

# Implementazione di Lector usando YAGO come KG di riferimento

**Mattia Iodice**

Email: mat.iodice1@stud.uniroma3.it

**Daniele Petrillo**

Email: dan.petrillo@stud.uniroma3.it

**Nicholas Tucci**

Email: nic.tucci@stud.uniroma3.it

Nella relazione viene presentata l'implementazione di Lector utilizzando YAGO come knowledge graph di riferimento. L'obiettivo prefissato era appunto l'implementazione del processo di estrazione di fatti utilizzato in Lector su un altro Knowledge Graph, valutando poi il risultato mediante il confronto con i risultati ottenuti su Freebase. Il documento è strutturato in sezioni, ciascuna focalizzata su un aspetto diverso del lavoro svolto. In particolare, dopo una breve descrizione delle caratteristiche fondamentali di YAGO, vengono descritti gli approcci seguiti nella fase di implementazione, insieme alla giustificazione delle scelte architeturali e alla presentazione degli strumenti utilizzati. In seguito è invece descritto l'approfondimento del precedente homework sullo studio dell'estrazione di fatti da frasi lista attraverso l'utilizzo di Stanford NLP Core. Nella sezione finale vengono quindi presentati i risultati ottenuti, fornendo un confronto diretto con l'esecuzione di Lector su Freebase.

## 1 YAGO

La prima parte del lavoro svolto ha riguardato lo studio dell'architettura di YAGO, insieme alla raccolta di dati circa le informazioni in esso contenute. Vengono quindi presentati i punti salienti emersi dall'indagine.

YAGO è un knowledge graph realizzato dal Max Planck Institute for Computer Science e deriva direttamente da Wikipedia, WordNet e GeoNames.

Dalla documentazione disponibile sulla pagina ufficiale di YAGO sono dichiarate **10 milioni di entità** e più di **120 milioni di fatti tra queste entità**. Inoltre per le entità sono individuate **circa 350 mila classi** (o tipi di appartenenza), mentre le **relazioni disponibili sono circa 100**.

Uno degli aspetti più rilevanti e originali di YAGO è da ricercare nel sistema di tassonomia usato per la rappresentazione gerarchica di entità e tipi, derivante da un utilizzo ibrido delle categorie di Wikipedia e Wordnet.

Ogni entità è essenzialmente un'istanza di una o più classi e la tassonomia risulta così strutturata:

1. *rdfs:Thing*: il nodo radice della tassonomia
2. classi di tassonomia provenienti da WordNet, che hanno

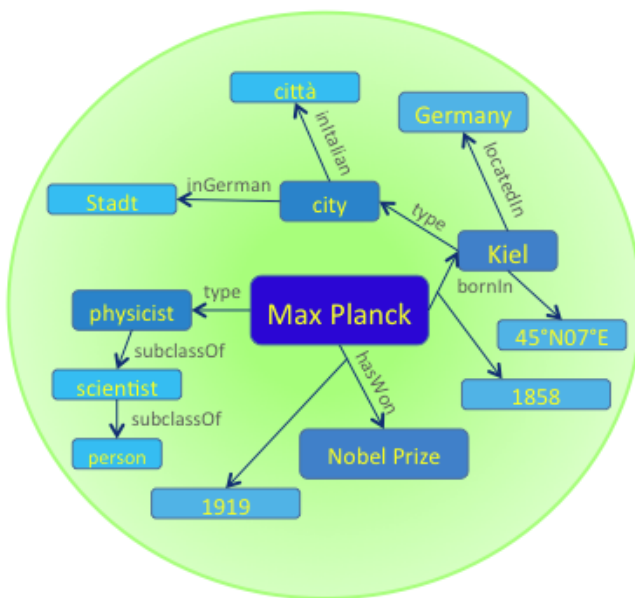


Fig. 1. Esempio della struttura di Yago. All'entità Max Planck sono associati due fatti che riguardano le relazioni *hasWon* e *bornIn*. All'entità Max Planck viene poi associata la *wikicat*, sottoclasse di *scientist*, a sua volta figlio di *person* nell'albero dei tipi.

un nome definito come *<wordnet\_XXX\_YYY>*, dove XXX è il nome del concetto rappresentato e YYY l'id del concetto all'interno di WordNet. Si ha quindi che le classi sono connesse a classi più generali tramite la relazione *rdfs:subClassOf*

3. classi derivate dalle categorie di Wikipedia, come ad esempio *<wikicategory\_American\_rock\_singers>*, che è derivata dalla categoria Wikipedia American Rock Singers. Queste classi sono connesse a quelle di WordNet sempre tramite la relazione *rdfs:subClassOf*
4. lo strato più basso della tassonomia è invece costituito dalle entità vere e proprie: oggetti del mondo reale, come fiumi persone, ecc. Ogni entità è poi connessa a classi di livello superiore tramite la relazione **rdf:type**.

Un esempio completo è il seguente:

Table 1. In tabella il confronto tra YAGO e Freebase in termini di fatti ed entit presenti

Knowledge Graph	# Fatti	# Entità
YAGO	120 M	10 M
Freebase	2.4 B	44 M

all'entità Ingrid Bergman sono associate le *wordnet* Woman, Female, Living Thing, Actress, Adult, Biographer e le *wikicat* Women, Swedish people, Jewish actors, 20th century Swedish actresses, 20th century Swedish people.

I dati utilizzati per il processo di implementazione di Lector provengono tutti dalla pagina ufficiale di YAGO e fanno riferimento ai file *yagoTaxonomy*, *yagoSchema*, *yagoTransitiveType* e *yagoFacts*.

*yagoTaxonomy* contiene l'intera tassonomia di YAGO ed è essenzialmente l'albero delle classi rappresentato sotto forma di file tsv.

*yagoTransitiveType* rappresenta invece l'associazione di ciascuna entità con le classi di cui è istanza.

Ciascuna relazione del knowledge graph è associata ad un particolare soggetto (domain) e impone un determinato oggetto (range): ciascuno di questi legami è definito in *yagoSchema*, mentre l'insieme dei fatti presenti è contenuto in *yagoFacts*.

Nell'utilizzare YAGO viene consigliato di usare il core dello stesso knowledge graph, perciò il processo è implementato sui file descritti in precedenza, che rappresentano proprio il nucleo del grafo. Abbiamo comunque ritenuto opportuno analizzare i dati presenti e abbiamo estratto diverse informazioni.

Dal file *yagoFacts.tsv* emerge che i fatti presenti sono **5.625.089**, mentre da *yagoTransitiveTypes.tsv* sono invece state contate **5.129.147** entità. Il dato che emerge è che si tratta di un grafo poco denso, in cui i nodi che hanno almeno un arco entrante o uno uscente sono **2.632.727**.

Un ulteriore aspetto interessante che abbiamo notato nell'analisi dei documenti messi a disposizione è che delle **100** relazioni attestate, solo **37** compaiono nei fatti presenti in *yagofacts*.

