

Estrazione automatica delle caratteristiche del personaggio d'opera attraverso pattern lessico-sintattici

¹Paolo Bonora, ²Angelo Pompilio

^{1,2}Università di Bologna, Italia

¹paolo.bonora@unibo.it

²angelo.pompilio@unibo.it

Abstract

L'articolo presenta la sperimentazione di regole basate su pattern sintattici per l'estrazione di informazioni relative alle relazioni interpersonali tra i personaggi del repertorio operistico del XVIII secolo. Lo studio illustra l'applicazione di questo approccio non supervisionato per individuare una specifica tipologia di relazioni descritte nelle didascalie che corredano gli elenchi di personaggi presenti nei libretti. I risultati dimostrano l'efficacia della soluzione proposta nell'estrarre relazioni definite in base a ontologie formali con una precisione tale da consentire la loro inclusione in una base di conoscenza di dominio senza supervisione. La sperimentazione si pone come primo contributo all'elaborazione di un modello formale di descrizione delle caratteristiche dei personaggi d'opera utile alla ricostruzione dei rispettivi profili. Data la dimensione del repertorio preso in esame e la numerosità dei suoi personaggi, l'impiego dell'analisi automatizzata offre allo studioso un efficace strumento a supporto dell'analisi critica delle fonti. L'identificazione delle caratteristiche e la ricostruzione della rete di relazioni è propedeutica all'utilizzo del personaggio come chiave di lettura della tradizione dei testi operistici.

The paper presents the experimentation of rules based on syntactic patterns for the extraction of interpersonal relationships between characters from the opera repertoire of the 18th century. The study illustrates the application of this unsupervised approach to identify this kind of relationships described within the captions accompanying the lists of characters in the librettos. The results demonstrate the effectiveness of the proposed solution in extracting relations defined through formal ontologies with a precision that allows them to be included in a domain knowledge base without further supervision. The experimentation contributes to the elaboration of a formal model for the description of the features of opera characters required for the reconstruction of their profiles. Given the size of the repertoire under examination and the number of characters, the use of automated analysis offers the researcher a useful tool to support the critical analysis of the sources. The ability to reconstruct the network of relationships and the

features of the characters is in turn preparatory to the use of the character as a dimension of analysis to reconstruct the complex tradition of this kind of texts.

Introduzione¹

Lo studio della tradizione dei testi dei libretti d'opera costringe a confrontarsi con una varietà significativa di problematiche. Oltre alle tradizionali difficoltà della filologia e bibliografia caratteristiche del testo poetico, il continuo rimaneggiamento di questi testi introduce le problematiche tipiche della filologia di testi a tradizione contaminata ([1]: 213). Il testo dei libretti è il risultato di una elaborazione strettamente connessa al processo di produzione dello spettacolo operistico. Un processo produttivo che, a partire dall'appalto della stagione affidato all'impresario, risponde ad un complesso sistema di convenzioni, prassi e rapporti tra committente, librettista, compositore e cantanti² tanto da rendere il libretto d'opera un testo instabile per principio ([6]: 34).

Nel descrivere il lavoro di Felice Romani, librettista tra i più attivi del primo Ottocento, Alessandro Roccatagliati individua, come premessa ineludibile perché il librettista possa procedere alla definizione di soggetto e personaggi dell'opera commissionatagli, la conoscenza degli interpreti che la porteranno in scena ([17]: 94). Questo ulteriore fattore di interdipendenza tra testo e contesto produttivo sottolinea la stretta corrispondenza tra le caratteristiche del personaggio e quelle del suo interprete ([17]: 88). Nel primo Ottocento, il periodo utile per la stesura di un libretto era pari a circa quaranta giorni, dalla scelta del soggetto alla consegna al committente ([17]: 67). Un tempo limitato in cui il processo compositivo partiva necessariamente dalla spoliatura di fonti preesistenti da cui ricavare nuovi soggetti o riprendere testi già impiegati³. I contenuti ricavati venivano adeguati, attraverso un lavoro a più mani, alle esigenze specifiche del nuovo allestimento⁴. Nella stesura di un nuovo dramma il nesso con gli antecedenti letterari è rilevante, così come l'influenza del contesto del singolo allestimento sulla definizione del profilo dei personaggi⁵.

1 I paragrafi relativi allo stato dell'arte, alla metodologia e alla sperimentazione sono stati curati da Paolo Bonora, l'introduzione e le conclusioni raccolgono il contributo di entrambi gli autori.

2 L'adesione alle 'convenzioni teatrali', ovvero l'adattamento alla piazza ed alla compagnia per cui è stata commissionata l'opera è spesso determinante: «Se l'impresario dispone di due prime donne, per dire, si potrà, anzi si dovrà fare un dramma fondato su una rivalità erotico-politica femminile, poniamo L'incoronazione di Poppea o Maria Stuarda» ([3]: 198).

3 «Il bricolage delle fonti è a volte capillare, quasi mai dichiarato» ([3]: 199) [corsivo mio].

4 Cfr. ([3]: 198).

5 Nella prefazione alla Fedra andata in scena a Milano nel 1820 (URI del libretto: <http://corago.unibo.it/resource/DOCUMENTI/DMBM18072>) Luigi Romanelli, riferendosi alla sua Fedra di Padova dello stesso anno (URI dell'opera: <http://corago.unibo.it/resource/BRANI/0000097987>), afferma: «Lo stesso argomento fu da me

Anche nel teatro di prosa coevo il testo drammatico acquisisce un ruolo centrale nel nuovo sistema produttivo teatrale. Qui l'interprete entra in relazione con il personaggio sullo schema del carattere⁶ ([10]: 64). La parte diventa quindi il materiale su cui, dove necessario, si può intervenire in funzione delle esigenze sceniche o dell'interprete. Interventi che vengono fatti nel rispetto della rigida codifica dei ruoli rispetto a cui ciascuna parte deve trovare una collocazione ([10]: 65). È così che nel passaggio dall'autore, al capocomico ed infine all'attore, il testo prima frammentato in parti, poi ricomposto sulla scena, assume forma e identità autonome. Anche nel teatro di parola il legame tra testo e personaggio sembra non sia dettato esclusivamente dalla sua funzione drammaturgica, ma anche da quanto i suoi tratti, le sue caratteristiche, risultino funzionali al processo produttivo che lo metterà in scena.

