



UNIVERSITÀ
DEGLI STUDI
DI PADOVA



DIPARTIMENTO DI INGEGNERIA
DELL'INFORMAZIONE

CORSO DI LAUREA IN INGEGNERIA INFORMATICA

**”Migrazione e standardizzazione di
terminologia multilingue da file Excel
a database relazionali: analisi dei dati
di Terminology Without Borders”**

Relatore: Prof. Giorgio di Nunzio Maria

Laureando: Yihui Zhu

ANNO ACCADEMICO 2022-2023

Data di laurea 25-09-2023

Sommario

Per un progetto come Terminology Without Borders, che si propone di immagazzinare e fornire crescenti quantità di dati strutturati come glossari multilingue, è vantaggioso disporre di un sistema di gestione dati efficiente. Attualmente le organizzazioni aderenti al progetto forniscono nuove risorse sottoforma di fogli di calcolo. Un database relazionale, in questo contesto, potrebbe portare numerosi vantaggi quali eliminare i duplicati, fornire una fonte centralizzata di dati, maggiore scalabilità, garantire consistenza e integrità dei dati e permettere interrogazioni complesse con il linguaggio SQL. Il presente elaborato si propone di illustrare un modo per migrare i dati da fogli di calcoli con schemi eterogenei esistenti verso un database relazionale in maniera automatizzata, applicato ad un sottoprogetto di Yourterm. A causa della struttura disomogenea tra vari fogli di uno stesso file e tra file di progetti diversi, è necessaria una analisi dei campi esistenti per creare uno schema onnicomprensivo, standardizzato, coerente e non ridondante. Infine, l'approccio di questo lavoro potrebbe essere riutilizzato per essere applicato ad altri progetti Yourterm, che presentano caratteristiche e problematiche simili.

Indice

1	Introduzione	3
1.1	Vantaggi dei DBMS rispetto ai fogli di calcolo per gestione database	3
1.2	Struttura esistente disomogenea	4
1.3	Obiettivi	8
2	Design database	9
2.1	Progettazione concettuale	10
2.2	Progettazione logica	12
2.3	Progettazione fisica	14
3	Migrazione da fogli di calcolo a database	17
3.1	Esportazione documenti Excel in formato CSV	17
3.2	Parsing	18
4	Risultati	22
4.1	Correttezza migrazione	22
4.2	Esempi di interrogazioni	22
5	Conclusioni	27
6	Appendice	29

1 Introduzione

Il progetto Terminology Without Borders del Parlamento europeo collabora con università e partner individuali esterni per raccogliere risorse terminologiche in modo da creare glossari e database multilingue che possano essere usati per comunicare con i cittadini in un linguaggio chiaro. I risultati ottenuti vengono inoltre usati per arricchire le risorse del sistema IATE con nuovi termini o aggiungere informazioni dettagliate a termini esistenti, successivamente alla validazione da parte di esperti di terminologia delle Istituzioni Europee. I domini che esso copre sono: medicina, ambiente e conservazione, cultura, diritti delle donne, giustizia e legalità, pesca e marittima, educazione, industria del cibo e gastronomia, IT e social media, economia e finanza, rispettivamente rappresentati da 10 progetti. Ciascuno di essi è composto da sotto-progetti. Le risorse sono ospitate sul sito di Yourterm yourterm.eu. Ogni sottoprogetto consiste in un insieme di pagine di fogli di calcolo compilati da vari collaboratori che sono stati accorpati in un unico documento Excel o Google Sheets. L'obiettivo dell'elaborato è mostrare un modo per migrare i dati contenuti in questi documenti verso un database gestito da un DBMS relazionale.

1.1 Vantaggi dei DBMS rispetto ai fogli di calcolo per gestione database

I fogli di calcolo sono utilizzati in statistica, contabilità, finanza ed in innumerevoli altri ambiti per l'analisi di dati o semplicemente per la conservazione di registri in formato tabulare. I programmi disponibili, per esempio Microsoft Excel e Google Sheets, presentano interfacce utenti grafiche intuitive che offrono ad utenti anche non-programmatori funzionalità avanzate come

la formattazione di celle, creazione di grafici, realizzazione di report, calcolo con formule.

I software per fogli elettronici hanno inoltre la capacità di gestire basi di dati relazionali [1]. In Google Sheets, per esempio, è possibile effettuare interrogazioni stile SQL con la funzione *QUERY* [2], simulare *join* di tabelle con la funzione *VLOOKUP* [3], si possono assegnare permessi agli utenti, forzare l'integrità dei dati con lo strumento *Data validation* e altre formule per verificare l'unicità degli attributi *UNIQUE* e *PRIMARY KEY* e garantire l'integrità referenziale per le chiavi esterne. Nonostante ciò non tutte le funzionalità e caratteristiche di un DBMS relazionale sono presenti e spesso molte implementazioni sono forzature inutilmente complesse rispetto ad una soluzione ad-hoc, d'altronde Microsoft stessa offre prodotti dedicati come Access o SQL Server.

Sebbene soluzioni come *Google Apps Script* permettano di creare API con cui realizzare applicazioni che si interfacciano a Google Sheets come data storage, quest'ultimo non è scalabile in quanto ammette fino ad un massimo di 10 milioni di celle [4], inoltre le prestazioni sono inferiori rispetto ad un DBMS come MySQL [5]. Per di più, affidarsi a soluzioni proprietarie così particolari rende difficoltosa la portabilità del database.

1.2 Struttura esistente disomogenea

I fogli elettronici dei vari sotto-progetti presentano una struttura di base comune, contengono sempre due tipi di pagine. Per esempio Climate Change.xlsx è composta da:

- una pagina chiamata *All languages* contenente, per ogni riga della tabella, i termini equivalenti di più lingue, come in tabella 1

EN	RO	ES	...
Oxidation	oxidare	Oxidación	...

Tabella 1: Esempio di record della tabella All languages in Climate Change.xlsx

- una pagina per ogni lingua, in cui ciascuna riga rappresenta un termine con le relative informazioni. In tabella 2 un esempio di record della pagina “IT”

colonna	valore
Source term	benzene
Target term*	benzolo
Term Reference*	https://www.treccani.it/...
Language ID	IT
Notes (Term)	
IATE ID	1566036
Definition*	liquido mobile d’odore etereo...
...	...

Tabella 2: Esempio di record della tabella IT in Climate Change.xlsx

Le regole per la compilazione delle tabelle e i significati delle colonne sono definite nelle linee guida IATE [6], inoltre esiste un documento Excel di riferimento con le istruzioni e contenente una struttura iniziale, di seguito riportati. Nota: i campi con l’asterisco(*) sono obbligatori.

- Source term*: insert the source term
- Language ID*: insert the two-letter code of the ISO 639-1 of the language of the target term, e.g. FR

- Target term*: insert the target term
- Part of speech*: noun, verb, adjective... for the correct value to enter (for each language), refer to the file named "pos_per_language".
- Note (Term): any relevant comments on the target term
- IATE: if available, the IATE ID corresponding to the term
- Term reference *: source in which the target term has been found
- Definition *: short and precise definition of the term, written in the language of the target term
- External cross reference (definition): web link to the source in which the definition has been found
- Source (definition)*: bibliographical citation (title/author/etc.) of the definition source
- Note (definition): any relevant information about the definition
- Synonym: two terms are synonyms if they share the same meaning and the same usage - they have to always be replaceable
- Quasi synonym: if two terms share the same meaning but not the same usage and therefore are not always replaceable, they are quasi-synonyms
- Common name: if the term has a common name, that is a popular variant (e.g. fever for pyrexia).
- Scientific name: if the term has a scientific correspondent (e.g. pyrexia for fever).
- Orthographic variant: any orthographic variant
- Acronym: any acronym

- Abbreviation: any abbreviation
- Domain*: domain of the term (e.g. "medicine") Note: on the FAIR-term WebApp this field is named with the ISO standard label as "subject field"
- Subdomain: specific sub-domain (e.g. oncology)
- Context*: sentence or short paragraph containing the term (it must be different from the definition, but it can be the same as the term reference)
- External cross reference (context): web link to the source in which the context has been found
- Source (context)*: bibliographical citation (title/author/etc.) of the context source
- Project*: name of the Your Term project (e.g. YourTerm MED)
- Subproject*: name of the Your Term subproject (e.g. "bipolar disorders")

