

# Relazione Settimanale - Sviluppo del Sistema di Fine-Tuning di BERT

Gruppo di Lavoro

1 aprile 2025

## Indice

<b>1</b>	<b>Introduzione</b>	<b>2</b>
<b>2</b>	<b>Pulizia e Preparazione dei Dati</b>	<b>2</b>
<b>3</b>	<b>Implementazione dei Test su Google Colab</b>	<b>2</b>
<b>4</b>	<b>Sviluppo di un Layer di Traduzione</b>	<b>3</b>
4.1	Obiettivi del Layer di Traduzione . . . . .	3
4.2	Versione 1.0 del Traduttore: Dettaglio del Codice . . . . .	3
4.3	Evoluzione del Workflow e Integrazione . . . . .	9
<b>5</b>	<b>Test Preliminari e Analisi dei Risultati</b>	<b>9</b>
<b>6</b>	<b>Conclusioni e Prospettive Future</b>	<b>10</b>

# 1 Introduzione

La presente relazione descrive le attività svolte durante la terza settimana di lavoro, incentrata sul perfezionamento della pulizia dei dati e sulla realizzazione dei primi test in ambiente Google Colab. In particolare, il focus principale è stato posto sull'implementazione iniziale del fine-tuning di BERT tramite la libreria `transformers` di Python e sull'elaborazione dei dati grezzi, al fine di predisporre un dataset bilanciato per i successivi esperimenti.

## 2 Pulizia e Preparazione dei Dati

Durante questa settimana abbiamo intensificato gli sforzi nella pulizia e organizzazione dei dati. Il lavoro si è articolato nelle seguenti fasi:

- **Estrazione e Aggregazione:** Dall'insieme dei file in formato `.txt` forniti dal Professor Ferretti, abbiamo estratto circa 2000 esempi etichettati come `jailbreak` e li abbiamo raccolti in un unico file. Per mantenere il bilanciamento del dataset, sono stati estratti altresì 2000 esempi etichettati come `confused` (intesi come non `jailbreak`).
- **Preliminare Validazione:** È importante notare che, in questa fase preliminare, non è stata ancora effettuata una verifica approfondita della correttezza delle etichette. Tale operazione verrà eseguita in iterazioni successive, poiché il primo approccio si concentra sul reperimento e l'aggregazione dei dati.
- **Utilizzo di Script Python:** Tutte le estrazioni e le operazioni di combinazione dei file sono state realizzate tramite script Python, garantendo così un elevato grado di automazione e riproducibilità del processo.

## 3 Implementazione dei Test su Google Colab

Un'altra parte fondamentale del lavoro di questa settimana ha riguardato l'implementazione dei primi test in Google Colab. Gli obiettivi di questi test sono:

- **Validazione dell'approccio:** Eseguire il fine-tuning di BERT su un dataset di dimensioni ridotte, al fine di monitorare le metriche di performance quali accuratezza, precisione e recall.
- **Analisi delle Risorse:** Identificare eventuali colli di bottiglia in termini di tempi di training e memoria GPU, predisponendo così il passaggio ad un modello più robusto da eseguire sul cluster del dipartimento.
- **Feedback Operativo:** I test hanno permesso di acquisire una prima esperienza sul workflow di training e di valutazione del modello, evidenziando il rischio di overfitting dovuto al testing con esempi già inclusi nel dataset di training.

## 4 Sviluppo di un Layer di Traduzione

Durante le operazioni di analisi dei dati è emersa la problematica della presenza di testi in lingue differenti dall'inglese. Nello specifico, una parte consistente del dataset conteneva risposte scritte in vietnamita, indonesiano, giapponese e turco. Per ovviare a questa problematica, abbiamo realizzato uno script Python sofisticato volto a tradurre in modo automatico i file contenenti gli esempi sia di `jailbreak` che di `no jailbreak`.

### 4.1 Obiettivi del Layer di Traduzione

Lo scopo del layer di traduzione è duplice:

1. **Standardizzazione dell'Input:** Garantire che tutte le risposte vengano convertite in lingua inglese, permettendo di utilizzare un unico modello (BERT) per il fine-tuning senza introdurre distorsioni semantiche.
2. **Compatibilità e Flessibilità:** Integrare strumenti di traduzione che spaziano dall'utilizzo di API esterne, come Google Traduttore, a soluzioni offline basate su librerie open source, per avere una copertura completa anche quando alcuni strumenti esterni si rifiutano di tradurre testi contenenti informazioni illegali.