In considerazione di questo stretto rapporto tra personaggio e testo drammatico si è voluto sperimentare l'impiego del personaggio come indicatore utile ad individuare eventuali relazioni di dipendenza tra le opere in cui compaiono. In base a questa ipotesi, la presenza di personaggi con lo stesso nome, o con caratteristiche affini, diventa un indizio di una possibile relazione tra i testi considerati. Diventa quindi importante identificare il personaggio non solo attraverso il nome, ma anche attraverso i tratti che lo caratterizzano, tra cui le relazioni che lo legano agli altri personaggi che compaiono nell'opera. Quest'ultime, infatti, sono uno dei cardini dell'intreccio drammaturgico. Ad esempio, il libretto *Laomedonte* (musica di Lorenzo Baseggio, senza indicazione del librettista, Venezia 1715) riprende sostanzialmente, comprese gran parte delle arie, il libretto de *L'Etearco* (testo di Silvio Stampiglia e musica di Giovanni Bononcini, Vienna 1707), cambia tutti i nomi dei personaggi seri ma ne conserva la caratterizzazione drammaturgica riportata nelle didascalie dei personaggi: Etearco/Laomedonte, "re di Asso in Creta"; Mirene/Cirene, "dama principale d'Asso"; Fronima/Dalinda, "figlia di" Etearco/Laomedonte; Aristeno/Eristono, "figlio di" Etearco/Laomedonte; Polinnesto/Polidoro, "re di Tera"; Temiso/Feraspe, "confidente di" Etearco/Laomedonte. I personaggi buffi Zelta e Delbo, che hanno un ruolo marginale nella vicenda, sono soppressi, ma il ruolo di Delbo, "servo di Polinnesto", è assunto da Liso, "capitano di Polidoro". In *Melinda* (testo di Giovanni Bertati, musica di Sebastiano Nasolini, Venezia 1798) il plot è in parte ripreso da *L'incanto superato* (testo dello stesso Giovanni Bertati, musica di Franz Xaver Süssmayr, Vienna 1793) ma il testo è completamente riscritto ed i nomi dei personaggi, solo in alcuni casi modificati, presentano affinità significative nella loro caratterizzazione: Falsirena/Melinda, "fata"; Carotta, "scudiero di Oliviero"; Firmina/Erminia, "damigella di Falsirena/Melinda"; Lidia, "damigella di" Falsirena/Melinda. Il personaggio

trattato in un Melodramma, che si rappresentò in Padova nell'occasione dell'ultima fiera del Santo. Ma il presente Melodramma, fuorché la sostanza del fatto, nulla ha di comune col primo né per la condotta, né per la versificazione; imperciocché la diversità del Teatro, la qualità degli Attori, ed altre circostanze esigevano, che la composizione fosse del tutto nuova».

6 Ci si riferisce qui alla definizione di carattere come «l'insieme di tratti psicologici, morali e fisici che contraddistinguono un personaggio teatrale, cioè la sua identità e personalità scenica. Ma, specie nel primo Ottocento, si considerava il carattere un segno, un'impronta, una qualità costante e fissa (che poteva coincidere con una passione dominante) rigidamente codificata.» ([10]; Dizionario dei Ruoli, voce: Carattere).

femminile di Logistilla, “altra fata, amica di Oliviero”, è sostituito dal personaggio maschile Malagigi, “compagno di Oliviero”, ma conserva lo stesso ruolo drammaturgico in relazione al cavaliere Oliviero.

L’uso delle caratteristiche fisiche, morali e sociali come elementi utili a ricostruire il rapporto tra testo e personaggio è stato oggetto d’attenzione della critica letteraria novecentesca. In Ferdinand De Saussure il personaggio dei grandi cicli leggendari è un simbolo di cui possono variare nell’ordine il nome, la relazione con gli altri personaggi, il carattere e le azioni ([19]: 185). Ben prima, Alessandro Piccolomini, nella dedica ad Antonio Cocco de *La sfera del mondo* nell’edizione del 1561, aveva già proposto una classificazione dei personaggi della commedia in base alle loro caratteristiche, relazioni e vicissitudini. Da questa Daniele Seragnoli ricava uno schema concettuale delle caratteristiche fondamentali individuate dal Piccolomini articolate in: relazioni di parentela, condizione economica, età, professione, indole e stati d’animo ([18]: 304). Da questo studio deriva l’idea di avviare la sperimentazione partendo proprio dalle relazioni di parentela. D’altro canto, secondo Vladimir Propp, gli attributi del personaggio delle favole costituiscono la parte mutevole che riveste il nucleo funzionale stabile. Questo spiega l’estrema varietà del genere a fronte di una sostanziale uniformità strutturale e tematica ([19]: 191). Volendo applicare questa lettura ai libretti, una estrazione sistematica delle caratteristiche dei personaggi consentirebbe di avere una valutazione quantitativa di questo fenomeno all’interno di testi fortemente connessi alle rispettive fonti letterarie. Se in Propp gli attributi e le caratteristiche superficiali del personaggio costituiscono le ‘dimensioni variabili’, a partire da Claude Lévi-Strauss questa distinzione viene superata riconducendo il personaggio a componente pienamente integrato all’interno della struttura narrativa del testo. Nelle successive analisi strutturaliste, a partire da Roland Barthes, sarà la nozione stessa di personaggio a risolversi in quella di soggetto agente, analizzato in base al suo ruolo nell’azione e non in relazione al suo essere ([19]: 195). Nel caso dei libretti va tenuto presente che il personaggio si posiziona al centro di quel delicato equilibrio tra referenzialità, mimesi del reale e puro fantastico che sostiene la ‘meccanica’ dell’universo drammaturgico del melodramma. È così che nel mondo di Don Giovanni lo spettatore è in grado di accettare che una statua presenti al dissoluto il conto dei suoi misfatti ([9]). Ecco quindi che, come primo passo, è necessario descrivere il personaggio ricavando dal testo le sue caratteristiche principali, o almeno quelle che l’autore del libretto ha voluto indicare come tali al pubblico.

Di norma nei libretti i personaggi vengono presentati allo spettatore attraverso le didascalie descrittive che corredano l’elenco delle *dramatis personae*. Queste descrizioni riportano in forma sintetica quei tratti del personaggio che consentono allo spettatore di riconoscerlo quando compare in scena e nello svolgersi dell’azione drammatica. Questi testi sono quindi una fonte primaria per ricavare quei caratteri che, a giudizio dell’autore, meglio identificano i personaggi della sua opera. Di norma il personaggio può essere descritto mediante uno o più attributi: la professione svolta, il suo rango sociale, il luogo di origine o la provenienza, le relazioni che lo

legano agli altri personaggi dell'opera, le caratteristiche d'animo o morali. Ad esempio, Cleonice⁷ sarà la “regina di Siria, amante corrisposta d'Alceste”, oppure Tirso⁸ il “villano affittuale de' beni d'Euristene in villa”.

