

ELENA MARTINELLI, VITA GARRAMMONE, FRANCESCA MORI,  
INCORONATA NOLÈ, FRANCA CAMERIERO, MATILDE MARTINO,  
GAETANO DI BELLO, GLORIA GAGLIARDI

## I correlati linguistici del deterioramento cognitivo: raccolta e analisi di un corpus di eloquio patologico di individui anziani lucani affetti da demenza

Linguistic correlates of cognitive decline: collection and analysis  
of a spoken corpus of patients with dementia living in Basilicata

The study aims to profile the linguistic-communicative disabilities related to dementia's onset in Italian. To this purpose, we recruited 40 subjects residing in Basilicata, an administrative region of Southern Italy not previously represented in this type of study. The enrolled cohort is balanced by sex and age: 20 healthy subjects and 20 subjects affected by different types of dementia (i.e., Alzheimer's disease, mixed dementia, frontotemporal dementia, and vascular dementia). After neuropsychological evaluation, we acquired approximately 9 hours of semi-spontaneous speech. Then, the corpus was transcribed and annotated to automatically extract a rich set of linguistic markers. Through statistical inference tests, we outline a degradation of the phonetic-acoustic, morphosyntactic/lexical, semantic, and syntax skills related to the disease. Finally, the qualitative analysis allows for better characterizing the linguistic profile of dementia, highlighting other traits such as local/global coherence and cohesion deficits, anomia's compensatory strategies, and repetitiveness/topic iteration.

*Keywords:* dementia, Italian, linguistic profiling, pathological speech, Clinical Linguistics.

### 1. *Caratteristiche linguistiche dell'eloquio nelle sindromi dementigene: lo stato dell'arte*

La demenza è una sindrome clinica caratterizzata dalla pervasiva compromissione delle funzioni cognitive e dalla conseguente perdita della personale autonomia funzionale. Si osserva in risposta all'insorgenza di un variegato spettro di patologie (es. malattie neurodegenerative, patologie cardiovascolari e/o metaboliche, ecc.) e costituisce la settima causa di morte a livello mondiale, la quarta nella popolazione ultrasettantacinquenne (Costa & Sinforiani, 2021).

Le stime epidemiologiche dell'incidenza e della mortalità riferite al disturbo nella popolazione anziana mostrano aumenti allarmanti: nel 2020 sono stati diagnosticati 50 milioni di casi, ne sono attesi 82 milioni nel 2030 e 152 milioni nel 2050 (ADI *et al.*, 2021). Tali previsioni, in combinazione con il dispendio di risorse economiche, socio-assistenziali e sanitarie richieste, motivano la definizione della demenza come "priorità mondiale di salute pubblica" (OMS, 2012).

In virtù di quanto appena delineato, negli ultimi anni la comunità scientifica ha dedicato notevole attenzione all'identificazione di "marcatori" in grado di supportare la diagnosi, come *biomarker* liquorali, indici demografici e cardiovascolari. Più di recente, ha destato crescente interesse l'individuazione di *marker* linguistici della patologia: analizzare la produzione verbale di una coorte di pazienti e individuare le disabilità linguistico-comunicative associabili alla condizione morbosa sottesa rappresentano procedure estremamente vantaggiose (perché meno invasive e più economiche) per progettare strumenti che possano essere di ausilio alle pratiche diagnostiche, di prevenzione e riabilitazione clinica, nonché per migliorare le conoscenze relative al funzionamento della competenza linguistica normofasica (Gagliardi, 2019).

A partire dagli anni Settanta del Novecento, studi condotti a livello internazionale hanno permesso di evidenziare e descrivere il deterioramento delle competenze linguistiche causato dall'insorgenza di una patologia dementigena (cfr. Boschi *et al.*, 2017): alto numero di disfluenze, alterazione delle caratteristiche spettro-acustiche e ritmiche della voce, impoverimento morfo-sintattico e lessicale, difficoltà di *word-retrieval*, *deficit* di pianificazione ed ecolalie a livello pragmatico e testuale rappresentano gli aspetti linguistici deficitari tipicamente indotti dai processi neurodegenerativi (anche se i *pattern* sindromici esibiti e la gravità delle disfunzioni linguistiche variano a seconda dello specifico quadro dementigeno al quale il soggetto è ascrivibile (Papagno & Bolognini, 2020; Marino, 2021).

Costruire campioni rappresentativi della popolazione clinica oggetto di indagine costituisce un essenziale prerequisito per delineare il profilo linguistico-comunicativo di una malattia: se per alcune lingue si dispone di adeguati *corpora* di eloquio dementigeno (es. TalkBank – DementiaBank, per la lingua inglese; Becker *et al.*, 2004), le risorse per la lingua italiana sono estremamente limitate. Escludendo gli studi condotti su un numero piuttosto ristretto di pazienti, i *corpora* di parlato dementigeno raccolti per l'italiano sono pochi, a maggior ragione se si considera il fatto che non siano confluiti in raccolte pubblicate o accessibili. In aggiunta, la notevole variazione diatopica interna al sistema linguistico italiano (Berruto, 2021) impone che vengano acquisiti campioni di linguaggio reclutando soggetti appartenenti a quante più regioni della penisola italiana, affinché si possa fornire una solida base alle generalizzazioni che scaturiscono dalle analisi dei *corpora* raccolti.