## 2 Implementazione di Lector

I dati riportati fanno riferimento all'esecuzione su un pc con le seguenti specifiche:

1. Intel Core i7-6700HQ
2. memoria fisica installata 8.00GB
3. SSD

### 2.1 Rimozione rumore e frasi non relazionali

Il processo di estrazione di fatti è stato implementato in Java, sebbene alcuni passaggi fondamentali siano stati sviluppati con l'ausilio di strumenti supplementari come

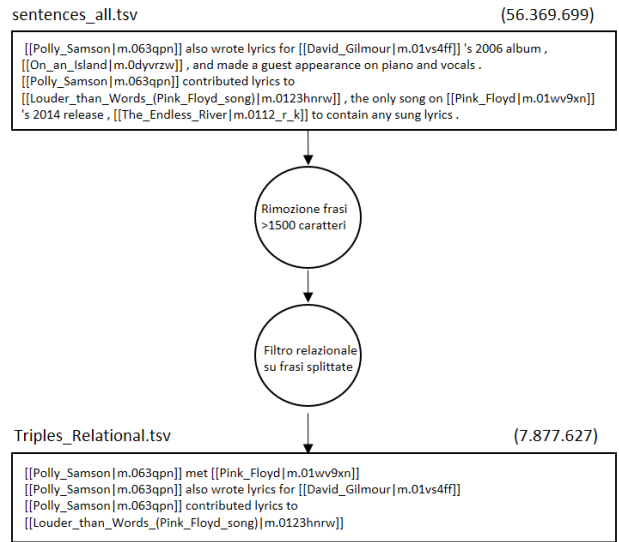


Fig. 2. Prefiltraggio del dump di Wikipedia

MySQL e script bash.

Analizzando il file annotato *sentences\_all.tsv* sono emerse alcune frasi particolari e problematiche per alcuni aspetti implementativi: si tratta di frasi spropositatamente lunghe e senza alcun significato, per lo più costituite da sequenze di entità. Abbiamo quindi filtrato il file escludendo le frasi più lunghe di 1500 caratteri. Seguendo i passaggi fondamentali del processo seguito da Lector abbiamo poi trasformato ciascuna frase in una sequenza di triple, costituite da frasi racchiuse tra due entità consecutive.

Ciascuna tripla è stata quindi data in pasto ad un filtro relazionale, con l'obiettivo di scartare tutte le frasi che non fossero relazionali, come mostrato in Figura 2. Il file ottenuto in output conteneva **7.631.515** triple distinte, e perciò ne consegue che il filtro ha scartato circa l'**85%** delle triple. Di queste **7.631.515** abbiamo contato quelle con soggetto e oggetto entrambi presenti nelle entità definite in *yagoTypes*: **5.480.597** triple. A questo punto abbiamo deciso di procedere attraverso quattro esecuzioni parallele:

1. la prima senza considerare né frasi lista né generalizzazione
2. la seconda effettuando la **generalizzazione** di alcune categorie di termini ricorrenti (lavori, numeri, stati, nazionalità)
3. la terza effettuando la **generalizzazione** parallelamente ad un pesante processo di **stopping**
4. la quarta, più che altro complementare alle precedenti, **normalizzando liste** di entità

Tale scelta ci ha consentito di valutare singolarmente i quattro approcci descritti, permettendoci poi di trarre conclusioni sull'effettivo beneficio apportato.

### 2.2 Associazione di frasi e relazioni senza stopping e generalizzazione

Nella prima fase si è quindi direttamente proceduto con l'associazione di frasi alle relazioni. Dovendo confrontare

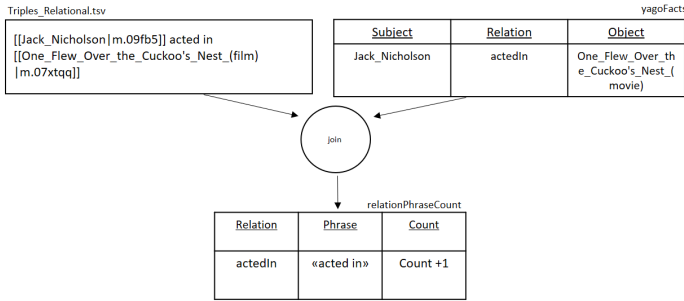


Fig. 3. Join tra le triple e i fatti

soggetto e oggetto di 7 milioni di triple con soggetto e oggetto di 5 milioni di fatti, abbiamo intuito che una buona soluzione in termini di efficienza e soprattutto di flessibilità potesse derivare dall'utilizzo di un DBMS. Infatti la logica SQL si presta particolarmente bene alle operazioni descritte nel paper: l'associazione di frasi alle relazioni non è altro che un join tra soggetto e oggetto di una tripla con soggetto e oggetto nei fatti presenti in YAGO. Abbiamo quindi innanzitutto definito la tabella *yagoFacts(sub,relation,ob)*, nella quale abbiamo importato i fatti dal .tsv *yagoFacts*. A questo punto per ciascuna tripla presente è stato sufficiente eseguire una query del tipo:

```
SELECT * FROM yagoFacts WHERE sub=subjT AND ob=objT
```

(1)

In questa fase è stato interessante osservare che:

1. le query che hanno restituito almeno un match sono state **475.062**. Questo dato è risultato particolarmente rilevante perché, se sottratto al numero di triple che hanno soggetto e oggetto entrambi appartenenti alle entità di YAGO<sup>1</sup>, indica il massimo numero di fatti estraibili: poco più di **5 milioni**.
2. ogni query restituisce un certo numero di relazioni, quelle appunto già esistenti tra soggetto e oggetto all'interno di YAGO. È emerso che la somma delle relazioni restituite è **547.601**. Ciò indica che, in media, per ogni coppia soggetto-oggetto ci sono **1.15** fatti.
3. il numero di frasi univoche presenti in *relPhraseCount*<sup>2</sup> è circa **85.000**, mentre le tuple presenti sono **111.514**

L'esecuzione di questo join ha richiesto circa 120 minuti.

### 2.3 Misurazione della forza di un'associazione entità-frase

La tabella *relPhraseCount* mette appunto in relazione ciascuna relazione con un certo numero di frasi e il numero di occorrenze. Per misurare quindi la forza dell'associazione tra relazione e frase abbiamo riutilizzato la formula proposta in Lector:

$$P(r_i|p) = \frac{c(p, r_i)}{\sum_{j \in R} c(p, r_j)} \quad (2)$$

<sup>1</sup>5.480.597

<sup>2</sup>Vedi 2.3 per definizione

Abbiamo quindi definito la nuova tabella *relPhraseScore(rel,phrase,score,probability)*, in cui inserire, appunto i valori di probabilità e score. Avendo già a disposizione la tabella *relPhraseCount* è stato possibile riempire la nuova tabella andando ad inserire in probability il risultato della query:

```
SELECT count FROM relPhraseCount WHERE rel = x AND phrase = y
SELECT sum(count) FROM relPhraseCount WHERE phrase = y
```

(3)

e nella colonna score:

$$score(p, r_i) = \log c(p, r_i) \cdot P(r_i|p). \quad (4)$$

Queste due operazioni sono risultate il collo di bottiglia dell'intero processo di estrazione di nuovi fatti. In particolare il tempo impiegato per popolare completamente la nuova tabella è stato di **4 h e 27 min**. Una considerazione su cui ci siamo soffermati è che si tratta però di operazioni che si prestano abbastanza bene alla **scalabilità orizzontale**.