In alcuni casi viene data qualche ulteriore informazione utile a contestualizzare il personaggio nell'azione drammatica come avviene per Semiramide⁹ che appare “in abito virile sotto nome di Nino re degli Assiri, amante di Scitalce, conosciuto ed amato da lei antecedentemente nella corte d'Egitto come Idreno”.



Figura 1: Pagina del libretto con l'elenco dei personaggi della Semiramide riconosciuta.

7 Cleonice, personaggio nel Demetrio, re di Siria, libretto di Pietro Metastasio, compositore Gioacchino Cocchi, andata in scena al King's Theatre in the Haymarket di Londra l'8 novembre 1757 (riferimento URI dell'opera: <http://corago.unibo.it/resource/BRANI/0000050181>).

8 Tirso compare in Demetrio tiranno, libretto di Aurelio Aureli, compositore Bernardo Sabatini, andato in scena al Teatro Nuovo di Piacenza nel 1694 (URI: <http://corago.unibo.it/resource/BRANI/0000050215>).

9 Personaggio eponimo della Semiramide riconosciuta di Pietro Metastasio, compositore Leonardo Vinci, rappresentata a Roma al Teatro delle Dame il 6 febbraio 1729 (URI: <http://corago.unibo.it/resource/BRANI/0000051966>).

In altri casi le informazioni sono strettamente legate al contesto in cui vengono presentate: Fausta è semplicemente indicata come “sua moglie”, in riferimento a Costantino, nell’omonimo libretto di Giovanni Albinoni¹⁰.

Queste descrizioni costituiscono il materiale da cui si è partiti per ricavare una descrizione strutturata del personaggio. Per ottenere questo scopo si è sperimentata la suddivisione del testo in unità informative elementari, sotto forma di sintagmi, corrispondenti alle sue caratteristiche indicate nella didascalia.

Di seguito viene illustrata l’estrazione di alcune di queste caratteristiche operata in forma automatizzata attraverso l’impiego di tecniche di Natural Language Processing (NLP). Il caso di studio è rappresentato dai personaggi dei drammi di Metastasio e Goldoni; il corpus è costituito dai testi delle didascalie ricavati dalla base di conoscenza (KB) Corago LOD¹¹ ([5]). Partendo da un’analisi morfosintattica dei contenuti, si sono sperimentate alcune euristiche di estrazione basate su pattern lessico-sintattici con l’obiettivo di estrarre le informazioni relative alle relazioni interpersonali che intercorrono tra i personaggi del repertorio preso in esame.

Stato dell’arte delle metodologie di analisi automatica dei contenuti

L’impiego di strumenti informatici per l’estrazione di dati strutturati a partire da informazioni presenti in testi espressi in linguaggio naturale ha ormai una storia più che trentennale ([12]: 796). L’insieme di tecniche che svolgono questo tipo di analisi dei testi costituiscono il campo di indagine dell’*information extraction* (IE). In questo ambito, l’annotazione morfosintattica, sotto forma di dipendenze funzionali¹² ([8]), e la lemmatizzazione dei testi consentono la costruzione di processi di estrazione basati sull’aderenza dei contenuti a specifici schemi sintattici di riferimento (*syntactic paths*) e termini in essi contenuti. La modalità di estrazione basata su *pattern* rientra nell’ambito delle tecniche cosiddette a bassa supervisione (*Lightly Supervised Approaches - LSA*) ([12]: 772). Le LSA sono applicazioni di estrazione dell’informazione che richiedono una ridotta supervisione dell’attore umano in fase di addestramento del sistema. La supervisione è infatti limitata alla sola definizione delle regole di estrazione definite in base allo specifico ambito di applicazione. Con l’introduzione delle *Universal Dependencies* (UD) è stata introdotta una

10 Costantino, libretto di Giovanni Kreglianovich Albinoni, compositore Joseph Hartmann Stuntz, rappresentata a Venezia al Teatro La Fenice l’8 febbraio 1820.

11 Il dataset è reperibile attraverso il DOI: <https://doi.org/10.5281/zenodo.3865867>

12 La proposta di una grammatica libera dal contesto (Context Free Grammar - CFG) in cui la frase viene rappresentata come una concatenazione di dipendenze sintattiche tra le parole che la compongono viene fatta risalire ai primi lavori di Noam Chomsky negli anni Cinquanta ([12]: 457).

tassonomia di dipendenze sintattiche multilingua ([7]) resa applicabile anche all'italiano grazie al *porting* delle funzioni di annotazione dello Stanford Core NLP ([2]).¹³

L'efficacia degli algoritmi per l'estrazione automatica di informazioni viene misurata in base alla capacità di reperire nuova informazione (*recall*) e alla precisione e aderenza al dominio di interesse dell'informazione estratta (*precision*). Nel nostro caso, il tipo di caratteristiche (*features*) e di relazioni che si vogliono individuare non sono necessariamente note a priori. Non è quindi possibile utilizzare criteri di estrazione basati su esempi attestati in un corpus. I sistemi per l'annotazione delle dipendenze sintattiche offrono una possibile soluzione ([13]). Analizzando le caratteristiche lessicali del testo e la sua articolazione sintattica è possibile individuare strutture ricorrenti nel discorso (*patterns*) in base alle quali definire le euristiche di estrazione. Accoppiando l'estrazione dei *pattern* a specifici lemmi utilizzati come criteri di restrizione è possibile raggiungere significativi livelli di *recall* mantenendo alta la *precision* pur avendo adottato un approccio non supervisionato dall'attore umano ([14]). L'applicazione di queste metodologie di estrazione dell'informazione al contesto dei libretti d'opera costituisce un primo tentativo del loro impiego per l'arricchimento della descrizione dei contenuti di questa particolare tipologia di testi.

