

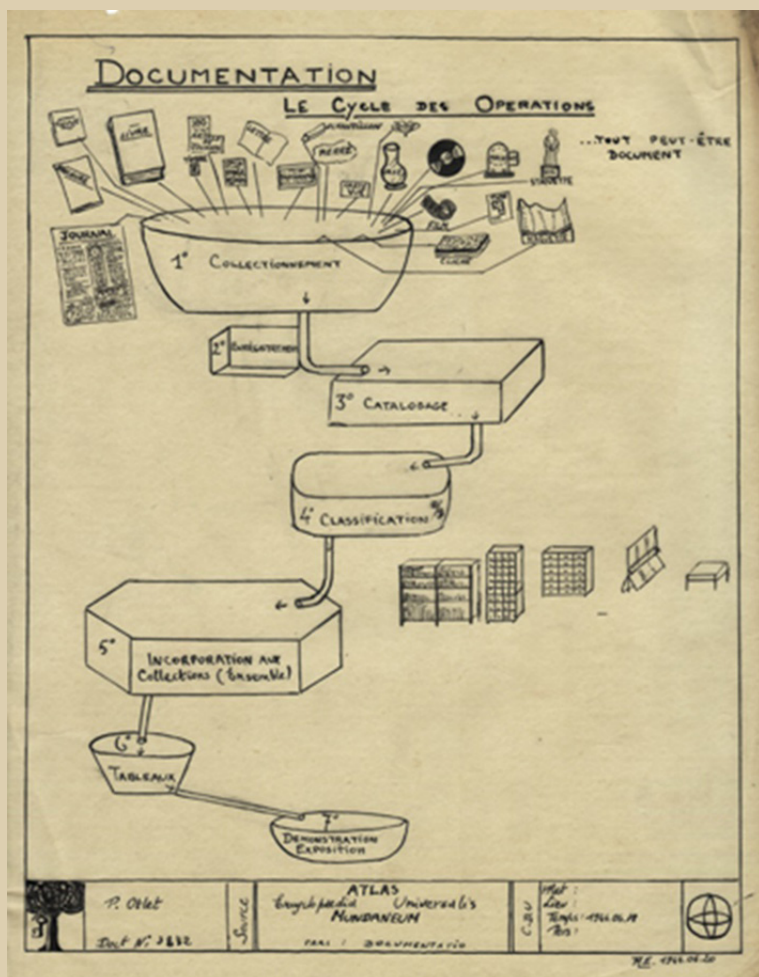
# AIDa informazioni

RIVISTA SEMESTRALE DI SCIENZE DELL'INFORMAZIONE

NUMERO 1-2

ANNO 42

GENNAIO-GIUGNO 2024



# **AIDA**informazioni

RIVISTA SEMESTRALE DI SCIENZE DELL'INFORMAZIONE

Fondata nel 1983 da Paolo Bisogno

**Proprietario della rivista:**

Università della Calabria

**Direttore Scientifico:**

Roberto Guarasci, *Università della Calabria*

**Direttore Responsabile:**

Fabrizia Flavia Sernia

**Comitato scientifico:**

Anna Rovella, *Università della Calabria*;

Maria Guercio, *Sapienza Università di Roma*;

Giovanni Adamo, *Consiglio Nazionale delle Ricerche* †;

Claudio Gnoli, *Università degli Studi di Pavia*;

Ferruccio Diozzi, *Centro Italiano Ricerche Aerospaziali*;

Gino Roncaglia, *Università della Toscana*;

Laurence Favier, *Université Charles-de-Gaulle Lille 3*;

Madjid Ihadjadene, *Université Vincennes-Saint-Denis Paris 8*;

Maria Mirabelli, *Università della Calabria*;

Agustín Vivas Moreno, *Universidad de Extremadura*;

Douglas Tudhope, *University of South Wales*;

Christian Galinski, *International Information Centre for Terminology*;

Béatrice Daille, *Université de Nantes*;

Alexander Murzaku, *College of Saint Elizabeth, USA*;

Federico Valacchi, *Università di Macerata*.

**Comitato di redazione:**

Antonietta Folino, *Università della Calabria*;

Erika Pasceri, *Università della Calabria*;

Maria Taverniti, *Consiglio Nazionale delle Ricerche*;

Maria Teresa Chiaravallotti, *Consiglio Nazionale delle Ricerche*;

Assunta Caruso, *Università della Calabria*;

Claudia Lanza, *Università della Calabria*.

**Segreteria di Redazione:**

Valeria Rovella, *Università della Calabria*

**Editrice:** Cacucci Editore S.a.s.

Via D. Nicolai, 39 – 70122 Bari (BA)

[www.cacuccieditore.it](http://www.cacuccieditore.it)

e-mail: [riviste@cacuccieditore.it](mailto:riviste@cacuccieditore.it)

Telefono 080/5214220



# AIDAinformazioni

## RIVISTA SEMESTRALE DI SCIENZE DELL'INFORMAZIONE

«AIDAinformazioni» è una rivista scientifica che pubblica articoli inerenti le Scienze dell'Informazione, la Documentazione, la Gestione Documentale e l'Organizzazione della Conoscenza. È stata fondata nel 1983 quale rivista ufficiale dell'Associazione Italiana di Documentazione Avanzata e nel febbraio 2014 è stata acquisita dal Laboratorio di Documentazione dell'Università della Calabria. La rivista si propone di promuovere studi interdisciplinari oltre che la cooperazione e il dialogo tra profili professionali aventi competenze diverse, ma interdipendenti. I contributi possono riguardare topics quali Documentazione, Scienze dell'informazione e della comunicazione, Scienze del testo e del documento, Organizzazione e Gestione della conoscenza, Terminologia, Statistica testuale e Linguistica computazionale e possono illustrare studi sperimentali in domini specialistici, casi di studio, aspetti e risultati metodologici conseguiti in attività di ricerca applicata, presentazioni dello stato dell'arte, ecc.

«AIDAinformazioni» è riconosciuta dall'ANVUR come rivista di Classe A per l'Area 11 – Settore 11/A4 e censita per le Aree 10 – Scienze dell'antichità, filologico-letterarie e storico-artistiche; 11 – Scienze storiche, filosofiche, pedagogiche e psicologiche; 12 – Scienze giuridiche; 14 – Scienze politiche e sociali, così come dall'ARES (Agence d'évaluation de la recherche et de l'enseignement supérieur) che la annovera tra le riviste scientifiche dell'ambito delle Scienze dell'Informazione e della Comunicazione. La rivista è, inoltre, indicizzata in: ACNP – Catalogo Italiano dei Periodici; BASE – Bielefeld Academic Search Engine; ERIH PLUS – European Reference Index for the Humanities and Social Sciences – EZB – Elektronische Zeitschriftenbibliothek – Universitätsbibliothek Regensburg; Gateway Bayern; KVK – Karlsruhe Virtual Catalog; Letteratura Professionale Italiana – Associazione Italiana Biblioteche; The Library Catalog of Georgetown University; SBN – Italian union catalogue; Summon™ – by SerialsSolutions; Ulrich's; UniCat – Union Catalogue of Belgian Libraries; Union Catalog of Canada; LIBRIS – Union Catalogue of Swedish Libraries; Worldcat.

I contributi sono valutati seguendo il sistema del *double blind peer review*: gli articoli ricevuti dal comitato scientifico sono inviati in forma anonima a due referee, selezionati sulla base della loro comprovata esperienza nei topics specifici del contributo in valutazione.

# AIDAinformazioni

Anno 42

N. 1-2 – gennaio-giugno 2024

CACUCCI  EDITORE  
BARI

---

PROPRIETÀ LETTERARIA RISERVATA

---

© 2024 Cacucci Editore – Bari

Via Nicolai, 39 – 70122 Bari – Tel. 080/5214220

<http://www.cacuccieditore.it> e-mail: [info@cacucci.it](mailto:info@cacucci.it)

Ai sensi della legge sui diritti d'Autore e del codice civile è vietata la riproduzione di questo libro o di parte di esso con qualsiasi mezzo, elettronico, meccanico, per mezzo di fotocopie, microfilms, registrazioni o altro, senza il consenso dell'autore e dell'editore.

