is a question of words chosen by the versifier for their form, we took into account non lemmatised data. The values of the index, very high, range between 0.60 (between *La Galerie du Palais* and *La Place Royale*) and 0.77 (between *Mélite* and *Les Femmes savantes*). The most striking result here is the distance displayed in the tree diagram (figure 2). Both author's plays are distributed very clearly into two distinct groups: all Corneille's comedies, including *Le menteur* and its sequel, are situated at the branching out of the upper part and all of Moliere's plays, including *Dom Garcie de Navarre*, his most singular play, occupy the opposite space, that is to say the lower half of the tree.

If the tree diagrams can be trusted, then it must be admitted that the index of intertextual distance indicates marked lexical differences between Corneille and Moliere.

The experiments reported in this article invalidate the conclusions of Cyril and Dominique Labbé and show that quantitative lexicology does not offer arguments in favour of Pierre Louÿs “intuitions”.

6. Bibliography

6.1. Corpus

CORNEILLE (Pierre): 1862-1868, *Œuvres*, new ed. by M. Ch. MARTY-LAVEAUX, t. 1-12 (Paris: Hachette).

CORNEILLE (Thomas): 1986, *Ariane*, text established, presented and annotated by J. TRUCHET, in *Théâtre du XVIIe siècle*, t. 2 (Paris: Gallimard, “Bibliothèque de la Pléiade”), pp. 907-966.

L'ESTOILE (Claude DE): 1977, *L'Intrigue des filous*, based on the 1st ed. (1648) by Roger GUICHEMERRE (Paris: H. Champion, “S.T.F.M.”), pp. 3-135.

LONGEPIERRE (Hilaire-Bernard DE ROQUELEYNE DE): 1981, *Électre*, text established and presented by T. TOBARI (Paris: A.-G. Nizet).

MOLIÈRE: 1873-1900, *Œuvres*, new ed. by Eugène DESPOIS, t. 1-14 (Paris: Hachette).

QUINAULT (Philippe): 1986, *Astrate, roi de Tyr*, text established, presented and annotated by J. TRUCHET, in *Théâtre du XVIIe siècle*, t. 2 (Paris: Gallimard, “Bibliothèque de la Pléiade”), pp. 1040-1100.

RACINE (Jean): 1865-1873, *Œuvres complètes*, new ed. by Paul MESNARD, t. 1-8 (Paris: Hachette).

REGNARD (Jean-François): 1820, *Le Distrait*, with *avertissements* by M. GARNIER, *Œuvres*, t. 2 (Paris: E. A. Lequien), pp. 282-395.

REGNARD (Jean-François): 1994, *Le Légataire universel*, text established, presented and annotated by Ch. MAZOUER, in *Le Légataire universel suivi de la Critique du Légataire* (Geneva: Droz, "T.L.F."), pp. 89-261.

6.2. Studies

BERNET (Charles): 1983, *Le Vocabulaire des tragédies de Jean Racine. Étude statistique* (Geneva: Slatkine-Champion).

BERNET (Charles): 2004, "Hasards de la rime", in PURNELLE (Gérald), FAIRON (Cédric) and DISTER (Anne), eds. *Le Poids des mots. Actes des septièmes Journées internationales d'analyse statistique des données textuelles*, vol. 1 (Louvain: PUL, 2004), pp. 148-159.

BOISSIER (Denis): 2004, *L'Affaire Molière. La Grande Supercherie littéraire* (Paris: J.-C. Godefroy).

BRUNET (Étienne): [2004], "Où l'on mesure la distance entre les distances", <http://www.revue-texto.net/Inedits/Brunet/Brunet_Distance.html#II>.

FERRAND (Franck): 2008, *L'Histoire interdite. Révélation sur l'Histoire de France* (Paris: Taillandier).

FORESTIER (Georges): 2003, "D'un vrai canular à une fausse découverte scientifique, à propos des travaux de Dominique et Cyril Labbé"; "Postscriptum (June 1, 2003)"; "Faux témoin ? ou falsification historique ? À propos des contrevérités contenues dans le livre de M. Labbé (July 1, 2003)"; "L'affaire Corneille-Molière, suite de l'histoire d'un canular qui a la vie dure", <<http://www.crht.org/>> [Web site of the Centre de Recherche sur l'Histoire du Théâtre].

