

# Risorse e sfide per la collazione automatica di testi medievali

Elena Spadini

L'utilizzo dell'informatica nelle pratiche ecdotiche ha una lunga storia, che ha inizio con il tentativo di automatizzare l'attività di collazione. A partire dagli anni '60, vengono identificate le operazioni di base di un programma di allineamento e riconosciuti gli aspetti più problematici della sua applicazione.<sup>1</sup> Questo contributo vuole offrire una panoramica dello stato dell'arte, con particolare attenzione alla collazione automatica di testi romanzati medievali.<sup>2</sup> In secondo luogo e brevemente, per completare l'analisi della situazione attuale, si propongono nuovi approcci ad alcuni degli annosi problemi già noti ai pionieri della filologia computazionale.

## *Collazione automatica e Gothenburg model*

La collazione consiste nel comparare due o più testi al fine di individuare le loro differenze. In termini filologici, ciò che si compara sono le lezioni dei testimoni e tramite la collazione si ottengono le varianti tra di essi.

Il procedimento definito 'collazione automatica' prevede l'utilizzo di un computer in una o più delle sue fasi: non si tratta quindi di delegare alla macchina l'intero procedimento, come vedremo, ma di avere la sua assistenza in vari momenti. In questo articolo useremo 'collazione automatica' nell'accezione qui esplicitata e per brevità, rispetto a formule più trasparenti ma che appesantirebbero il dettato, come 'collazione semi-automatica' o 'collazione assistita'.

<sup>1</sup>E. Nury, E. Spadini, «From giant despair to a new heaven: The early years of automatic collation», *it - Information Technology*, vol. 62, 2 (2020), DOI:10.1515/itit-2019-0047. <sup>2</sup> Alcune delle considerazioni presentate potranno essere valide anche per testi di tradizioni e epoche differenti.

Le operazioni costitutive di un programma di collazione automatica sono oggi ben definite. Chi voglia sviluppare un nuovo software può avvalersi di un modello e chi si accinga a studiare tali software può utilizzare lo stesso come griglia di analisi. Il modello in questione è quello definito a Göteborg,<sup>3</sup> nel 2009, dagli sviluppatori di Juxta e di CollateX e dai filologi del COST Action *Open Scholarly Communities on the Web* e del programma *Interedition*: esso è conosciuto come 'Gothenburg Model'.<sup>4</sup> Il modello stabilisce le cinque operazioni fondamentali che un programma di collazione automatica compie, seguendo il principio di progettazione della 'separation of concerns'. Ciò implica che le operazioni siano ripartite in sezioni separate dell'architettura software: ognuna di esse è addetta a un compito, non c'è sovrapposizione di inca-rici e, nel caso si voglia cambiare una parte, uno o più moduli possono essere sostituiti senza dover modificare l'intero ingegno.

Le operazioni fondamentali identificate dal modello di Göteborg sono la tokenizzazione, la normalizzazione, l'allineamento, l'analisi (o *feedback*) e l'*output*. Si passino brevemente in rassegna.

La tokenizzazione è il primo passo da compiere nella maggior parte dei trattamenti computazionali del testo e consiste nella divisione della sequenza di caratteri (lettere, punteggiatura, segni diacritici e spazi bianchi) in unità minime. La dimensione e la natura dell'unità dipendono dal trattamento da eseguire: per la collazione, così come per molte altre attività, l'unità minima o *token* sarà la parola.<sup>5</sup>

Sulla seconda operazione, la normalizzazione, torneremo in seguito. Basti dire che si tratta di quella fase in cui il testo viene preparato per l'allineamento: il grado di normalizzazione può variare, dalla neutralizzazione delle differenze tra maiuscole e minuscole o delle varianti

<sup>3</sup> Utilizziamo la forma svedese del nome, che non ha una traduzione italiana, tranne quando occorre nell'espressione inglese *Gothenburg model*.

<sup>4</sup> Cf. R. Haentjens Dekker, G. Middell, «Computer-Supported Collation with CollateX: Managing Textual Variance in an Environment with Varying Requirements», in *Supporting Digital Humanities*, Copenhagen, University of Copenhagen, 2011; R. Haentjens Dekker et al., «Computer-Supported Collation of Modern Manuscripts: CollateX and the Beckett Digital Manuscript Project», *Digital Scholarship in the Humanities*, 30/3 (2015), pp. 452-470; [http://www.interedition.eu/wiki/index.php/About\\_microservices](http://www.interedition.eu/wiki/index.php/About_microservices) e [http://wiki.tei-c.org/index.php/Textual\\_Variance](http://wiki.tei-c.org/index.php/Textual_Variance) (ultimo accesso 30 gennaio 2019).

<sup>5</sup> La complessità della tokenizzazione varia di lingua in lingua, ma generalmente bisogna considerare che gli spazi bianchi e la punteggiatura non sono sufficienti per dividere le parole (si pensi ad esempio alle abbreviazioni). Inoltre, diverse tokenizzazioni sono possibili sulla stessa sequenza: la punteggiatura, ad esempio, può essere eliminata, costituire un *token* a sé o un *token* insieme alla parola che precede.

ortografiche ad un livellamento più profondo, come ad esempio quello morfologico.

La terza operazione, l'allineamento, è il cuore del programma di collazione, la fase in cui il programma stabilisce quale *token* nel testo A corrisponde a quale *token* nel testo B.<sup>6</sup>

Il *feedback*, opzionale, permette al programma di correggere l'allineamento a partire dai risultati di un'analisi computazionale: ad esempio, il programma valuta la similarità tra i *token* e perfeziona l'allineamento in modo che *token* simili, seppure non uguali, corrispondano.<sup>7</sup>

L'ultima fase è quella dell'*output*, che consiste nella resa dei risultati in visualizzazioni quali la tabella e il grafo o in formati di codifica dei dati come XML o JSON.

L'architettura di CollateX e JuxtaCommons, i due programmi di collazione più adoperati negli ultimi anni, è basata sul modello di Göteborg.<sup>8</sup> Nello studio presentato in quest'articolo, è il primo ad essere usato: sia CollateX che Juxta sono software scritti nel linguaggio di programmazione Java, ma solo per CollateX è stata sviluppata parallelamente una versione in Python,<sup>9</sup> un linguaggio più diffuso nella ricerca in filologia computazionale; inoltre, CollateX gestisce per impostazione predefinita le forme normalizzate dei *token*, caratteristica che, come vedremo, risulta fondamentale per il trattamento dei testi medievali.<sup>10</sup>

### *La collazione automatica di testi medievali*

L'adozione della collazione automatica in filologia medievale è stata fre- nata da problemi concreti. Il primo di questi, forse il più importante, è la necessità di acquisire i testi, per poterli poi trattare e collazionare.<sup>11</sup>

<sup>6</sup> Faremo di seguito riferimento ai *token* come unità minime definite nella prima fase del programma di collazione; come abbiamo detto, essi corrispondono alle parole.