# Sommario

## Editoriale

ROBERTO GUARASCI, Quarant'anni di «AIDAinformazioni» (1984-2024) 9

## Contributo su invito

PIERO INNOCENTI, Que reste-t-il de nos amours? Arti tradizionali di un possibile moderno Trivio: Archivistica, Bibliografia, Documentazione 17

## Contributi

FRANCESCO AMATO, ANTONELLA BENVENUTO, SILVIA CENITI, MARIA TERESA CHIARAVALLOTTI, CLAUDIA LANZA, ERIKA PASCERI, Indagine e analisi comparativa delle procedure di codifica nella Terapia del Dolore in Italia 59

ANDREA CAPACCIONI, Sull'affidabilità delle bibliografie generate dai chatbot. Alcune considerazioni 81

ALESSANDRO MAISTO, La dimensione Testuale del Videogioco. Classificazione dei transcript dei videogiochi basata sul lessico 95

ALEXANDER MURZAKU, PONTISH YERAMYAN, CURT ANDERSON, STEVEN BUXBAUM, RUBEN DIAZ, MARIELLE LERNER, ARMENUI MINASYAN, HAZEL MITCHLEY, JODIE-ANN PENNANT, MIA SHANG, BRISA SPEIER BRITO, Discovering and documenting brilliance. A novel multimodal annotation method 117

CAMILLA NAPPI, Le document : miroir des enjeux institutionnels et de l'évolution linguistique des transitions écologique et énergétique en France (2007-2022) 145

## Testimonianze

PIERO CAVALERI, FERRUCCIO DIOZZI, La Documentazione in Italia e il ruolo di Paolo Bisogno: una rapida evoluzione 163

PAOLA GARGIULO, LUCIA MAFFEI, Domenico (Ingo) Bogliolo. Profilo di un navigatore nell'Information Science 169





# Contributi



# La dimensione Testuale del Videogioco

Classificazione dei *transcript* dei videogiochi basata sul lessico

Alessandro Maisto\*

**Abstract:** In this work, we explore the textual dimension of video games. Despite their pronounced visual and interactive characteristics, video games can be regarded as documents due to their narrative and communicative elements. Our research delves into this textual dimension to automatically generate rating tags associated with offensive language, violence, and the presence of drugs. We utilized a dictionary of English slang, compiled from various online sources and manually annotated with four categories: Slang, Violence, Drugs, and Discrimination. The resulting electronic dictionary facilitated the automatic assignment of the three rating tags with high precision. It has also been employed to classify video games based on their lexical content. The two classification tasks – by rating tags and by lexical dimension – could pave the way for an automatic warning system capable of analyzing the full textual dimension of a video game.

**Keywords:** Videogame transcript corpus, Automatic videogame rating, Bad language, Slang dictionary, Text analysis.

## 1. Introduzione

In questo lavoro proponiamo uno studio computazionale della dimensione testuale dei videogiochi, esplorando, in particolare, il potere classificatorio del linguaggio volgare, violento e discriminatorio associato allo slang presente nei dialoghi dei videogiochi in inglese. In particolare, si offrirà una analisi approfondita della dimensione testuale dei videogiochi, analizzata in ottica di classificazione automatica del rating di età tramite un approccio di estrazione terminologica, basato sullo slang.

Per definizione, il videogioco è un prodotto multimediale interattivo. Dal punto di vista puramente scientifico, le indagini sui videogiochi si sono concentrate sulla loro dimensione sociologica, sul design e sulla loro costruzione ed in parte anche sugli aspetti psicologici della loro fruizione. Tuttavia, in

---

\* Dipartimento di Scienze Politiche e della Comunicazione, Università degli Studi di Salerno, Fisciano (SA), Italia. amaisto@unisa.it

quanto prodotto multimediale, il videogame può essere considerato anche un documento (Aarseth 2004; Juul 2005), provvisto di una forte correlazione con immagini, suoni e altri elementi extra-linguistici, ma con una chiara e spiccata dimensione testuale.

Il testo, all'interno di un videogame, è diventato sempre più centrale man mano che la complessità di questi prodotti di intrattenimento (ma anche culturali) è andata crescendo. Volendo identificare i momenti di evoluzione del fenomeno videoludico negli ultimi 50 anni, cioè dalla loro nascita fino ad oggi, possiamo identificare diverse fasi durante le quali il testo ha assunto una certa importanza: benché i primi videogiochi fossero sostanzialmente costruiti grafici, dove l'intrattenimento era basato principalmente nello sviluppare una coordinazione occhio/mano tale da sconfiggere la macchina, già nei primi anni si svilupparono una quantità di giochi di natura puramente testuale. Questi giochi, definiti *avventure testuali* iniziarono a diffondersi negli anni Ottanta anche sulla scia del successo dei cosiddetti *giochi di ruolo* (GDR, o RPG dall'inglese *Role-Playing Game*) (Heritage 2021). Consistevano, sostanzialmente, nella descrizione testuale di una situazione e nella richiesta di un comando testuale specifico per risolvere la situazione.

Successivamente, questo tipo di gioco si è evoluto, trasformandosi in quelle che furono definite *avventure grafiche*, dove il testo era ancora un elemento centrale del gioco, ma veniva accompagnato da un meccanismo *punta e clicca* localizzato in un riquadro grafico (Marchiori et al. 2011).

Con lo sviluppo delle *console* e con l'avanzamento dei sistemi informatici ed il conseguente arricchimento dei comparti grafici, le componenti extra-linguistiche presero il sopravvento, ma senza poter mai escludere completamente la dimensione linguistica. Trattandosi di esperienze simulate, narrative e comunicative, infatti, il linguaggio non poteva che conservare la propria importanza, soprattutto in generi come il *gioco di ruolo*.

Ridotto alla propria dimensione testuale (i cosiddetti *transcript*) il videogame deve essere considerato un documento e può essere analizzato da un punto di vista linguistico attraverso una analisi computazionale.

In quanto documento, il videogame ha delle caratteristiche che lo distinguono da tutte le altre tipologie di documenti, inclusi quelli più simili legati, ad esempio, al linguaggio dei prodotti audiovisivi. Uno di questi elementi è, indubbiamente, l'interattività, intesa come la possibilità, da parte del giocatore, di scegliere una di varie possibili frasi in risposta ad un dialogo.

L'importanza della dimensione linguistica dei videogiochi è attestata anche dalla presenza, all'interno degli indicatori utilizzati per attribuire un *rating* di età minima del giocatore, di elementi di chiara matrice testuale. Uno di questi, ad esempio, riguarda il linguaggio volgare (*Bad Language*) come proposto dal sistema Europeo *Pan European Game Information* (PEGI), ma anche la discriminazione, che viene espressa soprattutto a livello linguistico o l'uso di Dro-

ghe. Tra questi indicatori, quello del linguaggio volgare è, inoltre, tra quelli più comunemente riscontrati assieme all'indicatore di violenza.

In questo contesto, dunque, possiamo sintetizzare gli obbiettivi dell'articolo come segue:

1. Presentare una risorsa in lingua inglese per lo studio della dimensione testuale dei videogame comprendente un dizionario elettronico dello slang ed un corpus composto da 51 *transcript* di videogiochi per un totale di oltre 700.000 token;
2. Proporre un approccio classificatorio automatico introduttivo che possa dimostrare l'efficacia dell'analisi linguistica nell'attribuzione di etichette di *warning* ai videogiochi.

Nel secondo paragrafo illustreremo brevemente il funzionamento del sistema PEGI ed i suoi principali difetti. Successivamente, nel paragrafo 3, parleremo degli studi effettuati sulla dimensione testuale del videogioco. Il paragrafo 4 descriverà nel dettaglio la risorsa sviluppata, partendo dal dizionario elettronico per poi arrivare al corpus. Il paragrafo 5 illustrerà la metodologia di classificazione adottata e l'esperimento effettuato sul corpus. In fine, nel paragrafo 6 concluderemo mostrando possibili sviluppi e implicazioni del nostro lavoro.

## 2. Il Sistema PEGI

I sistemi di valutazione dei videogiochi (*Video Game Rating Systems*) rappresentano un indicatore fondamentale per ragioni sociali, psicologiche, mediche e per aspetti legali. La multidisciplinarietà delle implicazioni dell'attribuzione di questo *rating* è attestata anche dalla varietà di settori di provenienza dei professionisti che si occupano di determinare le metodologie di valutazione. Le valutazioni, poi, sono di competenza dei produttori dei videogiochi che affidano a *tester* e ad esperti di giurisprudenza una eventuale revisione del *rating*.

Il Sistema di valutazione PEGI prevede, infatti, che i produttori propongano una lista di descrittori per il proprio prodotto sotto forma di risposte ad un questionario contenente 37 domande. Un sistema automatico individua nelle risposte la presenza di contenuti inappropriati proponendo l'aggiunta o la rimozione di alcuni indicatori. Alla fine di questo procedimento automatico avviene una fase di verifica manuale che può ulteriormente modificare le etichette attribuite al gioco (Felini 2015). Solo in poche occasioni, il software e i suoi *transcript* vengono effettivamente analizzati da un'audience specializzata e diversificata.

Il sistema PEGI utilizza una lista ristretta di etichette che include Linguaggio volgare, Paura, Sesso/nudità, Droghe, Discriminazione, Gioco d'azzardo, Violenza e Acquisti in gioco (dal 2018).

Le etichette di età sono numeriche e indicano l'età minima consigliata per un videogioco. Queste etichette sono cinque e includono 3, 7, 12, 16 e 18 anni.

