

**Scadenze:**

**Presentazione domanda e titolo:** entro il 22 Novembre

**Consegna riassunto:** dal 22 novembre al 6 dicembre

**Consegna tesi:** entro il 6 dicembre



**UNIVERSITÀ DEGLI STUDI DI MILANO**  
**FACOLTÀ DI SCIENZE E TECNOLOGIE**

*Corso di Laurea Magistrale in Informatica*

**TITOLO DA DEFINIRE**

**Relatore:** Prof. Alfio FERRARA

**Correlatore:** Dr. Francesco PERITI

**Autore:** Sebastiano Caccaro

**Matricola:** 958683

Anno Accademico 2020-2021

*dedicato a ...*

# Prefazione

hkjafgyruet.

## Organizzazione della tesi

La tesi è organizzata come segue:

- nel Capitolo 1 ....

# Ringraziamenti

asdjhgtry.

# Indice

	ii
Prefazione	iii
Ringraziamenti	iv
<b>1 Introduzione</b>	<b>1</b>
<b>2 Stato dell'Arte</b>	<b>2</b>
<b>3 Dataset e Perturbazione</b>	<b>3</b>
<b>4 Metodologia</b>	<b>4</b>
4.1 Introduzione . . . . .	4
4.2 Processo . . . . .	6
4.2.1 Panoramica Generale . . . . .	6
4.2.2 Error Detection . . . . .	6
4.2.3 Error Correction . . . . .	7
<b>5 Implementazione, Test e Risultati</b>	<b>10</b>
<b>6 Analisi dell'errore</b>	<b>11</b>

# Capitolo 1

## Introduzione

## Capitolo 2

### Stato dell'Arte

## Capitolo 3

# Dataset e Perturbazione



# Capitolo 4

## Metodologia

In questo capitolo è descritta la metodologia del sistema di correzione. Nella sezione 4.1 vengono descritti gli obiettivi del processo di correzione e le criticità che lo contraddistinguono. Nella sottosezione 4.2.1 è presente una panoramica generale del processo di correzione, e sono introdotte le varie fasi che lo compongono. Nella sottosezione 4.2.2 e sottosezione 4.2.3 sono descritte più nel dettaglio le fasi di error detection e correction che compongono il processo.

### 4.1 Introduzione

L'optical character recognition, o OCR, è una tecnologia tramite la quale è possibile estrarre e digitalizzare del testo da un'immagine. Più precisamente, per digitalizzare si intende ottenere del testo in una codifica leggibile da una macchina, come ad esempio ASCII o Unicode.

Provvedimenti per incrementare l'occupazione operaia, agevolando la costruzione di case per i lavoratori.

Prima di dare la parola al primo iscritto, mi permetto di far presente alla Camera che sono già stati presentati vari ordini del giorno e vari emendamenti, relativi a questo disegno di legge. Come la Camera sa, gli emendamenti, di regola, debbono esser presentati almeno ventiquattro ore prima della seduta in cui vengono discussi. Per eccezione, possono anche esser presentati senza rispettare questo termine, cioè nella seduta stessa, purché siano

→

Prima di dare la parola al primo iscritto, mi permetto di far presente alla Camera che sono già stati presentati vari ordini del giorno e vari emendamenti, relativi a questo disegno di legge. Come la Camera sa, gli emendamenti, di regola, debbono esser presentati almeno ventiquattro ore prima della seduta in cui vengono discussi.

Figura 1: Esempio di testo estratto tramite OCR

L'uso di sistemi OCR è particolarmente vantaggioso per l'archiviazione dei documenti. Infatti, l'estrazione del testo rende possibile eseguire sui documenti

operazioni di ricerca e di reperimento informazioni. Si pensi a un'operazione basilare come la ricerca di tutti i documenti che contengono una determinata parola: un task dispendioso e manuale che è quasi istantaneo se si ha a disposizione il testo digitalizzato.