È infatti sufficiente distribuire la tabella *relPhraseCount* su più macchine e operare concorrentemente su frasi distinte, aggregando poi i risultati parziali. Questo approccio permetterebbe di scalare facilmente e abbattere i tempi di esecuzione.

### 2.4 Estrazione di nuovi fatti e validazione

Al fine di scegliere le K=20 frasi più significative per ogni relazione è stato sufficiente eseguire la query che ci restituisse per ciascuna relazione le prime 20 frasi con probabilità maggiore di 0.5, ordinate per score decrescente.

A quel punto abbiamo definito una mappa *HashMap<String>, <Set<String>>* *phraseToRelations*. Attraverso questa mappa è stato possibile scansare il file iniziale, contenente le triple relazionali, per estrarre nuovi fatti: ogni volta che la frase contenuta nella tripla risulta contenuta in *phraseToRelations* viene aggiunto un nuovo fatto tra soggetto e oggetto della frase, corrispondente alle relazioni associate alla frase stessa.

I vincoli che abbiamo specificato nell'aggiunta del fatto sono:

1. sia soggetto che oggetto devono essere entità presenti in YAGO
2. non deve essere già presente un fatto corrispondente alla relazione individuata tra le due entità

Se per il secondo vincolo non c'è bisogno di particolari chiarimenti è bene spendere qualche riga per chiarire l'esigenza e le implicazioni del primo.

Al fine di validare i fatti estratti è infatti necessario che ciascuna entità compaia nel file *yagoTypes*. Questo requisito è necessario perché ciascuna relazione presente in YAGO lega un particolare tipo di soggetto ad un particolare tipo di oggetto. Vincoli che sono espressi in *yagoSchema*.

Proponendo un esempio concreto si ha che:

La relazione **actedIn** richiede un soggetto di tipo *actor* e

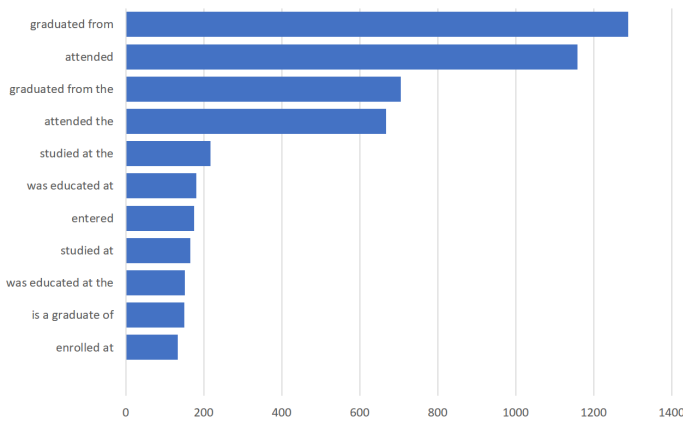


Fig. 4. Frasi pi ricorrenti per la relazione *graduatedFrom*

Table 2. Variazione del numero di fatti estratti al variare dei vincoli

Vincolo	# Fatti estratti	# Fatti validati
entit� in yagoFacts	340 K	190 K
entit� in yagoTypes	740 K	493 K
nessun vincolo	1.05 M	-

un oggetto di tipo *movie*. Supponendo che sia stato estratto il fatto `[[Cristiano_Ronaldo...]] played in [[Real_Madrid...]]`, associato alla relazione *actedIn*, la verifica dei tipi consentirebbe di scartare il fatto che, appunto, non rispetta il dominio della relazione. Ci siamo infatti accorti che in YAGO tutti i fatti presenti rispettano le imposizioni espresse in *yagoSchema* e perci  la validazione dei fatti   un processo indispensabile prima di poter inserire nuovi fatti nel knowledge graph.

Nonostante questi vincoli abbiamo comunque provato ad analizzare le variazioni delle quantit  di fatti estratti e validati modificando i requisiti:

1. le entit  devono essere gi  presenti in yagofacts
2. le entit  devono essere presenti in yagotypes
3. nessun vincolo sulle entit 

I risultati sono mostrati in Tabella 2 e come ci si poteva aspettare l'assenza di vincoli consente di estrarre un gran numero di fatti in pi .

In questa fase non c'  stato bisogno di caricare tutte le entit  in MySQL:   bastato leggere da tsv e memorizzare in un unico set le entit  univoche presenti (sia che provenissero dai fatti gi  presenti in YAGO sia che fossero stati presi da yagoType).

Per i tipi invece   stato necessario inserire i dati in MySQL, trattandosi di una quantit  di dati intorno ai 7 GB.

Abbiamo quindi definito la tabella *yagoTypes(entit ,tipo)*, e per ogni fatto estratto verificato che soggetto e oggetto avessero, tra i propri tipi, quelli richiesti dalla relazione in *yagoSchema*.

In un primo momento utilizzavamo un file *yagoTypes* sem-

Table 3. Fatti estratti e fatti validati con il processo di generalizzazione e con il processo ibrido generalizzazione e stopping

Approccio	# Fatti estratti	# Fatti validati
Generalizzazione	777 K	499 K
Generalizzazione e stopping	2.19 M	985 K

plificato, ma il contenuto non era per niente esaustivo, mancando spesso i tipi fondamentali delle entit . Questo implicava un'analisi non banale dei tipi, dovendo andare a percorrerne l'albero delle dipendenze fino alla radice e costringendoci a fare assunzioni pericolose.

Abbiamo quindi analizzato i documenti disponibili, fino a trovare il file *yagoTransitiveTypes*, che garantisce la chiusura transitiva delle entit  rispetto ai tipi. In esso ciascuna entit  risulta associata a tutti i suoi supertipi, fino alla radice *owl:Thing*, ed   quindi stato sufficiente fare un confronto 1:1 tra ciascun tipo dell'entit  e il tipo richiesto dalla relazione. In tutti e tre i casi il tempo di esecuzione per l'estrazione di fatti si   mantenuto al di sotto dei **18 minuti**, mentre quello per la validazione intorno ai **12**.

## 2.5 Generalizzazione e stopping

Dopo aver implementato questo primo processo di estrazione di nuovi fatti abbiamo ragionato sui modi di **generalizzazione delle frasi**, usando dei placeholder al posto di parole contenenti specifiche informazioni non essenziali per descrivere la relazione. Quindi, al di l  dei meccanismi proposti nel paper su Lector, abbiamo osservato i fatti estratti e cercato di individuare classi di termini in grado di generalizzare l'intero processo. Abbiamo quindi definito cinque dizionari:

1. nazionalit , *NAT*
2. stati, *STA*
3. professioni, *JOB*
4. numeri, *NUM*
5. numeri cardinali, *CAR*

Un esempio di una frase generalizzata   il seguente:

**input:** `[[Sidney_Olcott ...]] married actress [[Vale...]]`  
**output:** `[[Sidney_Olcott ...]] married JOB [[Vale...]]`

Abbiamo quindi ripetuto il processo di estrazione di nuovi fatti e di validazione basata su tipi. In Tabella 3 i risultati ottenuti.