Sottoporre i propri prodotti ad un *rating systems* come PEGI favorisce i produttori, che mostrano il proprio interesse per la salvaguardia dei clienti, evitando di rilasciare prodotti di dubbia moralità. Inoltre, trattandosi di un sistema europeo, si riduce il rischio di incorrere in sanzioni dai singoli stati dell'Unione. Infatti, nel mondo, esistono diversi organismi nazionali o sovranazionali che valutano l'impatto dei videogiochi e ne regolano l'utilizzo, e ciascuno di essi tende ad utilizzare indicatori diversi e a proporre punteggi diversi. Un videogioco come *The Sims 4* ad esempio, viene classificato da PEGI come 12+ a causa di "violenza non realistica tra umani", mentre è considerato 18+ da RARS, che è l'organismo che opera in Russia, *Matures (violence and sexual references)* da ACB (Australia) e *Teens (obscene humor, sexual themes, violence)* da ESRB (Nord America).

Tra le critiche principali a questi sistemi, tuttavia, possiamo annoverare il fatto che i *rating* sono attribuiti sulla base di una conoscenza parziale del prodotto. Diversi studiosi del settore hanno suggerito una revisione dei processi e dei criteri di valutazione, favorendo un approccio più dettagliato e preciso (Dogruel e Joeckel 2013; Felini 2015; Humphreys e Wang 2017).

Un approccio testuale, ad esempio, favorirebbe un'analisi completa e rapida dei *transcripts* portando alla luce possibili elementi importanti ai fini della valutazione del gioco sia dal punto di vista dall'uso di linguaggio volgare che da altri punti di vista (identificazione di indicatori testuali di violenza e nomi di droghe). L'uso di strumenti di questo tipo non dovrebbe sostituire un'analisi manuale del gioco, ma suggerire la presenza di tali indicatori al valutatore e favorire la loro ricerca nell'intero documento, e non solo in parziali campioni di testo.

### 3. Contesto e studi correlati

I principali studi sul linguaggio e sui videogiochi riguardano il loro uso nell'apprendimento di una lingua straniera (Klimova e Jaroslav 2017). In questo senso, però, numerosi studi si sono concentrati su giochi sviluppati appositamente per l'apprendimento e dunque per il cosiddetto *Game-Based Learning* (Wright, Betteridge, and Buckb 2006).

Esiste, tuttavia, un filone afferente alla Linguistica dei Corpora che considera il linguaggio del videogioco come base per studi linguistici e sociolinguistici di specifici fenomeni. Ne è un esempio il lavoro di Ray (2019) che ha analizza-

to le traduzioni della terminologia fantascientifica all'interno dei videogiochi. Fekete e Porkolab (2019) hanno utilizzato una metodologia computazionale per l'analisi dei toponimi nella serie di videogiochi *The Elder Scrolls*. In particolare, i due autori hanno esaminato oltre 1.000 nomi di luoghi inventati da un punto di vista morfologico ed onomastico.

Heritage (2021) ha, invece, studiato il problema del genere all'interno di diversi videogiochi, dimostrando l'importanza delle rappresentazioni di genere nel dominio e la varietà di sfumature che vengono attribuite a personaggi maschili e femminili in un corpus di oltre 300 mila token estratti da 10 diversi videogiochi.

Ensslin (2012), come anche Gee (2014), ha analizzato i videogiochi dal punto di vista dell'analisi del discorso. Ensslin, in particolare, ha dimostrato come questi prodotti multimediali funzionino come mezzo e obiettivo di comunicazione, favorendo la nascita di nuovi vocabolari, nuovi significati, generi testuali e pratiche discorsive.

Goorimoorthee (2019) ha analizzato gli accenti dei personaggi non giocatori dei videogiochi allo scopo di verificare l'assenza di diversità culturale nei giochi della *BioWare*, uno studio di produzione che ha fatto dell'inclusività sociale dei propri titoli un elemento distintivo.

A mettere in relazione slang e videogiochi c'è il lavoro di Balteiro (2019), che però analizza le conversazioni tra giocatori e non la dimensione testuale del videogioco. Ensslin e Finnegan (2019), invece si sono concentrati sull'uso del linguaggio volgare nelle cosiddette comunità "co-sit" (co-situate) di giocatori.

L'approccio proposto nel presente articolo si ispira fortemente al lavoro di Maisto et al. (2021), proponendo, tuttavia, una serie di modifiche che hanno condotto al miglioramento dei risultati ottenuti. In particolare, in Maisto et al. le etichette di *Bad Language*, *Violence* e *Drugs* erano estratte sulla base della frequenza relativa, mentre in questo lavoro proporremo un approccio più efficace basato esclusivamente sul conteggio dei termini. Infatti, l'idea alla base di queste modifiche è che, indipendentemente dalla dimensione del testo, la sola presenza di poche espressioni volgari, violente o discriminatorie possa rappresentare un indicatore sufficiente all'attribuzione dell'etichetta corrispondente. Sono stati aggiunti, inoltre, una serie di termini referenti alla sfera sessuale.

#### 4. Risorse linguistiche

Lo slang rappresenta una sfida nella lessicografia ed in particolare nello sviluppo di dizionari specifici. James B. McMillan (1978) considera il problema del lessico slang come un problema riguardante la sua stessa definizione. La varietà di etichette e l'impossibilità di suddividere lo slang in sottocategorie rende impossibile una sua analisi di tipo linguistico, sociolinguistico o psico-

linguistico. Ciononostante, negli anni ne sono state date varie definizioni e sono stati prodotti diversi studi sul tema.

Benché sia stato spesso considerato un fenomeno marginale o bollato come volgare (De Klerk 1990), lo slang ha una connotazione sociopsicologica che emerge nelle varie definizioni che ne sono state date nel corso degli anni. Citando Amari (2010), Walt Whitman nel 1885 definì questo registro linguistico come «un tentativo dell'umanità di fuggire dal puro letteralismo e di esprimersi senza limiti». Più recentemente, il *Concise Oxford Dictionary* (2006) ha definito lo slang come «un linguaggio informale più comunemente parlato che scritto che è tipicamente ristretto ad un particolare contesto o gruppo». L'*Historical Dictionary of American Slang* (Lighter 1994) sottolinea come il suo carattere irriverente faccia dello slang uno strumento di opposizione all'autorità e alla conformità linguistica, non per forza di natura sovversiva, ma utile ad indicare un'attitudine comune che permette il riconoscimento di un individuo all'interno di un particolare gruppo sociale (Amari 2010).

Partridge (1933) elencava diverse ragioni dell'utilizzo dello slang, tra cui il gioco o il divertimento, la sdrammatizzazione, un uso per differenziarsi, per essere pittoreschi, per brevità, per arricchire il linguaggio, per rendere più concreto l'astratto, per intimare con qualcuno, mostrando la propria appartenenza ad un gruppo e per comunicare nel codice segreto di quel gruppo. Per queste sue caratteristiche lo slang tende naturalmente ad esprimere concetti legati a criminalità, stupefacenti o sesso, utilizzando spesso un tono "indecoroso" (Dumas e Lighter 1978). Secondo Green (2011) oltre 5.000 parole appartenenti al lessico dello slang riguardano la criminalità, mentre quelle che identificano le droghe sono quasi 4.000. Numerosissimi sono anche i termini razzisti (circa 1.500, di cui la maggior parte riferita alle persone di colore o agli ebrei) e termini che riguardano la morte (oltre 1.000).

Lo slang, inoltre, è un tipo di linguaggio che fa dell'ambiguità, della metafora e del cambiamento di significato il proprio elemento distintivo (Mattiello 2005). Queste caratteristiche lo avvicinano al linguaggio dell'arte e della poesia (Amari 2010). Lo slang annovera termini di uso più o meno comune utilizzati in senso figurato e spesso dispregiativo. Questi usi nascono da una "esuberanza linguistica" che porta all'introduzione di nuove parole o di nuovi significati a parole già esistenti (Mencken 1967). Ne sono un esempio emblematico numerosi nomi di droghe (*coke* per *cocaina*; *baby* per *marijuana*, ecc.), relativi al sesso (*banana* per *pene*, *backdoor* per *ano*, *fruit* per *omosessuale*), alle bevande alcoliche (*jar*, *piss* o *belt*). Questa ambiguità, che rende oscuro, indeterminato o poco trasparente il significato di queste parole, deriva da meccanismi associativi complessi ed eterogenei, che spaziano dalla metonimia (*jar* per bevanda alcolica), alla metafora o da processi più oscuri che coinvolgono altri slang (*belt* è associato ad esempio al termine "a hit" che indica metaforicamente un cattivo post-sbornia) (Mattiello 2005).



Le caratteristiche qui citate, fanno dello slang un ottimo indicatore di elementi caratterizzanti alcune delle etichette utilizzate da PEGI e da altri sistemi di *rating* per classificare i videogiochi. Tuttavia, per poter rendere il lessico slang fruttuoso a questo scopo è necessario suddividerlo in diverse parti, ciascuna riguardante uno degli ambiti di valutazione del rating di un prodotto e ibridarlo con termini scientifici o di gergo estratti da domini specifici come quello delle droghe o quello militare.

*SsVdd (Slang-sex-Violence-drugs and discrimination)* è una risorsa lessicale terminologica costruita per la valutazione e la classificazione dei *transcript* dei videogiochi. La risorsa è stata sviluppata in lingua inglese e contiene oltre 6.000 termini monorematici e polirematici a cui sono stati associati quattro tag riguardanti le diverse categorie di appartenenza.