<sup>7</sup> È quello che succede attivando il parametro 'near match' della funzione 'collate' in CollateX, sul quale si tornerà di seguito.

<sup>8</sup> <https://collatex.net/doc/> e [http://juxtacommons.org/tech\\_info](http://juxtacommons.org/tech_info) (ultimo accesso, 30 gennaio 2019).

<sup>9</sup> <https://pypi.python.org/pypi/collatex> (ultimo accesso, 30 gennaio 2019).

<sup>10</sup> La scelta di CollateX è dovuta anche alla possibilità di avere un supporto immediato da parte del suo principale sviluppatore, Ronald Haentjens Dekker.

<sup>11</sup> La questione non è nuova: cf. ad esempio S. Waite, «Two programs for comparing texts», in *La pratique des ordinateurs dans la critique des textes*, Paris, Ed. du CNRS, 1979, pp. 241-244:

244. Il metodo del 'double keying' discusso da Waite è quello ritenuto anche in M. Piotrowski,

*Natural language processing for historical texts*, San Rafael, Morgan and Claypool, 2012, p. 48.

A seconda dell'opera, della sua lunghezza, della quantità di testimoni e del tipo di tradizione, l'acquisizione dei testi può richiedere un impegno diverso. Nonostante la ricerca sul riconoscimento delle scritture mano-scritte avvanzi rapidamente, chi si accinga oggi a lavorare su un'opera tradita da manoscritti medievali deve trascrivere, almeno in parte, manualmente.<sup>12</sup> Le letterature romanze del medioevo – proponendo una larga generalizzazione –, sono ricche in opere, largamente tradite e a tradizione attiva.<sup>13</sup> Trascrivere integralmente dieci o più manoscritti di un testo, magari lungo, a mano è un'operazione non solo onerosa, ma anche comunemente reputata inutile; mentre è in genere accettata l'idea di dedicare una fase del lavoro, anche non breve, alla collazione. Ci si potrebbe allora chiedere quale delle due attività sia impegnativa in termini di tempo. La risposta, anche in questo caso, dipende dalle variabili elencate prima, ovvero dalla quantità di testimoni e dalla natura della tradizione; la riflessione di Tara Andrews, qui riportata, può ad esempio non essere valida per un romanzo arturiano in prosa del Duecento: a causa della varianza diffusa e di diverso tipo, una nuova trascrizione non potrà essere approntata copiando e cambiando la trascrizione di un testimone simile.

[W]itness transcriptions need not take any more time than manual collation of texts, and can usually be made to take considerably less. A full text transcription of any witness may be made simply by copying and altering the transcription of a similar witness, a process that is akin to the creation of a spreadsheet of variants but simpler and easier in execution. There is a small risk that the readings of the similar text might influence the readings of the manuscript being transcribed, but that is more than offset by the fact that a full transcription removes the scholar's temptation to exclude a peculiar reading because it seems not to be worth the effort to set up a new variant location in a spreadsheet.<sup>14</sup>

<sup>12</sup> Se non altro per creare un 'training corpus' a partire dal quale la macchina possa poi imparare. Per il programma Transkribus, ad esempio, sono necessarie un centinaio di pagine di trascrizione per poter allenare il software e creare un modello, che si possa poi utilizzare sulle restanti pagine, cf. <https://transkribus.eu/Transkribus/> (ultimo accesso 30 gennaio 2019).

<sup>13</sup> La prima variabile non è rilevante per il lavoro sulla singola opera, ma assume importanza in una riflessione più ampia su un'agenda di ricerca: lo sforzo che si può dedicare ad una singola opera è inversamente proporzionale alla quantità di opere di una tradizione letteraria, facendo astrazione dell'importanza culturale relativa di ognuna di esse.

<sup>14</sup> T. Andrews, «The Third Way: Philology and Critical Edition in the Digital Age», *Variants*, 10 (2013), pp. 61-76: 67.

Il vantaggio, dunque, non è da cercare nel tempo risparmiato grazie all'utilizzo di metodi computazionali. Bisognerà piuttosto guardare nella direzione che Wilhelm Ott, autore del più longevo programma di edd- tica digitale, TuStep, proponeva nel 1973:

To sum up: by means of this new tools, which we have in electronic data processing, new and higher standards are imposed not only on the results of others sciences, but also on critical editions ... The question whether it is possible or not to save time and / or money by these methods is only of secondary importance. The expenses necessary for future critical editions may possibly be even higher than they have been in the past when these tools were not yet available.<sup>15</sup>

Parlare di più alti standard per l'edd- tica potrebbe equivalere a mettere in discussione i metodi non computazionali, e questo sarebbe un errore. Piuttosto, si può ritenere del proposito di Ott l'idea di nuovi standard che si inseriscono nel panorama scientifico. Per quanto riguarda la collazione, la novità e il vantaggio si troverebbero nel tipo di dati che vengono generati durante la *recensio*: non griglie di collazione ad uso privato, ma trascrizioni che permettono la riproducibilità del trattamento, divenuto esplicito, e che sono riutilizzabili. Le trascrizioni dei manoscritti possono ad esempio essere riprese per la creazione di un corpus di analisi linguistica; il programma di collazione può essere eseguito su quelle stesse trascrizioni a più riprese, tanto da colui che edita il testo quanto dal lettore.<sup>16</sup> Ci torneremo nelle conclusioni.

L'altro limite all'adozione di strumenti computazionali in filologia medievale è imposto dal tipo di variazione: la quantità di varianti che interessano aspetti formali della lingua, grafici e fonetici, è generalmente maggiore del numero di varianti di sostanza, che implicano un cambiamento di significato o di stile rilevante. Vero è che quest'ultimo criterio, quello della rilevanza, non si può in alcun modo generalizzare: nell'esercizio della filologia stemmatica, la variante rilevante è monogenetica, le altre poligenetiche; ma, come ricordano Leonardi e Morato:

[...] la distinction entre substance et forme et leur rapport avec la monogenèse ou la polygenèse ne peuvent ... être établis a priori, mais requièrent une évaluation

<sup>15</sup> W. Ott, «Computer Applications in Textual Criticism», in *The Computer and Literary Studies*, a cura di A.J. Aitken et al., Edinburgh, Edinburgh University Press, 1973, pp. 199-223: 222.