### 4.2 Versione 1.0 del Traduttore: Dettaglio del Codice

Di seguito viene riportata la versione 1.0 dello script del traduttore, sviluppata per essere il nucleo del futuro layer di traduzione. Questo script gestisce sia il download dei pacchetti per la traduzione offline sia la logica per suddividere il testo in blocchi traducibili, interfacciandosi con i vari strumenti disponibili.

```
1 import os
2 import time
3 from pathlib import Path
4 from tqdm import tqdm
5 import logging
6
7 # Configura logging
8 logging.basicConfig(
9     level=logging.INFO,
10    format='%(asctime)s - %(levelname)s - %(message)s',
11    handlers=[
12        logging.FileHandler("translation_log.txt"),
13        logging.StreamHandler()
14    ]
15 )
16 logger = logging.getLogger(__name__)
17
18 def get_translator(translator_type="google", to_lang="en"):
19     """Crea e restituisce il traduttore appropriato in base al tipo
20     specificato"""
21     if translator_type == "google":
22         try:
23             from deep_translator import GoogleTranslator
24             # Google Translator ha una funzione di auto-rilevamento
25             della lingua
26             # che funziona meglio con contenuti multilingue
27             return GoogleTranslator(source='auto', target=to_lang)
```

```

26     except ImportError:
27         logger.error("Per usare GoogleTranslator installa: pip
install deep-translator")
28         return None
29     elif translator_type == "offline":
30         try:
31             from argostranslate import package, translate
32             # Verifica se i pacchetti sono già installati
33             if not package.get_installed_packages():
34                 logger.info("Scaricamento e installazione dei pacchetti
di traduzione...")
35                 package.update_package_index()
36                 available_packages = package.get_available_packages()
37                 installed_count = 0
38                 for pkg in available_packages:
39                     if pkg.to_code == to_lang:
40                         try:
41                             pkg.install()
42                             installed_count += 1
43                             logger.info(f"Installato pacchetto: {pkg.
from_code} -> {pkg.to_code}")
44                         except Exception as e:
45                             logger.error(f"Errore installazione {pkg.
from_code}: {str(e)}")
46                 logger.info(f"Installati {installed_count} pacchetti di
traduzione verso {to_lang}")
47                 if installed_count == 0:
48                     logger.error(f"Nessun pacchetto trovato con target {
to_lang}")
49                 return None
50                 installed_languages = translate.get_installed_languages()
51                 target_lang = next((lang for lang in installed_languages if
lang.code == to_lang), None)
52                 if not target_lang:
53                     logger.error(f"Lingua target {to_lang} non disponibile")
54                     return None
55                 class MultilingualArgosTranslator:
56                     def __init__(self, installed_languages, target_lang_code
):
57                         self.target_lang_code = target_lang_code
58                         self.installed_languages = installed_languages
59                         self.models = {}
60                         for lang in installed_languages:
61                             if lang.code != target_lang_code:
62                                 translation = lang.get_translation(
target_lang_code)
63                             if translation:
64                                 self.models[lang.code] = translation
65                         if not self.models:
66                             logger.error("Nessun modello di traduzione
disponibile")
67                     def translate(self, text):
68                         import re
69                         if self._is_probably_english(text):
70                             return text
71                         detected_lang = self._detect_language(text)
72                         if detected_lang in self.models:

```

```

73         return self.models[detected_lang].translate(text
74     )
75     else:
76         best_translation = text
77         for lang_code, model in self.models.items():
78             try:
79                 translation = model.translate(text)
80                 if self._quality_score(translation) >
self._quality_score(best_translation):
81                     best_translation = translation
82             except:
83                 continue
84         return best_translation
85     def _is_probably_english(self, text):
86         common_english_words = {"the", "and", "is", "in", "
to", "of", "that", "for", "it", "with"}
87         words = text.lower().split()
88         if not words:
89             return True
90         english_count = sum(1 for word in words if word in
common_english_words)
91         return english_count / len(words) > 0.2
92     def _detect_language(self, text):
93         try:
94             import langdetect
95             return langdetect.detect(text)
96         except:
97             if any(ord(c) > 1000 for c in text):
98                 for candidate in ["ar", "ru", "zh", "ja", "
ko"]]:
99                     if candidate in self.models:
100                         return candidate
101                     return next(iter(self.models.keys())) if self.
models else "en"
102     def _quality_score(self, text):
103         latin_chars = sum(1 for c in text if 'a' <= c.lower
() <= 'z')
104         total_chars = max(1, len(text.strip()))
105         return latin_chars / total_chars
106     return MultilingualArgosTranslator(installed_languages,
to_lang)
107     except ImportError:
108         logger.error("Per usare il traduttore offline installa: pip
install argostranslate langdetect")
109         return None
110     elif translator_type == "combo":
111         google_translator = get_translator("google", to_lang)
112         if google_translator:
113             return google_translator
114         logger.warning("Google Translator non disponibile, uso
traduttore offline")
115         return get_translator("offline", to_lang)
116     else:
117         logger.error(f"Tipo di traduttore non supportato: {
translator_type}")
118         return None

```

```

119 def traduci_file_multilingue(file_input, file_output, etichetta_inizio,
120                               etichetta_fine,
121                               translator_type="google", to_lang="en",
122                               chunk_size=4500, delay=2):
123     """
124     Traduce un file di grandi dimensioni con supporto per testi
125     multilingue.
126
127     Parametri:
128     - file_input: percorso del file da tradurre
129     - file_output: percorso del file di output
130     - etichetta_inizio: tag che segna l'inizio di un blocco da tradurre
131     - etichetta_fine: tag che segna la fine di un blocco da tradurre
132     - translator_type: "google", "offline" o "combo"
133     - to_lang: lingua di destinazione (default "en" per inglese)
134     - chunk_size: dimensione massima in caratteri per ogni richiesta di
135     traduzione
136     - delay: ritardo in secondi tra le traduzioni
137     """
138     backup_file = f"{file_output}.progress"
139     temp_output = f"{file_output}.temp"
140     translator = get_translator(translator_type, to_lang)
141     if not translator:
142         logger.error("Impossibile inizializzare il traduttore")
143         return
144     blocchi_tradotti = []
145     ultimo_blocco_tradotto = -1
146     if os.path.exists(backup_file):
147         try:
148             with open(backup_file, 'r', encoding='utf-8') as f:
149                 ultimo_blocco_tradotto = int(f.read().strip())
150             if os.path.exists(temp_output):
151                 with open(temp_output, 'r', encoding='utf-8') as f:
152                     blocchi_tradotti = f.read().split("\n\n--- NUOVO
153     BLOCCO ---\n\n")
154                     blocchi_tradotti = [b for b in blocchi_tradotti if b
155     .strip()]
156             logger.info(f"Ripresa dalla traduzione: trovati {len(
157     blocchi_tradotti)} blocchi già tradotti")
158         except Exception as e:
159             logger.error(f"Errore durante il caricamento del progresso:
160     {e}")
161             ultimo_blocco_tradotto = -1
162             blocchi_tradotti = []
163     try:
164         with open(file_input, 'r', encoding='utf-8') as fin:
165             contenuto = fin.read()
166     except UnicodeDecodeError:
167         encodings = ['latin-1', 'iso-8859-1', 'cp1252']
168         for enc in encodings:
169             try:
170                 with open(file_input, 'r', encoding=enc) as fin:
171                     contenuto = fin.read()
172                     logger.info(f"File aperto con encoding: {enc}")
173                     break
174             except:
175                 continue
176     else:

```

```

170         logger.error("Impossibile aprire il file con gli encoding
supportati")
171         return
172     indice_inizio = 0
173     tutti_blocchi = []
174     while True:
175         inizio = contenuto.find(etichetta_inizio, indice_inizio)
176         if inizio == -1:
177             break
178         fine = contenuto.find(etichetta_fine, inizio)
179         if fine == -1:
180             break
181         fine_completa = fine + len(etichetta_fine)
182         blocco = contenuto[inizio:fine_completa]
183         tutti_blocchi.append(blocco)
184         indice_inizio = fine_completa
185     logger.info(f"Trovati {len(tutti_blocchi)} blocchi totali nel file")
186     for i, blocco in enumerate(tqdm(tutti_blocchi[ultimo_blocco_tradotto
+1:],
187                                     desc="Traduzione blocchi")):
188         indice_blocco = i + ultimo_blocco_tradotto + 1
189         try:
190             inizio_tag = blocco.find(etichetta_inizio)
191             fine_tag = blocco.rfind(etichetta_fine)
192             if inizio_tag != -1 and fine_tag != -1:
193                 inizio_contenuto = inizio_tag + len(etichetta_inizio)
194                 testo_da_tradurre = blocco[inizio_contenuto:fine_tag].
strip()
195                 if not testo_da_tradurre:
196                     blocco_tradotto = blocco
197                     blocchi_tradotti.append(blocco_tradotto)
198                     continue
199                 if len(testo_da_tradurre) > chunk_size:
200                     linee = testo_da_tradurre.split('\n')
201                     chunks = []
202                     chunk_corrente = []
203                     lunghezza_corrente = 0
204                     for linea in linee:
205                         if lunghezza_corrente + len(linea) > chunk_size
and chunk_corrente:
206                             chunks.append('\n'.join(chunk_corrente))
207                             chunk_corrente = []
208                             lunghezza_corrente = 0
209                             chunk_corrente.append(linea)
210                             lunghezza_corrente += len(linea) + 1
211                     if chunk_corrente:
212                         chunks.append('\n'.join(chunk_corrente))
213                     traduzioni = []
214                     for j, chunk in enumerate(chunks):
215                         try:
216                             if is_probably_english(chunk):
217                                 logger.info(f"  Chunk {j+1}/{len(chunks)}
} già in inglese, salto traduzione")
218                                 traduzioni.append(chunk)
219                             else:
220                                 traduzione_chunk = translator.translate(
chunk)
221                                 traduzioni.append(traduzione_chunk)

```

```

222         logger.info(f" Blocco {indice_blocco
+1}/{len(tutti_blocchi)}: Chunk {j+1}/{len(chunks)} tradotto")
223         time.sleep(delay)
224         except Exception as e:
225             logger.error(f"Errore durante la traduzione
del chunk {j+1}: {e}")
226             traduzioni.append(chunk)
227             testo_tradotto = '\n'.join(traduzioni)
228         else:
229             if is_probably_english(testo_da_tradurre):
230                 logger.info(f" Blocco {indice_blocco+1}/{len(
tutti_blocchi)} già in inglese, salto traduzione")
231                 testo_tradotto = testo_da_tradurre
232             else:
233                 testo_tradotto = translator.translate(
testo_da_tradurre)
234                 blocco_tradotto = f"{etichetta_inizio}\n{testo_tradotto
}\n{etichetta_fine}"
235                 blocchi_tradotti.append(blocco_tradotto)
236                 with open(backup_file, 'w', encoding='utf-8') as f:
237                     f.write(str(indice_blocco))
238                 with open(temp_output, 'w', encoding='utf-8') as f:
239                     f.write("\n\n--- NUOVO BLOCCO ---\n\n".join(
blocchi_tradotti))
240                 time.sleep(delay)
241             else:
242                 logger.error(f"Errore: tag non trovati nel blocco {
indice_blocco+1}")
243                 blocchi_tradotti.append(blocco)
244             except Exception as e:
245                 logger.error(f"Errore durante la traduzione del blocco {
indice_blocco+1}: {e}")
246                 blocchi_tradotti.append(blocco)
247                 with open(backup_file, 'w', encoding='utf-8') as f:
248                     f.write(str(indice_blocco))
249                 with open(temp_output, 'w', encoding='utf-8') as f:
250                     f.write("\n\n--- NUOVO BLOCCO ---\n\n".join(
blocchi_tradotti))
251                 with open(file_output, 'w', encoding='utf-8') as fout:
252                     fout.write("\n\n--- NUOVO BLOCCO ---\n\n".join(blocchi_tradotti)
)
253                 if len(blocchi_tradotti) == len(tutti_blocchi):
254                     try:
255                         if os.path.exists(backup_file):
256                             os.remove(backup_file)
257                         if os.path.exists(temp_output):
258                             os.remove(temp_output)
259                     except Exception as e:
260                         logger.warning(f"Avviso: impossibile rimuovere i file
temporanei: {e}")
261                 logger.info(f"{len(blocchi_tradotti)} blocchi tradotti e salvati in
'{file_output}'")
262
263 def is_probably_english(text):
264     """Determina se il testo è probabilmente già in inglese"""
265     common_english_words = {"the", "and", "is", "in", "to", "of", "that"
, "for", "it", "with",

