

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 Lanza

**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
1.3 Video real . . . . .	4
1.3.1 Scelta dei video real . . . . .	4
1.3.2 Processing . . . . .	5
1.4 Generazione dei video fake . . . . .	6
<b>2 Setting di acquisizione</b>	<b>7</b>
2.1 Modalità di acquisizione . . . . .	7
2.2 Estrazione delle feature . . . . .	7
2.2.1 Video . . . . .	7
2.2.2 Dati fisiologici . . . . .	7
2.2.3 Eye-tracking . . . . .	7
<b>3 Protocollo di acquisizione</b>	<b>8</b>
3.1 Stesura del protocollo . . . . .	8
3.2 Sviluppo dell'interfaccia . . . . .	8
3.3 Salvataggio dei dati raccolti . . . . .	8
<b>4 Analisi dei dati acquisiti</b>	<b>9</b>
<b>Conclusioni</b>	<b>10</b>
<b>Bibliografia</b>	<b>11</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 di questi ultimi, la qualità del lip-sync<sup>1</sup>, la qualità e la naturalezza della voce parlata generata, e il grado di realismo generale dato dai 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 movimenti dei muscoli facciali, e a malapena movimenti della testa, rendendo il

---

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

risultato finale poco convincente e innaturale. Inoltre, tutti gli avatar forniti dalla piattaforma per la generazione dei video sono chiaramente soggetti non reali, bensì generati a loro volta tramite IA.

### 1.2.2 Synthesia.io

Rispetto al precedente, Synthesia.io si mostra molto più capace. Fornisce avatar in mezzo busto, ed è in grado di generare movimenti del viso, della testa, e anche del corpo, producendo risultati decisamente più naturali di DupDub. Gli avatar forniti sembrano essere stati generati a partire da persone reali, ed il servizio offre la possibilità di generare avatar personali. I punti di debolezza individuati sono stati: la qualità del lip-sync e la qualità delle voci generate. In particolare, risultava frequente il disallineamento tra il movimento delle labbra dell'avatar e il suono della voce generata. La voce inoltre è stata valutata come poco espressiva e poco naturale. Vedremo in realtà come questi sono spesso i punti più deboli di questa tecnologia.

Nonostante questo, tale servizio sembrava un buon candidato per la ricerca, ma è stato scartato in base al piano offerto, in quanto offriva un servizio ad abbonamento basato su minuti.

### 1.2.3 HeyGen

Sin dal primo sguardo, HeyGen.com si è dimostrato essere al di sopra di tutti gli altri, offrendo anche la possibilità di generare video su sfondi reali, angolazioni diverse dello stesso avatar, e implementando movimenti del corpo avanzati come il movimento delle braccia e il gesticolamento delle mani. Gli avatar sono costruiti a partire da un video di riferimento del soggetto, il che li conferisce la possibilità di apprendere ed emulare i movimenti della persona inquadrata, producendo un risultato più naturale e realistico. Inoltre, la piattaforma si è dimostrata essere in costante evoluzione e sviluppo, arricchendo il suo catalogo di funzionalità, avatar e di voci durante il periodo di valutazione.

Per queste ragioni, tra le opzioni valutate, HeyGen è stato decretato come il migliore, in termini di qualità e naturalezza dei risultati prodotti, ed è stato quindi scelto come soluzione per la nostra ricerca. Un altro fattore che sicuramente ha giocato a suo favore è stato anche il piano offerto, il quale ci ha permesso di generare infiniti video durante il periodo di abbonamento, posto che questi fossero sotto i cinque minuti di durata.

**Profilo dei video fake** Si delinea così il tipo di video che siamo in grado di generare: gli avatar sono privi di sfondo, per cui i video fake sono video con sfondo bianco,

statici, privi di movimenti di macchina o cambi di inquadrature, e sono privi di animazioni o scritte che compaiono a corredo.

## 1.3 Video real

Identificata la piattaforma per la generazione dei video fake, è necessario procurarsi i video reali da utilizzare come riferimento per generare tali doppioni fake. Il tipo di video richiesto per i video real sono dei video brevi (durata inferiore ai 5 minuti), in modo che possano essere visti per intero. Fondamentale è trovare dei contenuti che non richiedano una particolare formazione pregressa per la loro comprensione, ma che non siano neanche banali, in modo di poter fare delle domande di comprensione a loro volta non banali.

### 1.3.1 Scelta dei video real

**Criterio di scelta** È stata utilizzata per la ricerca dei video la piattaforma YouTube, e il criterio di ricerca dei video reali è semplice: cercare i video più simili possibili ai video fake che siamo in grado di generare. Per tali ragioni, il criterio di scelta consiste in:

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

**Video trovati** Per la realizzazione di un pilot della ricerca, sono stati identificati quattro video che soddisfano i criteri stabiliti:

- *"How to make a GREAT impression - Presentation Tips"* di Expert Academy
- *"How to start a pitch or presentation"* di Dominic Colenso

---

<sup>2</sup>questo perché contenuti esplicitamente referenziati dallo speaker reale non sarebbero presenti nel corrispettivo video fake, rompendo l'illusione.



- *"How to start a presentation"* di Expert Academy
- *"How to Get Over Your Fear of Public Speaking"* di Expert Academy

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

### 1.3.2 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 ed 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 inoltre 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 dallo 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.

**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 che compare 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. Si effettua poi una maschera, che va a ritagliare il *clean plate*, in modo che, come una toppa, vada a coprire il testo in sovra-impressione, rimuovendolo dal video.

---

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



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.

## 1.4 Generazione dei video fake

# 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