

UNIVERSITÀ DEGLI STUDI DI MILANO

Facoltà di Scienze e Tecnologie

*Corso di Laurea in Informatica (L-31)*

ANALISI QUANTITATIVA E  
PERCETTIVA DI VIDEO CREATI CON  
GENERATORI AI

**Relatore:** Prof. Raffaella Lanzarotti

**Correlatore:** Prof. Andrea Gaggioli

Tesi di:

Federico COSCIA

Matricola: 977772

Anno Accademico 2023-2024

*a Celestina*



# Indice

<b>Introduzione</b>	<b>1</b>
<b>1 Generazione dei video fake</b>	<b>2</b>
1.1 Funzionamento . . . . .	2
1.2 Valutazione delle soluzioni disponibili . . . . .	2
1.2.1 DupDub . . . . .	2
1.2.2 Synthesia.io . . . . .	3
1.2.3 HeyGen . . . . .	3
1.2.4 La scelta . . . . .	4
1.3 Video generati da video scaricati dal web . . . . .	5
1.3.1 Criterio di ricerca . . . . .	5
1.3.2 Video trovati . . . . .	6
1.3.3 Processing . . . . .	6
1.3.4 Generazione dei video fake . . . . .	7
1.3.5 Scelta degli avatar . . . . .	8
1.3.6 Scelta delle voci . . . . .	8
1.3.7 Inserimento del testo . . . . .	10
1.3.8 Download dei risultati . . . . .	11
1.4 Video generati con attori . . . . .	12
1.4.1 Le limitazioni dell'approccio via web . . . . .	12
1.4.2 Scrittura dei testi . . . . .	13
1.4.3 Individuazione degli attori . . . . .	13
1.4.4 Acquisizione dei video . . . . .	14
1.4.5 Video processing . . . . .	15
1.4.6 Audio Processing . . . . .	16
1.4.7 Generazione dei video fake . . . . .	17
1.5 Pro e contro . . . . .	19
1.5.1 Video generati da video scaricati dal web . . . . .	19
1.5.2 Video generati con attori . . . . .	20
1.5.3 La scelta . . . . .	20

<b>2 Setting di acquisizione</b>	<b>22</b>
2.1 Modalità di acquisizione . . . . .	22
2.2 Estrazione delle feature . . . . .	22
2.2.1 Video . . . . .	22
2.2.2 Dati fisiologici . . . . .	22
2.2.3 Eye-tracking . . . . .	22
<b>3 Protocollo di acquisizione</b>	<b>23</b>
3.1 Stesura del protocollo . . . . .	23
3.2 Sviluppo dell’interfaccia . . . . .	23
3.3 Salvataggio dei dati raccolti . . . . .	23
<b>4 Analisi dei dati acquisiti</b>	<b>24</b>
<b>Conclusioni</b>	<b>25</b>
<b>Bibliografia</b>	<b>26</b>

# **Introduzione**

# Capitolo 1

## Generazione dei video fake

### 1.1 Funzionamento

### 1.2 Valutazione delle soluzioni disponibili

Per la generazione dei video fake sono stati valutati tre applicativi diversi, forniti come Software-as-a-Service (SaaS):

- DupDub.com
- Synthesia.io
- HeyGen.com

I criteri che sono stati valutati sono: la naturalezza dei movimenti generati, l'estensione dei movimenti generati, la possibilità di generare avatar personalizzati, la qualità del lip-sync<sup>1</sup>, la qualità e la naturalezza della voce parlata generata, e il grado di realismo generale dei video generati. Vediamo per ordine i punti di forza e di debolezza identificati di ognuno, e come si è pervenuti alla scelta finale.

#### 1.2.1 DupDub

DupDub si classifica come un prodotto "Talking-Photo". A partire da una fotografia di un persona, genera il movimento dei muscoli facciali e delle labbra per simulare il parlato. DupDub trova i suoi punti di forza nell'essere molto semplice, ma è stato valutato come troppo semplice per gli scopi di questa ricerca. La più grande limitazione è data dalla limitatezza dei movimenti, limitandosi appunto a generare solo i movimenti dei muscoli facciali, e a malapena movimenti della testa, rendendo

---

<sup>1</sup>sincronizzazione tra il movimento delle labbra di un soggetto e il suono delle parole pronunciate.

il risultato finale poco convincente e innaturale. Permette facilmente la creazione di un avatar personalizzato, a partire da un soggetto noto, ma i movimenti sono stati valutati come molto limitati e innaturali, risultando non idoneo per questa ricerca.

### 1.2.2 Synthesia.io

Rispetto al precedente, Synthesia.io fornisce avatar in mezzo busto, ed è in grado di generare movimenti del viso, della testa, e anche del corpo, producendo risultati più naturali di DupDub. Gli avatar forniti sono stati generati a partire da persone reali, ed il servizio offre la possibilità di generare dei propri avatar personalizzati. I punti di debolezza individuati sono stati: la qualità del lip-sync e la qualità delle voci generate. In particolare, è risultato frequente il disallineamento tra il movimento delle labbra dell'avatar e il suono della voce generato. La voce inoltre è stata valutata come poco espressiva e poco naturale.

Nonostante questo, tale servizio poteva essere un buon candidato per la ricerca, ma è stato scartato a causa del piano offerto, in quanto offriva un servizio ad abbonamento basato su minuti di video generati.

### 1.2.3 HeyGen

Sin dal primo sguardo, HeyGen.com si è dimostrato essere al di sopra di tutti gli altri. HeyGen si distingue dalle sue controparti supportando video con sfondi reali, e generando movimenti del corpo avanzati come il movimento delle braccia e il gesticolamento delle mani.