Dopo aver rieseguito Lector con la generalizzazione delle frasi abbiamo voluto rieseguire l'intero processo introducendo anche un modulo che eseguisse lo **stopping** delle frasi, unitamente al processo di generalizzazione. Abbiamo fatto uso della libreria open source Exude. I risultati sono messi a confronto con quelli della generalizzazione in Tabella 3.

Abbiamo scelto di non ripetere dall'inizio il processo di calcolo della probabilità e dello score, poiché il forte sospetto, confermato da diverse esecuzioni di testing, era che si rischiasse di introdurre troppo rumore. Essenzialmente quindi, dalla tabella *relPhraseScore*, costruita in precedenza sulle triple non processate in alcun modo viene effettuato il retrieve delle prime 30 frasi con probabilità maggiore di 0.5, ordinate per score. A queste frasi abbiamo quindi applicato il modulo di generalizzazione (e stopping). Il risultato è stato che diverse frasi sono state aggregate in una sola, sia nel processo di generalizzazione che in quello di generalizzazione e stopping. Abbiamo poi trasformato ciascuna frase nelle triple applicando il medesimo processo di generalizzazione. La nostra ipotesi, testata su un campione di 150K triple, è che costruire i valori di score e probabilità su frasi non processate assicuri una certa resistenza al rumore che altrimenti verrebbe introdotto.

In altre parole, in questo modo le frasi più rappresentative per ciascuna relazione dovrebbero essere più affidabili rispetto a quelle ottenute dal calcolo di score su frasi generalizzate o stoppage.

## 2.6 Normalizzazione di frasi lista

Ultimata la prima fase di implementazione di Lector abbiamo lavorato sul raffinamento del modulo realizzato nel precedente homework. Lo scopo del task era quello di migliorare l'estrazione di fatti da particolari tipologie di frasi, chiamate frasi lista.

Un esempio di questa tipologia di frasi è :

*Obama is a supporter of Chicago White Sox, Chicago Bears, and Pittsburgh Steelers.*

Da questo tipo di frasi Lector ha un approccio abbastanza limitato: riesce ad operare solo in presenza di “*such as*”, ricostruendo le frasi in modo parziale. Ad esempio, dalla frase precedente si vorrebbero estrarre tre fatti distinti:

*Obama is a supporter of Chicago White Sox*  
*Obama is a supporter of Chicago Bears*  
*Obama is a supporter of Pittsburgh Steelers*

Il nostro approccio si è basato quindi sull'utilizzo di *Stanford NLP Core*.

Questa libreria è in grado di analizzare una frase in termini di dipendenze sintattiche, consentendo di ricostruire l'analisi logica e i collegamenti logici tra parole all'interno della frase stessa.

Il principale svantaggio che deriva dall'utilizzo di questa libreria è il costo richiesto per l'analisi di ciascuna frase. Abbiamo infatti calcolato che in media il modulo riesce a calcolare le dipendenze di massimo 150 frasi al minuto. Parte del lavoro è stata quindi incentrata sulla realizzazione di un prefiltro che fosse in grado di individuare potenziali frasi lista, al fine di alleggerire il carico di lavoro a cui la libreria di Stanford viene sottoposta.

### 2.6.1 Prefiltraggio di frasi

La realizzazione del filtro è stata guidata da due esigenze particolari:

1. **efficienza**: tutte le 56 milioni di frasi presenti nel dump di Wikipedia devono essere passate al prefiltro. Non possono quindi essere tollerati tempi di esecuzione troppo lunghi per ciascuna singola frase.
2. **efficacia**: sostanzialmente si ambisce ad una precision molto alta, assicurandoci quindi di non andare a spendere esecuzioni inutili della libreria di Stanford.

Per sostenere l'efficienza abbiamo deciso di parallelizzare l'esecuzione del prefiltro attraverso l'utilizzo del *Lightweight Executable Framework*, utilizzando regex che richiedessero tempi di esecuzione ridotti. Lo speedup ottenuto è di circa 4 e il tempo di esecuzione inferiore ai 20 minuti.

Abbiamo quindi deciso di procedere secondo un approccio basato sull'attribuzione di un punteggio a ciascuna singola frase e sulla definizione di una soglia al di sotto della quale le frasi in input fossero giudicate come non-lista.

Il processo di attribuzione del punteggio si basa sulla ricerca di alcune caratteristiche strutturali all'interno della frase analizzata. In particolare:

**+0.9** se contiene una sequenza di entità dello stesso tipo (la verifica viene effettuata andando a verificare le classi di appartenenza di ciascuna entità su *yagoTransitiveType*).

**+0.5** se contiene almeno un'entità che segue un *list identifier*: una locuzione come *such as, like, both, ...*

**+0.35** se contiene un'entità che segue un *quantificatore*: *some, many, plenty, few, ...*

**+0.3** se nel primo sesto della lunghezza della frase è contenuta almeno un'entità

**+0.2** se sono contenuti dei numeri (numerici o letterali) seguiti da un numero di virgole pari al numero -1.

**+0.15** per ogni entità contenuta tra virgole

**+0.1** per ogni virgola presente

Prima di passare le frasi a questa componente però, attraverso tre espressioni regolari, è stato effettuato un consistente processo di filtraggio del file originale *sentences\_all.tsv*.

Innanzitutto sono state prese in considerazione soltanto le frasi composte da almeno 3 entità:

*egrep “. \*. \*. \*. \*”* (5)

Da questa regex sono state scartate circa **40 milioni** di frasi del file originale.

Ci siamo però resi conto che potrebbe non essere una scelta ottimale:

*[[Obama...]] is a supporter of Chicago White Sox, Chicago Bears, and [[Pittsburgh Steelers..]].*

Questa frase verrebbe scartata e non sarebbe trattata come frase lista, ragion per cui l'approccio semplice di estrazione di fatti non riuscirebbe ad estrarre nulla. Ad



ogni modo, la maggior parte di frasi lista individuate non é soggetta a questo tipo di vulnerabilit . Effettuato questo primo filtraggio abbiamo definito una nuova regex che, dalle **15.597.980** frasi rimanenti, individuasse quelle con una sequenza di entit  vicine.

`egrep "[a-zA-Z0-9\"][0,15][a-zA-Z0-9\"][0,15][a-zA-Z0-9\"][0,15]"` (6)

L'output   stato un file composto da **11.457.798** di frasi. Osservando la tipologia di frasi ottenute abbiamo quindi pensato di prendere soltanto quello che avessero una coppia di entit  separate da un "and":

`entit  = "\[a-zA-Z0-9_\]\/\&-\]*\m.[a-zA-Z0-9_\]*\N"` (7)

`egrep "'$entit '.{0,7} and .{0,7}'$entit '"` (8)

che ha restituito **3.009.416** di frasi.

A questo punto, quindi, **3.009.416** frasi sono state passate attraverso il modulo di attribuzione del punteggio, che ha scaricato tutte quelle al di sotto della soglia **1.65**. Alla fine le frasi candidate a contenere liste di entit , passate quindi allo Stanford NLP Core, sono state **302.533**. Vengono fornite, nell'Appendice A, alcuni esempi delle frasi che hanno passato l'intero processo di riconoscimento di frasi lista.