Purtroppo tra i vari fogli di sotto-progetti diversi e tra le stesse pagine di uno stesso documento esistono delle inconsistenze. Queste devono essere risolte ridefinendo uno schema ex novo. Inoltre molte colonne sono state compilate in maniera tale da richiedere la loro rielaborazione prima di importarli in un database. I problemi sono:

- schemi delle tabelle diversi: alcune pagine sono prive di certe colonne, altre hanno colonne aggiuntive non presenti nelle istruzioni. Spesso le colonne hanno ordini differenti e addirittura nomi diversi
- assenza di valori in campi obbligatori. A volte l'assenza di valore viene specificata con il simbolo "-"

- i dati sono stati inseriti sotto colonne sbagliate, a volte sistematicamente, a volte sono errori occasionali
- non esatta corrispondenza dei termini in “All languages” con gli equivalenti nelle altre pagine a causa di errori di compilazione, per esempio la capitalizzazione è diversa
- i campi come “Synonym”, “Quasi synonym”, “Abbreviation” e altri possono presentare più valori in una cella, ma sono separati con delimitatori diversi: a volte con il forward slash, a volte con la virgola, altre volte sono enumerati con cifre seguite dal punto
- Talvolta i campi “Target term” e le colonne delle pagine “All languages” presentano più valori che sono separati dai delimitatori in maniera inconsistente

In generale le celle mancano di formule per la verifica dei formati di dato inseriti. Creare un database con schemi uguali a quelli presenti nei fogli di calcolo è quindi impensabile. D'altronde anche considerando le intestazioni delle tabelle Excel come relazioni, queste presentano numerose dipendenze funzionali, perciò è necessario definire una nuova struttura che sia standard e più normalizzata.

1.3 Obiettivi

Il risultato finale della tesi è una base di dati relazionale gestita dal DBMS PostgreSQL adatta per un'applicazione multi utente che deve gestire operazioni di tipo inserimento, modifica ed eliminazione di termini, gestione corrispondenza di stessi termini in più lingue e occasionalmente stampa ed esportazione dei record, per esempio per importarli nel database IATE.

Il lavoro consiste in una fase di progettazione del database, in cui verranno illustrati gli schemi concettuale e relazionali ottenuti e le motivazioni dietro alle scelte progettuali, seguita da una fase di migrazione vera e propria dove si mostreranno le tecniche e gli script usati per estrarre dati dai fogli di calcolo ed inserirli nel database.

Il database e gli script sono stati usando i file appartenenti ai sottoprogetti “Climate change” e “Zoology” del progetto “ENVI” disponibili al link.

2 Design database

Il primo passo della migrazione è stato la progettazione della base di dati. In questa ricerca è stato adottato un metodo di progettazione graduale e metodico, con le seguenti fasi:

1. Nella fase di progettazione concettuale è stato realizzato un diagramma ER analizzando le strutture tabellari esistenti, per identificare tutte le informazioni necessarie da modellare
2. Lo schema concettuale ottenuto è stato tradotto in uno schema logico, su cui sono state apportate scelte di progetto atte a rendere il database più efficiente in termini di spazio occupato e di velocità di esecuzione di potenziali operazioni frequenti. In assenza di particolari requisiti, è stato deciso di definire degli schemi relazionali altamente normalizzati per garantire qualità al design
3. Nella progettazione fisica è stato prodotto il codice SQL per implementare gli schemi precedenti nel DBMS scelto, PostgreSQL

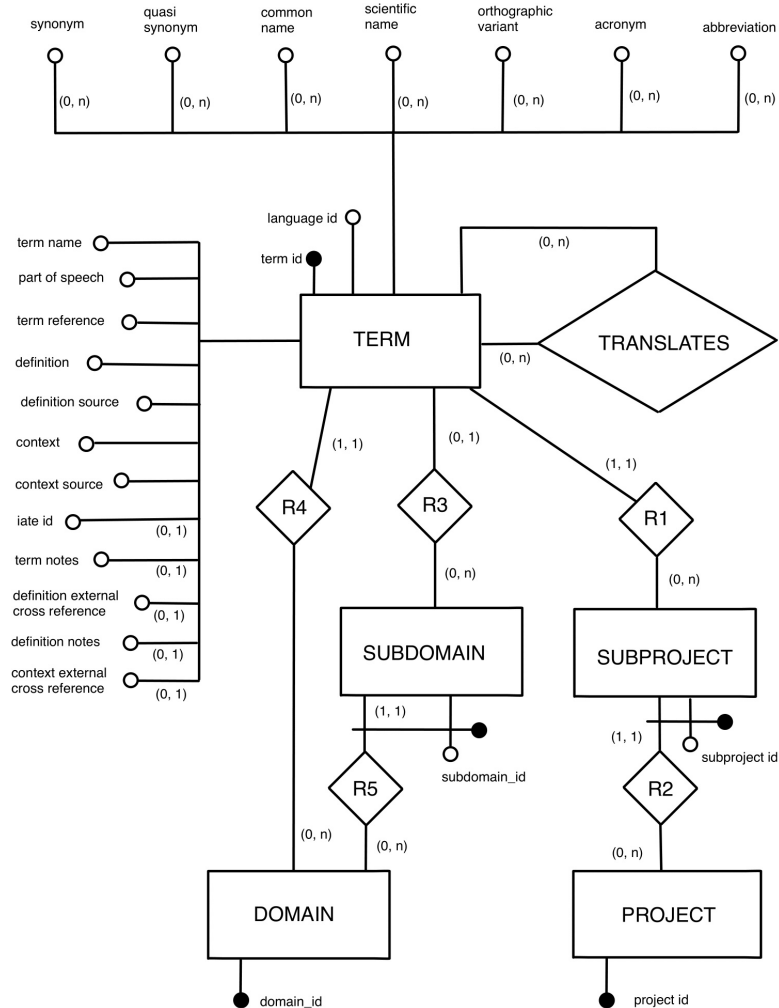


Figura 1: Diagramma ER del database

2.1 Progettazione concettuale

Prendendo come riferimento l'intestazione della tabella di esempio presente in Yourterm toolbox e descritta nella sezione 1.2, si sono rilevate le informazioni da modellare ed è stato ricavato il diagramma ER in fig 1. Eventuali campi aggiuntivi non standard introdotti in certe tabelle sono ignorati (come *Transaction*, *Responsibility*, *Name/Surname*, *Date* e *Animal*), perché sono dati non rilevanti per il progetto aggiunti su iniziativa delle singole univer-

sità e che non sono stati rimossi durante la fase di fusione dei file prima del caricamento sul sito `yourterm`. Per rendere la nomenclatura più chiara certe espressioni originali sono state rinominate riordinando le parole, sono stati rimossi caratteri speciali come asterischi e parentesi e i tutti caratteri sono stati convertiti in minuscolo: per esempio “External cross reference(definition)” è diventato “definition external cross reference”.

L’entità centrale è *TERM*, che rappresenta un termine, ovvero una riga della tabella del foglio di calcolo. Sebbene un termine possa essere identificato univocamente dalla combinazione degli attributi *term name*(*Target term* in origine) e *language id*, è stato deciso di introdurre un identificativo artificiale, *term id*, sottoforma di sequenza di numeri interi. L’approccio con chiave naturale non prevede che un termine possa avere più significati, almeno senza accumularli tutti in un’unica tupla, come in un dizionario, soluzione poco flessibile e che non rispetta la prima forma normale. La soluzione con identificativo artificiale, invece, permette l’esistenza di termini con lo stesso nome che rappresentano più concetti in contesti diversi, ed è anche l’approccio adottato da IATE [7], oltre ad essere potenzialmente più performante nelle operazioni di join in quanto viene usato solo un campo a numeri interi come chiave. I campi contrassegnati dall’asterisco nel file Excel di esempio denotano campi a compilazione obbligatoria, e sono stati resi attributi a cardinalità obbligatoria, e i rimanenti attributi sono stati resi facoltativi.

I campi *Subproject* e *Project* sono stati trasformati in entità separate, *SUBPROJECT* e *PROJECT*, per esprimere la dipendenza del primo dal secondo, cioè un *SUBPROJECT* appartiene ad uno e un solo *PROJECT*; attraverso la partecipazione di cardinalità uno a uno di *TERM* con la associazione *R1* e di *SUBPROJECT* con *R2* si può ottenere a quale *PROJECT*

appartiene un'istanza di TERM. Inoltre sia *PROJECT* che *SUBPROJECT* sono entità con un limitato numero di valori predefiniti ammissibili.