Metodologia

Nella sperimentazione è stata adottata una combinazione di diversi strumenti: l'annotazione delle dipendenze morfosintattiche; le euristiche di estrazione basate su *pattern* e l'estrazione vincolata alla semantica del lessico presente nelle porzioni di testo analizzate. L'obiettivo è riuscire ad estrarre informazione qualificata sulla base di criteri semantici applicati al grafo sintattico analizzato. Ad esempio, per Cherinto¹⁴, descritto come “figlio di Demofonte, amante di Creusa”, è necessario estrarre le due caratteristiche indicate: essere “figlio di Demofonte” e “amante di Creusa”.

Il primo passaggio è l'annotazione delle caratteristiche morfosintattiche del testo. Questo è il passaggio chiave che permette di «rendere esplicita, interpretabile ed esplorabile la struttura linguistica implicita nel testo» ([11]: 211). La sperimentazione ha applicato tre livelli di annotazione offerti dallo StanfordNLP Core: il *tagging* PoS, la lemmatizzazione e l'annotazione UD.

La prima assegna ad ogni parola del testo la rispettiva categoria grammaticale, ovvero la parte del discorso (Part of Speech - PoS) o categoria lessicale ([11]: 213). La seconda porta

13 Per approfondimenti sulla localizzazione dello stack Stanford CoreNLP per l'italiano si rinvia al sito web relativo al progetto Tint realizzato dalla Fondazione Bruno Kessler (<http://tint.fbk.eu/>; consultato il 10/02/2021).

14 Personaggio che compare nel Demofonte di Pietro Metastasio messo in musica da Antonio Ferrandini e andato in scena il 26 dicembre 1758 al Regio Ducale Teatro di Milano.

all'identificazione dell'esponente lessicale (o lemma) del singolo termine. La terza consente di esplicitare il grafo delle relazioni sintattiche che possono essere espresse attraverso triple RDF.

Quest'ultimo passaggio consente di definire i *pattern* di estrazione come interrogazioni espresse attraverso il linguaggio SPARQL.

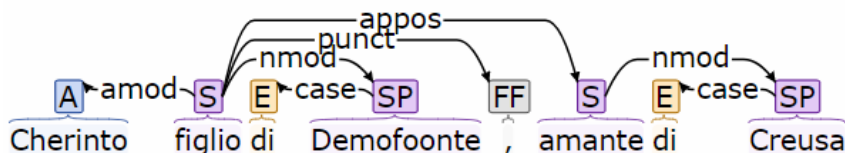


Figura 2: Esempio di annotazione sintattica UD.

La Figura 2 mostra il risultato prodotto dalla annotazione sintattica. Partendo dalla testa della frase, il sostantivo “figlio”, la prima relazione individuata è con il nome proprio del personaggio, “Cherinto”, identificato come modificatore aggettivale¹⁵ (“amod” è l’etichetta utilizzata per l’annotazione). A sua volta “figlio” è associato ad “amante” dalla relazione di apposizione (etichetta “appos”), ovvero di modificatore coordinato a “figlio” che ne estende la definizione. Ciascuno dei due modificatori “figlio” e “amante” hanno a loro volta dei modificatori: “Demofonte” specifica di chi è “figlio” Cherinto e “Creusa” di chi è amante. Considerando i modificatori aggettivali possiamo attribuire al personaggio Cherinto le due *features*: “figlio di Demofonte” e “amante di Creusa”.

Occorre qui notare che le didascalie, per loro natura, tendono ad avere una struttura sintattica vincolata dal fatto che il nome del personaggio¹⁶ è di norma la testa della frase. Questo deriva dal fatto che la didascalia, nell’impianto tipografico del libretto, è apposta al nome del personaggio nel relativo elenco. Le stesse sono state trascritte tenendo conto di quest’ordine e la frase sottoposta ad analisi è stata ricomposta utilizzando lo stesso criterio. Questa prassi compositiva, per quanto presente, non impedisce che possano presentarsi variazioni anche significative dell’albero sintattico, con la conseguente necessità di adattare la regola utile ad intercettarle.

Scorrendo la struttura sintattica della didascalia è quindi possibile segmentare i contenuti della descrizione del personaggio in sotto-unità informative. Ciascuna di queste può essere assunta come un attributo del personaggio o ulteriormente segmentata fino al raggiungimento della granularità corrispondente al singolo termine. L’individuazione delle porzioni di testo da

15 La guida alle relazioni sintattiche definite nelle UC versione 2.x è mantenuta online all’indirizzo: <https://universaldependencies.org/u/dep/> (consultato il 10/02/2021).

16 I nomi dei personaggi di ciascuna opera, ovvero quelli di cui si vogliono estrarre le relazioni, provengono dall’authority file dell’archivio Corago e sono ricavati dagli stessi libretti da cui derivano le didascalie. Ci si attende quindi una coerenza delle forme del nome all’interno dello stesso libretto. La normalizzazione delle varianti del nome del personaggio tra diverse edizioni della stessa opera, o tra testi diversi è invece ancora oggetto di studio ([4]).

considerare come attributi viene quindi effettuata in base a *pattern* sintattici espressi attraverso le UD.

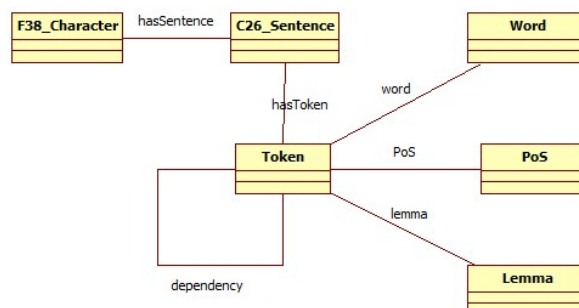


Figura 3: Modello di annotazione morfo-sintattica della descrizione del personaggio.

Il primo passaggio consiste nel produrre l'annotazione morfosintattica del testo. Il secondo consisterà nell'estrazione di sezioni le cui dipendenze corrispondono alla struttura del sintagma definita dal *pattern* di estrazione. Ottenuti i sintagmi si procede a valutare la coerenza semantica dei contenuti utilizzando un criterio che considera la semantica dei lemmi presenti nel sintagma estratto. Ad esempio, i sostantivi che fungono da modificatori aggettivali dovranno essere contenuti in una lista controllata di termini. In questo modo, se selezioniamo lemmi che esprimono una classe coerente di contenuti, ad esempio quelli di relazione parentale di primo grado (come "figlio", "figlia", "madre", "padre"), verranno estratti solo i sintagmi che descrivono la genealogia del personaggio.