GOUJON (Jean-Paul) et LEFRÈRE (Jean-Jacques): 2006, *Ôte-moi d'un doute... L'Énigme Corneille-Molière* (Paris: Fayard, « Littérature française »).

JADT-2004 Louvain-la-Neuve. CD included with *Le Poids des mots. Actes des septièmes Journées internationales d'analyse statistique des données textuelles* (Louvain: PUL).

KYLANDER (Britt-Marie): 1995, *Le Vocabulaire de Molière dans les comédies en alexandrins* (Göteborg: Acta universitatis Gothoburgensis).

LABBÉ (Dominique): 2003, *Corneille dans l'ombre de Molière. Histoire d'une découverte* (Paris-Brussels: Les Impressions nouvelles).

LABBÉ (Cyril) and LABBÉ (Dominique): 2001, "Inter-Textual Distance and Authorship Attribution. Corneille and Molière", *Journal of Quantitative Linguistics*, 8, 3, pp. 213-231.

LOUÏS (Pierre): 1938, *Broutilles recueillies par Frédéric Lachèvre (1870-1925)* [at the end of the volume: "Le problème Corneille-Molière vu par P. LouÏs (contribution au dossier définitif)"] (Paris: Frédéric Lachèvre).

LUONG (Xuan): 1994, "L'analyse arborée des données textuelles. Mode d'emploi", *Travaux du cercle linguistique de Nice*, 16.

MULLER (Charles): 1967, *Étude de statistique lexicale. Le Vocabulaire du théâtre de Pierre Corneille* (Paris: Larousse).

POULAILLE (Henry): 1951, *Tartuffe ou la comédie de l'Hypocrite*, presentation and preface by H. P. (Paris: Amiot-Dumont).

POULAILLE (Henry): 1957, *Corneille sous le masque de Molière* (Paris: Grasset).

VIPREY (Jean-Marie): [2003], "Morneille, Colière et messieurs Labbé", <<http://laseldi.univ-fcomte.fr/Archives/affaireMorneilleColiere/morneille.htm>>.

WOUTERS (Hippolyte) et VILLE DE GOYET (Cécile DE): 1990, *Molière ou l'auteur imaginaire* (Brussels: Complexe).

7. Appendices

<i>Comédies</i>						<i>Tragédies et pièces diverses</i>			
	C14	C15	ES01	Re01	Re02		Qu01	TC01	Lo01
C01	0,227	0,236	0,269	0,280	0,278	C02	0,283	0,287	0,300
C03	0,231	0,237	0,273	0,277	0,278	C07	0,293	0,297	0,293
C04	0,217	0,220	0,258	0,260	0,267	C09	0,258	0,274	0,297
C05	0,227	0,228	0,267	0,278	0,278	C10	0,261	0,292	0,301
C06	0,251	0,242	0,285	0,299	0,296	C11	0,245	0,278	0,304
C08	0,226	0,228	0,243	0,271	0,263	C12	0,222	0,239	0,292
C14		0,180	0,232	0,240	0,245	C13	0,255	0,283	0,307
C15	0,180		0,231	0,229	0,241	C16	0,219	0,254	0,288
ES01	0,232	0,231		0,257	0,251	C17	0,224	0,233	0,309
M01	0,205	0,206	0,224	0,221	0,211	C18	0,216	0,252	0,303
M02	0,215	0,211	0,231	0,219	0,213	C19	0,238	0,267	0,300
M03	0,280	0,273	0,296	0,295	0,289	C20	0,224	0,257	0,344
M04	0,223	0,217	0,237	0,222	0,227	C21	0,223	0,261	0,343
M05	0,248	0,248	0,264	0,258	0,248	C22	0,234	0,243	0,322
M06	0,226	0,217	0,228	0,228	0,222	C23	0,208	0,242	0,295
M07	0,252	0,243	0,256	0,260	0,254	C24	0,227	0,245	0,296
M08	0,242	0,232	0,256	0,233	0,240	C25	0,214	0,246	0,336