Anche campi *Domain* e *Subdomain* sono state rese entità separate, *DOMAIN* e *SUBDOMAIN*. Nelle regole di compilazione della tabella di riferimento, però, *Domain* è a compilazione obbligatoria mentre *Subdomain* non lo è, quindi è necessaria una ulteriore associazione tra *TERM* e *DOMAIN*.

La associazione *TRANSLATION* esprime le corrispondenze tra termini equivalenti in lingue diverse, specificate nelle tabelle *All languages*.

2.2 Progettazione logica

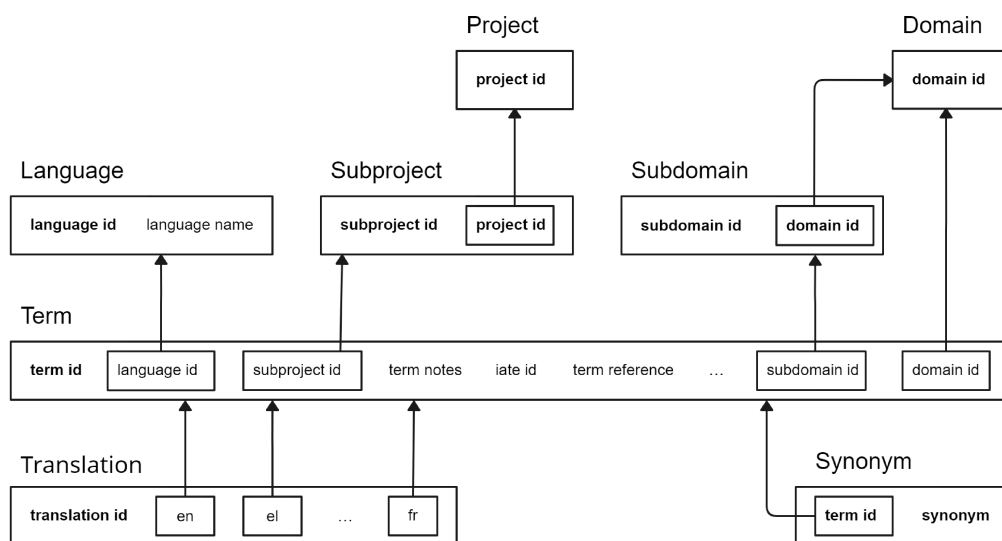


Figura 2: Schema relazionale del database. Le relazioni Acronym, Quasi synonym, Common name, Scientific name, Orthographic variant e Abbreviation sono simili a Synonym, non sono stati inseriti per motivi di spazio.

Traducendo il diagramma ottenuto nella fase precedente è stato ricavato lo schema relazionale in fig 2. Gli attributi multi-valore(*synonym*, *quasi sy-*

nonym, common name, scientific name, orthographic variant, acronym e abbreviation) sono stati tradotti in relazioni con vincoli di integrità referenziale con la relazione Term.

L'associazione ricorsiva *TRANSLATE* è stata trasformata nella relazione *Translation(translation id, en, el, ..., fr)* che ha come componenti, oltre ad una chiave surrogata *translation id*, una chiave esterna per ogni lingua con vincolo di integrità referenziale con Term.

Una variante più naturale e intuitiva sarebbe potuta essere per esempio *Translation(source term id, target term id)*, dove *source term id* e *target term id* sono coppie di termini equivalenti in lingue diverse. Purtroppo tale soluzione risulterebbe piuttosto dispendiosa in termini di storage in quanto presenta ridondanze: il numero di tuple totali per esprimere una traduzione di un termine richiede è pari al numero di disposizioni di n lingue a 2 a 2. Quindi per ogni traduzione, assumendo che il termine sia tradotto in tutte n le lingue disponibile, servirebbero $t = n \cdot (n - 1)$ tuple. Nei progetti sono state finora individuate 11 lingue, perciò potenzialmente sarebbero 110 tuple per traduzione.

La soluzione adottata, invece, è più efficiente in termini di spazio occupato in quanto richiede solo una tupla con $n + 1$ attributi per ogni traduzione, tranne i casi in cui un campo è composto da più valori separati da delimitatore. In tali situazioni si crea un numero di tuple pari al numero di combinazioni del prodotto cartesiano tra insiemi di termini. La chiave surrogata è necessaria in quanto i restanti campi sono frequentemente nulli a causa dell'assenza di corrispondenze in certe lingue.

Un compromesso sarebbe riservare il campo *source term id* per i soli termini in inglese, dato che nei documenti Excel ogni termine ha sempre il suo corrispondente in inglese. Si avrebbe potenzialmente un grande risparmio

di memoria evitando il grande numero di combinazioni tra ogni coppia di lingue o i numerosi valori nulli. Purtroppo questo approccio è poco flessibile in quanto implica l'obbligo della presenza di un termine in inglese prima di poter stabilire traduzioni.

L'attributo *language id* di *TERM* è stato reso una relazione a se stante. La motivazione è dovuta alla volontà di rendere possibile l'aggiunta di ulteriori attributi non esplicitamente richiesti o assolutamente necessari ma che potrebbero risultare utili, per esempio la denominazione per esteso della lingua *language name* potrebbe essere comoda in fase di realizzazione di un menù di selezione drop-down per la scelta del linguaggio in una ipotetica applicazione.

Lo schema ottenuto è quasi in terza forma normale: presenta solo una dipendenza funzionale, *subdomain id* \rightarrow *domain id*, non normalizzabile a causa del fatto che *subdomain_id* potrebbe essere ignoto per le caratteristiche del problema.

2.3 Progettazione fisica

È stato adottato il naming convention *snake_case*. le chiavi primarie hanno lo stesso nome delle tabelle con suffisso *_id* e campi che fungono da chiavi esterne sono generalmente denominati con i nomi degli attributi a cui si riferiscono.

A causa dell'imprevedibilità della lunghezza in caratteri della maggioranza dei campi, è stato scelto il tipo di dato *TEXT* per tutti gli attributi testuali, ad eccezione di quelli di lunghezza consistente come *project_id*, o quelli su cui è di interesse creare indici, come *term_name*, per cui sono stati scelti tipi *CHAR* o *VARCHAR*.

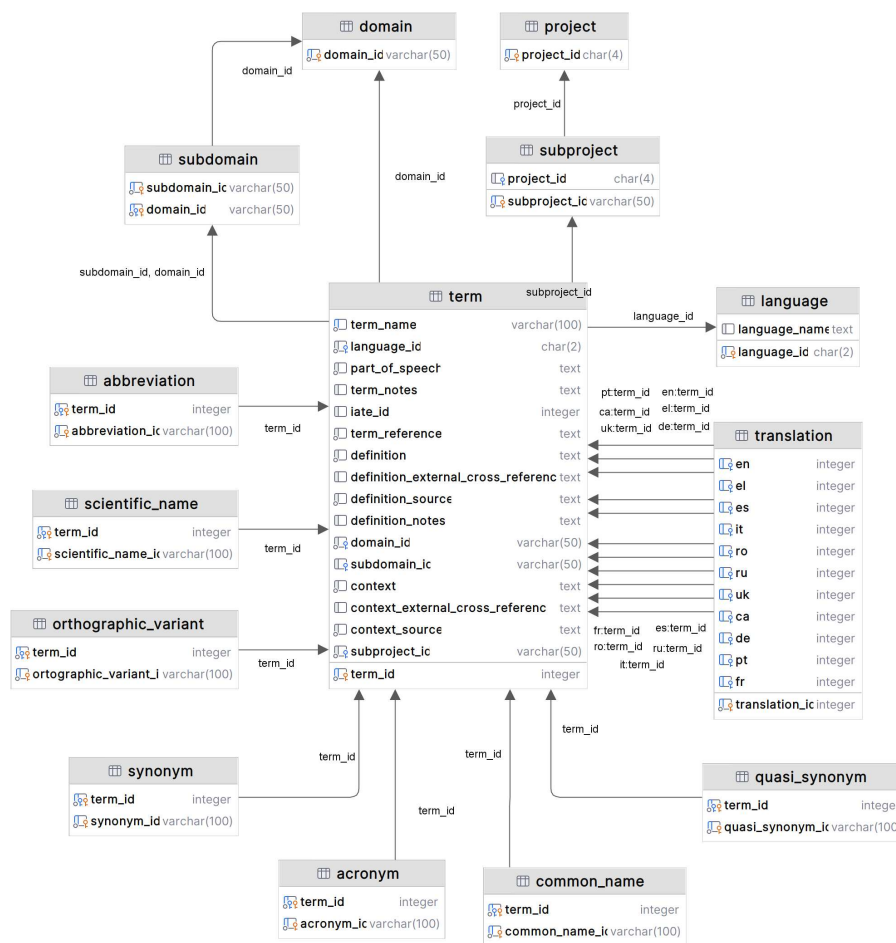


Figura 3: Schema fisico del database.