Come mostrato dalla tabella 1, la risorsa contiene le seguenti categorie di parole:

- *Slang (generico)*: consiste in una raccolta di termini volgari o generalmente associati al cosiddetto *bad language*. Al suo interno troviamo termini riguardanti insulti di ogni tipo o relativi a discriminazione, sesso e sessismo. Tra questi termini troviamo nomi come *asshole* (stronzo), *bollocks* (bugie), ma anche *dinkey* (pene), verbi come *chup* (disapprovare) o *uneat* (vomitare);
- *Sesso*: il tag Sex contiene parole con espliciti riferimenti ad atti sessuali come, ad esempio, *carnal knowledge* (conoscenza carnale), *orgy* (orgia) o *striptease*;
- *Violenza*: con questo tag sono identificati termini standard e gergali relazionati alla violenza, come, tra gli altri, i nomi di armi o di oggetti associati alla guerra. I termini identificati dal tag *violenza* includono parole come *kill* (uccidere), *ammunition* (munizioni), *assassin* (assassino), *hatred* (odio), *terrorism* (terrorismo);
- *Droghe*: contiene termini generici, commerciali, scientifici e gergali che indicano le droghe o il loro consumo. In questa sezione troviamo *methamphetamine*, *cocaine*, *pentobarbital* ma anche *sniff* (sniffare), *nose candy* (droghe inalate), *speedball*;
- *Discriminazione*: si tratta di termini che riguardano discriminazioni di genere o di natura etnica, regionale o religiosa. Troviamo termini come *burrhead* (persona di colore), *chinky* (di origine asiatica), *guido* (italo-americano), *shiksa* (donna attratta da persone di origine ebrea), *schizo* (schizofrenico), *wacko* (pazzo o drogato).

Categoria	Tag	Numero
Slang	Slang	3.795
Sesso	Sex	208
Indicatori lessicali di Violenza	Violence	635
Nomi di droghe	Drugs	533
Termini Discriminatori	Discrimination	1.094
TOTALE		6.122

Tabella 1: Suddivisione in categorie del dizionario sVdd.

La fonte dalla quale sono stati estratti i termini del dizionario *SsVdd* è, principalmente, il web, ed in particolare alcune pagine che riportano termini associati e neologismi dello slang. Le pagine esaminate sono *urbandictionary.com*, *theonlineSlangdictionary.com*, *peevish.co.uk/Slang*. Per i nomi di droghe, è stato necessario affidarsi alle definizioni di *wikipedia*. Molte definizioni sono state poi verificate nella terza edizione del *NTC's dictionary of American Slang and colloquial expressions* (Spears 2000).

Il dizionario necessita di un costante aggiornamento per poter conservare la propria efficacia pratica, dato che lo slang è un linguaggio in costante mutazione (Mattiello 2005; Benedict e Munro 1997; Eble 1996; Andersson e Trufgill 1992), che genera numerosi nuovi termini, li dimentica o ne modifica il significato in un breve lasso di tempo. Si tratta di un linguaggio estremamente sensibile ai mutamenti sociali, culturali, di costume e perfino tecnologici. Una variabilità simile è possibile trovarla anche nei nomi di droghe, dato che è possibile che, con il tempo, vengano prodotte nuove sostanze stupefacenti.

Come abbiamo detto, il dizionario include sia parole semplici che parole composte polirematiche. In particolare, le parole composte, o *MultiWord Expressions* (MWE) risultano essere più efficaci nella ricerca semantica delle espressioni di slang, perché sono generalmente meno ambigue (Elia 2002). Nei testi è possibile riscontrare determinati termini sia con uso ordinario e dunque privo di volgarità, sia nella loro accezione più volgare. Il termine *beaver* indica comunemente il castoro, ma in determinati contesti viene utilizzato come sinonimo di vagina. Allo stesso termine è possibile riferirsi anche tramite delle MWE che di norma sono ambigue, ma che poco probabilmente troveremo con alta frequenza nel loro significato ordinario. Un esempio è il termine *bearded oyster*. Tuttavia, nel linguaggio slang, molte MWE risultano comunque di elevata ambiguità e dunque di difficile interpretazione. Sempre relativa ai genitali femminili, infatti possiamo trovare il termine *promised land*, o riferito ad andare di corpo il termine *make a deposit*.

Questo fenomeno risulta chiaro anche ad una rapida analisi delle concordanze all'interno dei corpora dell'inglese. Il termine *bearded oyster* appare nel corpus *enTenTen21* (Jakubiček et al. 2013) 11 volte, di cui almeno 6 con il

significato di *vagina* («hint: a friend was on #21 and we hadn't even seen the Pussyfooters, the Camel Toe High Steppers, the Bearded Oysters», *thechicory.com*), in un caso con il significato di uomo barbuto («The father, a venerable, bearded oyster, of august appearance and solemn deportment, was much mortified that one of his family should happen to be so sickly», *gutenberg.ca*) e 4 in senso letterale («Add a dozen bearded oysters, that have been well rinsed in their own liquor», *chestofbooks.com*). Per *promised land* (66.492 occorrenze) o per *make a deposit* (circa 1.000 occorrenze), la stragrande maggioranza dei risultati riguarda il significato ordinario delle due sequenze di parole.

I quattro tag proposti all'interno del nostro dizionario non devono essere univoci, nel senso che ad ogni parola non deve essere associato un solo tag. I tag multipli possono derivare sia dalla ricchezza insita nel lessico dello slang e dunque nella polisemia che contraddistingue la maggior parte dei termini, ma anche dal fatto che una serie di termini definiti slang possono identificare crimini e criminali, droghe, forze dell'ordine, morte o assassinio, abuso e discriminazione razziale. Un esempio di entrata del nostro dizionario è il seguente:

*ass crack*;Slang  
*badonkadonk*;Slang+Sex  
*batty*;Slang+Discrimination

Nell'esempio sopra vediamo che il lemma è separato dalle informazioni ad esso relative da un punto e virgola. Ad ogni entrata del dizionario sono stati associati i cinque tag relativi alle classi individuate (*Slang*, *Sex*, *Violence*, *Drugs* e *Discrimination*).

Un termine, quindi, può ricevere un doppio tag sia per via della propria polisemia, sia, semplicemente, perché appartiene contemporaneamente alla categoria Slang e a quella Sex. *badonkadonk* ad esempio appartiene allo slang, ma indica contemporaneamente il sedere di una donna afroamericana, mentre *batty* può essere utilizzato sia per indicare un sedere che un omosessuale.

I tag multipli rappresentano il 2,25% del dizionario (140 parole) e riguardano principalmente i termini *slang* (128) che ricevono un secondo tag di tipo *sex* (27), *drugs* (17), *discrimination* (56) o *violence* (28). Ci sono poi 9 termini etichettati come *violence* che ricevono il tag *drugs* 5 volte, quello *discrimination* 3 volte ed una volta quello *sex*. Anche il tag *drugs* si trova in associazione con il tag *discrimination* in 3 occasioni.

Il dizionario elettronico è liberamente consultabile e scaricabile alla pagina (Github, n.d.a.).

I *transcript* selezionati per l'esperimento sono stati rintracciati sul web in varie piattaforme in cui utenti comuni e appassionati possono condividere i file contenenti i dialoghi o i testi dei videogiochi raccolti manualmente durante l'esperienza di gioco. Per le caratteristiche interattive dei videogiochi, infatti, non è facile estrarre ogni possibile dialogo, dato che spesso il giocatore

ha diverse possibilità che, come detto, generano risposte diverse e guidano la conversazione verso sponde diverse. Secondo Heritage (2021), esistono quattro metodi per estrarre i dialoghi dei videogiochi: utilizzare un software che estrae tutte le battute, giocare ogni possibile permutazione di testo ed annotarla, usare i *transcripts* creati dai fan o siti web specializzati che hanno utilizzato una delle prime tre opzioni.

I portali utilizzati per la costruzione del nostro corpus sono principalmente due: *transcripts.fandom.com* and *game-scripts.fandom.com*. Da questi due portali sono stati estratti 50 *transcripts* relativi a giochi pubblicati nel lasso temporale che va dal 1998 al 2020. Le dimensioni totali del corpus sono di 701.384 token.

La scelta dei titoli selezionati è stata effettuata sulla base di una serie di criteri non rigidi, tra cui figura, innanzitutto, la disponibilità del *transcript* al momento della raccolta. Verificata la disponibilità sulle piattaforme indicate sopra, i successivi criteri di scelta adottati sono stati i seguenti:

1. *Popolarità*: i titoli inseriti nel corpus hanno goduto o godono tutt'ora di buona popolarità o sono riconosciuti universalmente come capisaldi di determinati generi di videogiochi;
2. *Temi e Generi*: nella selezione effettuata si è cercato di includere titoli appartenenti a quanti più generi differenti, facendo però attenzione anche alle tematiche o alle diverse ambientazioni. La maggioranza dei giochi, tuttavia, appartiene ai generi in cui la dimensione testuale ha una maggiore importanza come gli RPG o i giochi di *Azione/Avventura*;
3. *Ratings*: in ultimo sono stati raccolti titoli con valutazioni differenti, dai giochi adatti a tutti (3+) a quelli esclusivamente per adulti (18). Vista la natura dell'esperimento che illustreremo nelle prossime pagine, tuttavia, abbiamo prediletto titoli che permettessero di far risaltare la presenza/assenza degli indicatori che cercheremo di evidenziare automaticamente.