Questa procedura di estrazione consente di identificare le porzioni del testo che sono portatrici di caratteristiche del personaggio che rientrano nella tipologia di nostro interesse. Per ottenere questo risultato è però necessario realizzare un modello descrittivo delle componenti del testo che consenta di operare contemporaneamente sui tre livelli di annotazione necessari: morfologico (PoS e lemmatizzazione), sintattico (annotazione delle dipendenze con le UD) e lessicale. Diventa così possibile definire criteri di estrazione che individuano le caratteristiche dei personaggi in base alla semantica del lessico contenuto nelle porzioni di testo analizzato.

Per questo studio i nomi dei personaggi, i testi delle didascalie ed i tre livelli di annotazione sono stati rappresentati attraverso triple RDF. I nomi dei personaggi, le didascalie ed i riferimenti al repertorio sono stati ricavati dal *dataset* Corago LOD accessibile attraverso l'identificativo DOI: <https://doi.org/10.6084/m9.figshare.7880489>. Le annotazioni, ricavate mediante la libreria Tint, sono state riportate in forma di triple nel dataset sperimentale.

Sperimentazione

La metodologia di estrazione descritta è stata applicata su un campione: i personaggi dei drammi di Pietro Metastasio¹⁷ e di Carlo Goldoni¹⁸. La scelta di questi due autori è motivata da tre fattori: l'appartenenza al medesimo periodo storico, ovvero la parte centrale del XVIII secolo¹⁹, da cui ci si attende una uniformità del lessico utilizzato nelle didascalie; la presenza di personaggi seri e comici; la consistenza quantitativa del campione esaminato: 1122 le opere su libretti di Metastasio, 333 di Goldoni. Nel complesso sono state considerate 1103 opere per un totale di circa 9000 personaggi. Di questi sono state analizzate le circa 6275 didascalie disponibili nel dataset Corago LOD da cui si è cercato di estrarre le relazioni interpersonali che intercorrono tra i personaggi di una stessa opera.

Dagli esempi riportati risulta evidente come queste relazioni abbiano un'articolazione a livello sintattico e una correlazione a livello terminologico. Il modello concettuale dell'annotazione proposto permette di coniugare i due piani. Le regole per l'estrazione di relazioni sono quindi composte da:

1. un *pattern* sintattico che definisce la struttura dei sintagmi che stabiliscono una relazione tra due personaggi;
2. un dizionario che restringa i contenuti del sintagma al dominio di interesse.

Il *pattern* ha lo scopo di identificare quei sintagmi che legano il personaggio 'soggetto' della relazione con il personaggio 'oggetto'. Nei nostri esempi, Cherinto (soggetto) "è figlio di" Demofoonte (oggetto). Il dizionario dei termini rilevanti, inteso come elenco di lemmi che definiscono il tipo di relazione, permette di discriminare la qualifica di "figlio" da quella di "amante". La prima viene quindi classificata come specifica relazione di parentela, la seconda come relazione interpersonale. Entrambe vengono attestate su un vocabolario formalizzato.

Nel nostro caso deve essere previsto un terzo criterio di estrazione che consenta di identificare il personaggio in posizione di 'oggetto' nella relazione come uno dei personaggi della stessa opera. In questo caso si dovrà verificare che il lemma corrispondente coincida con un nome presente nell'elenco dei personaggi della stessa opera.

Nella costruzione della regola di estrazione per ciascuna tipologia di relazione deve essere individuato l'insieme di *pattern* che rappresentano le diverse varianti di dipendenze sintattiche con le quali la relazione può essere espressa nei testi. A valle si specifica l'insieme dei termini la

17 Riferimento URI dell'autore: <http://corago.unibo.it/resource/RESPONS/Z00000337900>.

18 Riferimento URI dell'autore: <http://corago.unibo.it/resource/RESPONS/Z00000244400>.

19 La prima opera censita in archivio di Metastasio è la Angelica del 1720 (riferimento URI dell'opera: <http://corago.unibo.it/resource/BRANI/0000157674>), del Goldoni è L'uomo di mondo del 1728 (riferimento URI dell'opera: <http://corago.unibo.it/resource/BRANI/0001486902>). Le ultime opere del Goldoni risalgono al 1762, un anno prima della sua morte, e del Metastasio al 1781.

cui semantica corrisponde alla tipologia di relazione ricercata e la relativa posizione all'interno del pattern.

Operando su questi due insiemi, ci si aspetta che con l'estensione del primo si aumenterà la *recall* complessiva della regola, ovvero la capacità di selezionare un numero più vicino a tutti i sintagmi che specificano una relazione utile. Restringendo il numero di elementi del secondo insieme si aumenterà la *precision*, ovvero il controllo sulla coerenza semantica della relazione individuata.

L'individuazione della struttura dei *pattern* può essere fatta sia in modo deduttivo che empirico-induttivo. Nel primo caso si ipotizzano le catene di dipendenze che si reputano possano contenere relazioni utili, nel secondo si analizza un campione di relazioni note e se ne ricava la struttura. Nel nostro caso si è optato per il secondo approccio applicando alcune generalizzazioni dove opportuno²⁰.

Per la definizione degli elenchi dei termini che identificano il tipo di relazione si è ricorso all'utilizzo di dizionari²¹. Navigando la gerarchia di iponimi di un termine qualificante, si è costruito l'insieme di lemmi il cui significato rappresenta efficacemente la tipologia di relazione oggetto di interesse.

Il Listato 1 riporta un esempio di interrogazione SPARQL che estrae dal grafo annotato la relazione "figlio di Demofonte"²² per il personaggio di Cherinto.

Dal punto di vista implementativo, la regola di estrazione viene definita come una istruzione di tipo CONSTRUCT del linguaggio SPARQL. Questo tipo di istruzione restituisce un insieme di triple che asseriscono uno o più fatti (o assiomi) ricavati dalla base di conoscenza (blocco 0) in base ai criteri definiti nella corrispettiva clausola di WHERE. Questa è a sua volta articolata in più blocchi con funzioni specifiche.

Nell'esempio riportato, al blocco 1 viene assegnato il valore dell'URI del personaggio di cui si sta analizzando la didascalia; il vincolo può essere omesso per analizzare l'intero dataset alla ricerca della tipologia di relazione identificata dal blocco 6. Nel blocco 2 vengono estratte le variabili relative ai contenuti della didascalia. Nel blocco 3 vengono identificati i termini (sotto forma di *token*) che partecipano alla struttura del sintagma definita come *pattern* di relazioni UD tra di essi nel blocco 4. Il successivo blocco 5 estrae i lemmi corrispondenti ai *token* e nel blocco 6 viene imposto il filtro sul contenuto lessicale del *token* più significativo: in questo caso il sostantivo modificatore aggettivale. Nel blocco 7 viene imposto il vincolo che il *token* modificatore nominale corrisponda al nome di un personaggio della stessa opera.