## 2.6.2 Valutazione del prefiltro

Il prefiltro, per come   stato realizzato, in particolar modo per la scelta dell'attribuzione di un punteggio in base alle propriet  strutturali della frase e alla presenza di una soglia, si presta particolarmente bene all'uso di un classificatore.

Abbiamo per  scelto di procedere facendo un semplice tuning sui parametri e sull'attribuzione del punteggio in base ai risultati ottenuti. Sono quindi state etichettate 600 frasi, distinguendo le frasi lista da quelle non lista e dividendole quindi in training set e validation set. Il tuning sui parametri   stato effettuato manualmente e, nella migliore verifica sul validation set abbiamo ottenuto i seguenti risultati:

1. **Precision:**0.43
2. **Recall:**0.88
3. **Fallout:**0.18

## 2.6.3 Utilizzo di Stanford NLP Core

Una volta individuate le frasi candidate a contenere liste di fatti viene eseguita la libreria Stanford NLP Core su tale insieme di frasi. Ciascuna frase in input viene dunque analizzata dalla libreria, che   in grado di restituire la lista di dipendenze logiche all'interno delle frasi.

Le dipendenze rilevanti per le frasi lista sono risultate essere le *nmod*, che sono generalmente usate per modificatori nominali e per predicati con clausole.<sup>3</sup> Nella maggior parte di frasi lista, si presenta infatti un numero ristretto di tipologie di *nmod*: **nmod:such\_as**, **nmod:like**, **nmod:with** e **nmod:in** sono risultate le pi  ricorrenti.

Individuate quindi le dipendenze *nmod* all'interno della frase il nostro approccio si   basato sul riconoscimento dei

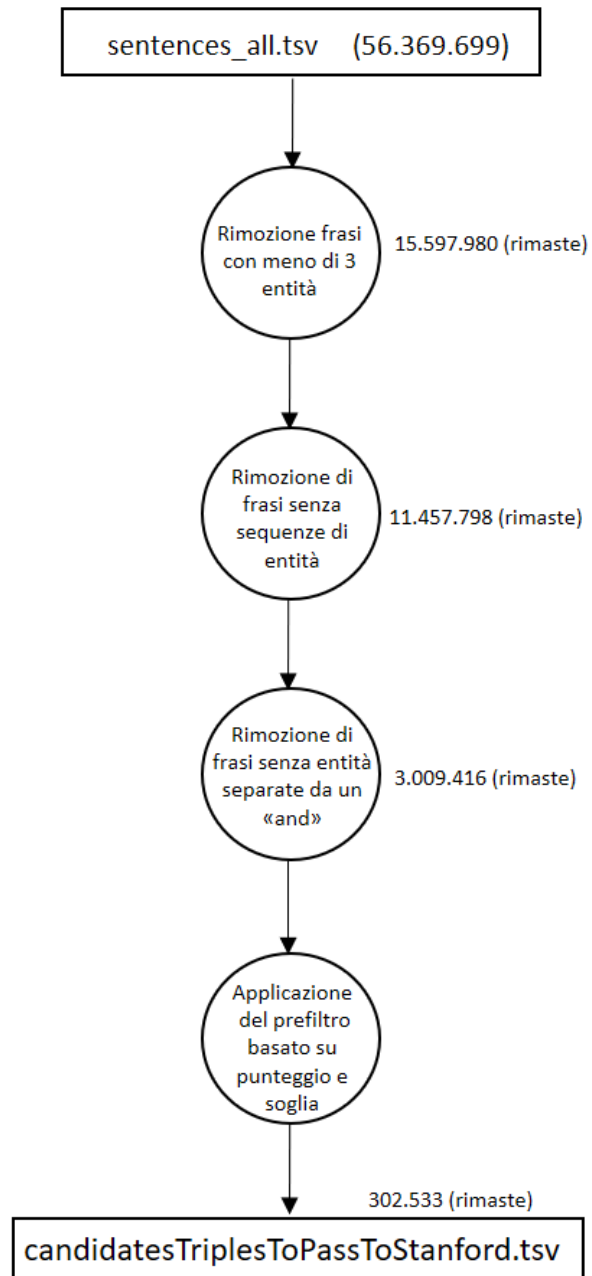


Fig. 5. Processo di prefiltraggio delle frasi lista

*governor* e dei *dependent* di ciascuna di esse.

Questi concetti si riferiscono alle estremit  della relazione e quindi alle parole coinvolte, come mostrato in Figura 6. A questo punto la frase viene suddivisa in sottoperiodi logici in base al numero di *nmod* rilevati, nella condizione in cui i *dependent* individuati siano almeno due. Essenzialmente si adotta questo processo perch  le *dependent* di ciascun *nmod* rappresentano le possibili liste di entit  contenute nella sequenza di fatti che si vuole estrarre.

La fase successiva consiste nell'individuazione del soggetto del sottoperiodo in esame. Si procede quindi con la navigazione delle dipendenze a partire dal *governor* comune degli *nmod*, cercando il ramo che porta alla dipendenza *nsubj*. Il *dependent* di *nsubj*   proprio il soggetto del sot-

<sup>3</sup><http://universaldependencies.org/docs/en/dep/>

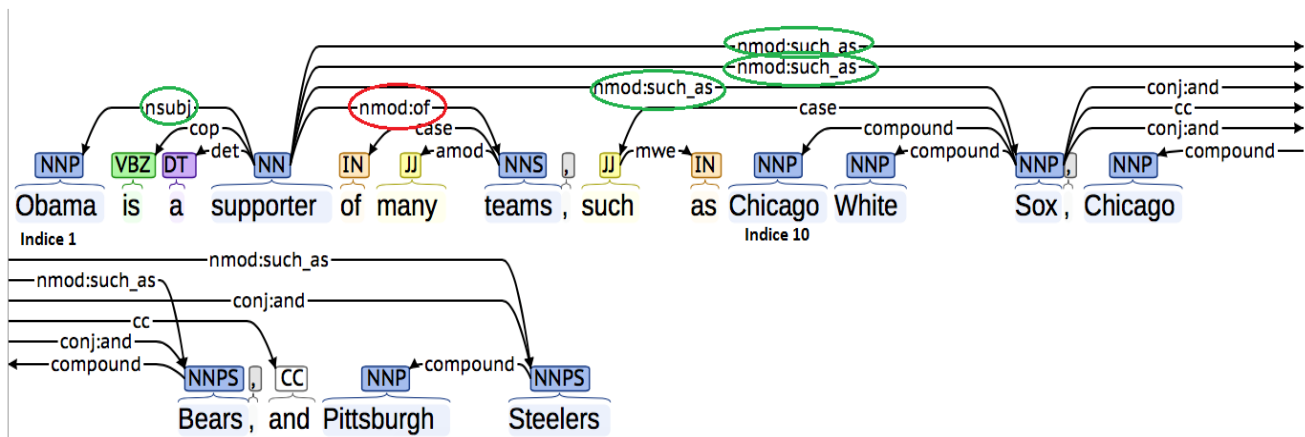


Fig. 6. Nella figura vengono mostrate le dipendenze individuate tramite la libreria di Stanford. Gli nmod rilevanti, cerchiati in verde, sono quelli che presentano diversi dependent. nmod:of viene quindi scartato in quanto presenta una sola dependent.

toperiodo.