Le date di pubblicazione dei titoli selezionati tendono a concentrarsi negli anni a partire dal 2010. Un buon numero di titoli (16) sono del primo decennio del 2000, mentre solo 2 titoli sono precedenti. Questa disparità è dovuta principalmente alla difficoltà di rintracciare i *transcript* dei videogiochi troppo datati, dovuta al fatto che le *fandom* utilizzate sono sorte solo in un momento successivo.

Per quanto riguarda la classificazione PEGI, abbiamo una predominanza di titoli per adulti (PEGI 18) seguita da titoli per adolescenti (PEGI 16) e bambini in età scolare (PEGI 7 e 7+). Solo pochi titoli sono riservati a bambini (PEGI 3) e preadolescenti (PEGI 12). In questo caso, si è scelto di privilegiare titoli più corposi dal punto di vista testuale, i cui *transcript* contenessero un numero di token abbastanza alto. La media dei token per *transcript* più alta,

infatti, riguarda i videogiochi per adulti (circa 20.000 token) seguita dai PEGI 16 (circa 17.000). Le altre categorie hanno medie estremamente più basse (intorno ai 4.000 token ciascuna), ma sono state inserite proprio come importante controprova dell'efficacia del modello su testi di dimensioni inferiori.

I generi scelti sono quelli più popolari e includono principalmente giochi di azione (19 titoli, 15.000 token in media), i giochi di ruolo (8 titoli, 27.000 token in media), i cosiddetti *picchiaduro* (3 titoli, 5.000 token in media), gli *sparatutto* (3, 9.000 token), *Hack and slash* (3, 5.000 token), *party e platform* (6, 3.000 token), gli *sportivi* (3, 3.000 token) e gli *horror* (6, 18.000 token).

Di seguito osserviamo in che modo si strutturano i *transcript*:

*[Dark lane. The main character wounded a goat-like demon. It lies on the ground.]*

*Starring GARCIA HOTSPUR*

*Garcia: The bullet train is here, hellmonkey.*

*Demon: As if you mortals can be saved with one squeeze. Kill me, and I shall be replaced by another and another... and another still. You cannot point that pet gun of yours at all demonkind. And wherever you are not looking is where the greatest threat shall be. One at a time, we shall seize the treasures of your his as spoils, and leave only emptiness and despair my last gifts to you.*

*Garcia: Just don't forget to wrap them, Puta Claus.*

*Demon: By the way, Hotspur... How is your dear sweet Paula? Is she "hanging in there"?*

*Garcia: Fuck you!*

*[He blows the demon head out with his gun and run away. He runs home and finds his girlfriend hanged.]*

L'esempio sopra è tratto dalla prima scena del videogioco *Shadow of the Damned*. Possiamo subito notare la presenza di alcune descrizioni dei principali elementi visivi della scena. Abbiamo poi un discorso diretto tra due personaggi. All'interno del discorso diretto, i termini di slang sono molti (*the bullet train, hellmonkey, Puta Claus, hanging in there, Fuck you*). Sia nelle descrizioni che nel discorso diretto, invece, è possibile trovare riferimenti alla violenza (*wounded, kill, gun, hanged*).

Trattandosi di dialoghi trascritti dagli utenti, nella maggior parte dei casi non sono presenti gli elementi interattivi, ma solo le linee di dialogo scelte dal giocatore durante la sua partita. Dunque, si tratta comunque di *transcript* parziali. Questa mancanza, tuttavia, è in parte mitigata dal fatto che in genere, gli sviluppatori tendono a rendere obbligatorie molte delle linee di dialogo riproponendole più volte durante la scena finché non vengono scelte. Visto questo tipo di approccio, i dialoghi effettivamente mancanti rappresentano una piccola parte della dimensione testuale del gioco. Il corpus è stato reso disponibile per il download alla pagina (Github, n.d.b). I diversi *transcript* sono stati inseriti in un unico file testuale (.txt) e separati dal seguente indicatore di fine documento: "--END--".

## 5. Estrazione dei termini

Come abbiamo detto, vista l'importanza della componente visiva, non è possibile stabilire automaticamente il rating di un videogioco basandosi esclusivamente sulla sua dimensione testuale. Tuttavia, attraverso l'estrazione di termini appartenenti al linguaggio dello slang e ad un lessico caratterizzato da specifiche emozioni, è possibile far emergere la presenza di indicatori legati alla presenza di linguaggio volgare, violenza ed uso di droghe.

L'esperimento che proponiamo in questo paragrafo vuole dimostrare come tali indicatori possano fungere da supporto ai sistemi di rating nell'attribuzione delle specifiche etichette di *warning*, andando ad approfondire più dettagliatamente la dimensione testuale di un videogioco spesso relegata ad attore secondario nelle fasi di analisi dei titoli.

La metodologia proposta è basilare ma si dimostra effettiva e meriterebbe un ulteriore approfondimento. Gli *step* che hanno portato al risultato finale sono così riassunti:

1. Estrazione dei termini del dizionario *SsVdd* e del dizionario *NRC*, relativamente ai tag *Anger*, *Fear* e *Disgust*, dai *transcript* dei videogiochi selezionati;
2. Formulazione di tre macro-etichette tramite l'accorpamento di diversi tag;
3. Calcolo delle frequenze relative per ciascuna macro-etichetta;
4. Calcolo dei margini di assegnazione dei tag;
5. Valutazione dei risultati.

Per quanto riguarda l'indicatore di violenza, benché il numero di termini associati a questa feature nel nostro dizionario siano un buon numero, si è deciso di integrare la risorsa con informazioni di tipo emozionale estratte da un dizionario delle emozioni spesso utilizzato in task di *emotion detection*. Il dizionario usato è l'*NRC Affect Intensity Lexicon* (Mohammad e Turney 2013), disponibile gratuitamente per il download sul web. Questo dizionario contiene oltre 10.000 entrate lessicali a cui sono stati attribuiti 8 punteggi riguardanti otto emozioni basilari. I punteggi di associazione tra le parole e le emozioni di base sono stati assegnati tramite un questionario sottoposto ad un elevato numero di soggetti umani tramite una campagna di *crowdsourcing*. Da questo ampio dizionario abbiamo deciso di estrarre solo i termini che presentavano una associazione con tre emozioni in particolare: rabbia (1.515 parole), paura (1.798 parole) e disgusto (1.108 parole).

Dal corpus sono stati estratti un totale di circa 37.000 termini (circa il 5,2% del corpus) di cui, tuttavia, un buon numero appartenente all'*Affect Intensity Lexicon*. Questi ultimi sono stati aggiunti al conteggio delle parole relative alla violenza. Se raggruppiamo le occorrenze dei termini per i quattro tag, Slang,



Violenza, Droghe, possiamo notare che oltre 32.000 parole sono etichettate come Violenza, circa 5.000 appartengono allo Slang, e 800 a Droghe. Se prendiamo in esame i type estratti, cioè le singole occorrenze di ogni parola, abbiamo circa 650 parole estratte per Slang e Violenza e circa 110 per Droghe.

Osservando i termini estratti, ci accorgiamo, però, della presenza di innumerevoli parole ambigue presenti in ciascuna categoria. Tra le parole afferenti al dizionario delle Droghe troviamo parole come *blow* (105 occorrenze), *idiot* (46), *green* (46) tra le più presenti. Troviamo, con un minor numero di occorrenze anche parole non ambigue come *medicine* (23), *acid* (16), *cigarette* (9), *pills* (6), *morphine* (2). I termini che abbiamo definito Slang presentano, a loro volta, una simile incidenza di parole ambigue, benché, al primo posto, troviamo la parola *fucking* con 452 occorrenze. Per quanto riguarda i termini che indicano Violenza, non sembra che l'aggiunta delle circa 4.500 parole del *lessico delle emozioni* abbia portato ad un aumento del numero di termini estratti. Il numero di type estratti per violenza (661), infatti, è molto simile al totale di type afferenti al tag Slang (640). Tuttavia, il numero dei singoli token estratti è estremamente maggiore per violenza. Questo effetto è stato in parte influenzato dal fatto che i termini del *lessico delle emozioni* sono, a volte, molto frequenti, ma è dovuto, principalmente, all'alta frequenza di termini violenti nei giochi selezionati. Le parole *dead*, *kill*, *fight* e *soldier* contano rispettivamente 658, 955, 593 e 382 occorrenze e sono tra le parole con maggior frequenza del corpus. Se uniamo le parole *dead*, *death* e *die* otteniamo un totale di 1.340 occorrenze.