20 Un esempio di applicazione dello stesso approccio per la definizione di regole di estrazione basate su una combinazione di lessico e sintassi è reperibile in [16].

21 Come dizionario per i lemmi è stato utilizzato il synset Babelnet ([15]).

22 In questo caso rappresentata formalmente dalla proprietà fhkb:isSonOf.

```
#Interrogazione SPARQL per l'estrazione della relazione padre/figlio tra personaggi
CONSTRUCT
{
  #BLOCCO 0: costruzione della tripla che specifica la relazione tra i personaggi
  ?characterSbjURI <http://www.cs.man.ac.uk/~stevensr/ontology/fhkb.owl#isSonOf> ?characterObjURI.
}
WHERE
{
  #BLOCCO 1: identificazione del personaggio in posizione di 'soggetto' [OMETTERE per ricercare su tutto il dataset]
  BIND(<http://corago.unibo.it/resource/PERSONAG/0000050806> as ?characterSbjURI)

  #BLOCCO 2: estrazione didascalia del personaggio
  ?characterSbjURI <http://corago.unibo.it/sm/hasSentence> ?sent.
  ?sent rdfs:label ?lblSent.
  ?characterSbjURI ^<http://corago.unibo.it/sm/CNLP1_refers_to_character> ?1.
  ?1 <http://universaldependencies.org/u/dep#word> ?characterName.

  #BLOCCO 3: estrazione TOKEN cui applicare il pattern
  ?sent <http://universaldependencies.org/u/dep#hasToken> ?1.
  ?sent <http://universaldependencies.org/u/dep#hasToken> ?2.
  ?sent <http://universaldependencies.org/u/dep#hasToken> ?3.

  #BLOCCO 4: definizione del pattern sintattico -> inverse(ud:amod)/ud:nmod/ud:case con PoS(ud:amod)=S
  OPTIONAL { ?1 ^<http://universaldependencies.org/u/dep#amod> ?2.}
  ?2 <http://universaldependencies.org/u/dep#nmod> ?3.
  ?3 <http://universaldependencies.org/u/dep#case> ?3c.
  ?2 <http://universaldependencies.org/u/dep#POS> "S". # restrizione ai soli SOSTANTIVI
```

Listato 1: Interrogazione SPARQL per l'estrazione della relazione figlio di.

In termini più generali le clausole di estrazione devono prevedere tre sezioni fondamentali, qui rappresentate dai blocchi 4, 6 e 7, ciascuno dei quali specifica una delle tre tipologie di criteri ipotizzate come necessarie ad estrarre una relazione. Il blocco 4 specifica la struttura del *pattern* sintattico che i sintagmi devono avere: il 'criterio sintattico'. Il blocco 6 specifica il contenuto

lessicale del nodo centrale del grafo del sintagma: il ‘criterio semantico’. Il blocco 7 applica il vincolo di appartenenza dell’oggetto della relazione al contesto il ‘criterio di dominio’.

Il ‘criterio sintattico’ specifica la struttura che i sintagmi devono rispettare sotto forma di grafo di relazioni tra i termini che li compongono. Questa viene formalmente definita utilizzando le proprietà che rappresentano le dipendenze UD. Alla struttura delle dipendenze sintattiche possono essere aggiunti ulteriori criteri morfologici del sintagma. Nel caso specifico si noti che nel blocco 4 è previsto che il *token* in posizione 2 deve essere annotato come PoS di tipo ‘S’ (sostantivi).

Il ‘criterio semantico’, necessario per garantire la coerenza della relazione asserita nel sintagma identificato attraverso il ‘criterio sintattico’, viene definito mediante un insieme controllato di termini. In questo caso la regola prevede che uno o più *lemmi* del sintagma devono essere presenti nell’insieme dato. Formalmente il vincolo è espresso attraverso l’istruzione FILTER che valuta un’espressione booleana di tipo rdfs:IN (valore 1, ..., valore n). Nel caso del blocco 6 l’insieme dei valori considerati è limitato al singolo elemento ‘figlio’. Questa restrizione è legata alla semantica molto specifica della relazione fhkb:isSonOf che la regola mira ad estrarre. La proprietà, infatti, asserisce una relazione tra un genitore ed un figlio maschio. Sarebbe quindi possibile aumentare la *recall* estendendo l’insieme dei termini con l’aggiunta di altri termini come “figliolo” o “rampollo”, ma non “figlia”.

Nel nostro esempio, ai criteri di estrazione relativi al contenuto della descrizione è necessario aggiungere un vincolo relativo al contesto di dominio. Il blocco 7 utilizza un costrutto del linguaggio che esegue un’interrogazione su un grafo esterno per verificare che il personaggio identificato come secondo termine della relazione appartenga alla stessa opera²³.

A partire dall’esempio mostrato, è stata sperimentata la possibilità di analizzare le didascalie attraverso l’uso di procedure di IE in grado di considerare la struttura (sintassi), i contenuti (lessico) e applicare vincoli di dominio (repertorio e fonti documentali) per produrre informazione qualificata che potrà essere inserita direttamente nella base di conoscenza Corago LOD²⁴.

Seguendo la tecnica di estrazione descritta, la sperimentazione è stata orientata all’ estrazione delle qualifiche più frequentemente attribuite ai personaggi del campione. Il *pattern* è stato disegnato in modo da selezionare i modificatori nominali direttamente riferiti al personaggio descritto dalla didascalia.

23 Nel caso in analisi la funzione attinge al repertorio delle opere del dataset Corago LOD. Il grafo di riferimento su cui opera la regola, essendo uno strumento intermedio di analisi, è tenuto separato dal dataset principale.

24 Il dataset delle relazioni estratte in formato RDF N-Triple è reperibile attraverso il DOI:
<https://doi.org/10.6084/m9.figshare.14748531>

Applicata sul campione preso in esame, la regola ha estratto 67 occorrenze di relazione padre/figlio. La Tabella 1 riporta le 11 occorrenze distinte in base al nome dei personaggi. Nome personaggio soggetto	Relazione	Nome personaggio associato
Arbace	figlio di	Artabano
Argeno	figlio di	Ernesto
Cherinti	figlio di	Demofonte
Cherinto	figlio di	Demofonte
Demetrio	figlio di	Antigono
Everardo	figlio di	Gualtiero
Germondo	figlio di	Alarico
Idelberto	figlio di	Berengario
Learco	figlio di	Eurinome
Sammete	figlio di	Amasi
Siroe	figlio di	Cosroe

Tabella 1: Nomi di personaggi associati da una relazione figlio di.