È stato prodotto il codice SQL per implementare gli schemi del database progettato (9).

Per migliorare le prestazioni delle operazioni di join tra *term* e *translation* sono stati creati indici sulle tutte le chiavi esterne della tabella *translation*. Prevedendo frequenti ricerche sui nomi dei termini è stato inoltre creato un indice sulla colonna *term_name* della tabella *term*.

Prima di effettuare la migrazione è necessario inserire manualmente le tuple relative ai linguaggi, progetti, sotto-progetti, domini e sotto-domini.

La procedura per estrarre rapidamente tali valori univoci usando Excel è:

1. selezionare il menù *Dati*
2. dalla sezione *Ordinazione* cliccare *Avanzate*
3. selezionare la colonna di interesse
4. spuntare il check-box *Copia univoca dei record*

Per imporre l'inserimento corretto di codici lingua definiti dallo standard ISO 639-1 è stato applicato il vincolo di integrità referenziale tra il campo *language_id* di *term* e la tabella *language*, in questo modo solo i valori presenti in quest'ultima sono ammessi.

Alternativamente si sarebbe potuto usare un vincolo di controllo sulla colonna *language_id*:

```
1 CREATE TABLE term(  
2     language_id CHAR(2) NOT NULL,  
3     ...  
4     CONSTRAINT ADMITTED_LANGUAGES CHECK(language_id IN  
5     ('en', 'el', 'es', 'it', 'ro', 'ru', 'uk', 'ca', 'de', 'pt', 'fr'))  
6 );
```

Un'ulteriore soluzione sarebbe stata la definizione di un *DOMAIN*, ma il vincolo viene usato solo nella tabella *term* e perciò non è da preferire ad un vincolo di controllo su colonna.

```
1 CREATE DOMAIN language AS CHAR(2)  
2 CHECK(VALUE IN('en', 'el', 'es', 'it', 'ro', 'ru', 'uk', 'ca', 'de', 'pt', '  
fr'));
```

Con la soluzione adottata risulta più immediato aggiungere nuove lingue e la modifica alle lingue esistenti si propaga in modo automatico sui record della tabella “term”, non è necessario infatti eliminare il vecchio *CONSTRAINT/DOMAIN*, crearne un nuovo ed aggiornare manualmente i valori.

3 Migrazione da fogli di calcolo a database

Dopo che il database è stato creato, sono stati trasformati i fogli di calcolo in file interpretabili da uno script di parsing, in questo caso è stato scelto il formato CSV, anche perché facilmente esportabile da Excel o Google Sheets. La codifica dei caratteri è UTF-8, considerato che si stanno importando testi scritti in lingue europee. I due programmi dedicati al parsing e all'inserimento in database sono stati scritti in modo tale da elaborare i documenti CSV ottenuti e risolvere i problemi elencati in sezione 1.2.

3.1 Esportazione documenti Excel in formato CSV

Sono state scaricate le copie dei file *Climate Change.xlsx* e *Zoology.xlsx* da ENVI projects. Usando il software Microsoft Excel per manipolare i documenti scaricati. Sono state eseguite le seguenti modifiche:

1. Si sono eliminate le colonne vuote delle tabelle insieme a quelle non appartenenti allo schema logico
2. Sono state eliminate le righe finali vuote inutili
3. Se sono stati presenti errori di compilazione sistematici, come per esempio nella pagina *DE* di *Zoology*, si è intervenuti manualmente, in questo caso spostando i valori della colonna *Orthographic variant* in *Part of speech*, dove si dovrebbero trovare correttamente
4. Le singole pagine di un file sono state esportate in formato CSV UTF-8
5. le colonne delle tabelle sono state rinominate con i nomi definiti nello schema fisico. In particolare i campi *Synonym*, *Quasi synonym*,

Common name, Scientific name, Orthographic variant, Acronym e Abbreviation sono stati rinominati con i nomi delle loro rispettive tabelle, mentre i restanti attributi seguono i nomi degli attributi della tabella *term*. Le pagine *All languages* seguono la denominazione della tabella *translation*.

I file risultanti sono puliti da campi o righe inutili ed condividono intestazioni con nomi standard. Il campo *Source term* non viene mai utilizzato perché l'informazione è già presente nella pagina *All languages*, quindi viene ignorato.

3.2 Parsing

Il parsing dei file ottenuti viene effettuato da due script. Uno si occupa di analizzare le pagine “*All languages*” contenente la traduzioni, l'altro si occupa di tutte le restanti pagine contenenti i termini. Gli script completi sono stati scritti in Python e sono disponibili in appendice 6. Effettuano sia il parsing dei file CSV usando la libreria Pandas che l'inserimento nel database usando l'adapter Python Psycopg2 per PostgreSQL.

```
1 CLASS Term(term_name, language_id, part_of_speech, ...)
2 -- synonym, quasi_synonym, abbreviation, ...
3 CLASS Variant(term_name, variant)
4
5 PARSE_TERMS(filename, language, subproject, domain)
6 dataframe <- read_csv(filename)
7
8 -- output
9 terms <- []
10 variations <- {'synonym': [], 'quasi_synonym': [], ...}
11
12 FOR entry IN dataframe DO
13     -- creazione di un termine per ogni riga di dataframe e assegnargli i
14     suoi valori
```

```

14 term <- create_empty_term()
15 FOR field IN fields(term) DO
16     IF field IN {'language_id', 'subproject_id', 'domain_id'} SKIP
17     value <- entry[field]
18     value <- trim(value)
19     IF exists(value) THEN
20         IF field == 'iate_id' THEN term.iate_id <- to_integer(value)
21         ELSE IF field == 'subdomain_id' THEN term.subdomain_id <-
lowercase(value)
22         ELSE term.field <- value
23
24     -- se record vuoto non inserirlo in database
25     IF is_empty(term) THEN CONTINUE
26
27     -- valori costanti assegnati da argomenti e non da file
28     term.language_id <- language
29     term.subproject_id <- subproject
30     term.domain_id <- domain
31
32     -- term name potrebbe avere valori separati
33     new_terms <- split_and_trim(term.term_name, ';', ',', '/')
34     new_terms <- lowercase(new_terms)
35     terms <- [terms, new_term]
36
37     -- associazione di term alle sue varianti
38     FOR key IN keys(variant) DO
39         value <- entry[key]
40         IF exists(value) THEN
41             FOR term IN new_terms DO
42                 FOR variant IN split_and_trim(value, ',', '/', '[0-9].') DO
43                     variants[key] <- [variants[key], Variant(term, variant)]
44
45 insert_into_db(terms)
46 insert_into_db(variations)

```

Listing 1: Pseudocodice script di parsing file di termini

```

1 CLASS Translation(en, el, ... fr)
2
3 PARSE_TRANSLATIONS(filename)
4 dataframe <- read_csv(filename)
5

```

```

6 -- output
7 translations = []
8
9 FOR entry IN dataframe DO
10     translation <- create_empty_translation()
11     FOR field IN fields(translation) DO
12         value <- entry[field]
13         IF exists(value) THEN
14             value <- trim(value)
15             value <- lowercase(value)
16             translation.field <- value
17
18 -- se i campi contengono valori multipli
19 -- sets contiene una lista di liste
20 sets <- []
21 FOR field IN fields(translation) DO
22     set <- split_and_trim(translation.field, ';', ',', '/')
23     sets <- [sets, set]
24 combinations <- cartesian_product(sets)
25
26 -- creare una traduzione per ogni combinazione di valori
27 FOR combination IN combinations DO
28     translation <- create_translation(combination)
29     translations <- [translations, translation]
30
31 insert_into_db(translations)

```

Listing 2: Pseudocodice script di parsing file di traduzioni

L'esecuzione dello script per i termini richiede come parametri da riga di comando, oltre al file da analizzare, il codice del linguaggio, il nome del sottoprogetto e il nome del dominio a cui appartengono i termini del file, che vengono assegnati alle nuove tuple ignorando valori presenti all'interno. Questo per garantire l'inserimento esatto di *language_id* nel formato ISO 639-1 e degli stessi *subproject_id* e *domain_id* presenti già nelle tabelle *subproject* e *domain*; inoltre sono attributi obbligatori costanti in uno stesso file che potrebbero non essere stati compilati, oppure compilati non uniformemente

causando violazioni di vincoli. La colonna *subdomain_id* spesso contiene gli stessi valori ma con capitalizzazione inconsistente, perciò sono stati resi tutti in minuscolo per renderli compatibili con i valori univoci inseriti nella tabella *subdomain*.

```
1 python termsParser.py csv/ZoologyEN.csv en "Zoology" Environment
```

Listing 3: Esempio di comando per eseguire lo script

Frequentemente i campi obbligatori non sono stati compilati, spesso perché i file CSV mancano interamente di alcune colonne; in questi casi agli attributi *NOT NULL* viene assegnata una stringa vuota, mentre agli attributi senza tale vincolo viene assegnato il valore *NULL* per indicare l'assenza di valori. Per rispettare le linee guida servirebbe l'intervento manuale per la compilazione dei campi mancanti.