Per la generazione dei video, HeyGen offre due soluzioni:

1. Generazione con i modelli di avatar e di voce forniti dalla piattaforma
2. Creazione di un avatar personalizzato, che a partire da un video di riferimento clona l'aspetto e la voce della persona raffigurata, mantenendo lo sfondo raffigurato

Selezionato il metodo di generazione, si inserisce un testo di riferimento, a partire dal quale verrà generato il video fake.

#### **Generazione con i modelli di avatar e di voce forniti dalla piattaforma**

La piattaforma fornisce un catalogo di avatar e di modelli di voce già pronti per l'utilizzo. Con questi è possibile generare un video fake utilizzando soltanto un testo di riferimento. È il metodo più veloce per la generazione di video fake, poiché non ha bisogno di un video di riferimento per la creazione di un avatar ad-hoc, e permette di iniziare a generare video immediatamente. Si è, d'altra parte, limitati dall'aspetto

degli avatar forniti dalla piattaforma. Gli avatar forniti sono privi di sfondo, per cui il video generato presenta l'avatar al centro dell'inquadratura, posto su uno sfondo bianco. Per la generazione è necessario identificare l'avatar che si intende utilizzare, ed identificare il modello di voce più adatto all'avatar scelto, tra quelli forniti dalla piattaforma. La piattaforma fornisce modelli di voce compatibili con tutte le lingue del mondo, ma tra tutti i modelli forniti, i modelli in lingua inglese sono i più naturali.

### **Generazione con avatar personalizzato**

È possibile creare un avatar personalizzato a partire da un video di riferimento. A partire da tale video, la piattaforma identifica la persona raffigurata nel video, e ne crea un suo avatar. L'avatar creato non è privo di sfondo, bensì è inserito nello stesso sfondo in cui è stato registrato il video originale, aumentando il grado di realismo del video prodotto. Nella creazione di questo avatar è anche clonata la voce del soggetto rappresentato, per cui per la generazione dei video fake verrà utilizzata la voce della persona raffigurata, eliminando il problema di dover scegliere il modello di voce più adatto. Inoltre, la piattaforma si è dimostrata in grado di apprendere bene l'inflessione e l'accento della persona raffigurata, producendo risultati naturali indipendentemente dalla lingua parlata. Con questa soluzione è per cui possibile usare anche video in lingua italiana senza compromettere la qualità del risultato prodotto. C'è solo un dettaglio di cui tener conto, per la creazione di un avatar personalizzato è necessario il consenso esplicito in formato video della persona raffigurata che acconsente verbalmente l'utilizzo della sua immagine per la creazione di un avatar sulla piattaforma.

#### **1.2.4 La scelta**

Per queste ragioni, tra le opzioni valutate, HeyGen è stato valutato come il più adatto, in termini di qualità e naturalezza dei risultati prodotti, ed è stato quindi scelto come soluzione per questa ricerca. Sono state valutate entrambe le opzioni offerte da HeyGen per la generazione dei video fake, i cui approcci vengono approfonditi nella prossima sezione. Un altro fattore che sicuramente ha giocato a suo favore è stato anche il piano offerto, il quale permette di generare infiniti video durante il periodo di abbonamento, posto che questi siano sotto i cinque minuti di durata.

### **Profilo dei video fake**

Si delinea così il tipo di video che siamo in grado di generare: video raffiguranti una persona che parla, inquadrata a mezzo busto, privi di movimenti di macchina o cambi di inquadrature, e privi di animazioni o scritte che compaiono a corredo. Il video può

essere a sfondo bianco (generazione con avatar forniti dalla piattaforma) o con uno sfondo reale (generazione con avatar personalizzato).

Viste le due possibili soluzioni per la generazione dei video fake, sono state valutate due soluzioni diverse per l'acquisizione dei video reali di riferimento:

- Generazione di video fake a partire da video scaricati dal web, facendo utilizzo degli avatar già forniti dalla piattaforma per la generazione dei video fake, così facendo si ottengono soggetti diversi tra video real e fake
- Generazione di video fake a partire da video registrati con attori, realizzando avatar personalizzati così da avere lo stesso soggetto tra video real e fake

Vediamo ora i dettagli di entrambe le soluzioni, riportando l'approccio seguito, e valutando infine i pro e i contro di ogni soluzione.

## 1.3 Video generati da video scaricati dal web

### 1.3.1 Criterio di ricerca

È stata utilizzata per la ricerca dei video la piattaforma YouTube. Il criterio di ricerca usato è stato: cercare i video più simili possibili ai video fake che siamo in grado di generare, così da minimizzare le differenze tra video reali e video fake. Minimizzando le differenze tra video reali e video fake massimizziamo le possibilità che diverse percezioni dei video visualizzati siano dovute solo alla natura (reale o fittizia) dei video visualizzati e non ad altri dettagli come ambientazione, soggetto, etc. Per tali ragioni, sono stati cercati video:

- frontali, con un soggetto al centro su sfondo bianco
- con nessun movimento di macchina o cambi di inquadrature
- autodescrittivi, in altre parole non vengono utilizzate immagini, slide o grafici di supporto che vengono esplicitamente referenziati dallo speaker<sup>2</sup>
- con il minor numero di scritte o immagini che compaiono a corredo, preferibilmente nessuna
- con i sottotitoli preferibilmente inseriti a mano dall'autore del video, in modo da poter scaricare il copione associato al video più facilmente, per la generazione del video associato
- preferibilmente di durata inferiore ai cinque minuti

---

<sup>2</sup>questo perché contenuti esplicitamente referenziati dallo speaker reale non sarebbero presenti nel corrispettivo video fake, creando un'incongruenza e rendendo il video fake inefficace.