<sup>16</sup> Per far sì che la riproducibilità sia reale, bisognerà mettere a disposizione tutti i pre- e post-trattamenti applicati ai dati, insieme ai testi e ai programmi.

empirique des mécanismes de transmission actifs pour chaque genre littéraire et chaque domaine linguistique et chronologique, si ce n'est pour chaque texte et chaque manuscrit.<sup>17</sup>

In ogni caso, la collazione automatica, effettuata sulle trascrizioni dei testi, produrrà una selva di varianti minori, ovvero una massa di informazioni di poco o scarso rilievo ai fini della *recensio*. Lo constatano i filologi che vi hanno fatto ricorso:

[t]he main problem in the automatic comparison of different manuscripts of medieval texts is not so much to identify textual variants, but to distinguish between important and unimportant variants.<sup>18</sup>

La questione non si pone solo per la collazione, ma per buona parte dei trattamenti computazionali applicati a testi medievali. In un articolo del 2012, Karina van Dalen-Oskam applica metodi stilometrici di attribuzione autoriale ad un corpus di manoscritti della *Scolastica* di Jacob van Maerlant per studiare l'attitudine dei diversi copisti verso il testo; il procedimento è però ostacolato da differenze insignificanti:

the measurements on the Middle Dutch text are indeed highly influenced by irrelevant spelling differences. If we want to abstract from spelling, we need to do the measurements on normalised text, or on lemmatised text.<sup>19</sup>

Le soluzioni proposte da van Dalen-Oskam, la normalizzazione o la lemmatizzazione, sono quelle più adottate per ridurre il disordine di lingue non standardizzate ad un ordine che ne permetta la gestione tramite strumenti computazionali, generalmente pensati per lingue moderne con

<sup>17</sup> L. Leonardi, N. Morato, «L'édition du cycle de Guiron le Courtois. Établissement du texte et surface linguistique», in *Le Cycle de Guiron le Courtois. Prolégomènes à l'édition intégrale du corpus*, Paris, Garnier, 2018, pp. 453-509: 472. Il ricco articolo di Leonardi e Morato insiste sulla necessità di una linea di divisione mobile tra forma e sostanza, tra monogenesi e poligenesi, linea che va spostata a seconda del corpus in esame: non si può ridurre d'ufficio la poligenesi ai fenomeni grafico-fonetici, al contrario «l'objectif est d'identifier les phénomènes d'ordre morpho-syntaxique, lexical, discursif, qui se comportent comme des phénomènes polygénétiques, ou qui sont à tout le moins tendanciellement ou potentiellement polygénétiques» *ivi*, p. 475.

<sup>18</sup> L. Zeevaert, «Easy Tools to Get to Grips with Linguistic Variation in the Manuscripts of *Njáls Saga*», *Digital Medievalist*, 10/1 (2015), p. 50.

<sup>19</sup> K. van Dalen-Oskam, «The Secret Life of Scribes. Exploring Fifteen Manuscripts of Jacob van Maerlant's *Scolastica* (1271)», *Literary and Linguistic Computing*, 27/4 (2012), pp. 355-72: p. 359.

forme canoniche. Michael Piotrowski identifica le seguenti caratteristiche delle lingue pre-moderne, che ne impediscono il trattamento automatico: variazione diacronica, variazione sincronica e incertezza della trascrizione, in quanto risultato di un'attività soggetta ad errori. Per ovviare almeno ai primi due problemi, in linguistica computazionale si tende a normalizzare la lingua:

in NLP [Natural Language Processing]<sup>20</sup> canonicalization and normalization usually imply a mapping to a modern form. Spelling modernization would possibly be a more precise term for this approach.<sup>21</sup>

Una volta stabilite le corrispondenze tra la lingua storica e quella moderna si possono usare metodi, strumenti e risorse disponibili per la seconda e applicarli alla prima. Per identificare le corrispondenze, Piotrowski menziona la canonizzazione assoluta, che prevede l'utilizzo di una lista di forme corrispondenti, rigorosa ma onerosa da implementare. Altri approcci discussi, detti di canonizzazione relativa, mirano alla creazione di forme canoniche diverse da quelle della lingua moderna e includono metodi che tengono conto della distanza di *edit*<sup>22</sup> e della similitudine fonetica tra le varianti.<sup>23</sup>

Nella storia della collazione automatica diversi modi di gestire la varianza di 'spelling' (inteso come separazione delle parole, grafia, fonetica, abbreviazioni e punteggiatura) si sono avvicinati. La scelta non è stata generalmente quella della modernizzazione, ma di altri tipi di normalizzazione, equivalenti agli approcci di canonizzazione relativa discussi da Piotrowski. Tutti vengono effettuati prima dell'allineamento vero e proprio e corrispondono dunque alla seconda fase (in alcuni casi alla prima e alla seconda) del modello di Göteborg, quella della normalizzazione. Tutti sono procedimenti di preparazione del testo, ai fini di ottenere risultati ottimali nella fase successiva: non solo riducendo

<sup>20</sup> L'inglese si può rendere in italiano con 'Trattamento automatico del linguaggio' o 'Elaborazione del linguaggio naturale'.

<sup>21</sup> Piotrowski, *Natural language processing for historical texts*, p. 69 (corsivo nell'originale).

<sup>22</sup> La 'edit distance' tra due *token* indica la quantità di operazioni (aggiunte, cancellazioni e sostituzioni) sulle lettere necessarie per passare da una parola all'altra; generalmente, una soglia di 'edit distance' è definita, sotto alla quale i due *token* vengono allineati, dunque considerati corrispondenti. Ad esempio, la 'edit distance' tra 'legemmo' e 'leggemmo' (aggiunta di una lettera) è 1, così come quella tra 'scrisse' e 'scripse' (sostituzione di una lettera).

<sup>23</sup> Piotrowski, *Natural language processing for historical texts*, Capitolo 6.3.2.

il numero di varianti poco significative, ma anche limitando gli errori di allineamento che le varianti formali possono indurre. È importante notare che la normalizzazione e la preparazione del testo per l'allineamento non hanno conseguenze sulla trascrizione, ovvero sui dati raccolti dal filologo: la trascrizione rimane intatta e se ne crea una copia di lavoro da modificare per l'allestimento della collazione. Tramite meccanismi di codifica, le due copie possono convivere nello stesso file, come due facce della stessa medaglia: per la lettura della trascrizione se ne utilizzerà una e per l'allineamento l'altra. Il programma di collazione dovrebbe dare la possibilità di passare dall'una all'altra ad ogni momento; ciò implica, ad esempio, che i risultati dell'allineamento eseguito sulla copia normalizzata possano essere visualizzati secondo la lettera della versione originale.