A questo punto del processing sostanzialmente sono state individuate le entità coinvolte nel fatto da estrarre, il passo successivo consiste dunque nel riuscire a localizzare la frase caratterizzante il fatto stesso.

Vengono quindi richiesti gli indici corrispondenti al soggetto e alla prima entità coinvolta nella lista. Tali indici corrispondono alla posizione della parola all'interno della frase.

Prendendo in considerazione l'esempio fornito nella Figura 6 avremmo che gli indici in questione sono 1, quello di "Obama", e 10, "Chicago".

Preso la frase delimitata dagli indici individuati si prosegue con il Pos Tagging, eseguito ancora dalla libreria di Stanford. L'obiettivo di questa fase è l'individuazione dei verbi e delle congiunzioni subordinate<sup>4</sup>.

Queste operazioni sono risultate particolarmente rilevanti perché le frasi individuate vengono poi state utilizzate per la ricerca di nuovi fatti.

I principali problemi riscontrati sono i seguenti:

1. prendere tutta la frase compresa tra gli indici, escludendo magari soltanto i list identifiers, come ad esempio "such as" o "like", compromette la possibilità di generalizzare la frase.  
Poiché le frasi estratte vengono successivamente passate al filtro relazionale si rischierebbe di perdere tutte le frasi eccessivamente lunghe.
2. prendere soltanto il verbo all'interno del periodo rischia invece di stravolgerne il significato, con la possibilità concreta di inserire fatti sbagliati nel knowledge graph.

Abbiamo quindi cercato dei compromessi tra queste soluzioni opposte andando ad includere anche le preposizioni e le congiunzioni subordinate presenti tra il verbo e la prima entità della lista. Con riferimento ancora alla Figura 6 otterremmo come output:

*Obama is a supporter of Chicago White Sox*

*Obama is a supporter of Chicago Bears*

*Obama is a supporter of Pittsburgh Steelers*

Su questo tipo di frasi elementari e ben definite i risultati sono più che soddisfacenti, tuttavia all'aumentare della complessità del periodo, l'output peggiora sensibilmente. In particolare, si verifica spesso che la frase estratta ha un significato diverso o parziale rispetto a quella originale. Al termine del processing le frasi estratte vengono quindi passate al filtro relazionale, per poi seguire il processo di estrazione di fatti presentato in precedenza. I tempi di esecuzione di questa componente sono risultati estremamente lunghi e per questo è senza dubbio necessario scalare orizzontalmente, splittando il file ed eseguendo il calcolo di dipendenze su più calcolatori. Anche la valutazione del processo di estrazione di fatti da frase lista, di conseguenza, è stata effettuata su un campione di **10.000** frasi passate alla libreria di Stanford.

### 3 Valutazione dei fatti estratti

Al termine del processo di implementazione abbiamo quindi valutato i fatti estratti mediante tecniche diverse:

1. fatti estratti senza generalizzazione di frasi e senza preprocessing
2. fatti estratti mediante generalizzazione
3. fatti estratti mediante generalizzazione e stopping
4. fatti estratti dalle frasi lista

Abbiamo quindi preso in considerazione le medesime relazioni usate per la valutazione dei nuovi fatti estratti su FreeBase, impostando inoltre la probabilità minima di 0.5 e  $K=20$ . Per la tecnica di estrazione senza preprocessing abbiamo valutato 470 fatti, per quelle basate sul preprocessing 240 ciascuna e per l'analisi delle frasi lista circa 280. Abbiamo inoltre adottato una soluzione di crowdsourcing, rivolgendoci a quattro nostri colleghi. Il costo di tale operazione è stato quantificato in un caffè ogni 80 fatti valutati. Sebbene, paragonando il carico di lavoro e il costo relativo alle offerte presenti su Amazon Mechanical Turk ci siamo accorti che la ricompensa era troppo elevata, abbiamo comunque accettato il compromesso, trattandosi di persone qualificate.

<sup>4</sup>[http://www.ling.upenn.edu/courses/Fall\\_2003/ling001/penn\\_treebank\\_pos.html](http://www.ling.upenn.edu/courses/Fall_2003/ling001/penn_treebank_pos.html)

Table 4. In tabella i risultati ottenuti applicando Lector su **FreeBase**

Relazione	# Fatti in FreeBase	# nuovi fatti	# fatti valutati	# accuratezza
birthPlace	662.192	57.140	347	88,9%
deathPlace	178.849	18.458	104	80,6%
nationality	584.792	50.234	290	95,6%
team	145.080	49.809	286	96,5%
almaMater	378.043	46.342	286	98,3%
spouse	130.425	14.939	97	31,6%
child	141.860	3.149	50	38,8%
award	98.625	1.934	50	96,4%
party	65.300	3.684	50	94,5%

Table 5. In tabella i risultati ottenuti applicando Lector su **YAGO** senza generalizzazione e stopping

Relazione	# Fatti in YAGO	# nuovi fatti	# fatti valutati	# accuratezza
wasBorn	280.942	57.837	40	90,4%
diedIn	94.728	22.203	40	81,6%
isCitizenOf	36.253	28.763	40	96,3%
playsFor	525.358	23.716	40	87,9%
graduatedFrom	51.478	50.229	40	96,3%
isMarriedTo	34.219	16.119	40	35,1%
hasChild	42.330	14.617	40	25,4%
hasWonPrize	120.465	25.736	50	93,1%
isPoliticianOf	32.774	630	50	82,5%
actedIn	126.236	2.385	50	96,7%
<b>Totale fatti estratti ( contando tutte le relazioni)</b>		493.279		



Table 6. In tabella i risultati ottenuti applicando Lector su **YAGO** con generalizzazione

Relazione	# Fatti in YAGO	# nuovi fatti	# fatti valutati	# accuratezza
wasBorn	280.942	57.456	60	90,2%
diedIn	94.728	22.044	60	79,8%
isCitizenOf	36.253	28.568	-	-
playsFor	525.358	23.653	-	-
graduatedFrom	51.478	50.056	60	87,7%
isMarriedTo	34.219	15.842	-	-
hasChild	42.330	14.345	-	-
hasWonPrize	120.465	25.564	60	92,1%
isPoliticianOf	32.774	1.757	-	-
actedIn	126.236	2.374	-	-
<b>Totale fatti estratti</b>		498.941		

Table 7. In tabella i risultati ottenuti applicando Lector su **YAGO** con generalizzazione e stopping

Relazione	# Fatti in YAGO	# nuovi fatti	# fatti valutati	# accuratezza
wasBorn	280.942	58.547	60	89,8%
diedIn	94.728	26.411	60	72,6%
isCitizenOf	36.253	31.397	-	-
playsFor	525.358	96.002	-	-
graduatedFrom	51.478	55.194	60	85,6%
isMarriedTo	34.219	19.038	-	-
hasChild	42.330	59.475	-	-
hasWonPrize	120.465	36.309	60	86,4%
isPoliticianOf	32.774	7.619	-	-
actedIn	126.236	2.739	-	-
<b>Totale fatti estratti</b>		985.597		