### 1.3.2 Video trovati

Durante il periodo di ricerca, sono stati trovati quattro video che soddisfano i criteri stabiliti:

- *"How to make a GREAT impression - Presentation Tips"* di Expert Academy (<https://youtu.be/lZg6H0WqPVY>)
- *"How to start a pitch or presentation"* di Dominic Colenso (<https://youtu.be/P2LwuF7zn9c>)
- *"How to start a presentation"* di Expert Academy (<https://youtu.be/LrjlW00kkws>)
- *"How to Get Over Your Fear of Public Speaking"* di Expert Academy (<https://youtu.be/So3Z93hEPDk>)

I video sono stati scaricati utilizzando il tool open source `yt-dlp` (<https://github.com/yt-dlp/yt-dlp>).

### 1.3.3 Processing

I video individuati non corrispondevano tutti perfettamente alle specifiche richieste, per cui per poterli integrare nella ricerca è stato necessario fare del pre-processing.

#### Trimming

Tutti i video individuati presentavano un introduzione e una coda al video, con musiche, scritte o elementi animati. I video individuati sono per cui stati tagliati, in modo da eliminare gli elementi non utili al nostro studio, e mantenere solamente la parte di video parlata. Per il trimming dei video è stato utilizzato il tool open source gratuito `ffmpeg` (<https://www.ffmpeg.org>), così da favorire un'operazione veloce e priva di operazioni di re-encoding ove possibile.

#### Pulizia dello sfondo

Alcuni dei video individuati presentavano alcuni elementi grafici a comparsa durante la parte parlata del video, come grafici o piccole scritte. Questo è stato valutato come accettabile visto che tali elementi non venivano referenziati esplicitamente dal speaker, e comparivano solo in sovrapposizione dello sfondo.<sup>3</sup> Questo ha permesso la rimozione di tali elementi aggiuntivi tramite una semplice operazione di video-editing, detta mascheramento.

---

<sup>3</sup>Ricordiamo che tutti i video individuati presentano uno sfondo bianco uniforme, che non cambia nel tempo.

**Mascheramento** Si identifica un fotogramma dell’immagine dove non vi sono elementi a coprire la parte dello sfondo interessata, e si salva tale fotogramma come file a parte. Questo fotogramma ”pulito” è detto *clean plate*. Dal momento che lo sfondo è statico, ovvero non cambia nel tempo, il clean plate funge da copia pulita dell’immagine, che possiamo utilizzare per coprire qualunque elemento in sovrapposizione dello sfondo. Con un qualunque programma di editing, si sovrappone il clean plate alla porzione temporale di video in cui compare l’elemento da rimuovere, ad esempio una scritta, e si effettua poi una maschera, che va a ritagliare il clean plate. Come una toppa, il clean plate mascherato copre il testo in sovra-impressione, rimuovendolo dal video. È possibile vedere un esempio di questa operazione in Figura 1.



Figura 1: Una operazione di mascheramento con clean plate in Adobe Premiere Pro

### Estrazione del testo

Per poter generare i doppioni fake, è stato estratto il testo associato al parlato presente nei video individuati. È stato utilizzato il sito web gratuito <https://downsub.com> per scaricare i sottotitoli già forniti da YouTube. La maggior parte dei video presentavano dei sottotitoli ufficiali, ovvero inseriti direttamente dagli autori dei video. Per gli altri, sono stati scaricati i sottotitoli generati automaticamente da YouTube, utilizzando quindi di fatto il motore SpeechToText integrato di YouTube.

In ogni caso, tutti i sottotitoli scaricati sono stati poi revisionati a mano per eliminare refusi, errori di battitura o di trascrizione, e per eliminare elementi non parlati o associati alle parti di video che sono state tagliate via. Questi file di sottotitolo sono tutto il necessario per generare i video fake.

#### 1.3.4 Generazione dei video fake

Dal momento che non è possibile ottenere il consenso esplicito dei soggetti rappresentati per la generazione di un avatar personalizzato, è necessario ricorrere agli avatar già forniti dalla piattaforma HeyGen per la generazione dei video fake. Tale processo prevede la selezione di un avatar tra quelli forniti dalla piattaforma, la selezione

di una voce tra i modelli TextToSpeech disponibili per generare il parlato, e infine l'inserimento del testo di riferimento.

Per ognuno dei video real individuati sono stati generati due video fake, uno con un avatar di genere maschile e uno con un avatar di genere femminile. Per la generazione di un video fake è stata seguita la seguente procedura, per ogni video real:

1. Scelta di un avatar
2. Scelta del modello di voce più adatto all'avatar scelto
3. Se non è stata trovata una coppia avatar-voce convincente tornare al passo 1 passando al prossimo avatar
4. Inserimento del testo estratto dal video real
5. Fine-tuning del testo per migliorare intonazione, pronuncia e pause
6. Revisione del risultato, ripetere il passo 5 se necessario
7. Ripetizione del processo con un avatar del genere opposto

### 1.3.5 Scelta degli avatar

Il punto di partenza per la generazione di un video fake è la scelta dell'avatar da utilizzare, ovvero la persona che verrà animata per realizzare il video parlato. La piattaforma HeyGen mette a disposizione una sua selezione di avatar proprietari, disponibili a tutti gli utenti del servizio, per realizzare i video fake. Gli avatar sono figure di persone a mezzo busto o in primo piano, prive di sfondo. È possibile vedere in Figura 2 un esempio ridotto della schermata di selezione degli avatar forniti da HeyGen. Tra gli avatar sono disponibili look molto variegati, tra cui figure in abiti formali, completi, in camice, abiti da lavoro, abiti casual, etc. Per il nostro studio, sono stati considerati avatar con un look semi-formale o casual.