Prima di passare in rassegna meccanismi specifici di normalizzazione, si trattino brevemente quelli più comuni, che sono rimasti quasi intatti negli ultimi cinquant'anni, seppur per motivi diversi. Nel 1966, Dom Jacques Froger scriveva:

Avec la petite machine dont nous nous sommes servis, les moyens d'exprimer un texte se réduisaient à ceci :

— des lettres exclusivement capitales, et par conséquent aussi les chiffres romains, puisqu'ils sont formés de lettres ;

— les chiffres arabes ;

— le point pour toute ponctuation ;

— l'espace pour séparer les mots, les chiffres et de façon générale des éléments quelconques du texte.

Dans des limites aussi étroites, on ne peut mettre à profit la différence entre majuscules et minuscules, entre romain et italique, entre lettres grasses ou maigres etc. On ne peut pas non plus faire usage de lettres qui portent des accents, puisqu'on ne dispose que de capitales toutes nues.<sup>24</sup>

Le limitazioni dell'*hardware* negli anni '60 imponevano l'utilizzo di convenzioni, come il segno \$ per indicare una maiuscola o una virgola separata da spazi per il ritorno a capo. Froger decide di conservare la differenza tra maiuscole e minuscole, ma considera irrilevanti gli accenti e la punteggiatura.<sup>25</sup> Nel tempo sono venuti meno i limiti tecnici per

<sup>24</sup> Dom J. Froger, «La collation des manuscrits à la machine électronique», *Bulletin d'information de l'Institut de Recherche et d'Histoire des Textes*, 13 (1964-1966), pp. 135- 171: 144.

<sup>25</sup> Ivi, p. 145.

la rappresentazione di accenti, punteggiatura e maiuscole,<sup>26</sup> ma i programmi di collazione hanno continuato a dare la possibilità di neutralizzare queste componenti del testo: a seconda della lingua, del corpus e del fine della collazione, gli accenti, la punteggiatura e la distinzione tra maiuscole e minuscole possono essere conservati o rimossi; nei differenti casi, la migliore strategia può essere individuata empiricamente tramite prove sullo stesso corpus con parametri diversi. Gli esempi rimontano agli anni '70: «[t]he editor may specify in the STRIP parameter which punctuation marks and special symbols the module is to ignore», scrive Penny Gilbert nel 1979, presentando il sistema Collate per le edizioni di testi in prosa (da non confondere con quello elaborato da Peter Robinson dieci anni dopo).<sup>27</sup> Questi meccanismi di preparazione possono valere per tutti i testi, e non solo per quelli in lingue ad alto tasso di variabilità come quelle medievali.

Altri tipi di normalizzazione, che rispondono alle caratteristiche degli specifici *corpora* in questione, sono stati utilizzati nel tempo. Tra i primi casi di normalizzazione *ad hoc* ricordiamo il lavoro di Georgette Silva e Harold Love, che utilizzano come caso di studio il *Prelude* di Wordsworth, in cui prima della comparazione tra due versi «[i]n order to avoid the most common types of variation, such as liv'd/lived, or ;/, or capitalized/non-capitalized words, the subroutine first eliminates letters "e" and "s", the \$-sign, blanks, and all punctuation marks from both lines».<sup>28</sup> Un approccio simile è quello di Waite: «if one of the texts systematically represents consonantal I and U by J and V while the other uses J and V for the same purpose, this change can readily be allowed for in the programs, so that only significant differences are marked».<sup>29</sup> In entrambi i casi, il filologo identifica i tratti formali non significativi e li elimina dal testo, oppure istruisce il programma per ignorarli, in modo da ridurre il rumore nei risultati.

Il programma che introduce meccanismi più sofisticati per gestire la varianza formale è Collate, sviluppato da Peter Robinson insieme a diversi

<sup>26</sup> Per la codifica dei caratteri, valgono le specificazioni di Unicode, cf. <http://www.unicode.org/> (ultimo accesso 30 gennaio 2019).

<sup>27</sup> P. Gilbert, «The preparation of prose-text editions with the "Collate" system», in *La pratique des ordinateurs dans la critique des textes*, Paris, Ed. du CNRS, 1979, pp. 245-254: 247.

<sup>28</sup> G. Silva, H. Love, «The identification of text variants by computer», *Information Storage and Retrieval*, 5/3, (1969), pp. 89-108: 93. Il segno \$ viene usato dopo la lettera per indicare che si tratta di una maiuscola; cf. ad esempio v. 90: «T\$O GRATIFY THE PALATE WITH REPASTS» ivi, p. 94.

<sup>29</sup> Waite, «Two programs for comparing texts», p. 242.

collaboratori,<sup>30</sup> inizialmente concepito per l'edizione di un testo islandese medievale.<sup>31</sup> Il programma prevede tre modi per gestire la varianza formale: un meccanismo di normalizzazione che interviene prima dell'allineamento e due procedimenti che avvengono al momento dell'allineamento o dopo. Questi ultimi sono il cosiddetto 'fuzzy match', che consente di allineare *token* simili ma non uguali;<sup>32</sup> e la funzione 'defvars', che fa ricorso ad un file esterno in cui sono registrati *token* graficamente dissimili ma equivalenti, come, ad esempio, gli avverbi di negazione 'mie' e 'pas' in antico francese.<sup>33</sup> La normalizzazione che avviene prima dell'allineamento, invece, prevede l'utilizzo di un dizionario di varianti da scartare, approccio utilizzato ancora oggi. Il dizionario è costituito da una serie di righe: per ognuna di esse, la prima parola è la forma normalizzata, seguita da tutte le forme corrispondenti nei manoscritti. Per preparare il testo alla collazione, si crea una copia normalizzata dei testi, controllando ogni parola nel dizionario e, se la parola è presente, sostituendo la forma originale con quella normalizzata. Il testo prodotto risulterà un'aberrazione filologica: ma si tratta, come detto, di una copia di lavoro allestita per la collazione; le corrispondenze tra forme originali e normalizzate sono presenti ad ogni momento e possono facilmente essere recuperate in ogni momento del processo. Un'altra questione apparen-

<sup>30</sup> P. Robinson, «The Collation and Textual Criticism of Icelandic Manuscripts (1): Collation», *Literary and Linguistic Computing*, 4/2 (1989), pp. 99-105. Annunciato nella mailing-list ('discussion group' o 'international seminar') *Humanist* nel 1991, lo sviluppo di Collate ha dato vita a tre distribuzioni principali del software; le funzionalità e le caratteristiche qui discusse sono dunque frutto di un processo durato vent'anni. Cf. P. Robinson, «Rationale and Implementation of the Collation System», in *The Miller's Tale on CD-ROM. The Canterbury Tales Project*, Saskatchewan, Scholarly Digital Editions, 2004; <https://dhhumanist.org/Archives/Virginia/v04/1240.html>; <https://scholarlydigitaleditions.blogspot.com/2014/09/the-history-of-collate.html>; <https://scholarlydigitaleditions.blogspot.com/2014/09/collate-2-and-design-for-its-successor.html> (ultimo accesso 30 gennaio 2019).