M09	0,263	0,247	0,268	0,255	0,250	C26	0,216	0,239	0,334
M10	0,292	0,289	0,284	0,275	0,266	C27	0,219	0,238	0,341
M11	0,252	0,233	0,270	0,241	0,255	C28	0,209	0,231	0,360
M12	0,297	0,289	0,281	0,269	0,271	C29	0,229	0,251	0,328
M13	0,252	0,256	0,247	0,253	0,252	C30	0,208	0,228	0,348
M14	0,292	0,279	0,286	0,268	0,268	C31	0,204	0,227	0,354
M15	0,257	0,244	0,263	0,246	0,242	C32	0,202	0,209	0,346
M16	0,292	0,282	0,280	0,275	0,268	Qu01		0,218	0,325
M17	0,282	0,279	0,279	0,282	0,281	R01	0,239	0,275	0,291
M18	0,294	0,282	0,291	0,269	0,274	R02	0,259	0,282	0,302
M19	0,269	0,263	0,266	0,251	0,253	R03	0,252	0,242	0,298
M20	0,260	0,250	0,269	0,238	0,252	R05	0,258	0,272	0,310
M21	0,286	0,279	0,288	0,260	0,258	R06	0,265	0,259	0,316
R04	0,296	0,293	0,283	0,260	0,264	R07	0,256	0,249	0,308
Re01	0,240	0,229	0,257		0,206	R08	0,235	0,257	0,276
Re02	0,245	0,241	0,251	0,206		R09	0,266	0,292	0,265
						R10	0,282	0,292	0,263
						R11	0,357	0,380	0,306
						R12	0,343	0,379	0,293
						TC01	0,218		0,344
						Lo01	0,325	0,344	

Table 1. Intertextual distance

Abbreviations - C01 : *Mélite* ; C02 : *Clitandre* ; C03 : *La Veuve* ; C04 : *La Galerie du Palais* ; C05 : *La Suivante* ; C06 : *La Place Royale* ; C07 : *Médée* ; C08 : *L'Illusion comique* ; C09 : *Le Cid* ; C10 : *Horace* ; C11 : *Cinna ou la Clémence d'Auguste* ; C12 : *Polyeucte* ; C13 : *La Mort de Pompée* ; C14 : *Le menteur* ; C15 : *La Suite du menteur* ; C16 : *Rodogune* ; C17 : *Théodore* ; C18 : *Héraclius* ; C19 : *Andromède* ; C20 : *Don Sanche d'Aragon* ; C21 : *Nicomède* ; C22 : *Pertharite* ; C23 : *OEdipe* ; C24 : *La Toison d'Or* ; C25 : *Sertorius* ; C26 : *Sophonisbe* ; C27 : *Othon* ; C28 : *Agésilas* ; C29 : *Attila* ; C30 : *Tite et Bérénice* ; C31 : *Pulchérie* ; C32 : *Suréna* ; Es01 : *L'Intrigue des filous* ; Lo01 : *Électre* ; M01 : *L'Étourdi ou les Contretemps* ; M02 : *Le Dépit amoureux* ; M03 : *Dom Garcie de Navarre ou le Prince jaloux* ; M04 : *L'École des maris* ; M05 : *Les Fâcheux* ; M06 : *L'École des femmes* ; M07 : *La Princesse d'Élide* ; M08 : *Tartuffe ou l'Imposteur* ; M09 : *Dom Juan ou le Festin de pierre* ; M10 : *L'Amour médecin* ; M11 : *Le Misanthrope ou l'Atrabilaire amoureux* ; M12 : *Le Médecin malgré lui* ; M13 : *Amphitryon* ; M14 : *George Dandin ou le Mari confondu* ; M15 : *L'Avare ou l'École du mensonge* ; M16 : *Monsieur de Pourceaugnac* ; M17 : *Les Amants magnifiques* ; M18 : *Le Bourgeois gentilhomme* ; M19 : *Les Fourberies de Scapin* ; M20 : *Les Femmes savantes* ; M21 : *Le Malade imaginaire* ; Qu01 : *Astrate, roi de Tyr* ; R01 : *La Thébaine* ; R02 : *Alexandre le Grand* ; R03 : *Andromaque* ; R04 : *Les Plaideurs* ; R05 : *Britannicus* ; R06 : *Bérénice* ; R07 : *Bajazet* ; R08 : *Mithridate* ; R09 : *Iphigénie* ; R10 : *Phèdre* ; R11 : *Esther* ; R12 : *Athalie* ; Re01 : *Le Distrait* ; Re02 : *Le Légataire universel* ; TC01 : *Ariane*.