La piattaforma offre anche la possibilità di realizzare un proprio avatar, a propria immagine e somiglianza, ma non è stato possibile nel nostro studio usufruire di questa feature, avendo utilizzato come video real video di terzi.<sup>4</sup>

### 1.3.6 Scelta delle voci

Come anticipato, la scelta degli avatar non è stata fatta in modo indipendente, ma è stata fatta in funzione dei modelli di voce forniti dalla piattaforma HeyGen. Difatti, anche se il video generato è visivamente impeccabile, una voce innaturale o

---

<sup>4</sup>È richiesto il consenso esplicito del soggetto rappresentato per realizzare un avatar a sua immagine.

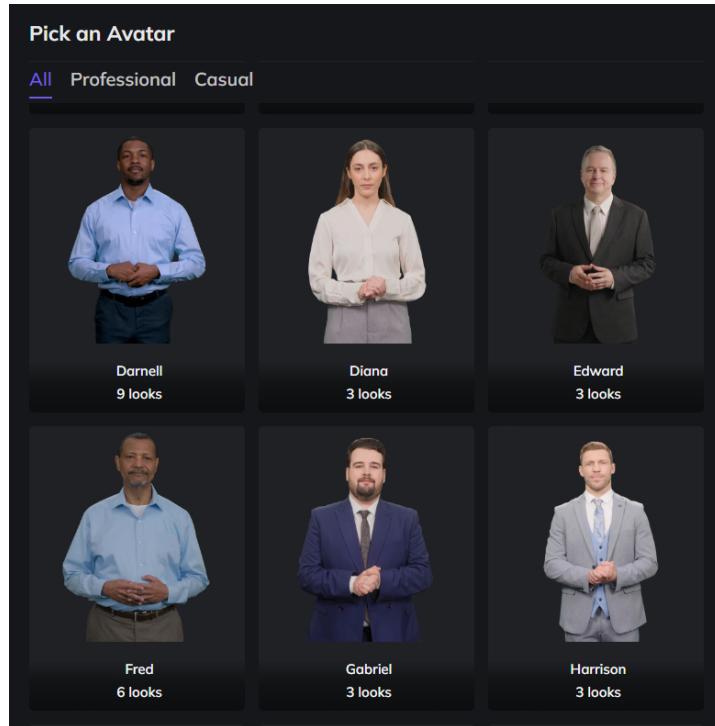


Figura 2: Esempio di schermata di selezione dell'avatar della piattaforma HeyGen

non calzante all'avatar selezionato è in grado di rompere completamente l'illusione, risultando chiaramente artificiale, o può essere un fattore di distrazione, in grado di impedire la fruizione normale del contenuto. Riconosciamo come il giudizio di una proprietà come una voce "calzante al soggetto identificato" può essere fortemente personale, così come anche fortemente umana, e meriterebbe uno studio approfondito a parte. Per i nostri scopi, la scelta è stata guidata dal giudizio umano.

### Filtraggio tramite categorie di voci

La piattaforma mette a disposizione un catalogo di voci molto variegato, suddiviso per categorie. Le categorie fornite sono visibili in Figura 3, e sono: genere (Maschio, Femmina), età (Child, Young adult, Middle-aged, Old), e "use case" (Conversazionale, Pubblicità e social, Informativo ed educativo, Narrativo). Sono state innanzitutto filtrate le voci selezionando il genere appropriato e la fascia di età appropriata per l'avatar selezionato. In aggiunta, sono state favoreggiate voci categorizzate come a scopo "Informativo ed educativo", ma se necessario sono state valutate anche voci con altri use-case.

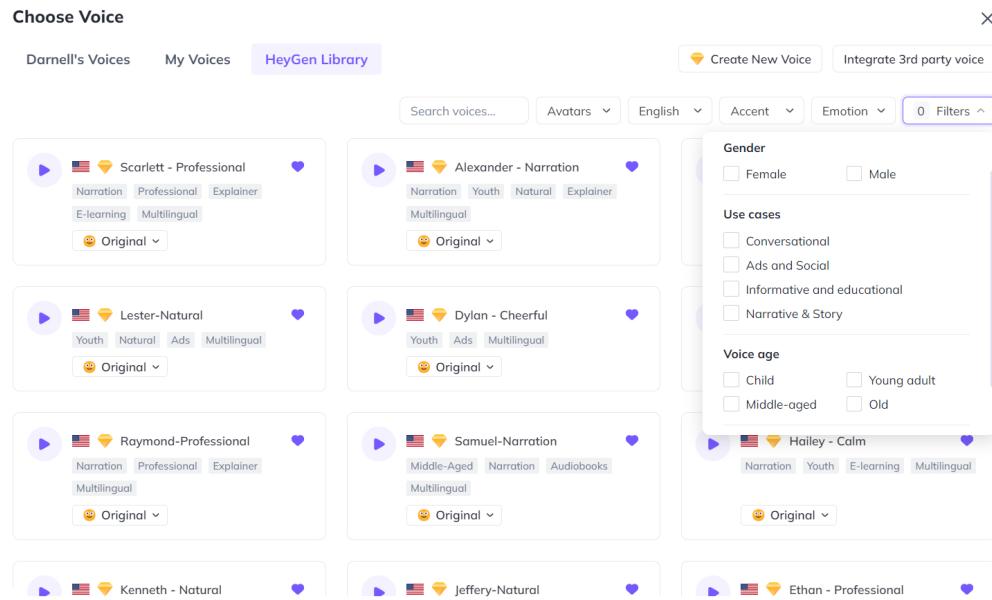


Figura 3: Schermata di selezione della voce sulla piattaforma HeyGen.