<sup>31</sup> Le varianti grafiche nel testo studiato sono davvero molte: «[t]he anarchy is such that I estimate that there are over fourteen million possible spellings of "FjolsviSr", the name of one of the protagonists of Fj [Fjölsvinnsmál]. Thankfully, we get only 97 different spellings, in the 243 occurrences of the name in the 37 manuscripts in which it occurs» Robinson, «The Collation and Textual Criticism of Icelandic Manuscripts», p. 100.

<sup>32</sup> Il parametro fuzzy match di Collate corrisponde concettualmente al 'near match' di CollateX. Il secondo è basato sulla 'edit distance' tra i due *token* (cf. nota 20). Per il fuzzy match, il meccanismo è diverso; cf. Robinson, «The Collation and Textual Criticism of Icelandic Manuscripts», p. 103: «For FUZZYMATCH I devised a mathematical formula in which the variables are the number of letters in the two words, the letters themselves and their order».

<sup>33</sup> L'esempio proposto da Robinson è quello degli avverbi di congiunzione 'ok' e 'en' in islandese.

temente spinosa è quella delle forme normalizzate: per una lingua non standardizzata, non è facile, né corretto, scegliere una forma canonica. In effetti, come ricorda Robinson, «the normalisation adopted is quite unimportant, as it is only to be used for comparing readings between manuscripts. All that matters is that it be consistent».<sup>34</sup> L'importante non sarà dunque stabilire quale forma ritenere come canonica, ad esempio 'scripse' o 'scrisse'; ma assicurarsi che entrambe rimandino ad una forma canonica arbitraria ('scripse' o 'scrisse' o il disegno di una penna) e che possano essere quindi riconosciute dal programma come equivalenti.

Questo tipo di normalizzazione tramite dizionario, estremamente preciso ma anche oneroso da implementare, è quello che oggi viene considerato standard. E che, insieme al costo della trascrizione, frena l'adozione della collazione automatica, come ribadito in un articolo recente: «l'attuabilità del piano di lavoro ... è contrastata dalle difficoltà insite nella codifica di un numero così ampio di testimoni, il cui testo deve essere standardizzato sufficientemente da essere collazionato da una macchina virtuale».<sup>35</sup>

### *Nuove risorse e nuove sfide*

In questa sezione vedremo come le difficoltà discusse sopra possono essere almeno in parte sormontate grazie all'utilizzo di nuove risorse, occupandoci dello stato dell'arte recente per quanto riguarda la normalizzazione e la gestione delle varianti nella collazione di testi medievali.

La sfida da affrontare con nuove risorse è quella di mettere insieme le due componenti problematiche della collazione di testi medievali, ovvero la normalizzazione dei testi per ottenere un allineamento corretto e la grande quantità di varianti formali, e di utilizzare la prima per controllare la seconda: distinguere, quindi, categorie di varianti attraverso l'informazione usata per la normalizzazione. Il procedimento sembra tautologico – le varianti formali sono quelle definite durante la fase di normalizzazione come varianti formali –, ma diventa vantaggioso se le varianti sono identificate tramite un'annotazione linguistica automatica, come proposto di seguito.<sup>36</sup>

<sup>34</sup> Robinson, «The Collation and Textual Criticism of Icelandic Manuscripts (1): Collation», p. 100.

<sup>35</sup> S. Bertone, et al., «Verso il trattamento automatico della tradizione manoscritta del liber di Catullo», *Umanistica Digitale*, 2/3 (2018).

<sup>36</sup> Procedimenti simili sono studiati da Jean-Baptiste Camps e Lucence Ing, che ringrazio per aver messo a disposizione il loro lavoro in corso. Cf. J.-B. Camps, L. Ing,

Per proporre un esempio, dovremo innanzitutto specificare cosa intendiamo per varianti formali: includeremo qui solamente le varianti grafiche ('scripse' vs 'scrisse') e fonetiche ('bon' vs 'buon'). Prima di passare all'utilizzo dell'annotazione linguistica automatica, presentiamo il procedimento, attraverso un esempio in cui la normalizzazione è effettuata manualmente.

Consideriamo il verso «soli eravamo e senza alcun sospetto» (Dante, *Commedia*, ed. Petrocchi, Inf. V, 129) nei tre manoscritti di mano di Giovanni Boccaccio, che conservano «sança alcun sospecto» (Toledano) e «sença alcun sospecto» (Riccardiano e Chigiano).<sup>37</sup> Per collazionare i testi, essi andranno normalizzati; definiremo quindi un dizionario di arbitrarie forme canoniche, in cui alle forme 'sança' e 'sença' corrisponde la forma canonica 'SENZA': le copie di lavoro per la collazione registreranno dunque «SENZA alcun sospecto» e «SENZA alcun sospecto». Il programma potrà allinearli correttamente e non riconoscerà nessuna differenza. In questo modo abbiamo ottenuto l'allineamento ottimale ed escluso la varianza grafica. Possiamo ora recuperare l'informazione usata per la normalizzazione e creare delle semplici regole: quando due o più *token* allineati hanno la stessa forma canonica ma una forma originale diversa, i *token* sono varianti formali; quando hanno diverse forme canoniche e diverse forme originali, i *token* sono varianti di sostanza; quando hanno la stessa forma canonica e originale, i *token* sono invariati. I *token* 'sança' e 'sença' saranno quindi riconosciuti come varianti formali, perché la loro forma canonica 'SENZA' coincide nei testimoni ma le forme originali differiscono. L'applicazione delle regole appena descritte avviene dopo l'allineamento, durante la quarta fase, quella del *feedback*, ed è un esempio di analisi che serve ad arricchire l'*output*.

*Collation assistée par ordinateur de témoins de textes en ancien français*. 19 ottobre 2018,

<<https://halshs.archives-ouvertes.fr/halshs-02023936/document>>; L. Ing, *Outils numériques pour corpus diachronique: quelques disparitions lexicales en œuvre dans le micro-corpus du Lancelot en prose*, mémoire de master 2 «Humanités numériques et computationnelles», dir. J.-B. Camps e F. Duval, Université Paris Sciences & Lettres, 2018.

<sup>37</sup> Per uno studio dei manoscritti boccacciani della *Commedia*, cf. S. Tempestini, «Boccaccio copista e interprete della "Commedia"». La "Commedia" nei codici Toledano 104.6, Riccardiano 10135, Chigiano L VI 213: alcuni dati sulla variantistica», in *Intorno a Boccaccio / Boccaccio e dintorni. Atti del Seminario internazionale di studi (Certoaldo Alta, Casa di Giovanni Boccaccio, 9 settembre 2015)*, a cura di S. Zamponi, Firenze, University Press, 2015, pp. 89-107; S. Tempestini, «Boccaccio copista e editore della "Commedia"». Per un'analisi della variantistica», *Critica del testo*, 21/2 (2019). Il corpus è parzialmente interrogabile all'indirizzo <http://boccacciocommedia.unil.ch/> (ultimo accesso 30 gennaio 2019); cf. S. Tempestini, E. Spadini, «Querying Variants: Boccaccio's "Commedia" and Data-Models», *Digital Medievalist*, XII (2019).