A causa dell'inconsistente uso delle lettere maiuscole e delle minuscole nel campo *Target term*, è stato deciso di convertirne tutte le occorrenze in minuscolo. In questo modo si sono evitati i casi in cui un non si riusciva a stabilire la traduzione di un termine perché il valore nella tabella *All languages* di una riga è diverso da quello in *Target term* in termini di capitalizzazione. Talvolta però, questo è dovuto a errori di battitura che sono risolvibili attraverso l'intervento manuale e sono stati ignorati dallo script.

Quando il campo *Target term* presenta più valori separati da un delimitatore come “;”, “,” o “/”, per esempio in numerose righe della pagina *ES* di *Zoology.xlsx*, vengono create più tuple tutte con gli stessi valori.

4 Risultati

4.1 Correttezza migrazione

Completata la migrazione dei due sotto-progetti analizzati, la tabella *term* contiene 1681 righe di cui 72, contati dallo script, sono valori aggiuntivi che nei documenti originali condividono una stessa cella con altri da cui sono separati da delimitatori. Utilizzando la formula *CONTA.VALORI* di Excel è stato contato un numero di righe totale originale pari a 1609, perciò tutte le occorrenze sono state trasferite.

La tabella *translation* contiene 643 righe di cui 123, contati dallo script, sono gli elementi aggiuntivi dovuti alle combinazioni di traduzioni per i casi in cui i campi presentavano più valori separati da delimitatori. Altre 66 righe sono i duplicati dovuti alla presenza di termini con *term_name* e *language_id* uguali, il cui numero è stato trovato con la query (8). Nei documenti originali le righe sono 454, perciò tutte le traduzioni sono state importate.

4.2 Esempi di interrogazioni

È possibile riottenere la rappresentazione completa di tutti termini con tutti gli attributi originari eseguendo la seguente query

```
1 SELECT term_name AS "term name",
2         language_name AS "language name",
3         part_of_speech AS "part of speech",
4         term_notes AS "note(term)",
5         iate_id AS IATE,
6         term_reference AS "term reference",
7         definition,
8         definition_external_cross_reference AS "external cross reference(
9         definition)",
10        definition_source AS "source(definition)",
11        definition_notes AS "note(definition)",
12        context,
```

```

12     context_external_cross_reference AS "external cross reference(context
13 )",
14     context_source AS "source(context)",
15     synonym_id AS synoynm,
16     quasi_synonym_id AS "quasi synonym",
17     common_name_id AS "common name",
18     scientific_name_id AS "scientific name",
19     ortographic_variant_id,
20     acronym_id,
21     abbreviation_id,
22     domain.domain_id,
23     subdomain.subdomain_id,
24     subproject.subproject_id,
25     project.project_id
26 FROM term
27 LEFT JOIN language ON term.language_id = language.language_id
28 LEFT JOIN synonym ON term.term_id = synonym.term_id
29 LEFT JOIN quasi_synonym ON term.term_id = quasi_synonym.term_id
30 LEFT JOIN common_name ON term.term_id = common_name.term_id
31 LEFT JOIN scientific_name ON term.term_id = scientific_name.term_id
32 LEFT JOIN orthographic_variant ON term.term_id = orthographic_variant.
33 term_id
34 LEFT JOIN acronym ON term.term_id = acronym.term_id
35 LEFT JOIN abbreviation ON term.term_id = abbreviation.term_id
36 LEFT JOIN domain ON term.domain_id = domain.domain_id
37 LEFT JOIN subdomain ON term.subdomain_id = subdomain.domain_id
38 LEFT JOIN subproject ON term.subproject_id = subproject.subproject_id
39 LEFT JOIN project ON subproject.project_id = project.project_id;

```

Le traduzioni si ottengono con

```

1 SELECT en.term_name, el.term_name, es.term_name,
2        it.term_name, ro.term_name, ru.term_name,
3        uk.term_name, ca.term_name, de.term_name,
4        pt.term_name, fr.term_name
5 FROM translation
6 LEFT JOIN term as en on translation.en = en.term_id
7 LEFT JOIN term as el on translation.el = el.term_id
8 LEFT JOIN term as es on translation.es = es.term_id
9 LEFT JOIN term as it on translation.it = it.term_id
10 LEFT JOIN term as ro on translation.ro = ro.term_id
11 LEFT JOIN term as ru on translation.ru = ru.term_id

```



```

12 LEFT JOIN term as uk on translation.uk = uk.term_id
13 LEFT JOIN term as ca on translation.ca = ca.term_id
14 LEFT JOIN term as de on translation.de = de.term_id
15 LEFT JOIN term as pt on translation.pt = pt.term_id
16 LEFT JOIN term as fr on translation.fr = fr.term_id

```

Siccome l'associazione *TRANSLATE* è stata trasformata nella relazione *Translation*(*translation id*, *en*, *el*, ..., *fr*)(1) e non nella più intuitiva *Translation*(*source term id*, *target term id*)(2), le operazioni che coinvolgono la tabella *translation* richiedono più attenzione.

Per esempio, se l'operazione fosse: "Trovare tutte le traduzioni del termine 'xxx' dalla lingua 'yy'", la query potrebbe essere:

```

1 WITH trans AS(
2     SELECT en, el, es, it, ro, ru, uk, ca, de, pt, fr
3     FROM term JOIN translation ON term.term_id = translation.en
4     WHERE term_name = 'younger dryas')
5 SELECT *
6 FROM term
7 WHERE term_id IN (
8     SELECT en FROM trans UNION
9     SELECT el FROM trans UNION
10    SELECT es FROM trans UNION
11    SELECT it FROM trans UNION
12    SELECT ro FROM trans UNION
13    SELECT ru FROM trans UNION
14    SELECT uk FROM trans UNION
15    SELECT ca FROM trans UNION
16    SELECT de FROM trans UNION
17    SELECT pt FROM trans UNION
18    SELECT fr FROM trans);

```

Listing 4: Query 1

La stessa query, se si fosse adottato lo schema (2), appare più immediata:

```

1 SELECT term.*
2 FROM term JOIN translation ON term.term_id = translation.target_term_id
3 WHERE source_term_id IN (
4     SELECT term_id

```

```

5     FROM term
6     WHERE term_name = 'aditiv' AND language_id = 'ro'
7 );

```

Listing 5: Query 2

Per comparazione è stato implementato lo schema (2) nella variante con *source_term_id* che contiene solo termini in inglese. Di conseguenza le query che eseguono la stessa operazione si distinguono per i casi in cui il termine da cercare sia inglese o non.

```

1 WITH source_terms AS(
2     SELECT term_id
3     FROM term
4     WHERE term_name = 'younger dryas' AND language_id = 'en')
5
6 SELECT term.*
7 FROM term JOIN translation ON term.term_id = translation.target_term_id
8 WHERE source_term_id IN (SELECT * FROM source_terms)
9 UNION
10 SELECT term.*
11 FROM term
12 WHERE term_id IN (SELECT * FROM source_terms);

```

Listing 6: Query 3. Seleziona le traduzioni in tutte le lingue di un certo termine inglese

```

1 WITH source_terms AS(
2     SELECT source_term_id
3     FROM term JOIN translation ON term_id = target_term_id
4     WHERE term_name = 'dryas recente' AND language_id = 'it')
5
6 SELECT term.*
7 FROM term JOIN translation ON term.term_id = translation.target_term_id
8 WHERE source_term_id IN (SELECT * FROM source_terms)
9 UNION
10 SELECT term.*
11 FROM term

```

```
12 WHERE term_id IN (SELECT * FROM source_terms);
```

Listing 7: Query 4. Seleziona le traduzioni in tutte le lingue di un certo termine solo se non inglese

Preponendo il comando *EXPLAINANALYZE* alle query si ottengono i query plan. Effettuando la stessa query 20 volte in seguito ad una query iniziale di warmup e mediando i risultati del campo *actual time* si sono ottenute informazioni che forniscono indicativamente indizi sulle prestazioni dei diversi schemi e l'importanza degli indici.