### Processo di selezione

Isolate le possibili voci candidate, è stato generato un video per ogni voce. È stata selezionata poi, tra le candidate, la voce che, visionando il video generato, al giudizio umano è parsa più naturale e convincente con l'avatar selezionato. Se nessuna voce delle voci provate tra quelle fornite dalla piattaforma HeyGen è risultata convincente, l'avatar è stato scartato.<sup>5</sup>

### Lingua

Tutti i video sono stati generati in lingua inglese poiché, sebbene HeyGen fornisca modelli di voci italiane, questi al tempo della ricerca erano limitati in numero e di qualità fortemente limitata rispetto alle controparti anglosassoni.

#### 1.3.7 Inserimento del testo

L'ultimo passaggio per generare un video fake è l'inserimento del testo da far esporre all'avatar. Nel nostro caso, si tratta del testo estratto dai video real, come spiegato in

---

<sup>5</sup>C'è da notare come con il tempo la piattaforma si è evoluta, e al tempo della scrittura di questo documento, HeyGen fornisce insieme agli avatar una pre-selezione di voci adatte all'avatar selezionato. Questo processo risulterebbe per cui molto semplificato.

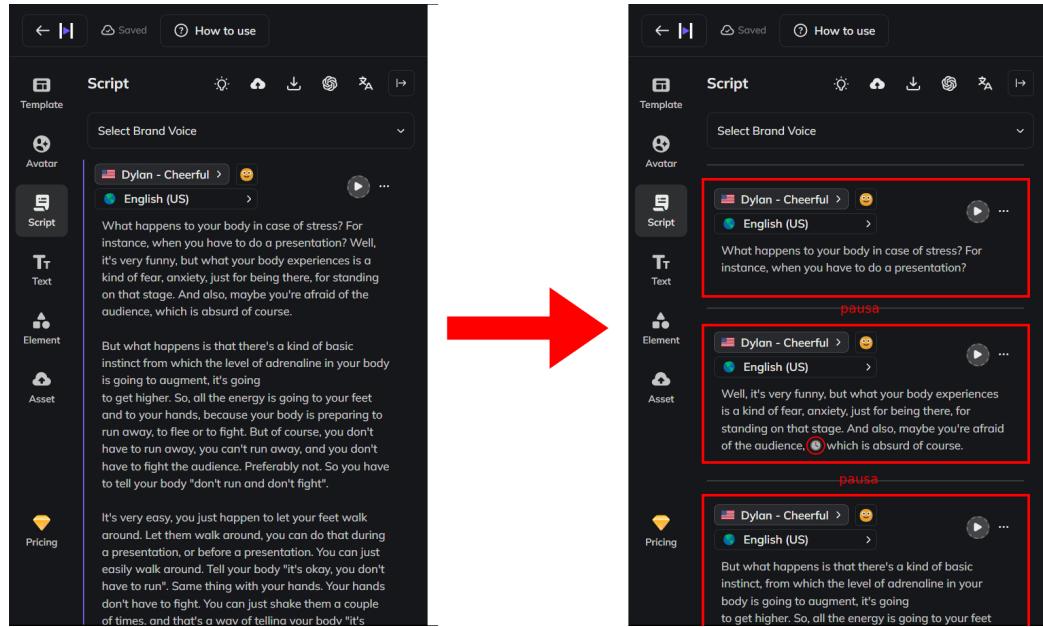


Figura 4: Separazione del testo in paragrafi per introdurre pause naturali.

1.3.3. I punti chiave di questo passaggio sono: l'introduzione di pause per un flusso naturale del discorso, e la specifica di particolari pronunce, ove necessario.

### Introduzione di pause

La piattaforma HeyGen permette di introdurre pause nel discorso in modo naturale separando il testo in "paragrafi". I paragrafi sono blocchi indipendenti di testo a partire dai quali è generata la voce. Tra un paragrafo e l'altro viene inserita automaticamente una piccola pausa, permettendo un flusso naturale del discorso. È possibile vedere un esempio di separazione in paragrafi in Figura 4.

### Revisione del risultato

A partire dal testo inserito viene generato l'audio della voce, che farà da guida per la generazione dei movimenti dell'avatar, come spiegato in 1.1. Prima di avviare la generazione del video è possibile generare un'anteprima della voce. Se non si identificano problemi di pausa o di pronuncia, si fa partire la generazione del video.

#### 1.3.8 Download dei risultati

Una volta generati i video, questi sono visualizzabili sulla piattaforma e scaricabili gratuitamente in formato 720p o 1080p. I video sono stati scaricati in formato 1080p.

## 1.4 Video generati con attori

### 1.4.1 Le limitazioni dell'approccio via web

L'approccio appena presentato è molto semplice e conveniente, in quanto non richiede la registrazione di video appositi per poter realizzare dei video fake, ma presenta delle forti limitazioni. Innanzitutto, il contenuto.

#### Complessità del contenuto

A livello contenutistico, si è limitati dai video che si è in grado di trovare in rete che soddisfano i requisiti richiesti (sfondo bianco, nessun cambio di inquadratura, ecc.). Nonostante la vastità della piattaforma YouTube, i criteri richiesti sono molto specifici, per cui i video trovati sono stati valutati molto limitati, ed in particolare il contenuto, seppur educativo, è stato valutato come di semplice comprensione. Il protocollo sperimentale stabilito prevede la somministrazione di alcune domande di comprensione sui contenuti affrontati, ma se questi sono troppo semplici, la somministrazione di tali domande diventa inefficace nel valutare l'efficacia dell'esperienza. I contenuti esposti nei video devono avere il giusto livello di complessità per essere comprensibili dallo spettatore, ma stimolarne l'attenzione, mettendo alla prova le sue capacità di comprensione.