**Errori** I sistemi OCR presentano però alcuni problemi che possono influenzare negativamente le performance nel reperimento di informazioni [1][2]. Si veda, ad esempio, il testo estratto in Figura 1: seppur la maggior parte delle parole sono state riconosciute correttamente, è possibile notare alcuni errori. Generalmente la presenza di tali errori è dovuta a piccole imprecisioni nell'immagine iniziale: si pensi ad esempio un granello di polvere che può essere scambiato per un punto, o ad una "n" battuta male che viene scambiata per la sequenza "ii". Si possono distinguere le seguenti categorie di errore:

- **Errore di token.** Si ha un errore di token aggiunte, rimozioni o sostituzioni di caratteri occorrono all'interno di una parola senza però dividerla o eliminarla. Questo tipo di errore si divide in due ulteriori sottocategorie:
  - Non-word error (NW), quando la parola risultante non è presente in un dato vocabolario.
  - Real-word error (RW), quando la parola risultante è presente in un dato vocabolario.
- **Errore di segmentazione (SG).** Si ha un errore di segmentazione quando aggiunte, rimozioni o sostituzioni di caratteri occorrono in maniera tale da dividere o unire parole. Sono considerati errori di segmentazione anche aggiunte di caratteri e segni di punteggiatura spuri all'interno del testo.

---

Lo scopo del processo di OCR-Post Processing è quello di correggere e minimizzare gli errori introdotti dall'acquisizione di testo da immagini. Più in generale, data una frase contenente degli errori, lo scopo del processo di correzione è quello produrre in output la stessa frase senza errori. Nel fare ciò è inoltre necessario assicurarsi di non introdurne di nuovi.

La metodologia di correzione sviluppata, una volta identificati gli errori, fa uso del BERT Masked Language Modeling per produrre una serie di candidati per la correzione. Un approccio simile è utilizzato in [3], dove una combinazione di BERT e FastText riesce a produrre candidati corretti nel 70% dei casi. L'approccio appena citato però non include le fasi di error detection e scelta del candidato corretto, che sono invece implementate nella metodologia proposta e applicate al dataset descritto nel Capitolo 3.

## 4.2 Processo

### 4.2.1 Panoramica Generale

Il processo di correzione si articola in più fasi, alcune delle quali si ripetono più volte. È possibile distinguere i seguenti passaggi:

- **Error detection:** in questa fase vengono individuati e contrassegnati gli errori all'interno della frase, se presenti.
- **Error correction:** in questa fase si tenta la correzione degli errori individuati in precedenza.

Come mostrato in Figura 2, queste due fasi sono ripetute più volte. Ciò serve per sfruttare al massimo le caratteristiche del sistema di Error Correction, che si basa sul BERT Masked Language Modeling. Come spiegato più approfonditamente nella sottosezione 4.2.3, questa funzione fa uso del contesto della frase e dell'intorno della parola da correggere per proporre dei candidati per la correzione da effettuare. Il sistema di correzione, specialmente in caso di frasi contenenti molti errori, potrebbe non essere in grado di correggerli tutti. È però possibile che, dopo aver corretto alcuni errori, il sistema sia in grado di correggerne altri grazie al maggior contesto che le correzioni hanno portato all'interno della frase.

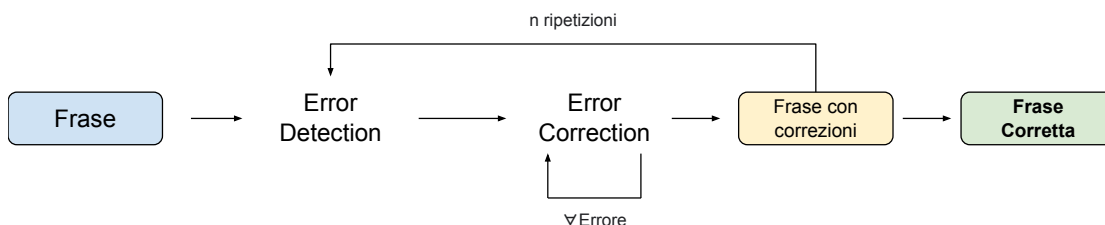


Figura 2: Schema riassuntivo della Metodologia

### 4.2.2 Error Detection

Lo scopo dell'error detection è quello di contrassegnare gli errori presenti all'interno della frase per la successiva fase di error correction. Dato che l'error correction corregge gli errori a livello di token, vengono contrassegnati per la fase successiva tutti i token contenenti errori. A tale scopo, è necessario in primis tokenizzare la frase di partenza. Successivamente tutti i token contenenti errori vengono marcati per la correzione.