Calcolando la *Term Frequency/Inverse Document Frequency* (TF/IDF) delle parole nei testi, abbiamo un'ulteriore conferma dell'incidenza della terminologia violenta e spesso militare che caratterizza i videogiochi selezionati. Il TF/IDF rappresenta un valore di peso di una parola che rapporta la sua frequenza in un determinato testo alla frequenza della stessa parola in tutti i testi. Più essa appare nelle diverse sezioni del corpus, meno efficace sarà nel determinare le caratteristiche distintive di una specifica sezione. Abbiamo estratto le 100 parole con il valore di TF/IDF più alto per ogni testo e le abbiamo contate per verificare in quanti testi questi termini avessero un peso determinante. Una parola come *blood*, appartenente alla categoria violenza, è tra le 100 parole più importanti in più della metà dei giochi. Lo stesso vale per *guards*. Poco dietro troviamo termini come *shooting* e *destroyed*.

Al termine dell'estrazione, i punteggi dei vari dizionari sono stati combinati in modo da creare un indice che riguardasse tre indicatori nello specifico:

- *Indice di linguaggio volgare*: calcolato combinando i valori di occorrenza delle parole di Slang, Discriminazione e Sesso;
- *Indice di violenza*: associando i termini violenti a quelli estratti dal *lessico delle emozioni* e standardizzando per la dimensione del testo;
- *Indice di droghe*: contenenti i termini del dizionario delle droghe.

Nonostante la presenza di numerosi termini ambigui, i risultati ottenuti sembrano confermare l'efficacia di questa prima estrazione nell'identificare i tre indicatori all'interno dei giochi selezionati. Nella tabella 3 possiamo osservare alcuni dei valori limite ottenuti sulla base della frequenza relativa dei termini dei tre indicatori.

Indicatore	Valore Minimo	Valore Massimo	Media
Linguaggio volgare	<10 (Ninja Gaiden, Medieval)	>300 (Cyberpunk 2077)	85.27
Violenza	0.05 (Disney Infinity 3.0)	0.15 (Mortal Kombat vs DC Universe)	0.09
Droghe	1 (Toy Story 3, Cars 3: Driven to win)	55 (Cyberpunk 2077)	12.8

Tabella 3: Punteggi minimi, massimi e media dei tre indici.

Tra i giochi privi di linguaggio volgare troviamo *Ninja Gaiden*. Si tratta di un gioco di genere *hack and slash* dei primi anni 2000, etichettato come PEGI 16+ ed ESRB M a causa della violenza dei contenuti. Il tono dei dialoghi, tuttavia, non sembra essere caratterizzato da volgarità o slang.

???: *"I am Maith. You will never get pass me."*

Ryu: *"Then I must defeat you!"*

Maith: *"You are as bold as your father. But he is a much better swordsman."*

Ryu: *"You know my father?"*

Maith: *"Come and fight, young Hayabusa!"*

(GAMEPLAY)

Ryu: *"You killed my father."*

Maith: *«Killed? It is true that we fought. But your father is alive.»*

Ryu: *«Lair!»*

Maith: *«No, it is not a lie. If you proceed further, you will see him. But it will be the last thing you see.»*

Il tono della conversazione è teso, ma non si incorre mai in espressioni volgari. Ne è un esempio l'uso dell'espressione *liar*, preferita a numerose possibili alternative slang (*bullshitter, fibber, storyteller*).

Il secondo caso di videogioco privo di espressioni volgari è *Medieval*, un gioco del 1998 contraddistinto da un linguaggio di stile medievale e dunque caratterizzato proprio dall'assenza di uno slang moderno.

First Mate: *Captain, I thought you ought to know we have a stowaway on board, I've told the men to scour the decks for him.*

Captain: *Good, I want that scurvy dog dangling from yon yardarm by his bowels!*

First Mate: *Is that really necessary? Couldn't we just give him a good tongue lashing*



*and drop him off at the next port?*

Possiamo notare un elemento quantomeno gergale come *scurvy dog*, che appartiene però ad un linguaggio “piratesco”, soprattutto se accompagnato da sequenze come *dangling from yon yardarm by his bowels*. Si tratta dunque di un linguaggio fortemente stereotipato su elementi tipici della narrativa di avventura del ‘900.

Il gioco che presenta un numero maggiore di termini volgari è, invece, *Cyberpunk 2077*. Anche in questo caso, la valutazione sembra abbastanza corretta, trattandosi di un videogioco per maggiorenti ambientato tra la malavita di una città immaginaria del prossimo futuro.

*Stout: Now listen close. This piece of shit, Antony Gilchrist is he your contact? Is he the one who leaked intel on the convoy?*

*V: That guy? Never seen him before.*

*Militech Guard: Checks out.*

*V: Listen, I know where the transport is. I can help you. Just want a favor in return.*

*Stout: Hmm...*

*Gilchrist: I told you! I fucking told you! I'm not the mole! Jesus Christ!*

*Stout: Shut him up.*

*Gilchrist: Unhand me now before I-- ungh!*

L'intero *transcript*, così come l'esempio sopra, risulta caratterizzato da una fortissima presenza di slang e parolacce come *piece of shit* o *fucking*.

Per quanto riguarda l'indice di violenza, troviamo un punteggio molto basso per un gioco di corse di macchine come *Cars 3*. Benché ridotti al minimo, i termini che indicano combattimenti o scontri sono presenti anche in giochi come *Super Mario Party* o *Toy Story 3*. In ambito automobilistico, invece, la presenza di parole che indicano violenza è molto bassa.

Sono, al contrario, molto frequenti i termini violenti in uno dei videogiochi di lotta famosi proprio per la propria crudezza e per le scene sanguinose, *Mortal Kombat vs DC Universe*. Benché in questa versione del 2008 il franchise abbia introdotto alcuni eroi della *DC Comics*, l'universo fumettistico di Batman e Superman, bisogna ricordare come il primo *Mortal Kombat* fece così tanto scalpore al proprio rilascio che il caso fu portato al Senato degli Stati Uniti.

*DARK KAHN: You seek to destroy me. Good...*

*[Everyone begins to hold their heads in pain as the rage surges through them.]*

*DARK KAHN: Give in to your rage. Unleash your aggression. FIGHT!!!*

*[The two forces, with the exception of Raiden and Superman, charge at each other. Just before everyone collides, the scene cuts to darkness. It is replaced by pictures of the ensuing brawl. Liu Kang throws a flying kick at Green Lantern. Jax and Lex Luthor pound away at each other. Wonder Woman blocks a high kick from Kitana. Scorpion throws his spear as Flash dives in with his fist. Baraka and Deathstroke draw their blades...]*

Anche in questo breve dialogo, accompagnato da una descrizione delle immagini, è facile notare la costante presenza di termini legati al combattimento (*destroy, fight, kick, spear, blades*), al dolore (*pain*), e alla rabbia (*rage*).

In fine, abbiamo l'indice di Droghe, dove il videogioco *Cyberpunk 2077* ottiene il punteggio più alto. Non si tratta di un punteggio altissimo, considerato che la frequenza dei nomi di droghe è di 55 su circa 73 mila token; tuttavia, è il punteggio più alto di un videogioco, ed è coerente con quello che sappiamo della trama. Il videogioco è ambientato in un futuro distopico di tipo *cyberpunk* dove le persone si fanno impiantare arti e organi cibernetici per potenziare le proprie capacità, rischiando di perdere il contatto con la realtà. I nomi di droga individuati riguardano più farmaci e altre sostanze che droghe vere e proprie (*analgesic, painkiller*). I videogiochi per bambini e quelli *fantasy* o ambientati nel medioevo, invece, sono caratterizzati dalla completa assenza di queste parole. Oltre a *Toy Story 3* e *Cars 3*, tra giochi in cui le droghe sono assenti troviamo anche *Dark Souls II: Scholar of the First Sin* e *Assassin's Creed: Bloodlines*.

Una volta estratte le occorrenze dei termini afferenti al dizionario *SsVdd*, abbiamo valutato automaticamente la presenza dei tre indicatori nei videogiochi, verificando poi l'etichetta attribuita ad ogni gioco dai due principali sistemi di valutazione, PEGI e ESRB. In particolare, abbiamo stabilito tre livelli di rischio per ogni videogame:

- Rischio Zero: videogiochi privi di *warnings* di tipo *violence, bad language* o *Drugs*;
- Rischio Uno: giochi a cui è stata attribuita una delle tre etichette prese in esame;
- Rischio Due: videogiochi in cui l'etichetta è preceduta dai termini *insane/intense/strong*.