L'esempio è servito a mostrare il procedimento, ma se la normalizzazione deve essere effettuata manualmente l'intero procedimento perde di interesse perché estremamente oneroso; di fatto, poi, le varianti formali sarebbero identificate in fase di normalizzazione. Il vantaggio risiede invece nell'utilizzare strumenti computazionali per la normalizzazione, in modo da automatizzare, almeno in parte, questo passaggio. Come detto, la normalizzazione serve a stabilire una forma canonica arbitraria: essa può essere costituita dall'equivalente moderno, da un'immagine, ma anche dalle proprietà linguistiche della parola; l'utilizzo del lemma e della categoria grammaticale come forma di normalizzazione è stato già applicato con successo.<sup>38</sup> Inoltre, si possono oggi utilizzare risorse che sono andate sviluppandosi negli ultimi decenni per l'analisi linguistica delle lingue medievali. Si tratta in particolare di modelli da usare per la lemmatizzazione e l'annotazione morfologica, disponibili ad oggi per varie lingue.<sup>39</sup> Per il francese antico, su cui ci concentreremo qui, due modelli per il software TreeTagger<sup>40</sup> sono a disposizione, basati sulle risorse del *Nouveau Corpus d'Amsterdam* e della *Base de français médiéval*, riuniti nel *Medieval French Language Toolkit*.<sup>41</sup>

Il procedimento rimane lo stesso rispetto a quello già descritto, salvo per l'utilizzo di strumenti automatici per l'annotazione linguistica durante la fase di preparazione del testo. Vediamolo in un altro esempio.

Immaginiamo tre manoscritti, che inizino in questo modo: A 'Grant fu la ioie'; B 'Or dist li contes que grans fu la joie', C 'Grans fu la joie'.

<sup>38</sup> Cf. van Dalen-Oskam, «The Secret Life of Scribes».

<sup>39</sup> Sono discusse risorse per arabo, cinese, olandese, inglese, francese, tedesco, latino, greco antico, portoghese e lingue nordiche in Piotrowski, *Natural language processing for historical texts*, Capitolo 8. Per lo spagnolo, cf. G. Boleda, «Extending the tool, or how to annotate historical language varieties», in *Proceedings of the 5th ACL-HLT Workshop on Language Technology for Cultural Heritage, Social Sciences, and Humanities*, 2011, pp. 1-9; J. Porta, et al., «Edit transducers for spelling variation in Old Spanish», in *Proceedings of the workshop on computational historical linguistics at NODALIDA 2013; May 22-24; 2013; Oslo; Norway. NEALT Proceedings Series 18*, Linköping University Electronic Press, 2013, pp. 70-79.

<sup>40</sup> H. Schmid, «Probabilistic Part-of-Speech Tagging Using Decision Trees», in *Proceedings of International Conference on New Methods in Language Processing*, 1994. Il programma è disponibile online all'indirizzo <http://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/>. Il programma ha un corrispondente Python wrapper, che ne può facilitare l'integrazione con CollateX: <https://pypi.org/project/treetaggerwrapper/> (ultimo accesso, 30 gennaio 2019).

<sup>41</sup> <https://github.com/sheiden/Medieval-French-Language-Toolkit>. Il modello della *Base de français médiéval* è disponibile all'indirizzo <http://bfm.ens-lyon.fr/spip.php?article324> (ultimo accesso 30 gennaio 2019).

Senza nessuna forma di normalizzazione, CollateX produce il risultato riprodotto nella TAB. 1.<sup>42</sup>

TABELLA 1

Risultati della collazione senza normalizzazione ('var' = varianti; 'inv' = invarianti).

A	B	C	D
Grant	Or	Grans	<i>var</i>
-	dist	-	<i>var</i>
-	li	-	<i>var</i>
-	contes	-	<i>var</i>
-	que	-	<i>var</i>
-	grans	-	<i>var</i>
fu	fu	fu	<i>inv</i>
la	la	la	<i>inv</i>
ioie	joie	joie	<i>var</i>

Se invece neutralizziamo la differenza tra maiuscole e minuscole, si risolve la posizione di 'Grans' in C, ma l'allineamento rimane poco soddisfacente per A, come mostrato nella TAB. 2, perché il programma non identifica la vicinanza tra 'grant' e 'grans'.

TABELLA 2

Risultati della collazione dopo la neutralizzazione di maiuscole e minuscole ('var' = varianti; 'inv' = invarianti).

A	B	C	D
Grant	Or	-	<i>var</i>
-	dist	-	<i>var</i>
-	li	-	<i>var</i>

<sup>42</sup> I risultati riprodotti in questa e nelle successive tabelle sono ottenuti con CollateX 2.1.3 (Python). Si noti che la demo di CollateX disponibile all'indirizzo <https://collatex>.

-	contes	-	<i>var</i>
-	que	-	<i>var</i>
-	grans	Grans	<i>var</i>
fu	fu	fu	<i>inv</i>
la	la	la	<i>inv</i>
ioie	joie	joie	<i>var</i>

L'analisi linguistica automatica, invece, riconoscerà che 'Grant', 'grans' e 'Grans' sono aggettivi e che il lemma corrispondente è 'grant'; così come che 'ioie' e 'joie' sono sostantivi, il cui lemma corrispondente è 'joie'. La collazione si effettuerà quindi sui *token* normalizzati, ovvero sull'anno- tazione linguistica, come mostrato nella TAB. 3. I risultati saranno resi utilizzando invece le forme originali: si ottiene quindi la TAB. 4.

TABELLA 3

Risultati della collazione dopo la normalizzazione ('var' = varianti; 'inv' = invariati): token normalizzati.

A	B	C	D
-	avverbio, <i>or</i>	-	<i>var</i>
-	verbo, <i>dire</i>	-	<i>var</i>
-	articolo, <i>il</i>	-	<i>var</i>
-	sostantivo, <i>conte</i>	-	<i>var</i>
-	congiunzione, <i>que</i>	-	<i>var</i>
aggettivo, <i>grant</i>	aggettivo, <i>grant</i>	aggettivo, <i>grant</i>	<i>inv</i>
verbo, <i>estre</i>	verbo, <i>estre</i>	verbo, <i>estre</i>	<i>inv</i>
articolo, <i>il</i>	articolo, <i>il</i>	articolo, <i>il</i>	<i>inv</i>
sostantivo, <i>joie</i>	sostantivo, <i>joie</i>	sostantivo, <i>joie</i>	<i>inv</i>

net/demo/ (ultimo accesso 30 gennaio 2019) utilizza invece la versione Java, che prevede la neutralizzazione delle maiuscole come normalizzazione di default.