Quanto appena spiegato è riportato nello schema in Figura 3.

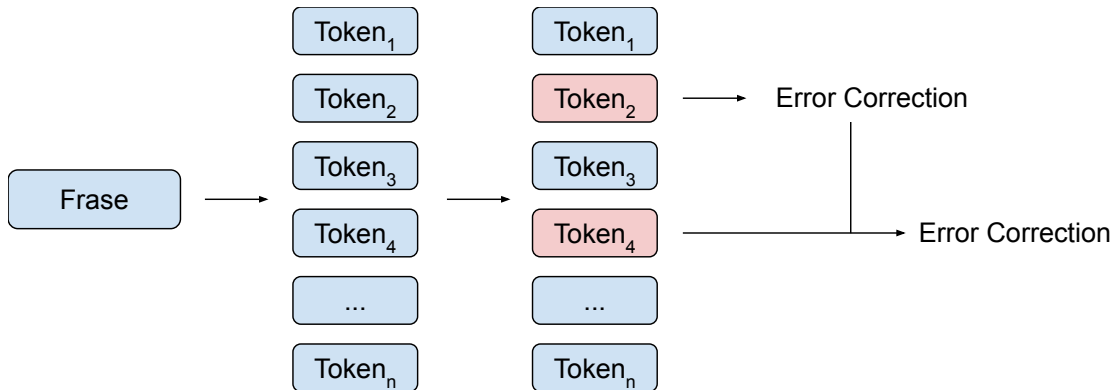


Figura 3: Schema del processo di error detection

### 4.2.3 Error Correction

Il processo di error correction si basa sul BERT Masked Language Modeling. Data una qualsiasi frase, è possibile sostituire una parola con la stringa *[MASK]*, detta maschera. Dando in input la frase mascherata al modello BERT, esso produrrà una serie di parole (da qui in poi candidati) associate alla loro probabilità di corrispondere al token mascherato.

**Esempio** Data la frase

*"che assistono ragazze in difficoltà, le persone soie e abbandonate, gli ammalati e gli anziani."*

la parola "*soie*" sottolineata è stata individuata come errore. È quindi necessario mascherarla, per dare la frase in input al modello BERT. La frase diventa dunque:

*"che assistono ragazze in difficoltà, le persone [MASK] e abbandonate, gli ammalati e gli anziani."*

BERT produce quindi una lista di candidati, di cui sono riportati solo i primi 5:

- "*sole*" con probabilità 0.42
- "*anziane*" con probabilità 0.28
- "*povere*" con probabilità 0.08
- "*care*" con probabilità 0.03
- "*disabili*" con probabilità 0.01

Bisogna sottolineare come la parola originale sia trasparente al modello BERT. Ciò significa che i candidati prodotti dal modello sono del tutto indipendenti dalla parola originale, e sono inferite unicamente dal contesto derivato dal resto della frase.

La fase di error correction inizia con i token contrassegnati come errore nella fase precedente. Si procede mascherando il primo token errato all'interno della frase. Siccome BERT necessita di una frase, e non di un'insieme di token, è necessario detokinizzare la frase. Fatto ciò è possibile produrre i candidati per la correzione, come mostrato in Figura 5.