Table 8. In tabella i risultati ottenuti applicando Lector su **YAGO** con estrazione di fatti da frasi lista

Relazione	# Fatti in YAGO	# nuovi fatti	# fatti valutati	# accuratezza
wasBorn	280.942	-	-	-
diedIn	94.728	-	-	-
isCitizenOf	36.253	-	-	-
playsFor	525.358	4	4	100%
graduatedFrom	51.478	1	1	100%
isMarriedTo	34.219	9	9	0%
hasChild	42.330	6	6	0%
hasWonPrize	120.465	-	-	-
isPoliticianOf	32.774	-	-	-
actedIn	126.236	28	28	100%
hasNeighbor	558	128	110	98,1%
worksAt	10.254	25	25	100%
<b>Totale fatti estratti</b>		278		

**su 10.000 frasi in input a Stanford**

#### 4 Conclusioni e sviluppi futuri

Dall'analisi delle precedenti tabelle sono emerse le seguenti considerazioni:

1. I fatti estratti senza preprocessing delle frasi (Tabella 5) sono circa 493 K. Ci aspettavamo un numero maggiore, essendo YAGO un knowledge graph meno denso di Freebase. Osservando anche le quantità di fatti estratte su Freebase i valori risultano piuttosto simili.
2. Il preprocessing delle frasi tramite stopping e generalizzazione consente di estrarre quasi il doppio di fatti. Per alcune relazioni, quelle mostrate nelle Tabelle 6 e 7, potrebbe rivelarsi una scelta vincente. Essendo però una quantità di dati considerevolmente maggiore è necessaria una valutazione su un numero di fatti più grande.
3. Lo studio delle frasi lista ha portato notevoli miglioramenti nel numero di fatti estratti per le relazioni **worksAt** e **actedIn**. Il collo di bottiglia deriva dai tempi di esecuzione della libreria di Stanford.

Osservando il rapporto tra i fatti estratti e il numero di frasi passate al modulo di estrazione (tramite libreria Stanford) è lecito domandarsi se abbia senso effettuare l'analisi delle frasi per tutte le relazioni. Un buon approccio sarebbe quello di eseguire l'estrazione di fatti da frasi lista soltanto per alcune relazioni, magari quelle che hanno un riscontro positivo in fase di valutazione.

In ultima analisi quindi, la fase di estrazione di fatti da frasi lista, andrebbe effettuata solo dopo il processo di individuazione delle frasi più significative per ogni relazione.

## APPENDICE A

Alcune delle frasi che passano le regex di prefiltraggio nel processo di estrazione di fatti da frasi lista.

World\_Series\_Cricketm.0198kr 5 To clarify the legal implications , [[Frank\_Packer|m.0321kg]] backed a challenge to the TCCB in the [[High\_Court\_of\_Justice|m.02rc3p]] by three of his players : [[Tony\_Greig|m.03t\_kh]], [[Mike\_Procter|m.054s4d]] and [[John\_Snow\_(cricketer)|m.01gkl2]] .

Seminole m.0198dl 5 The unified [[Seminole|m.0198dl]] spoke two languages : [[Creek\_language|m.04mj8v]] and [[Mikasuki\_language|m.08xx8f]] , two among the [[Muskogean\_languages|m.04ml67]] family .

Seminole m.0198dl 4 [[Seminole|m.0198dl]] tribes generally follow Christianity , both [[Protestantism|m.05sfs]] and [[Roman\_Catholicism|m.02vxy\_]] , and their traditional Native religion , which is expressed through the stomp dance and the [[Green\_Corn\_Ceremony|m.089pdj]] held at their ceremonial grounds .

Vespers m.0199r7 5 The psalms and hymns of the [[Vespers|m.0199r7]] service have attracted the interest of many composers , including [[Claudio\_Monteverdi|m.01vlj]], [[Antonio\_Vivaldi|m.0pth]], [[Wolfgang\_Amadeus\_Mozart|m.082db]] , and [[Anton\_Bruckner|m.0c05p]] .

Canary\_Wharf m.019bwp16 Around 105,000 people work in [[Canary\_Wharf|m.019bwp]] and it is home to the world or European headquarters of numerous major banks , professional services firms and media organisations including [[Barclays|m.05t8c5]] , [[Citigroup|m.01dfb6]] , [[Clifford\_Chance|m.0515b7]] , [[Credit\_Suisse|m.03g776]] , [[Infosys|m.01k8x6]] , [[Fitch\_Ratings|m.04grk2]] , [[HSBC|m.01vq1f]] , [[J.P.\_Morgan\_&\_Co|m.04f75w7]] , [[KPMG|m.0k2gt]] , [[MetLife|m.03kt1t]] , [[Moody's|m.03jhzk]] , Morgan Stanley , [[Royal\_Bank\_of\_Canada|m.02wjhf]] , [[Skadden,\_Arps,\_Slate,\_Meagher\_&\_Flom|m.023nw1]] , [[State\_Street\_Corporation|m.06nlq]] and [[Thomson\_Reuters|m.02r5t91]] .

Summer\_Triangle m.0198k8 7 The [[Summer\_Triangle|m.0198k8]] is an astronomical asterism involving an imaginary triangle drawn on the northern hemisphere 's celestial sphere , with its defining vertices at [[Altair|m.013f4]] , [[Deneb|m.022\_r2]] , and [[Vega|m.080jt]] , the brightest stars in the three constellations of [[Aquila\_(constellation)|m.01npxs]] , [[Cygnus\_(constellation)|m.01wt5]] , and [[Lyra|m.04lsj]] , respectively .

Mia\_Farrow m.01cpqk 10 [[Moses\_Farrow|m.0c155dc]] was in a relationship with actor-director [[Woody\_Allen|m.081lh]] from 1980 to 1992 and appeared in twelve of his thirteen films over that period , including [[Zelig|m.01gp1v]] , [[Broadway\_Danny\_Rose|m.07bxl4]] , [[The\_Purple\_Rose\_of\_Cairo|m.05jf85]] , [[Hannah\_and\_Her\_Sisters|m.0sxns]] , [[Radio\_Days|m.07bxqz]] , [[Crimes\_and\_Misdemeanors|m.04smd]] , [[Alice\_(1990\_film)|m.02drn4]] and [[Husbands\_and\_Wives|m.03m4mj]] .

Umkhonto\_we\_Sizwem.019c5l 5 On 11 July 1963 , 19 ANC and [[Umkhonto\_we\_Sizwe|m.019c5l]] leaders , including [[Arthur\_Goldreich|m.03j\_nt]] and [[Walter\_Sisulu|m.01gf5z]] , were arrested at [[Liliesleaf\_Farm|m.051vk7g]] , [[Rivonia|m.06x59z]] .