La scelta di attribuire un'etichetta ad un videogioco è stata effettuata stabilendo un margine di assegnazione. Questo margine, benché differente per i tre indicatori, è uniforme per tutti i videogiochi selezionati. La tabella 4 presenta i valori di assegnazione ottimizzati per ogni categoria. I punteggi sono stati stabiliti tramite una sperimentazione che ha preso il via dal calcolo dei quantili su percentuali del 33 e del 67%. Una fase di ottimizzazione del modello, effettuata testando un gran numero di valori, ha permesso di modificare i quantili per ottenere risultati migliori. In particolare, è stato necessario aumentare di 10 punti il quantile del livello 1 per *Bad Language*, e duplicare i valori per *Drugs*. I quantili di *Violence*, al contrario sono stati modificati in maniera contrastante, aumentando il livello uno e diminuendo il due.

Indice	Livello Due	Livello Uno
<b>Bad Language</b>	>80	>60
<b>Violence</b>	>0.105	>0.05
<b>Drugs</b>	>30	>22

Tabella 4: Valori soglia per l'attribuzione delle etichette.

Come abbiamo detto, per l'indicatore di violenza, vista la presenza di un enorme numero di termini generici estratti grazie al lessico emozionale, sono stati impiegati i valori di frequenza relativa. L'*Affect Intensity Lexicon*, infatti, contiene un numero molto alto di parole generiche etichettata con emozioni specifiche come rabbia, paura o tristezza. Con l'etichetta *Anger*, ad esempio, possiamo trovare parole come *brutality* (con un punteggio di 0.964), ma anche *flagrant*, *impermeable*, *clamor* o *intense* (0.455). Visto il numero di occorrenze estratte tramite questi termini, si è preferito utilizzare il punteggio di frequenza per diminuire il loro impatto in giochi caratterizzati da *transcript* più grandi.

La tabella 5 presenta i valori di precisione ottenuta nell'identificazione delle etichette per i videogiochi selezionati.

	Bad Language	Violence	Drugs	Media
Valutazione del livello di rischio	0.62	0.54	0.66	0.61
Valutazione della presenza dell'indicatore	0.78	0.86	0.86	0.83

Tabella 5: Valutazione dei risultati.

Prendendo in considerazione il livello di rischio (zero, uno e due) delle etichette, il sistema è in grado di riconoscere la presenza di un *warning* con una precisione media del 61%. Benché questo risultato non sia trascurabile, l'obiettivo, come abbiamo detto, non è attribuire un'etichetta, ma suggerire la presenza di eventuali indicatori e favorire un controllo manuale del software prima del suo rilascio. Prendendo in considerazione la capacità del sistema di predire la sola presenza/assenza di un *warning*, il punteggio di precisione è di gran lunga superiore (83%).

Nella tabella 7 possiamo osservare le differenze nella valutazione dei videogiochi in cui gli indicatori sono presenti (Positivi) e quelli in cui non lo sono (Negativi).

	SsVdd		
	Precisione Positivi	Precisione Negativi	F-score
Bad Language	0.83	0.73	0.78
Violence	0.98	0.14	0.92
Drugs	0.64	0.94	0.72

Tabella 6: Precisione del modello per Positivi e Negativi.

Il dato più significativo è legato alla precisione sui negativi dell'indicatore Violenza. Il punteggio estremamente basso ottenuto può essere dovuto all'e-

strazione dei termini generici presenti nell'*Affect Intensity Lexicon*, che induce il modello a considerare come positivi la maggior parte dei testi analizzati. Seppur in maniera meno marcata, la dimensione ridotta del tag *Drugs* produce un risultato inverso nella rispettiva etichetta. Il punteggio dell'etichetta *Bad Language*, invece, risulta più equilibrato, ma genera un F-score inferiore a quello di *Violence*.

Alla luce dell'avvento dei *Large Language Models* (LLM) è necessario proporre un confronto con le capacità di questi modelli. I LLM, o IA generative, sono dei costrutti software, basati sulla tecnologia delle reti neurali, addestrati su enormi quantità di dati linguistici ed in grado di comprendere e produrre il linguaggio in maniera quasi del tutto indistinguibile da un essere umano. Questi modelli possono essere considerati un'evoluzione dei modelli di Semantica Distribuzionale di tipo predittivo (Lenci 2023). A differenza di questi, però, i LLM sono addestrati su una quantità di dati estremamente maggiore (GPT-3, ad esempio, era addestrato su un corpus di 500 miliardi di parole), e con una complessità della rete neurale molto superiore.

Per il confronto, abbiamo utilizzato ChatGPT, basato sul modello *GPT-4* (aprile 2024) sottoponendo al *prompt* la seguente query: «Given the 51 texts in the file, each separated by the --END-- sequence, estimate the presence of drugs in each sequence and present the results in a table with the following columns: A: number of the sequence; B: Description of content; C presence of drugs (high, moderate, low)».

ChatGPT ha prodotto le tabelle richieste affermando, nel caso della violenza e del linguaggio volgare, che «this table summarizes the thematic content of [violence, ranging from discussions and threats to physical confrontations and strategic actions involving violence, across the sequences in the document | explicit language based on the dialogue and themes present in the text]». Per quanto riguarda l'uso di droghe, pur modificando la *query*, abbiamo ricevuto in risposta una semplice estrazione di elementi lessicali o una assegnazione random dell'etichetta. I risultati ottenuti da ChatGPT con gli altri due indicatori, tuttavia, sono di gran lunga inferiori alle aspettative, riuscendo a raggiungere un valore massimo di precisione per la violenza di 0.31, considerando due distinti livelli, e di 0.68 considerando solo la presenza, e per il *Bad Language* rispettivamente di 0.27 e 0.45. L'impressione, tuttavia è che, sottoponendo all'analisi di ChatGPT ogni *transcript* singolarmente, i risultati possano essere migliori.

Al contrario, chiedendo al software di indicare un possibile rating di età per ogni videogioco abbiamo ottenuto una risposta estremamente positiva, con una correlazione tra i rating automatici e quelli di PEGI dello 0.89 (correlazione di Pearson). In questo caso, ChatGPT ha prodotto la risposta basandosi su valutazioni interne che fanno affidamento alle sue specifiche conoscenze linguistiche di natura semantica, applicate al testo. Dato che nel corpus non

sono presenti i titoli dei singoli videogiochi, probabilmente ChatGPT non ha fatto uso di conoscenze pregresse su tali titoli, e possiamo considerare valido questo secondo esperimento.

## 6. Conclusioni

In questo articolo abbiamo presentato una risorsa elettronica formata da un dizionario elettronico del linguaggio slang ed un corpora composto da *transcript* di videogiochi in inglese, ed un sistema per classificare automaticamente i *transcript* dei videogiochi. Pur trattandosi di prodotti multimediali interattivi, i videogiochi possono essere considerati dei testi per la loro natura comunicativa e narrativa. L'uso di questi prodotti culturali, tuttavia, è regolato da specifiche etichette che allertano sui rischi di esposizione dei minori a contenuti non idonei a determinate fasce d'età. Le agenzie di *rating*, tuttavia, non sempre concordano nel classificare i videogiochi, perché le loro analisi prendono in considerazione elementi parziali del prodotto.

Il sistema di classificazione proposto è in grado di allertare le agenzie su eventuali rischi legati alle etichette relative a linguaggio volgare, violenza e uso di droghe, sulla base del dizionario elettronico dello slang. Il modello proposto è stato testato sul corpus dei *transcript* che include 50 videogiochi pubblicati tra gli anni 90 ed il 2020. I risultati sono incoraggianti, in quanto con questa semplice metodologia è possibile identificare correttamente circa l'80% dei videogiochi in cui uno dei tre indicatori è presente.

Benché l'attribuzione di etichette di *warning* rappresenti una possibile applicazione pratica per le risorse elettroniche presentate, l'analisi della dimensione testuale dei videogiochi può proporsi come *step* fondamentale per lo studio di fenomeni sociali o psicologici. Vista la natura sociale dello slang, l'estrazione di questi termini potrebbe fungere da base per uno studio più dettagliato della natura sociologica e sociolinguistica dei videogiochi più di tendenza, portando alla luce, ad esempio, il desiderio di evasione sociale dei giocatori. Allo stesso modo, da un'analisi approfondita dei dati, anche su scala diacronica, si potrebbe portare alla luce l'evoluzione di fenomeni come la discriminazione razziale o sessuale, analizzando l'influenza che questo media ha sulle giovani generazioni di giocatori. La risorsa che abbiamo presentato, pertanto, sarà a disposizione di ricercatori di diverse discipline, interessati a svariate tematiche legate al linguaggio utilizzato nei videogiochi.

In un possibile prosieguo di questo progetto ci proponiamo di lavorare ad un sistema per la disambiguazione dei termini slang. Ciò potrebbe essere possibile sia lavorando sulla sintassi della frase, sia analizzando i termini tramite un modello di Semantica Distribuzionale di tipo contestuale, cioè in grado di carpire il significato di un termine polisemico dal contesto in cui appare.