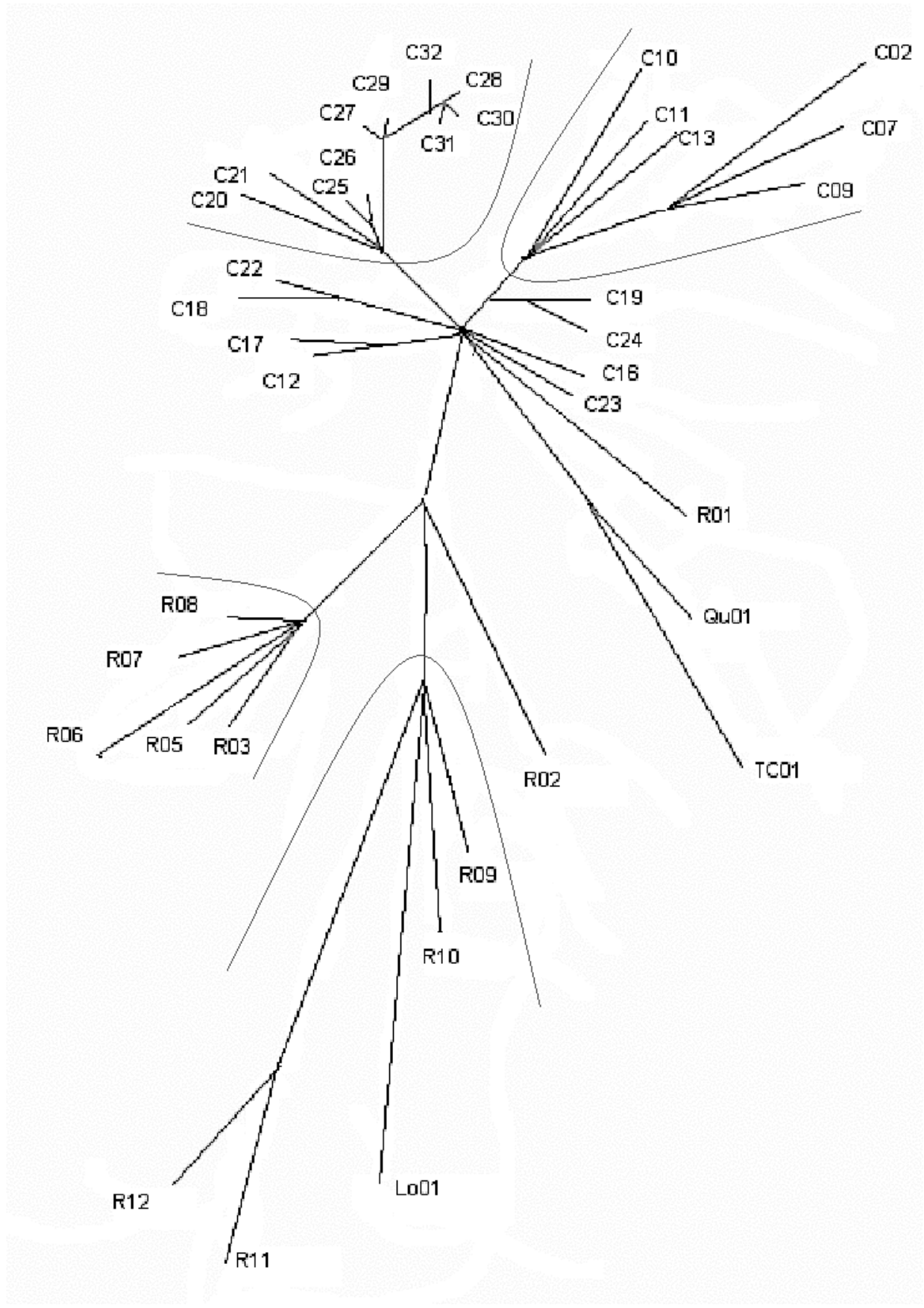


Figure 1. Tragedies and other plays
 Tree diagram by Xuan Luong (Université de Nice)