#### Differenza di soggetti tra video reali e video fake

Un'altra limitazione, più tecnica, è la presenza di soggetti diversi tra un video reale e il video fake associato. Come già detto, individuare i video sul web riduce i tempi di acquisizione dei video, ma ci costringe a usare gli avatar forniti dalla piattaforma HeyGen per la generazione dei video fake. Non è possibile creare degli avatar personalizzati (quindi a immagine e somiglianza) dei video reali, poiché sarebbe necessario il consenso esplicito formato video dei soggetti raffigurati, di cui ovviamente non si può avere a disposizione, essendo i video stati trovati in rete. Bisogna, per cui, utilizzare gli avatar forniti dalla piattaforma HeyGen per realizzare i video. I video reali e i video fake presentano così soggetti diversi.

Seppure entrambi i video sono privi di sfondo, non hanno cambi di inquadratura, immagini o testi di corredo, ecc., non è esattamente vero che l'unica differenza tra i due tipi di video è la loro natura (reale o fittizio). Tra due video che presentano lo stesso argomento con le stesse identiche parole, se i due presentatori sono diversi vi può essere una differenza di percezione nello spettatore, in termini di livello di attenzione, interesse verso l'argomento, e fattore di comprensione degli argomenti esposti. Questo introdurrebbe una variabile esterno che potrebbe influenzare e sporcare i risultati

ottenuti, creando differenze tra i dati non dipendenti dal fenomeno che si intende studiare.

### La soluzione

Entrambi questi problemi trovano una soluzione comune: registrare dei video ad-hoc con attori. Registrando dei video è possibile avere il controllo sul livello di qualità, complessità e sull'argomento trattato nei contenuti esposti, scrivendo un copione preciso da seguire. Inoltre, è possibile fare uso degli avatar personalizzati, garantendo di avere lo stesso soggetto tra video real e video fake, poiché possiamo ottenere il consenso esplicito per la creazione degli avatar di chi si presta come soggetto per il video.

#### 1.4.2 Scrittura dei testi

La realizzazione dei video con attori parte dalla scrittura dei testi che verranno esposti. Come previsto dal protocollo sperimentale, sono stati realizzati due testi, inerenti a due argomenti diversi, di media lunghezza. Questo progetto è stato svolto in collaborazione con l'Università Cattolica del Sacro Cuore, e i testi necessari per la realizzazione dei video sono stati scritti da un tesista dell'Università Cattolica, Matteo Scarinzi, come parte del suo progetto di tesi, seguito dal Prof. Andrea Gaggioli, docente della facoltà di Psicologia della Cattolica.

#### 1.4.3 Individuazione degli attori

Una volta scritti i testi, è stato necessario individuare gli attori disposti a prestare il loro aspetto e la loro voce per registrare i video reali, a partire dai quali vengono realizzati i video fake.

Sono stati cercati un uomo e una donna, di età simile, preferibilmente con già qualche esperienza nell'esporre, spiegare, e parlare in pubblico, per facilitare le operazioni di ripresa dei video. È stato valutato come importante che i due attori avessero età simile per evitare che una chiara differenza di età potesse avere un'influenza diversa sull'esperienza di visione dei video, in termini di comprensione e grado di attenzione e coinvolgimento.

La ricerca è stata effettuata dal tesista della Cattolica Matteo Scarinzi all'interno del corpo docenti e dottorandi della facoltà di Psicologia della Cattolica, seguendo l'intuizione che docenti e dottorandi abbiano buona esperienza nell'esporre argomenti ad alta voce e parlare al pubblico. Sono stati individuati come attori:

- Michele Paleologo, per la parte maschile, dottorando della facoltà di Psicologia dell'Università Cattolica del Sacro Cuore

- Marta Pizzolante, per la parte femminile, dottoranda della facoltà di Psicologia dell’Università Cattolica del Sacro Cuore

#### 1.4.4 Acquisizione dei video

I video sono stati registrati presso la sede dell’Università Cattolica del Sacro Cuore in Via Buonarroti, 30 a Milano, nell’aula privata del Prof. Andrea Gaggioli, docente di Filosofia della Cattolica e collaboratore di questo progetto.

I video sono stati registrati con luce artificiale, così da avere il pieno controllo della luce in scena, e non dipendere dalle condizioni meteo, come il movimento delle nuvole e del sole per la quantità di luce presente in scena durante i video e tra un video e l’altro. Per compensare la luminosità ridotta delle luci della stanza è stata utilizzata una ring light<sup>6</sup>, posta davanti la telecamera, così da illuminare gli attori.

Per la acquisizione dei video è di fondamentale importanza che i testi scritti siano seguiti parola per parola dagli attori, per evitare che differenze tra diverse modalità di esposizione possano influenzare i risultati ottenuti. Per garantire ciò senza che gli attori debbano imparare a memoria i testi scritti agli attori, è stato acquistato e utilizzato un teleprompter<sup>7</sup>. Per la registrazione dell’audio è stato utilizzato un microfono lavalier wireless. Il setup utilizzato per la registrazione è visibile in Figura 6.

Nonostante la capacità del servizio di generazione dei video fake di replicare lo sfondo nei video catturati, i video sono stati registrati di fronte a uno sfondo bianco, per evitare la presenza di possibili elementi di distrazione nello sfondo. È possibile vedere in Figura 5 un esempio della sessione di registrazione.

#### Video di riferimento per la creazione degli avatar

Per la realizzazione dei video fake è stato registrato per ogni attore un video a parte, richiesto dalla piattaforma HeyGen per la creazione dell’avatar personalizzato (Instant Avatar). La piattaforma richiede un video di breve-media durata, 2-3 minuti, dove si può parlare di qualsiasi argomento. Il video è utilizzato come punto di partenza per creare l’avatar, estraendo da tale video l’aspetto del soggetto, dello sfondo, così come i movimenti del soggetto e il suono della sua voce. Tale video è stato registrato per ogni attore nello stesso identico setup, così che l’avatar realizzato fosse visivamente identico ai corrispettivi video reali. È stato necessario realizzare un video a parte poiché è richiesto dalla piattaforma HeyGen che in tale video vengano inserite frequenti pause tra una frase e l’altra. Questo porta a un flusso del discorso molto innaturale, non

---

<sup>6</sup>Lampada a forma di anello a luce bianca, molto luminosa, che permette attraverso la sua forma caratteristica di illuminare uniformemente il viso senza creare ombre dure sul volto.