```



```

266         "this", "on", "are", "as", "was", "by", "be",
        "have", "you", "not"}
267 words = text.lower().split()
268 if not words:
269     return True
270 if len(words) < 5:
271     return False
272 english_count = sum(1 for word in words if word in
common_english_words)
273 return english_count / len(words) > 0.15
274
275 # Esempio di utilizzo
276 if __name__ == "__main__":
277     traduci_file_multilingue(
278         '2490_Confused.txt',
279         '2490_Confused_Tradotto.txt',
280         '###CONFUSED###',
281         '###ENDCONFUSED###',
282         translator_type="google",
283         to_lang="en",
284         chunk_size=4500,
285         delay=2
286     )

```

Listing 1: Versione 1.0 del traduttore

### 4.3 Evoluzione del Workflow e Integrazione

Una volta verificato il corretto funzionamento dello script di traduzione, abbiamo testato la sua efficacia passando i file di esempio per `jailbreak` e `no jailbreak` per tre volte consecutive. Quest’approccio ha portato a due ulteriori sviluppi:

- **Adozione del Formato JSON:** Considerata la necessità di gestire file di grandi dimensioni e di integrare agevolmente i dati con i nostri script Python, abbiamo convertito i file iniziali in formato JSON. Tale formato si rivela più adatto rispetto al CSV perché permette una migliore strutturazione e manipolazione dei dati.
- **Preparazione del Dataset Finale:** Successivamente, abbiamo sviluppato un ulteriore script per combinare i due file JSON, alternando gli esempi etichettati (1 per esempi `jailbreak` e 0 per esempi `no-jailbreak`). Inoltre, è stato prelevato un sottoinsieme di 300 esempi per classe da utilizzare come dataset di test, sebbene tale operazione si sia poi rivelata non ottimale in quanto comportava l’inclusione di esempi già presenti nel dataset di training.

## 5 Test Preliminari e Analisi dei Risultati

Una volta completata la fase di preparazione del dataset, sono stati eseguiti i primi test di fine-tuning di BERT su Google Colab, configurati per l’uso di GPU. Le metriche ottenute, seppur solo indicative, hanno mostrato risultati sorprendentemente ottimi. Tuttavia, si sospetta che il modello stia andando in overfitting, poiché i test sono stati effettuati su esempi già visti durante la fase di training. A seguito di questi test preliminari, abbiamo avuto modo di interfacciarci con il professore per ottenere un feedback mirato e definire i prossimi step della ricerca:

- La realizzazione di un dataset di training più pulito ed affidabile.
- La definizione di un dataset di test separato per evitare sovrapposizioni e ridurre il rischio di overfitting.

## 6 Conclusioni e Prospettive Future

Riassumendo, durante questa settimana sono stati compiuti importanti passi avanti, che includono:

- L'approfondimento del processo di pulizia dei dati e l'aggregazione dei file grezzi in dataset bilanciati.
- L'implementazione dei primi test in ambiente Google Colab per il fine-tuning iniziale di BERT, con una valutazione preliminare delle performance del modello.
- Lo sviluppo e il collaudo di un complesso script Python destinato alla traduzione automatica dei testi, che costituirà il nucleo del futuro layer di traduzione.
- La conversione in formato JSON dei file di dataset e la successiva combinazione per la creazione di un dataset finale alternato, sebbene con alcuni limiti riscontrati nella fase di test.

In prospettiva, nei prossimi step ci concentreremo su:

- Il perfezionamento del processo di scrematura del dataset per eliminare dati ridondanti e non validi.
- L'ottimizzazione del modello di fine-tuning per evitare l'overfitting, tramite una revisione attenta dei dataset di training e test.
- L'implementazione di ulteriori funzionalità nel layer di traduzione per garantire una gestione ancora più robusta delle diverse lingue.

## Note Finali

Il lavoro di questa settimana ha fornito una solida base per le fasi successive del progetto, evidenziando sia i punti di forza del nostro approccio, sia le criticità da affrontare nelle prossime iterazioni. L'integrazione dei vari script Python e l'adozione di un formato dati più strutturato rappresentano elementi fondamentali per garantire la scalabilità e la riproducibilità del sistema.