	Query 1	Query 1 senza indici	Query 3
startup cost	17.29	181.30	23.72
total cost	92.48	256.50	23.76
actual startup time[ms]	0.05	0.522	0.081
actual total time[ms]	0.062	0.533	0.085

Tabella 3: Comparativa prestazioni query diverse per stessa operazione

Una grande parte del miglioramento del tempo di esecuzione di Query 1 è dovuto alla presenza dell'indice sulla colonna *term_name*, in quanto il costo maggiore è dovuto alla scansione sequenziale sulla tabella “term” per cercare il nome del termine.

Query 1 (4) necessita di un *index scan* su *term* e di un *nested loop join* in meno rispetto a Query 3 (6), ed è quindi più efficiente. Lo spazio disco occupato, però, è superiore. La tabella con schema (1) occupa 304kB e contiene 643 righe, quella con schema (2) occupa 232kB con 1195 record. Il divario aumenta nei casi in cui la maggior parte delle traduzioni coinvolga un numero ridotto di lingue.

5 Conclusioni

Nel seguente elaborato è stato illustrato il processo che ha portato alla definizione di uno schema database adatto a rappresentare le terminologie con le relative traduzioni del progetto Yourterm e la creazione di procedimenti e script per la migrazione automatizzata dei dati preesistenti da fogli di calcolo a un database relazionale gestito da DBMS.

È stato seguito un procedimento di progettazione metodica del database, partendo dallo definire uno schema concettuale, a sua volta tradotto in uno schema logico, fino ad arrivare alla implementazione fisica in linguaggio SQL, adottando scelte progettuali e implementazioni che rendessero il database altamente normalizzato e performante.

Nonostante i fogli di calcolo su cui si sono svolte le analisi fossero compilati in maniera non sempre rigorosa e senza particolari standard, lo script di parsing è stato in grado di estrarre con successo, attraverso elaborazioni, la maggior parte dei dati per rappresentarli in maniera significativa nel database. Purtroppo, i valori del campo *Part of speech* presentano una eccessiva eterogeneità in termini di capitalizzazione, linguaggio, delimitatori e varietà tale che non si è riusciti a importarli in maniera strutturata. Ci si è limitati dunque a copiare i valori sottoforma di testo in un attributo di tipo *TEXT* nella tabella *term*. Sarebbe necessario intervenire manualmente per aggiornare i record con dei valori standard.

Sebbene la progettazione del database e dei due programmi di estrapolazione si sia focalizzata sull'elaborazione di solo due sotto-progetti, *Zoology* e *Climate Change*, del progetto ENVI di Yourterm, il materiale prodotto durante questa ricerca può essere base di ulteriori modifiche e affinamenti futuri per trasferire i restanti progetti, d'altronde tutti i documenti Excel presenti sul sito condividono una struttura comune. In tal caso sarà possibi-

le utilizzare il database come back-end per una applicazione di gestione per le terminologie più centralizzata ed efficace del sistema attuale, e gli script di parsing potranno comunque essere utili per migrare altre collezioni di termini se si volesse continuare ad usare fogli di calcolo come metodo di inserimento dati.

Riferimenti bibliografici

- [1] J. Tyszkiewicz, “Spreadsheet as a relational database engine,” in *Proceedings of the 2010 ACM SIGMOD International Conference on Management of data*, pp. 195–206, 2010.
- [2] Google, “Query function - google docs editors help,” 2023. Last accessed 19 September 2023.
- [3] L. Inc., “Workaround for using sql join on google sheets,” 2023. Last accessed 19 September 2023.
- [4] Google, “Files you can store in google drive - google drive help,” 2023. Last accessed 19 September 2023.
- [5] L. J. Ekanayake, D. Ihalage, and S. P. Abyesundara, “Performance evaluation of google spreadsheet over rdbms through cloud scripting algorithms,” in *2021 International Conference on Computer Communication and Informatics (ICCCI)*, pp. 1–7, IEEE, 2021.
- [6] T. C. Unit, “Toolbox,” 2023. Last accessed 19 September 2023.
- [7] IATE, “Iate faq,” 2023. Last build date 22 August 2023 09:39.

6 Appendice

```
1 WITH duplicates AS
2   (SELECT term_name, language_id, COUNT(*) AS duplicate_count
3   FROM term
4   GROUP BY term_name, language_id
5   HAVING COUNT(*) > 1)
6
7 SELECT SUM(num) AS "total duplicates", COUNT(*) AS "duplicated rows" FROM
8   (SELECT en.term_name, el.term_name, es.term_name,
9         it.term_name, ro.term_name, ru.term_name,
10        uk.term_name, ca.term_name, de.term_name,
11        pt.term_name, fr.term_name, count(*) AS num
12   FROM translation
13     LEFT JOIN term AS en ON translation.en = en.term_id
14     LEFT JOIN term AS el ON translation.el = el.term_id
15     LEFT JOIN term AS es ON translation.es = es.term_id
16     LEFT JOIN term AS it ON translation.it = it.term_id
17     LEFT JOIN term AS ro ON translation.ro = ro.term_id
18     LEFT JOIN term AS ru ON translation.ru = ru.term_id
19     LEFT JOIN term AS uk ON translation.uk = uk.term_id
20     LEFT JOIN term AS ca ON translation.ca = ca.term_id
21     LEFT JOIN term AS de ON translation.de = de.term_id
22     LEFT JOIN term AS pt ON translation.pt = pt.term_id
23     LEFT JOIN term AS fr ON translation.fr = fr.term_id
24   WHERE
25     (en.term_name, en.language_id) IN (SELECT term_name, language_id
26   FROM duplicates) OR
27     (el.term_name, el.language_id) IN (SELECT term_name, language_id
28   FROM duplicates) OR
29     (es.term_name, es.language_id) IN (SELECT term_name, language_id
30   FROM duplicates) OR
31     (it.term_name, it.language_id) IN (SELECT term_name, language_id
32   FROM duplicates) OR
33     (ro.term_name, ro.language_id) IN (SELECT term_name, language_id
34   FROM duplicates) OR
35     (ru.term_name, ru.language_id) IN (SELECT term_name, language_id
36   FROM duplicates) OR
37     (uk.term_name, uk.language_id) IN (SELECT term_name, language_id
38   FROM duplicates) OR
39     (ca.term_name, ca.language_id) IN (SELECT term_name, language_id
```

```

FROM duplicates) OR
33     (de.term_name, de.language_id) IN (SELECT term_name, language_id
FROM duplicates) OR
34     (pt.term_name, pt.language_id) IN (SELECT term_name, language_id
FROM duplicates) OR
35     (fr.term_name, fr.language_id) IN (SELECT term_name, language_id
FROM duplicates)
36 GROUP BY (en.term_name, el.term_name, es.term_name,
37           it.term_name, ro.term_name, ru.term_name,
38           uk.term_name, ca.term_name, de.term_name,
39           pt.term_name, fr.term_name)
40 ) AS duplicated_rows

```

Listing 8: Query per conteggio di traduzioni duplicate

```

1 CREATE TABLE language(
2     language_id CHAR(2) PRIMARY KEY,
3     language_name TEXT
4 );
5
6 CREATE TABLE project(
7     project_id CHAR(4) PRIMARY KEY
8 );
9
10 CREATE TABLE subproject(
11     subproject_id VARCHAR(50) PRIMARY KEY,
12     project_id CHAR(4),
13     FOREIGN KEY (project_id) REFERENCES project(project_id)
14         ON UPDATE CASCADE
15         ON DELETE NO ACTION
16 );
17
18 CREATE TABLE domain(
19     domain_id VARCHAR(50) PRIMARY KEY
20 );
21
22 CREATE TABLE subdomain(
23     subdomain_id VARCHAR(50),
24     domain_id VARCHAR(50),
25     PRIMARY KEY(subdomain_id, domain_id),
26     FOREIGN KEY (domain_id) REFERENCES domain(domain_id)
27         ON UPDATE CASCADE