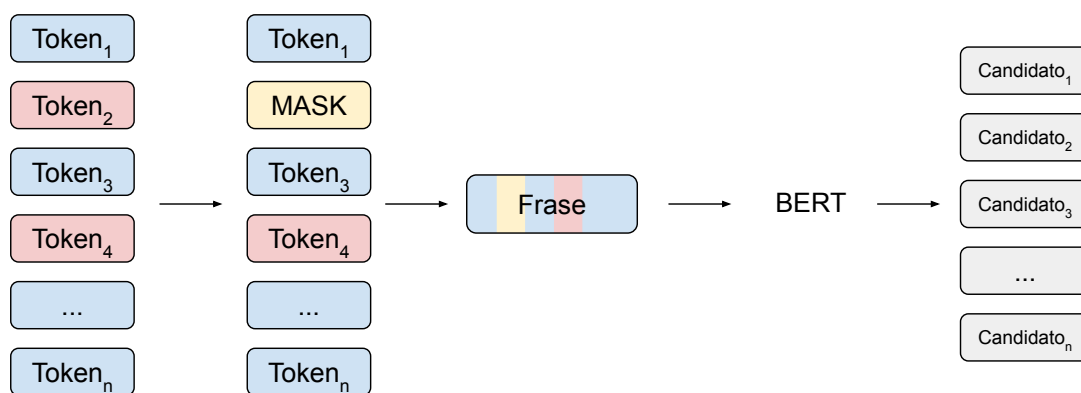


Figura 4: Schema del processo di generazione dei candidati

Fra i vari candidati proposti da BERT ne viene scelto solamente uno, che viene determinato sulla base di alcuni criteri che tengono conto anche del token originale contenente l'errore. Può però accadere che nemmeno il candidato sia la correzione adatta: si pensi al caso in cui l'error detection contrassegna un token corretto come errore, o ad una frase in cui bert non produce la parola corretta fra i candidati. Ogni candidato scelto deve quindi passare un ulteriore filtro, che ha lo scopo di distinguere queste evenienze. Nel caso il candidato non passi il filtro, il sistema ignora la correzione; in caso contrario, la correzione viene sostituita al token mascherato. Quanto appena descritto è rappresentato nello schema in Figura 5.

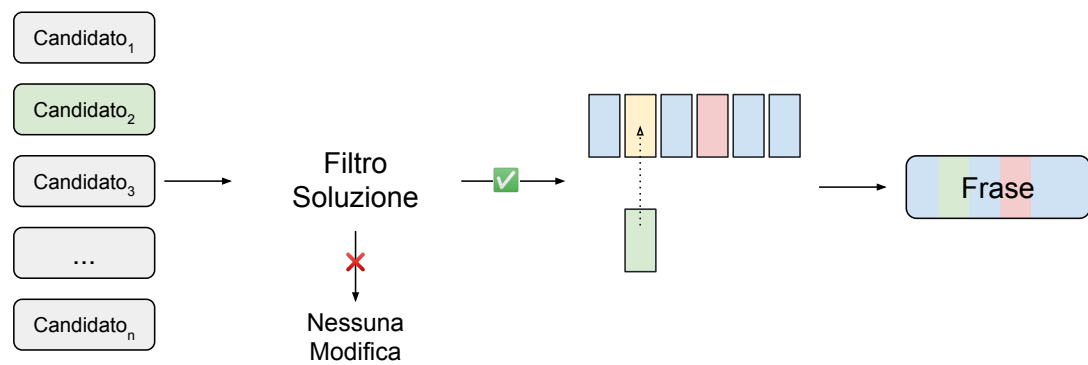


Figura 5: Schema del processo di scelta dei candidati

Il processo appena descritto si ripete per ogni errore contrassegnato durante la fase di error detection. Una volta completata la correzione dell'ultimo token errato, la fase di error correction può dirsi conclusa.

## Capitolo 5

### Implementazione, Test e Risultati

## Capitolo 6

### Analisi dell'errore



# Bibliografia

- [1] Guillaume Chiron, Antoine Doucet, Mickaël Coustaty, Muriel Visani, and Jean-Philippe Moreux. Impact of ocr errors on the use of digital libraries: towards a better access to information. In *2017 ACM/IEEE Joint Conference on Digital Libraries (JCDL)*, pages 1–4. IEEE, 2017.
- [2] Stephen Mutuvi, Antoine Doucet, Moses Odeo, and Adam Jatowt. Evaluating the impact of ocr errors on topic modeling. In *International Conference on Asian Digital Libraries*, pages 3–14. Springer, 2018.
- [3] Mahdi Hajiali, Jorge Ramón Fonseca Cacho, and Kazem Taghva. Generating correction candidates for ocr errors using bert language model and fasttext subword embeddings. In *Intelligent Computing*, pages 1045–1053. Springer, 2022.