## Riferimenti bibliografici

- Aarseth, Espen J. 2004. "Quest Games as Post-Narrative Discourse Espen Aarseth." In *Narrative across media: The languages of storytelling*, edited by Marie-Laure Ryan, 361-76. Lincoln: University of Nebraska Press.
- Amari, Jorgen. 2010. "Slang lexicography and the problem of defining Slang." In *The fifth international conference on historical lexicography and lexicology*, 1-11. University of Oxford.
- Andersson, Lars-Gunnar, and Peter Trudgill. 1992. *Bad Language*. London: Penguin Books Downes.
- Balteiro, Isabel. 2019. "Lexical and morphological devices in gamer language in Fora." In *Approaches to videogame discourse: Lexis, interaction, textuality*, edited by Astrid Ensslin and Isabel Balteiro, 39-57. Bloomsbury Academic. <https://doi.org/10.5040/9781501338489.0008>.
- Benedict, Jennifer, and Pamela Munro. 1997. *U.C.L.A. Slang 3: [a dictionary by 25 U.C.L.A. students]*, 3-28. Westwood: Department of Linguistics, University of California.
- De Klerk, Vivian. 1990. "Slang: A male domain?" *Sex Roles* 22: 589-606.
- Dogruel, Leyla, and Sven Joeckel. 2013. "Video Game Rating Systems in the US and Europe: Comparing Their Outcomes." *International Communication Gazette* 75(7): 672-92. <https://doi.org/10.1177/1748048513482539>.
- Dumas, Bethany K., and Jonathan Lighter. 1978. "Is Slang a word for linguists?" *American Speech* 53(1): 5-17.
- Eble, Connie C. 1996. *Slang & sociability: In-group language among college students*. The University of North Carolina Press.
- Elia, Annibale. 2002. "Discorso scientifico e linguaggio settoriale. Un esempio di analisi lessico-grammaticale di un testo neuro-biologico." *Simboli, linguaggi e contesti* 2: 71-85. Carocci. <https://www.iris.unisa.it/handle/11386/1061667>.
- Ensslin, Astrid, and John Finnegan. 2019. "Bad language and bro-up cooperation in co-sit gaming." *Approaches to videogame Discourse: Lexis, interaction, textuality*, 139-56. Bloomsbury Publishing.
- Ensslin, Astrid. 2012. *The language of gaming*. Basingstoke: Palgrave Macmillan.
- Fekete, Tamás, and Ádám Porkoláb. 2019. "From Arkngthand to Wretched Squalor: Fictional place-names in universe." *ICAME Journal* 43(1): 23-58.
- Felini, Damiano. 2015. "Beyond Today's Video Game Rating Systems: A Critical Approach to PEGI and ESRB, and Proposed Improvements." *Games and Culture* 10(1): 106-22. <https://doi.org/10.1177/1555412014560192>.

- Gee, James Paul. 2014. *Unified discourse analysis: Language, reality, virtual worlds and video games*. Routledge.
- Github. n.d.a. "Slang and Videogames." Consultato il 15 maggio 2014. <https://github.com/amatusNLP/Slang-and-Videogames/blob/main/Ss-Vdd.txt>.
- Github. n.d.b. "Slang-and-Videogames." Consultato il 15 maggio 2014. <https://github.com/amatusNLP/Slang-and-Videogames>.
- Goorimoorthee, Tejasvi, Adrianna Csipo, Shelby Carleton, and Astrid Ensslin. 2019. "Language Ideologies in Videogame Discourse: Forms of Sociophonetic Othering in Accented Character Speech." In *Approaches to Videogame Discourse: Lexis, Interaction, Textuality*, edited by Astrid Ensslin and Isabel Balteiro, 269-87. London: Bloomsbury Academic. <http://dx.doi.org/10.5040/9781501338489.0020>.
- Green, Jonathon, 2011. "Some thoughts on Slang." *Revista Alicantina de Estudios Ingleses* 24: 153-71. <https://doi.org/10.14198/raei.2011.24.06>
- Heritage, Frazer. 2021. "Language, Gender, and Videogiochi." In *Language, Gender and Videogiochi*, edited by Frazer Heritage, 27-61. Cham: Springer International Publishing. [https://doi.org/10.1007/978-3-030-74398-7\\_2](https://doi.org/10.1007/978-3-030-74398-7_2).
- Humphreys, Ashlee, and Rebecca Jen-Hui Wang. 2018. "Automated text analysis for consumer research." *Journal of Consumer Research* 44(6): 1274-1306.
- Jakubíček, Miloš, Adam Kilgarriff, Vojtěch Kovář, Pavel Rychlý, and Vít Suchomel. 2013. "The TenTen corpus family." In *7th international corpus linguistics conference CL*, 125-27.
- Juul, Jesper. 2011. *Half-real: Video games between real rules and fictional worlds*. MIT press.
- Klimova, Blanka, and Kacet Jaroslav. 2017. "Efficacy of computer games on language learning." *Turkish Online Journal of Educational Technology-TOJET* 16(4): 19-26.
- Lenci, Alessandro, 2023. *Understanding natural language understanding systems. A critical analysis*. <https://arxiv.org/pdf/2303.04229>.
- Lighter, Jonathan E. 1994. *Historical Dictionary of American Slang*. Random House. Oxford University Press.
- Maisto, Alessandro, Giandomenico Martorelli, Antonietta Paone, and Serena Pelosi. 2021. "Extracting video games rating labels from transcript files." *Internet of Things* 16: 100439.



- Marchiori, Eugenio J., Ángel Del Blanco, Javier Torrente, Iván Martínez-Ortiz, and Baltasar Fernández-Manjón. 2011. "A visual language for the creation of narrative educational games." *Journal of Visual Languages & Computing* 22(6): 443-52.
- Mattiello, Elisa. 2005. "The pervasiveness of Slang in standard and non-standard English." *Mots Palabras Words* 6: 7-41.
- McMillan, James B. 1978. "American Lexicology 1942-1973." *American Speech* 53(2): 141-63.
- Mencken Henry Louis. 1967. "American Slang, Its origin and history." In *The American Language*, edited by Raven I. McDavid and David W. Maurer, 4th edition with supplements, 702-61. New York, Alfred A Knopf.
- Mohammad, Saif M., and Peter D. Turney. 2013. "Crowdsourcing A Word-Emotion Association Lexicon." *Computational Intelligence* 29(3): 436-65. <https://doi.org/10.1111/j.1467-8640.2012.00460.x>.
- Partridge, Eric. 1933. *Slang Today and Yesterday*. London: Routledge and Kegan Paul.
- Ray, Alice. 2019. "Playing with the language of the future: The localization of science-fiction terms in videogames." In *Approaches to videogame discourse: Lexis, interaction, textuality*, edited by Astrid Ensslin and Isabel Balteiro, 87-115. Bloomsbury Academic.
- Spears, Richard A. 2000. *NTC's dictionary of American Slang and colloquial expressions*. NTC Publishing Group. <http://thuvienso.bvu.edu.vn/handle/TVDHBRVT/14818>.
- Wright, Andrew, David Betteridge, and Michael Buckby. 2006. *Games for language learning*. Cambridge University Press.





# AIDAinformazioni

Rivista semestrale di Scienze dell'Informazione

Anno 42

N. 1-2 – gennaio-giugno 2024

## Editoriale

ROBERTO GUARASCI

*Quaran'anni di «AIDAinformazioni»  
(1984-2024)*

## Contributi su invito

PIERO INNOCENTI

*Que reste-t-il de nos amours? Arti  
tradizionali di un possibile moderno Trivio:  
Archivistica, Bibliografia, Documentazione*

## Contributi

FRANCESCO AMATO, ANTONELLA  
BENVENUTO, SILVIA CENITI, MARIA  
TERESA CHIARAVALLOTI, CLAUDIA  
LANZA, ERIKA PASCERI

*Indagine e analisi comparativa delle  
procedure di codifica nella Terapia del  
Dolore in Italia*

ANDREA CAPACCIONI

*Sull'affidabilità delle bibliografie generate  
dai chatbot. Alcune considerazioni*

ALESSANDRO MAISTO

*La dimensione Testuale del Videogioco.  
Classificazione dei transcript dei  
videogiochi basata sul lessico*

ALEXANDER MURZAKU, PONTISH  
YERAMYAN, CURT ANDERSON, STEVEN  
BUXBAUM, RUBEN DIAZ, MARIELLE  
LERNER, ARMENUI MINASYAN, HAZEL  
MITCHLEY, JODIE-ANN PENNANT, MIA  
SHANG, BRISA SPEIER BRITO

*Discovering and documenting brilliance.  
A novel multimodal annotation method*

CAMILLA NAPPI

*Le document : miroir des enjeux  
institutionnels et de l'évolution linguistique  
des transitions écologique et énergétique en  
France (2007-2022)*

## Testimonianze

PIERO CAVALERI, FERRUCCIO DIOZZI

*La Documentazione in Italia e il ruolo di  
Paolo Bisogno: una rapida evoluzione*

PAOLA GARGIULO, LUCIA MAFFEI

*Domenico (Ingo) Bogliolo. Profilo di un  
navigatore nell'Information Science*



mundaneum

In copertina

Disegno di Paul Otlet, Collections Mundaneum, centre d'Archives, Mons (Belgique).

ISBN 979-12-5965-407-6

ISSN 1121-0095



9 791259 654076



9 770112 100950