<sup>7</sup>Un meccanismo a specchio che permette, a chi si pone davanti alla telecamera, di leggere il testo senza distogliere lo sguardo dalla camera.



Figura 5: Sessioni di registrazione dei video reali con attori, con: (a) Michele Paleologo nella parte maschile, (b) Marta Pizzolante nella parte femminile.

rendendo possibile adoperare i video reali per la realizzazione dell'avatar, richiedendo la registrazione di un video dedicato a parte. Tale video è stato processato allo stesso modo di tutti gli altri, così che fosse visivamente e qualitativamente uguale ai corrispettivi video reali.

#### 1.4.5 Video processing

I video registrati sono stati processati per migliorare la qualità dell'immagine. Sono stati applicati un filtro di riduzione del rumore e un aumento della luminosità dell'immagine. Per l'elaborazione dei video è stato utilizzato il programma gratuito Da Vinci Resolve.

##### Riduzione del rumore

I video sono stati registrati in condizioni di luminosità ristretta, utilizzando la luce artificiale dell'aula a disposizione e una ring light come luce di supporto. Per questo motivo, l'immagine presentava del leggero rumore, che è stato rimosso tramite un filtro di riduzione del rumore, incluso nel programma utilizzato. Il filtro di riduzione del rumore è stato controbilanciato con un filtro di sharpening<sup>8</sup>, per evitare la perdita di dettagli dovuta alla riduzione del rumore. È possibile vedere un esempio del risultato di tale processo in Figura 7<sup>9</sup>.

---

<sup>8</sup>Un aumento della nitidezza dell'immagine, per rendere i dettagli più evidenti.

<sup>9</sup>Si tratta di un esempio molto sottile, e vista la dovuta compressione dell'immagine, necessaria per l'inserimento delle figure in questo documento, e la sua successiva stampa, potrebbe non essere facilmente apprezzabile la differenza tra le due immagini.



Figura 6: Setup di registrazione dei video fake, con videocamera, teleprompter e ring light.

### Color Grading

L'immagine è stata migliorata con una correzione del bilanciamento del bianco, per correggere i colori, un leggero aumento della saturazione e della nitidezza, per correggere il profilo neutro della fotocamera utilizzata per le riprese, ed è stata aumentata la luminosità dell'immagine. È possibile vedere un esempio risultato finale in Figura 8.

I video sono stati esportati in formato .mp4, con codifica H.264 con un target di bitrate di 2048 kb/s, a 23.976 fps.

### 1.4.6 Audio Processing

L'audio registrato è stato normalizzato a un valore di loudness di -23 LUFS ( $\pm 0.5$  LU) integrati (I), secondo lo standard Europeo EBU R128 di distribuzione audio in ambito broadcast. Inoltre, l'audio è stato pulito con un filtro di riduzione del rumore di fondo, ed è stata applicata una curva di equalizzazione (Figura 9) e una leggera compressione per rendere l'ascolto più stabile, naturale e gradevole.

L'audio è stato convertito in formato AAC mono a 128 kb/s per poter essere inserito nel contenitore .mp4.



Figura 7: Esempio di filtro di riduzione del rumore dell’immagine nei video con attori.



Figura 8: Esempio della correzione dell’immagine effettuata su video con attori.

#### 1.4.7 Generazione dei video fake

La generazione dei video fake risulta molto semplificata nel caso di video realizzati con attori. Una volta ottenuto il consenso esplicito degli attori che si sono prestati per la realizzazione dei video, tramite la piattaforma HeyGen è possibile effettuare la creazione degli avatar personalizzati (Instant Avatar). Gli Instant Avatar creano una copia virtuale sia dell’aspetto che della voce degli attori, replicando in tutto e per tutto le fattezze del video originale. Non è più necessario, quindi, dover selezionare l’avatar e la voce più adatti per il video in questione.

**Creazione degli Instant Avatar** Per la creazione degli Instant Avatar, come già anticipato, è stato usato un video a parte, registrato per ogni attore, realizzato seguendo le linee guida indicate dalla piattaforma per la creazione di un Instant Avatar:

- Durata di 2-5 minuti

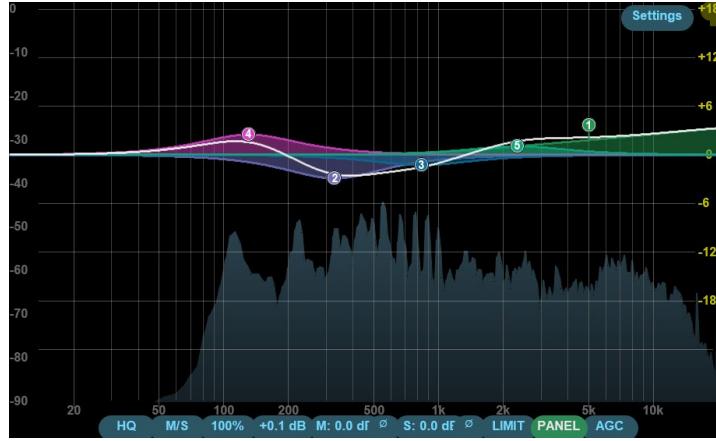


Figura 9: Curva di equalizzazione dell’audio registrato.