TABELLA 4

Risultati della collazione dopo la normalizzazione ('var' = varianti; 'inv' = invariati): token originali.

A	B	C	D
-	Or	-	<i>var</i>
-	dist	-	<i>var</i>
-	li	-	<i>var</i>
-	contes	-	<i>var</i>
-	que	-	<i>var</i>
Grant	grans	Grans	<i>inv</i>
fu	fu	fu	<i>inv</i>
la	la	la	<i>inv</i>
ioie	joie	joie	<i>inv</i>

Si noterà che si è passati da due a quattro righe invariati. Non solo, quindi, 'Grant', 'grans' e 'Grans' sono stati correttamente allineati, ma, come 'ioie' e 'joie', non sono più considerati varianti. Possiamo a questo punto applicare le regole enunciate sopra e ottenere il risultato rappresentato nella TAB. 5.<sup>43</sup>

TABELLA 5

Risultati della collazione dopo la normalizzazione e la distinzione delle varianti ('var sost' = varianti di sostanza; 'var form' = varianti di forma; 'inv' = invariati).

A	B	C	D
-	Or	-	<i>var sost</i>
-	dist	-	<i>var sost</i>
-	li	-	<i>var sost</i>

<sup>43</sup> Una visualizzazione a mezzo di colori e altri espedienti grafici permetterebbe di cogliere le categorie espresse nell'ultima colonna a destra a colpo d'occhio. Export dei risultati in altri formati, come JSON e XML/TEI, sono configurabili.

-	contes	-	<i>var sost</i>
-	que	-	<i>var sost</i>
Grant	grans	Grans	<i>var form</i> <sup>44</sup>
fu	fu	fu	<i>inv</i>
la	la	la	<i>inv</i>
ioie	joie	joie	<i>var form</i>

Un altro breve esempio, tratto dal *Lancelot en prose*, prevede tre manoscritti che leggono rispettivamente: «A ge te conois mielz que tu ne conois moi»; B «Artu je te conois miels que tu ne fes moi», C «Artus je te conois mult miaus que tu ne fas moi».

La collazione senza nessun tipo di normalizzazione produce i risul-  
tati, per molti versi insoddisfacenti, presentati nella TAB. 6.

TABELLA 6

Risultati della collazione senza normalizzazione (*var* = varianti; *inv* = invarianti).

A	B	C	D
ge	Artu	Artus	<i>var</i>
-	je	je	<i>var</i>
te	te	te	<i>inv</i>
conois	conois	conois	<i>inv</i>
mielz	miels	mult	<i>var</i>
-	-	miaus	<i>var</i>
que	que	que	<i>inv</i>
tu	tu	tu	<i>inv</i>
ne	ne	ne	<i>inv</i>

<sup>44</sup> L'aggettivo *grant* al femminile, seppur caso soggetto, può non mutare la terminazione, dunque possiamo considerare la differenza di tipo fonetico e non morfologico.

conois	fes	fas	<i>var</i>
moi	moi	moi	<i>inv</i>

Reiterando il procedimento descritto sopra, ad ogni *token* viene assegnato un lemma e una parte del discorso, che faranno ufficio di forma canonica per l'allineamento, producendo infine i risultati visualizzati nella TAB. 7.

TABELLA 7

Risultati della collazione dopo la normalizzazione e la distinzione delle varianti (*var sost* = varianti di sostanza; *var form* = varianti di forma; *inv* = invarianti).

A	B	C	D
	Artu	Artus	<i>var sost</i>
ge	je	je	<i>var form</i>
te	te	te	<i>inv</i>
conois	conois	conois	<i>inv</i>
-	-	mult	<i>var sost</i>
mielz	miels	miaus	<i>var form</i>
que	que	que	<i>inv</i>
tu	tu	tu	<i>inv</i>
ne	ne	ne	<i>inv</i>
conois	fes	fas	<i>var sost</i>
moi	moi	moi	<i>inv</i>

L'approccio descritto tramite questi esempi permette di risolvere alcuni degli aspetti problematici della collazione di testi medievali, grazie alla sua capacità di limitare gli effetti dell'instabilità linguistica sui risultati dell'allineamento e alla possibilità di riutilizzare l'informazione per gestire la grande quantità di varianti tramite categorie.

La facilità di implementare questo tipo di soluzione dipende in parte dall'abilità del programma di gestire le informazioni per ogni *token* e, in particolare, di tenere in memoria la forma originale insieme alla forma normalizzata. CollateX offre un'architettura ideale, dato che per ogni *token* sono previste due proprietà per dare spazio alle forme canoniche e originali. L'efficienza del procedimento riposa invece soprattutto sull'accuratezza dell'analisi linguistica: se i valori per il software TreeTagger possono essere superiori al 95%,<sup>45</sup> essi possono diminuire nel caso delle lingue medievali proprio a causa dell'estesa varianza ortografica. Ulteriori ricerche e l'applicazione a corpora diversi serviranno per misurare l'efficienza delle proposte qui delineate.<sup>46</sup>

### Conclusioni

In questo articolo si fornisce un panorama critico di approcci alla collazione automatica, menzionando implementazioni del passato e soluzioni recenti. Per quanto riguarda i testi medievali, si è visto come l'instabilità linguistica – principalmente grafica e fonetica, ma non solo – mette a repentaglio la possibilità di usare strumenti automatici per la collazione. Gestire la diffusa varianza formale, declinata a seconda degli interessi di ricerca, è importante per ottenere un allineamento corretto e limitare il 'rumore' nei risultati.

Tra i meccanismi passati in rassegna, figura il 'near' o 'fuzzy-match'. La necessità di trovare il giusto parametro, che potrebbe variare caso per caso anche all'interno dello stesso testo, rende però questo tipo di approccio difficilmente applicabile su larga scala.

L'altra possibilità esplorata, e che acquista interesse alla luce degli sviluppi nel trattamento automatico di lingue storiche, è quella della normalizzazione. Essa può essere di diversa natura e permette dunque una grande flessibilità: ogni elemento del linguaggio può essere oggetto di normalizzazione, tanto i tratti fonetici quanto gli aspetti semantici, ad

<sup>45</sup> Cf. H. Schmid, «Improvements in Part-of-Speech Tagging with an Application to German», in *Proceedings of the ACL SIGDAT-Workshop*, Dublin, 1995.