Per estrarre le relazioni interpersonali in forma strutturata è stata costruita una regola basata sul modello descritto nel Listato 1. Ad esempio: il personaggio Farnaspe²⁵, descritto come “principe parto, amico e tributario d'Osroa, amante e promesso sposo d'Emirena”, deve essere messo in relazione con Emirena²⁶, descritta come “prigioniera d'Adriano, amante di Farnaspe”. Il risultato atteso dalla procedura di estrazione è una tripla RDF che asserisca: <Farnaspe> relationship:lifePartnerOf <Emirena>. È stata definita una regola di estrazione per ciascuna tipologia di relazione composta da un pattern sintattico e da un dizionario dei termini che esprimono la semantica della relazione. Si noti che, a parità di pattern sintattico, è sufficiente

25 Riferimento URI del personaggio: <http://corago.unibo.it/resource/PERSONAG/0000437343>

26 Riferimento URI del personaggio: <http://corago.unibo.it/resource/PERSONAG/0000437342>

operare sul criterio lessicale per ricavare relazioni diverse. Ad esempio, intervenendo sul blocco 6 del Listato 1 sostituendo il lemma “amante” a “figlio”, la regola intercetta 505 relazioni riferibili al tipo relationship:lifePartnerOf²⁷. La semantica della relazione estratta è stata formalmente espressa mediante due ontologie formali: la Relationship Ontology (<https://vocab.org/relationship/>) e la Family History Knowledge Base ([20]).

L'analisi condotta utilizzando la metodologia descritta ha consentito di estrarre dal campione dei personaggi dei drammi di Metastasio e Goldoni le relazioni riportate nella Tabella 2.

LEMMMA	RELAZIONE	CONTEGGIO LEMMA	RELAZIONI ESTRATTE	RECALL	CORRETTA	PRECISION
amante	relationship:lifePartnerOf	2.902	505	17,40	496	98
amico/a	relationship:friendOf	706	113	16,01	95	84
padre	fhkb:isFatherOf	471	95	20,17	95	100
sorella	fhkb:isSisterOf	469	237	50,53	237	100
figlio	fhkb:isSonOf	391	98	25,06	98	100
figlia	fhkb:isDaughterOf	277	84	30,32	84	100
moglie	fhkb:isWifeOf	135	78	57,78	78	100
fratello	fhkb:isBrotherOf	71	14	19,72	14	100
madre	fhkb:isMotherOf	43	21	48,84	21	100
nipote	fhkb:isNephewOf	35	27	77,14	27	100
zia	fhkb:isAuntOf	7	1	14,29	1	100

Tabella 2: Relazioni interpersonali estratte.

Il confronto tra il conteggio delle occorrenze dei lemmi e il numero di relazioni individuate evidenzia una significativa distanza: per esempio il lemma “amante” registra 2901²⁸ occorrenze

²⁷ L'insieme delle query SPARQL sviluppate per la sperimentazione è reperibile attraverso il DOI: <https://doi.org/10.6084/m9.figshare.14743239>

²⁸ Va sottolineato che la presenza del lemma all'interno della didascalia è una misura indiretta della potenziale presenza della relazione ricercata. Infatti, per poterla ricavare è necessario che siano presenti

contro 505 relazioni. Questa distanza, nel complesso, porta ad una *recall* media del 34,30%. Per contro, la *precision* sul campione considerato risulta ottimale²⁹. Per individuare in modo puntuale le relazioni binarie tra i personaggi è stato necessario costruire regole definite da criteri particolarmente selettivi. Questa condizione incide in modo negativo sulla *recall* e di conseguenza sullo F-score³⁰ che risulta pari a 50,86%. D’altro canto, l’obiettivo primario della sperimentazione è l’estrazione di informazione qualificata senza necessità di successiva validazione. In questa prospettiva si è scelto di privilegiare la *precision* a garanzia dell’attendibilità dell’informazione estratta.

L’impiego di *pattern* calibrati rispetto alla *precision* ha consentito di ottenere il risultato desiderato, ma allo stesso tempo non ha intercettato tutta l’informazione utile. Ad esempio, nel caso di “Idelberto figlio di Berengario e di Matilde”, il *pattern* ha intercettato la relazione tra Idelberto e Berengario ma ha omissso la relazione con la madre Matilde. Esempi affini sono: “Sandro villano, amante di Menghina, poi di Cecca” o “Pandolfo padre di Cecilia e Dorina”; le prime relazioni sono state individuate, le successive invece sono state omesse. Più complessa è la ricostruzione dell’albero genealogico nel caso di: “Giuda fratello di Giuseppe e Beniamino, figliuoli di Giacobbe e di Lia”. L’algoritmo indentifica Giuda come fratello di Giuseppe e figlio di Giacobbe grazie alla lemmatizzazione delle parole, ma non individua la relazione tra Giuda, fratello di Beniamino e Lia madre di Giuda e Beniamino. In questi casi il *pattern* dovrebbe essere esteso in modo da poter trattare correttamente la paratassi tra le diverse parti del discorso. Questo risultato può essere conseguito sia con l’impiego di più *pattern* in parallelo per intercettare più strutture sintattiche al fine di migliorare la *recall*, sia con l’estensione del vocabolario dei termini che identificano le relazioni analizzate per incrementare la *recall*.

di tutti i termini della stessa. Non necessariamente l’oggetto corrisponde ad un personaggio dell’opera come nel caso di Linco: “figlio d’Egitto, amante d’Ipermestra” (<http://corago.unibo.it/resource/PERSONAG/0001185313>). Trattandosi comunque quantomeno di una condizione necessaria, è stata utilizzata come termine per il calcolo della *recall* anche se potenzialmente peggiorativo.