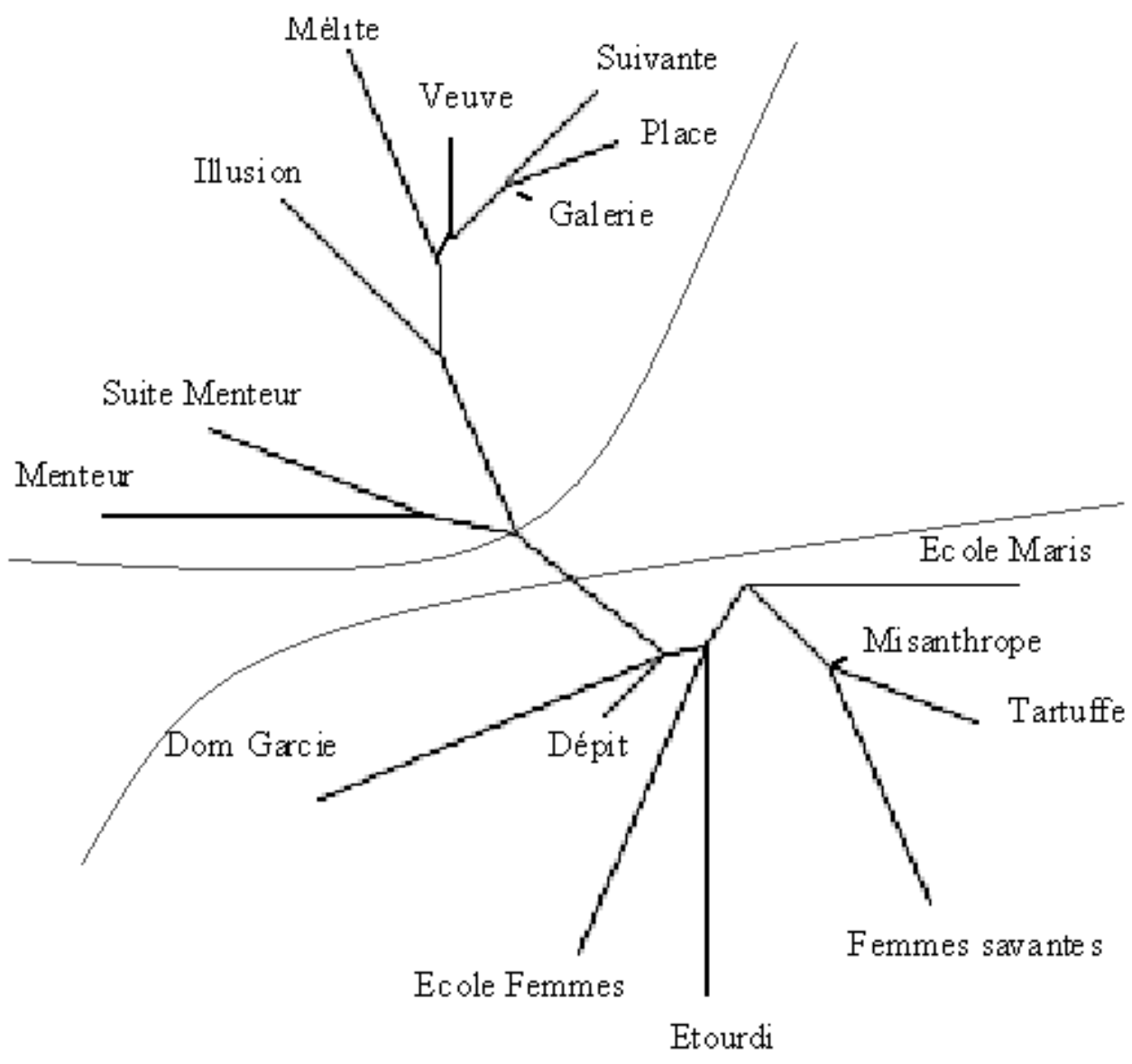


Figure 2. Rhyming word-forms – Corneille and Moliere – 16 comedies
 Tree diagram by Xuan Luong (Université de Nice)