<sup>46</sup> Insieme a Jean-Baptiste Camps e Lucence Ing stiamo sviluppando un'implementazione dell'approccio qui presentato; i primi risultati si trovano in: J.-B. Camps,

L. Ing e E. Spadini, «Collating Medieval Vernacular Texts. Aligning Witnesses, Classifying Variants», in *Digital Humanities Conference 2019 (Utrecht)*,

<https://hal.archives-ouvertes.fr/hal-02268348>. Cf. anche

<https://github.com/CondorCompPhil/falcon>. Lavori preliminari si trovano in

[https://github.com/elsesdn/collation\\_spelling](https://github.com/elsesdn/collation_spelling).

esempio tramite la neutralizzazione dei sinonimi. La normalizzazione può essere effettuata manualmente, come è avvenuto fino ad oggi, o automaticamente, come ormai è possibile. Nel primo caso, uno sforzo considerevole deve essere dedicato alla creazione di un dizionario di corrispondenze. Nel secondo caso, l'utilizzo di strumenti di trattamento della lingua per normalizzare e annotare il testo riduce l'impegno richiesto e permette un ampio spettro di analisi: l'annotazione descritta nell'esempio qui sopra si riduce al lemma e alla categoria grammaticale, che permettono di identificare varianti grafiche e fonetiche, applicando semplici regole; un'annotazione più ricca, che comprendesse anche dati morfologici e semantico-lessicali, permetterebbe di creare nuove categorie di varianti corrispondenti.<sup>47</sup> Bisognerà senz'altro tenere a mente che gli strumenti per l'analisi linguistica della lingue medievali sono ancora in corso di sviluppo: si potrebbe immaginare un circolo virtuoso in cui i testi annotati e manualmente controllati possano servire come 'training corpus' per migliorare le risorse linguistiche, e così via.

L'utilizzo di una serie di software, uno dopo l'altro - in una 'pipeline' -, necessita poi di attenzione dal punto di vista dell'interoperabilità dei formati. L'utilizzo di standard e la semplificazione sono dunque fondamentali: formati come XML, CSV e TXT, a seconda della struttura dei dati, corrispondono ad entrambi i requisiti.

In questo stato dell'arte si sono discussi brevemente i limiti all'adozione della collazione automatica. Esistono però anche dei vantaggi, da considerare in relazione al circolo virtuoso appena evocato e al cambiamento di paradigma menzionato da Ott: la collazione automatica è riproducibile ed esplicita. Quest'ultimo punto si riferisce alla normalizzazione mentale implicita che applica chi collaziona manualmente, registrando solo le varianti significative; la collazione automatica, tramite le varie operazioni descritte nel modello di Göteborg, rende invece esplicito ogni passaggio. Ognuno di essi può essere inoltre verificato, nel caso in cui si volesse effettuare un controllo, ma anche per cambiare approccio in corso di rotta: basterebbe rieseguire la collazione, con parametri diversi. I testi delle trascrizioni, inoltre, possono essere riutilizzati, come detto, per finalità diverse.

Infine, consideriamo le osservazioni di Robinson che, in termini quasi provocatori, sottolinea:

<sup>47</sup> La proposta è abbozzata in Tempestini, Spadini, «Querying Variants: Boccaccio's 'Commedia' and Data-Models». Nella stessa direzione si muovono anche le ricerche di Camps e Ing, che grazie all'utilizzo del software Pandora per l'annotazione linguistica ottengono risultati completi di analisi morfologica, da poter utilizzare per l'identificazione di varianti (cfr. nota 34).

you simply cannot judge what might or might not be significant until you have seen every word of every manuscript. The best guide to what was significant was not theory but the manuscripts themselves.<sup>48</sup>

In un già citato articolo, Leonardi e Morato riflettono sulla distinzione tra i fenomeni di superficie linguistica e i fenomeni testuali e tra le metodologie per trattare gli uni e gli altri nella pratica ecdotica; gli autori fanno riferimento all'edizione Segre del *Bestiaire d'amours*, in particolare all'identificazione di categorie di varianti tramite liste di occorrenze, e aggiungono:

C'est ce type d'analyse qui devrait être élargie et approfondie, de façon à rendre moins aléatoire la distinction dont il est question. Or, la seule voie pour tenter de mieux définir ce que nous attribuons à la surface linguistique est de recourir à l'expérience de la recensio, à savoir aux résultats de la collation et aux indications qui émergent de celle-ci concernant la monogenèse et la distribution des variantes.<sup>49</sup>

Seppur gli autori non si riferiscono alla collazione automatica, procedimenti come quelli qui presentati si iscrivono nella direzione tracciata: lo studio della tradizione dei testi può trarre vantaggio da risultati della *recensio* espliciti, riproducibili e riutilizzabili secondo il paradigma delineato dalla filologia computazionale.

#### ABSTRACT

Collation is one of the first philological activities for which the use of computers has been considered. Already during the 60's, algorithms for semi-automatic collation came to light. This article pursues their development, focusing on a specific type of materials: medieval vernacular texts. The linguistic resources available nowadays for these materials allow to design a collation pipeline, following the architecture proposed by the Gothenburg model: the text of each witness is annotated with linguistic information; the alignment is made on the lemma, in order to neutralize the orthographic variation; eventually, the linguistic annotation is used to identify categories of variants. While most of the steps in this pipeline have been conceived earlier in the history of semi-automatic collation, they were partially carried out manually; the possibility of automatize them might influence the adoption of computers for the collation of medieval vernacular texts.

<sup>48</sup> Robinson, «The Collation and Textual Criticism of Icelandic Manuscripts», p. 101.

<sup>49</sup> Leonardi, Morato, «L'éditio du cycle de Guiron le Courtois», p. 472.

### *Keywords*

Collation, textual criticism, digital philology.

### ABSTRACT

La collazione è una delle prime attività filologiche per le quali è stato proposto l'uso del computer. Fin dagli anni '60 assistiamo al fiorire di algoritmi per la collazione detta semi-automatica. In questo articolo se ne segue la storia, con particolare attenzione alle problematiche dei testi medievali in volgare. Le risorse linguistiche ad oggi disponibili per alcune lingue rendono possibile la progettazione di una *pipeline* computazionale, basata sul modello di Gotheborg: il testo di ogni testimone viene annotato con informazioni linguistiche; l'allineamento si effettua sui lemmi, per neutralizzare la diffusa varianza ortografica; infine, le informazioni linguistiche sono recuperate per categorizzare le varianti. Sebbene la maggior parte dei trattamenti nella *pipeline* siano stati adottati in precedenza nella storia della collazione semi-automatica, un passo avanti è rappresentato dalla proposta di una maggiore automatizzazione qui avanzata grazie alle risorse oggi disponibili.

### *Keywords*

Collazione, filologia, filologia computazionale.