29 I risultati rispetto alle prime due tipologie di relazione sono influenzati dalla presenza di modificatori che alterano la semantica del lemma utilizzato come criterio di estrazione. In particolare la precisione è influenzata per relazione *relationship:friendOf* dal caso di Zopiro “falso amico di Radamisto” per il *pattern* individua una relazione di “amicizia”, sulla relazione *relationship:lifePartnerOf* i casi di Learco “amante ricsusato d’ Issipile” e Rodope “amante ingannata di Learco”. In entrambi i casi relazioni in realtà negate dal profilo drammaturgico dei rispettivi personaggi. Questa considerazione, di natura interpretativa, ha suggerito di considerarle come casi negativi.

30 Misura utilizzata per stabilire l’accuratezza dei modelli calcolata come la media armonica di *precision* e *recall*.

Conclusioni e sviluppi ulteriori

Nel complesso la sperimentazione, volutamente circoscritta ad una tipologia specifica di relazioni, ha evidenziato la fattibilità dell'approccio basato su regole che uniscano criteri sintattici a criteri semantici. Grazie alla metodologia descritta è stato possibile ricavare in modo non supervisionato 1273 relazioni interpersonali da un corpus di 6275 didascalie. In particolare, si è verificata la possibilità di raggiungere una precisione tale da consentire l'introduzione dell'informazione estratta nella base di conoscenza senza ulteriore validazione del dato estratto. D'altra parte, questo risultato è stato raggiunto a scapito della capacità delle regole di estrazione di intercettare tutte le relazioni effettivamente presenti nelle didascalie. Sulla base di questi primi esiti è possibile individuare possibili affinamenti della procedura attraverso due direttrici. La prima dovrà verificare la possibilità di aumentare l'indice di *recall* operando una estensione dei *pattern* sintattici impiegati per selezionare le parti di testo. La seconda dovrà verificare l'efficacia della soluzione nell'estrarre relazioni di natura diversa operando sui dizionari impiegati per definire il criterio semantico. I risultati ottenuti in base all'approccio basato su regole presentato dovranno essere confrontati con un equivalente basato sul *machine learning* supervisionato per valutare quali dei due risultino più efficaci nel contesto considerato.

References

- [1] Accorsi, Maria Grazia. 1989. «Problemi testuali dei libretti d'opera fra Sei e Settecento». *Giornale storico della letteratura italiana*, 1989.
- [2] Aproso, Alessio Palmero, e Giovanni Moretti. 2016. «Italy Goes to Stanford: A Collection of CoreNLP Modules for Italian». *ArXiv:1609.06204 [Cs]*, settembre. <http://arxiv.org/abs/1609.06204>
- [3] Bianconi, Lorenzo. 2017. «Il libretto d'opera». In *Musica*. Istituto della Enciclopedia Italiana.
- [4] Bonora, Paolo. 2020. «Impiego del Web Semantico per lo sviluppo e la consultazione di archivi musicali. Un caso di studio sulla storia e la documentazione del melodramma italiano: l'archivio Corago». PhD Thesis, Alma Mater Studiorum - Università di Bologna. <http://amsdottorato.unibo.it/9174/>
- [5] Bonora, Paolo, e Angelo Pompilio. 2021. «Corago in LOD: the debut of an Opera repository into the Linked Data arena». *JLIS.it* 12 (2): 54–72.
- [6] Coletti, Vittorio. 2017. *Da Monteverdi a Puccini*. Einaudi.
- [7] De Marneffe, Marie-Catherine, Timothy Dozat, Natalia Silveira, Katri Haverinen, Filip Ginter, Joakim Nivre, e Christopher Manning. 2014. «Universal Stanford Dependencies: A Cross-Linguistic Typology». *LREC* 14.
- [8] De Marneffe, Marie-Catherine, e Christopher D Manning. 2008. «Stanford Typed Dependencies Manual». Stanford University.
- [9] Elam, Keir. 1988. *Semiotica del teatro*. Bologna: Il Mulino.

- [10] Jandelli, Cristina. 2002. *I ruoli nel teatro italiano tra Otto e Novecento*. Firenze: Le lettere.
- [11] Lenci, Alessandro, Simonetta Montemagni, e Vito Pirrelli. 2005. *Testo e computer*. Roma: Carocci.
- [12] Martin, James H, e Daniel Jurafsky. 2009. *Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition*. Prentice Hall series in artificial intelligence. Upper Saddle River, NJ: Prentice Hall, Pearson Education International.
- [13] Mintz, Mike, Steven Bills, Rion Snow, e Dan Jurafsky. 2009. «Distant Supervision for Relation Extraction without Labeled Data». In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 2 - ACL-IJCNLP '09*, 2:1003. Suntec, Singapore: Association for Computational Linguistics.
<https://doi.org/10.3115/1690219.1690287>
- [14] Moro, Andrea, Hong Li, Sebastian Krause, Feiyu Xu, Roberto Navigli, e Hans Uszkoreit. 2013. «Semantic Rule Filtering for Web-Scale Relation Extraction». In *Advanced Information Systems Engineering*, a cura di Camille Salinesi, Moira C. Norrie, e Óscar Pastor, 7908:347–62. Berlin-Heidelberg: Springer. https://doi.org/10.1007/978-3-642-41335-3_22
- [15] Navigli, Roberto, e Simone Paolo Ponzetto. 2012. «BabelNet: The Automatic Construction, Evaluation and Application of a Wide-Coverage Multilingual Semantic Network». *Artificial Intelligence* 193: 217–50.
- [16] Poria, Soujanya, Erik Cambria, Lun-Wei Ku, Chen Gui, e Alexander Gelbukh. 2014. «A Rule-Based Approach to Aspect Extraction from Product Reviews». In *Proceedings of the Second Workshop on Natural Language Processing for Social Media (SocialNLP)*, 28–37. Dublin, Ireland: Association for Computational Linguistics and Dublin City University.
<https://doi.org/10.3115/v1/W14-5905>
- [17] Roccatagliati, Alessandro. 1996. *Felice Romani librettista*. Quaderni di Musica/realità. Lucca: Libreria musicale italiana.
- [18] Seragnoli, Daniele. 1987. «La struttura del personaggio e della fabula». In *Il teatro italiano nel rinascimento*, 297–317. Problemi e prospettive. Serie di musica e spettacolo. Bologna: Il Mulino.
- [19] Stara, Arrigo. 2004. *L'avventura del personaggio*. Firenze: Le Monnier università.
- [20] Stevens, Robert, Nicolas Matentzoglou, Uli Sattler, e Margaret Stevens. 2014. «A Family History Knowledge Base in OWL 2». *Informal Proceedings of the 3rd International Workshop on OWL Reasoner Evaluation (ORE 2014)*, 6.