- Frequenti pause tra una frase e l’altra
- Si può parlare di qualsiasi argomento, ma cercando di esagerare con le emozioni, così da donare espressività all’avatar
- Uso di movimenti delle braccia generici, evitando gesti di indicazione

I video di riferimento registrati sono venuti di circa 3 minuti. Caricato il video di riferimento, la piattaforma si occupa di tutto il resto, e in una decina di minuti l’avatar è pronto per essere utilizzato.

### Creazione dei video

La creazione dei video fake è a sua volta semplificata. Creati i due Instant Avatar associati ai due attori, sono stati presi i due testi, scritti dal tesista della Cattolica Matteo Scarinzi, e sono stati caricati sulla piattaforma per la generazione dei video.

Come già visto in 1.3.7, i testi sono stati separati in sezioni, così da introdurre naturalmente delle pause nel discorso generato. Inoltre, è stata aggiunta della punteggiatura aggiuntiva e sono state modificate alcune parole dei testi, esplicitandone l’accento, in modo da correggere alcuni problemi di intonazione e pronuncia che gli avatar hanno incontrato durante la generazione del testo. Ad esempio, la parola “prosodia” è stata scritta esplicitando l’accento, “prosodìa”, così da correggerne la pronuncia, dal momento che la voce generata pronunciava la parola con l’accento sbagliato. La qualità e la precisione dell’intonazione del discorso generato dipendono fortemente dal discorso di riferimento che è stato utilizzato per la generazione dell’avatar. Seguendo l’esempio precedente, è evidente che i modelli di voce generati sono stati generati da un discorso che non conteneva la parola “prosodia” al suo interno,

lasciando al modello il compito di ricavarne la pronuncia, non garantendo ricavi quella corretta, rendendo necessario la specifica a mano in caso di errore.

### Download del risultato

Sono stati realizzati un video fake per ogni attore (maschio e femmina), e per ogni attore un video per ognuno dei due testi realizzati, per un totale di quattro video fake, associati ai quattro video reali registrati. I video fake sono stati scaricati in formato 1080p e 25fps.

## 1.5 Pro e contro

In questa sezione analizziamo i pro e i contro dei due approcci valutati per la generazione dei video:

- Video generati da video scaricati dal web (1.3)
- Video generati con attori (1.4)

### 1.5.1 Video generati da video scaricati dal web

Questo approccio è consistito nel trovare dei video validi da utilizzare come video reali sul web, e utilizzare gli avatar forniti dalla piattaforma HeyGen per creare dei video fake con lo stesso contenuto.

#### Pro

- Non è necessario scrivere dei testi appositi
- Non è necessario registrare da sé dei video, per cui non serve cercare attori, recuperare l'attrezzatura necessaria, e investire il tempo richiesto per la registrazione e la produzione dei video registrati
- Ottenimento dei video reali veloce

#### Contro

- Non si ha controllo sulla qualità, argomento, lunghezza e complessità del contenuto esposto
- La ricerca dei video può richiedere tempo e fornire risultati insoddisfacenti a seconda dei criteri di selezione richiesti e la disponibilità della piattaforma di ricerca video utilizzata

- I video reali e fake presentano soggetti diversi, che dipendono dagli avatar forniti dalla piattaforma scelta, il che può influenzare l'esperienza visiva introducendo una variabile esterna indesiderata
- Creazione dei video fake più lenta, vista la necessità di selezione dell'avatar e della voce più adatti

### 1.5.2 Video generati con attori

Questo approccio è consistito nella scrittura dei testi da esporre, la ricerca di due attori (maschio e femmina), la registrazione dei video reali, la creazione di avatar personalizzati e l'utilizzo dei testi scritti per la generazione dei video fake.

#### Pro

- I video reali e fake sono visivamente identici, eliminando qualsiasi variabile esterna, e l'unica differenza tra i due video è la loro natura (reale o fittizia)
- Si ha il totale controllo sul contenuto esposto, in termini di qualità, complessità, lunghezza e argomento, essendo questo scritto appositamente
- Generazione dei video fake molto più veloce, in quanto non serve selezionare l'avatar e la voce più adatta, essendo l'aspetto e la voce clonati dalla piattaforma di generazione dei video

#### Contro

- Acquisizione dei video reali molto più lenta, in quanto richiede la scrittura dei contenuti, la ricerca di attori per la registrazione dei video, e la registrazione e post-produzione dei video stessi
- La qualità dei video fake realizzati dipende dalla qualità dei video registrati, in termini di qualità video/audio e in termini di qualità dell'avatar personalizzato creato

### 1.5.3 La scelta

Visti i pro e i contro di ogni approccio, è stato scelto per la ricerca il secondo approccio, e quindi sono stati utilizzati i video realizzati con attori. Nonostante il tempo non indifferente richiesto dalla scrittura, acquisizione e produzione dei video reali, è stato valutato di estrema importanza avere lo stesso soggetto tra video reale e video fake, così da eliminare qualunque tipo di bias dovuto a soggetti diversi, e poter essere sicuri

che eventuali risultati diversi ottenuti da i due tipi di video (reali o fake) possano essere causati solo da tale differenza, e non da fattori esterni.

# **Capitolo 2**

## **Setting di acquisizione**

**2.1 Modalità di acquisizione**

**2.2 Estrazione delle feature**

**2.2.1 Video**

**2.2.2 Dati fisiologici**

**2.2.3 Eye-tracking**

# **Capitolo 3**

## **Protocollo di acquisizione**

**3.1 Stesura del protocollo**

**3.2 Sviluppo dell'interfaccia**

**3.3 Salvataggio dei dati raccolti**

# **Capitolo 4**

## **Analisi dei dati acquisiti**

# **Conclusioni**

# Bibliografia

# **Ringraziamenti**