45th\_Annual\_Grammy\_Awards m.019bk0 8 [[Norah\_Jones|m.0197tq]] and her song [[Don't\_Know\_Why|m.0cg72g]] were the main recipients of the night , running away with a grand total of six Grammys , including all four major awards : [[Grammy\_Award\_for\_Record\_of\_the\_Year|m.01bgqh]] , [[Grammy\_Award\_for\_Album\_of\_the\_Year|m.01by1l]] , [[Grammy\_Award\_for\_Song\_of\_the\_Year|m.0c4z8]] , and [[Grammy\_Award\_for\_Best\_New\_Artist|m.01c427]] , plus [[Grammy\_Award\_for\_Best\_Female\_Pop\_Vocal\_Performance|m.01c99j]] and [[Grammy\_Award\_for\_Best\_Pop\_Vocal\_Album|m.03qbnj]] .

Spiritual\_(music) m.019c24 6 [[The\_Gospel\_Train|m.047rshj]] , [[Song\_of\_the\_Free|m.04mygf9]] , and [[Swing\_Low,\_Sweet\_Chariot|m.03vvnz]] are likewise supposed to contain veiled references to the [[Underground\_Railroad|m.07x8t]] , and many sources assert that [[Follow\_the\_Drinking\_Gourd|m.0119syxv]] contained a coded map to the [[Underground\_Railroad|m.07x8t]] .

Balafon m.019c0x 10 [[Balafon|m.019c0x]] is a [[Manding\_languages|m.03007g]] name , but variations exist across [[West\_Africa|m.0hqjh]] , including the balangi in [[Sierra\_Leone|m.06s\_2]] and the gyl of the [[Dagara\_people|m.019b\_w]] , [[Lobi\_people|m.036y0m]] and [[Gurunsi|m.0277vdq]] from [[Ghana|m.035dk]] , [[Burkina\_Faso|m.01699]] and [[Ivory\_Coast|m.0fv4v]] .

## APPENDICE B

### Alcuni dei fatti estratti dalle frasi lista

#### Frase iniziale:

Mia\_Farrow m.01cpqk 10 [[Moses\_Farrow|m.0c155dc]] was in a relationship with actor-director [[Woody\_Allen|m.081lh]] from 1980 to 1992 and appeared in twelve of his thirteen films over that period , including [[Zelig|m.01gp1v]] , [[Broadway\_Danny\_Rose|m.07bxl4]] , [[The\_Purple\_Rose\_of\_Cairo|m.05jf85]] , [[Hannah\_and\_Her\_Sisters|m.0sxns]] , [[Radio\_Days|m.07bxqz]] , [[Crimes\_and\_Misdemeanors|m.04smdd]] , [[Alice\_(1990\_film)|m.02drn4]] and [[Husbands\_and\_Wives|m.03m4mj]] .

#### Fatti estratti:

[[Moses_Farrow m.0c155dc]]	appeared in	[[Broadway_Danny_Rose m.07bxl4]]	actedIn
[[Moses_Farrow m.0c155dc]]	appeared in	[[Zelig m.01gp1v]]	actedIn
[[Moses_Farrow m.0c155dc]]	appeared in	[[The_Purple_Rose_of_Cairo m.05jf85]]	actedIn
[[Moses_Farrow m.0c155dc]]	appeared in	[[Broadway_Danny_Rose m.07bxl4]]	actedIn
[[Moses_Farrow m.0c155dc]]	appeared in	[[Zelig m.01gp1v]]	actedIn
[[Moses_Farrow m.0c155dc]]	appeared in	[[The_Purple_Rose_of_Cairo m.05jf85]]	actedIn

#### Frase iniziale:

Vito\_Aconci m.01gfd1 8 [[Vito\_Aconci|m.01gfd1]] has taught at many institutions , including the [[Nova\_Scotia\_College\_of\_Art\_and\_Design|m.02x60g]] , Halifax ; [[California\_Institute\_of\_the\_Arts|m.015zyd]] , Valencia ; [[Cooper\_Union|m.01p7x7]] ; [[School\_of\_the\_Art\_Institute\_of\_Chicago|m.0bjqh]] ; Yale University ; [[University\_of\_Iowa|m.01j\_9c]] , [[Pratt\_Institute|m.03bmmc]] ; and the [[Parsons\_School\_of\_Design|m.08htt0]] .

#### Fatti estratti:

[[Vito_Aconci m.01gfd1]]	has taught at	[[Nova_Scotia_College_of_Art_and_Design m.02x60g]]	worksAt
[[Vito_Aconci m.01gfd1]]	has taught at	[[California_Institute_of_the_Arts m.015zyd]]	worksAt
[[Vito_Aconci m.01gfd1]]	has taught at	[[Cooper_Union m.01p7x7]]	worksAt
[[Vito_Aconci m.01gfd1]]	has taught at	[[School_of_the_Art_Institute_of_Chicago m.0bjqh]]	worksAt
[[Vito_Aconci m.01gfd1]]	has taught at	[[University_of_Iowa m.01j_9c]]	worksAt
[[Vito_Aconci m.01gfd1]]	has taught at	[[Pratt_Institute m.03bmmc]]	worksAt
[[Vito_Aconci m.01gfd1]]	has taught at	[[Parsons_School_of_Design m.08htt0]]	worksAt

Frase iniziale:

The [[New\_England\_town|m.099ltc]] is bordered by [[Lynnfield,\_Massachusetts|m.01m27v]] to the north , [[Lynn,\_Massachusetts|m.0tz1x]] to the east , [[Revere,\_Massachusetts|m.0t\_xb]] to the south , and [[Melrose,\_Massachusetts|m.0t\_5b]] and [[Wakefield,\_Massachusetts|m.01m2pk]] to the west , in [[Middlesex\_County,\_Massachusetts|m.0k3k1]]

Fatti estratti:

[[New_England_town m.099ltc]]	is bordered by	[[Wakefield,_Massachusetts m.01m2pk]]	hasNeighbor
[[New_England_town m.099ltc]]	is bordered by	[[Melrose,_Massachusetts m.0t_5b]]	hasNeighbor
[[New_England_town m.099ltc]]	is bordered by	[[Middlesex_County,_Massachusetts m.0k3k1]]	hasNeighbor
[[New_England_town m.099ltc]]	is bordered by	[[Revere,_Massachusetts m.0t_xb]]	hasNeighbor
[[New_England_town m.099ltc]]	is bordered by	[[Lynnfield,_Massachusetts m.01m27v]]	hasNeighbor
[[New_England_town m.099ltc]]	is bordered by	[[Lynn,_Massachusetts m.0tz1x]]	hasNeighbor

Frase iniziale:

Hans\_Bronsart\_von\_Schellendorff m.01w3kr 7 [[Hans\_Bronsart\_von\_Schellendorff|m.01w3kr]] went to [[Weimar|m.0cpyv]] in 1853 where [[Hans\_Bronsart\_von\_Schellendorff|m.01w3kr]] met [[Franz\_Liszt|m.0hqgp]] and became familiar with all the musicians in [[Franz\_Liszt|m.0hqgp]] 's circle at the time , including [[Hector\_Berlioz|m.0f00f]] and [[Johannes\_Brahms|m.0459z]] .

Fatti estratti:

[[Hans_Bronsart_von_Schellendorff m.01w3kr]]	became familiar with	[[Hector_Berlioz m.0f00f]]	hasAcademicAdvisor
[[Hans_Bronsart_von_Schellendorff m.01w3kr]]	became familiar with	[[Johannes_Brahms m.0459z]]	hasAcademicAdvisor