```

```

28         ON DELETE NO ACTION
29 );
30
31 CREATE TABLE term(
32     term_id SERIAL PRIMARY KEY,
33     term_name VARCHAR(100) NOT NULL,
34     language_id CHAR(2) NOT NULL,
35     part_of_speech TEXT NOT NULL,
36     term_notes TEXT,
37     iate_id INTEGER,
38     term_reference TEXT NOT NULL,
39     definition TEXT NOT NULL,
40     definition_external_cross_reference TEXT,
41     definition_source TEXT NOT NULL,
42     definition_notes TEXT,
43     domain_id VARCHAR(50) NOT NULL,
44     subdomain_id VARCHAR(50),
45     context TEXT NOT NULL,
46     context_external_cross_reference TEXT,
47     context_source TEXT NOT NULL,
48     subproject_id VARCHAR(50) NOT NULL,
49     FOREIGN KEY (subproject_id) REFERENCES subproject(subproject_id)
50         ON UPDATE CASCADE
51         ON DELETE NO ACTION,
52     FOREIGN KEY (language_id) REFERENCES language(language_id)
53         ON UPDATE CASCADE
54         ON DELETE NO ACTION,
55     FOREIGN KEY (domain_id) REFERENCES domain(domain_id)
56         ON UPDATE CASCADE
57         ON DELETE NO ACTION,
58     FOREIGN KEY (subdomain_id, domain_id) REFERENCES subdomain(subdomain_id,
59         domain_id)
60         ON UPDATE CASCADE
61         ON DELETE NO ACTION
62 );
63
64 CREATE TABLE translation(
65     translation_id SERIAL PRIMARY KEY,
66     en INTEGER,
67     el INTEGER,
68     es INTEGER,

```



```

68     it INTEGER,
69     ro INTEGER,
70     ru INTEGER,
71     uk INTEGER,
72     ca INTEGER,
73     de INTEGER,
74     pt INTEGER,
75     fr INTEGER,
76     FOREIGN KEY (en) REFERENCES term(term_id)
77         ON DELETE SET NULL,
78     FOREIGN KEY (el) REFERENCES term(term_id)
79         ON DELETE SET NULL,
80     FOREIGN KEY (es) REFERENCES term(term_id)
81         ON DELETE SET NULL,
82     FOREIGN KEY (it) REFERENCES term(term_id)
83         ON DELETE SET NULL,
84     FOREIGN KEY (ro) REFERENCES term(term_id)
85         ON DELETE SET NULL,
86     FOREIGN KEY (ru) REFERENCES term(term_id)
87         ON DELETE SET NULL,
88     FOREIGN KEY (uk) REFERENCES term(term_id)
89         ON DELETE SET NULL,
90     FOREIGN KEY (ca) REFERENCES term(term_id)
91         ON DELETE SET NULL,
92     FOREIGN KEY (de) REFERENCES term(term_id)
93         ON DELETE SET NULL,
94     FOREIGN KEY (pt) REFERENCES term(term_id)
95         ON DELETE SET NULL,
96     FOREIGN KEY (fr) REFERENCES term(term_id)
97         ON DELETE SET NULL
98 );
99
100 CREATE TABLE synonym(
101     term_id    INTEGER,
102     synonym_id VARCHAR(100),
103     PRIMARY KEY (term_id, synonym_id),
104     FOREIGN KEY (term_id) REFERENCES term(term_id)
105         ON UPDATE CASCADE
106         ON DELETE CASCADE
107 );
108

```

```

109 CREATE TABLE quasi_synonym(
110     term_id            INTEGER,
111     quasi_synonym_id  VARCHAR(100),
112     PRIMARY KEY (term_id, quasi_synonym_id),
113     FOREIGN KEY (term_id) REFERENCES term(term_id)
114         ON UPDATE CASCADE
115         ON DELETE CASCADE
116 );
117
118 CREATE TABLE common_name(
119     term_id            INTEGER,
120     common_name_id    VARCHAR(100),
121     PRIMARY KEY (term_id, common_name_id),
122     FOREIGN KEY (term_id) REFERENCES term(term_id)
123         ON UPDATE CASCADE
124         ON DELETE CASCADE
125 );
126
127 CREATE TABLE scientific_name(
128     term_id            INTEGER,
129     scientific_name_id VARCHAR(100),
130     PRIMARY KEY (term_id, scientific_name_id),
131     FOREIGN KEY (term_id) REFERENCES term(term_id)
132         ON UPDATE CASCADE
133         ON DELETE CASCADE
134 );
135
136 CREATE TABLE orthographic_variant(
137     term_id            INTEGER,
138     ortographic_variant_id VARCHAR(100),
139     PRIMARY KEY (term_id, ortographic_variant_id),
140     FOREIGN KEY (term_id) REFERENCES term(term_id)
141         ON UPDATE CASCADE
142         ON DELETE CASCADE
143 );
144
145 CREATE TABLE acronym(
146     term_id            INTEGER,
147     acronym_id         VARCHAR(100),
148     PRIMARY KEY (term_id, acronym_id),
149     FOREIGN KEY (term_id) REFERENCES term(term_id)

```

```

150         ON UPDATE CASCADE
151         ON DELETE CASCADE
152 );
153
154 CREATE TABLE abbreviation(
155     term_id          INTEGER,
156     abbreviation_id VARCHAR(100),
157     PRIMARY KEY (term_id, abbreviation_id),
158     FOREIGN KEY (term_id) REFERENCES term(term_id)
159         ON UPDATE CASCADE
160         ON DELETE CASCADE
161 );
162
163 CREATE INDEX term_term_name_index
164     ON term(term_name);
165 CREATE INDEX translation_en_index
166     ON translation(en);
167 CREATE INDEX translation_el_index
168     ON translation(el);
169 CREATE INDEX translation_es_index
170     ON translation(es);
171 CREATE INDEX translation_it_index
172     ON translation(it);
173 CREATE INDEX translation_ro_index
174     ON translation(ro);
175 CREATE INDEX translation_ru_index
176     ON translation(ru);
177 CREATE INDEX translation_uk_index
178     ON translation(uk);
179 CREATE INDEX translation_ca_index
180     ON translation(ca);
181 CREATE INDEX translation_de_index
182     ON translation(de);
183 CREATE INDEX translation_pt_index
184     ON translation(pt);
185 CREATE INDEX translation_fr_index
186     ON translation(fr);
187
188 INSERT INTO language VALUES('en', 'English'),
189                             ('el', 'Greek'),
190                             ('es', 'Spanish'),

```

```

191         ('it', 'Italian'),
192         ('ro', 'Romanian'),
193         ('ru', 'Russian'),
194         ('uk', 'Ukrainian'),
195         ('ca', 'Catalan'),
196         ('de', 'German'),
197         ('pt', 'Portuguese'),
198         ('fr', 'French');
199
200 INSERT INTO project VALUES('ENVI');
201 INSERT INTO subproject VALUES('Climate change', 'ENVI');
202 INSERT INTO subproject VALUES('Zoology', 'ENVI');
203
204 INSERT INTO domain VALUES('Environment');
205 INSERT INTO subdomain VALUES('climate change', 'Environment'),
206         ('biochemistry', 'Environment'),
207         ('pollution', 'Environment'),
208         ('chemistry', 'Environment'),
209         ('biology', 'Environment'),
210         ('zoology', 'Environment'),
211         ('sea', 'Environment'),
212         ('green finance', 'Environment'),
213         ('atmosphere', 'Environment'),
214         ('degradation of the environment', 'Environment'),
215         ('international relations', 'Environment'),
216         ('radiation', 'Environment'),
217         ('construction and town planning', 'Environment'),
218         ('finance', 'Environment'),
219         ('meteorology', 'Environment'),
220         ('processes', 'Environment'),
221         ('forestry', 'Environment'),
222         ('natural enviroment', 'Environment'),
223         ('water', 'Environment'),
224         ('water pollution', 'Environment'),
225         ('earth sciences', 'Environment');

```

Listing 9: comandi SQL per creazione database. Alcuni valori di subdomain sono stati omessi causa caratteri particolari non renderizzabili

```
1 """
```

```

2 parses one file at a time. Arguments required: 1. path to the file 2. ISO
   639-1 language code 3. subproject 4. domain
3 For example: python termsParser.py csv/ClimateChangeEN EN "Climate change"
   Environment
4 """
5
6 import pandas as pd
7 from dataclasses import dataclass, fields
8 import sys
9 import psycopg2
10 import re
11 import copy
12
13 # splits concatenated_string on separator and strips the resulting values of
   whitespaces
14 def split(concatenated_string, separator):
15     split_result = re.split(separator, re.sub(r'\s+', ' ',
   concatenated_string))
16     return [item.strip() for item in split_result if item.strip()]
17
18 @dataclass
19 class Term():
20     term_name: str = None
21     language_id: str = None
22     part_of_speech: str = ''
23     term_notes: str = None
24     iate_id: str = None
25     term_reference: str = ''
26     definition: str = ''
27     definition_external_cross_reference: str = None
28     definition_source: str = ''
29     definition_notes: str = None
30     domain_id: str = ''
31     subdomain_id: str = None
32     context: str = ''
33     context_external_cross_reference: str = None
34     context_source: str = ''
35     subproject_id: str = ''
36
37 # generic class for synonyms, quasi_synonyms, common_names...
38 @dataclass

```

```

39 class Variant():
40     term_name: str = None
41     variant: str = None
42
43 filename = sys.argv[1]
44 language = sys.argv[2]
45 subproject = sys.argv[3]
46 domain = sys.argv[4]
47
48 # returns a list of terms and a dictionary of lists containing variants of
49 # the terms
50 def parseTable(dataframe):
51     terms = []
52     variants = {'synonym': [], 'quasi_synonym': [], 'common_name': [], '
53 scientific_name': [], 'orthographic_variant': [], 'acronym': [], '
54 abbreviation': []}
55
56 for index, entry in dataframe.iterrows():
57     term = Term()
58
59     # loop through all possible columns of Term
60     for field in fields(Term):
61         # special values given by program arguments
62         if field.name in ['language_id', 'subproject_id', 'domain_id']:
63             continue
64
65         # if dataframe does not have the column then the default value
66         # for term remains
67         try: value = entry[field.name]
68         except KeyError: continue
69
70         # if value is nan or '-' then term retains the default
71         if pd.isna(value) and value != '-':
72             # remove leading and trailing whitespaces and newlines
73             value = re.sub(r'\s+', ' ', str(value)).strip()
74             # if column is 'iate_id' then convert first to integer
75             if field.name == 'iate_id':
76                 setattr(term, field.name, str(int(float(value)))) #
77                 # cannot convert float string to int directly
78             elif field.name == 'subdomain_id':
79                 setattr(term, field.name, value.lower())

```

```

74         else:
75             setattr(term, field.name, value)
76
77         # if term_name was null then don't even insert it into database
78         if pd.isna(term.term_name): continue
79
80         term.language_id = language
81         term.subproject_id = subproject
82         term.domain_id = domain
83
84         # if term_name were multiple values separated by ';' create multiple
terms and append to terms
85         new_terms = []
86         for item in split(term.term_name, ';|,|/'):
87             new_term = copy.copy(term)
88             new_term.term_name = item.lower()
89             new_terms.append(new_term)
90         if len(new_terms) > 1:
91             print(f'line {index} has {len(new_terms)} separated values')
92         terms.extend(new_terms)
93
94         # parsing variants' fields
95         for key in variants.keys():
96             try: value = entry[key]
97             except KeyError: continue
98             if pd.notna(value):
99                 # split variants separated by commas, forward slashes and
digit+period
100                 for term in new_terms:
101                     for variant in split(value, ',|/[0-9]\.'):
102                         variants[key].append(Variant(term.term_name, variant
))
103
104         return terms, variants
105
106 df = pd.read_csv(filename)
107
108 # replaces single quotes with two single quotes so postgresql escapes them
109 stringify = lambda x: '"' + x.replace('"', '""') + '"'
110
111 # insert parsed data into database

```

```

112 with psycopg2.connect('dbname=yourterm user=postgres password=password') as
    conn:
113     with conn.cursor() as cur:
114         terms, variants = parseTable(df)
115         for term in terms:
116             values = ', '.join(['NULL' if getattr(term, field.name) is None
117 else stringify(getattr(term, field.name)) for field in fields(Term)])
118             try:
119                 cur.execute('INSERT INTO term VALUES(DEFAULT, {})'.format(
120 values))
121                 conn.commit()
122             except psycopg2.Error as err:
123                 print(err)
124                 conn.rollback()
125                 continue
126         for key, variant_list in variants.items():
127             for variant in variant_list:
128                 try:
129                     cur.execute('SELECT term_id from term WHERE term_name =
130 {} AND language_id = \'{}\'''.format(stringify(variant.term_name),
131 language))
132                     term_ids = cur.fetchall()
133                     for term_id in term_ids:
134                         cur.execute('INSERT INTO {} VALUES({}, {})'.format(
135 key, term_id[0], stringify(variant.variant)))
136                         conn.commit()
137                 except psycopg2.Error as err:
138                     print(err)
139                     conn.rollback()
140                     continue

```

Listing 10: Script per parsing file di termini

```

1 """
2 parses 'all languages' csv file. Argument required: path to file.
3 For example: python translationsParser.py csv/ClimateChangeAllLanguages.csv
4 """
5
6 import pandas as pd
7 import sys
8 from dataclasses import dataclass, fields
9 import psycopg2

```



```

10 import re
11 from itertools import product
12 from math import prod
13
14 @dataclass
15 class Translation():
16     en: str = None
17     el: str = None
18     es: str = None
19     it: str = None
20     ro: str = None
21     ru: str = None
22     uk: str = None
23     ca: str = None
24     de: str = None
25     pt: str = None
26     fr: str = None
27
28 # splits concatenated_string on separator and strips the resulting values of
    whitespaces
29 def split(concatenated_string, separator):
30     split_result = re.split(separator, re.sub(r'\s+', ' ',
    concatenated_string))
31     return [item.strip() for item in split_result if item.strip()]
32
33 def parseTable(dataframe):
34     for index, entry in dataframe.iterrows():
35         translation = Translation()
36         for field in fields(Translation):
37             try: value = entry[field.name]
38             except KeyError: continue
39             if pd.isna(value):
40                 value = re.sub(r'\s+', ' ', value).strip().lower()
41                 setattr(translation, field.name, value)
42
43         # case fields are ; separated values
44         values = [split(str(getattr(translation, field.name)), ';|,|/')] for
    field in fields(Translation)]
45         combination_number = prod([len(value) for value in values])
46         if combination_number > 1: print(f'line {index} has {
    combination_number} combinations')

```

```

47     for combination in product(*values):
48         yield Translation(*combination)
49
50 # returns the list of all term_id that have the term_name == term and the
    specified language
51 # if there is no correspondence, [None] is returned
52 def generate_term_ids(cur, term, language):
53     term_ids = []
54     if term is not None:
55         cur.execute('SELECT term_id from term WHERE term_name = {} AND
    language_id = \'{ }\'.format(stringify(term), language))
56         for id in cur.fetchall():
57             term_ids.append(id[0])
58     if len(term_ids) == 0:
59         term_ids.append(None)
60     return term_ids
61
62 filename = sys.argv[1]
63
64 df = pd.read_csv(filename)
65
66 # replaces single quotes with two single quotes so postgresql escapes them
67 stringify = lambda x: '"' + x.replace("'", "'') + '"'
68
69 with psycopg2.connect('dbname=yourterm user=postgres password=password') as
    conn:
70     with conn.cursor() as cur:
71         for translation in parseTable(df):
72             try:
73                 # ids contains all the lists corresponding to the term_ids
    of each term in a translation
74                 ids = []
75                 for field in fields(Translation):
76                     ids.append(generate_term_ids(cur, getattr(translation,
    field.name), field.name))
77                 for combination in product(*ids):
78                     if not all([item is None for item in combination]):
79                         values = ', '.join(['NULL' if id is None else str(id
    ) for id in combination])
80                         cur.execute('INSERT INTO translation VALUES(DEFAULT,
    { })'.format(values))

```

```
81         conn.commit()
82     except psycopg2.Error as err:
83         conn.rollback()
84         print(err)
85         continue
```

Listing 11: Script per parsing file di traduzioni