I progetti più rilevanti in tale campo censiscono le produzioni linguistiche semi-spontanee di parlanti emiliano-romagnoli (progetto OPLON; Calzà *et al.*, 2021), campani (corpus CIPP-ma; Dovetto *et al.* 2022) e calabresi (corpus CIPP-mci; Dovetto *et al.* 2021). Al gruppo Anchise, invece, si deve la costruzione del primo *corpus* di parlato dementigeno spontaneo (Vigorelli, 2016; Benvenuti *et al.*, 2020): il *Corpus* Anchise 320, infatti, raccoglie le trascrizioni degli scambi conversazionali avvenuti in contesto ecologico tra operatori sanitari e soggetti affetti da demenza provenienti da molteplici regioni italiane (prevalentemente, però, dalla Lombardia).

## 2. Razionale dello studio

Lo studio si propone di confrontare a livello fonetico-acustico, morfosintattico, semantico, sintattico e pragmatico l'eloquio di soggetti italo-foni affetti da demenza e quello di soggetti coetanei cognitivamente integri, al fine di evidenziare divergenze quantitative/qualitative funzionali alla profilazione delle caratteristiche linguistico-comunicative del deterioramento cognitivo per la lingua italiana.

Al fine di garantire una più affidabile generalizzabilità ai risultati ottenuti, nonché per disporre di materiale prezioso per futuri confronti in merito alla potenziale incidenza della variazione diatopica dell'italiano sui profili linguistici dei pazienti, sono stati reclutati soggetti nati e residenti in Basilicata, una regione non precedentemente rappresentata in studi di tale natura.

Il disegno sperimentale ha previsto l'acquisizione di campioni di linguaggio semi-spontaneo, appartenenti a due gruppi di soggetti adulti/anziani (> 60 anni), reclutati in seguito ad una valutazione neuropsicologica realizzata mediante la somministrazione di classici test carta e matita: *Mini-Mental State Examination* -MMSE (Magni *et al.*, 1996), *MOntreal Cognitive Assessment* – MoCA (Conti *et al.*, 2015) e due test di fluenza verbale, fonemica (Caltagirone *et al.*, 1995; Carlesimo *et al.*, 1995, 1996) e semantica (Spinnler & Tognoni, 1987).

Il campione risulta in definitiva costituito da 40 individui nati e residenti in Basilicata, di cui 24 di sesso femminile e 16 di sesso maschile, distinti in due gruppi (bilanciati per sesso ed età):

- Gruppo dei pazienti (GP): 20 soggetti (12 F, 8 M; età:  $81 \pm 6.9$ ), con diagnosi conclamata di demenza e in cura presso la Residenza Sanitaria Assistenziale Opera Salute – Don Uva (sede di Potenza, PZ). Il gruppo è più finemente suddivisibile, sulla base della diagnosi ricevuta, in:
  - MA: 9 soggetti (8F, 1 M) con diagnosi di malattia di Alzheimer;
  - DM: 2 soggetti (1 F, 1M) con diagnosi di demenza mista;
  - DNS: 5 soggetti (2 F, 3 M) con diagnosi di demenza non ulteriormente specificata;
  - DV: 3 soggetti (M) con demenza vascolare;
  - DF-T: 1 soggetto (F) con demenza frontotemporale.
- Gruppo di controllo (GC): 20 soggetti (12 F e 8 M; età:  $81 \pm 6.3$ ), con stato cognitivo nella norma.

Il gruppo dei pazienti è stato reclutato presso la Residenza Sanitaria Assistenziale “Universo Salute – Opera Don Uva” (PZ), durante un arco temporale di circa quattro mesi (gennaio 2022 – aprile 2022).

La successiva tabella riassume i criteri di reclutamento del campione.

Tabella 1 - *Criteri di reclutamento dei partecipanti allo studio*

CRITERI DI RECLUTAMENTO	
GRUPPO P	GRUPPO C
Età > 60 anni	Età > 60 anni
Esposizione linguistica monolingue	Esposizione linguistica monolingue
L1 italiano	L1 italiano
Diagnosi clinica di demenza	Assenza di deficit neurologici / sensoriali
<i>Punteggi ai test carta-matita:</i>	<i>Punteggi ai test carta-matita:</i>
- MMSE < 22	- MMSE ≥ 22
- MoCA ≤ 19,262	- MoCA > 19,262
- Fluenza fonemica < 17,35	- Fluenza fonemica ≥ 17,35
- Fluenza semantica < 7,25	- Fluenza semantica ≥ 7,25

Il progetto di ricerca è stato approvato dal Comitato di Bioetica dell'Alma Mater Studiorum, Università di Bologna (Prot. n. 72032 del 29/01/2022).

Per tutti i pazienti è stata ottenuta l'autorizzazione da parte dei relativi *caregiver*/ tutori legali, in virtù del loro stato di deterioramento cognitivo; i soggetti di controllo, invece, hanno firmato personalmente i moduli di consenso informato per la partecipazione allo studio e di consenso per il trattamento dei dati personali.

I dati anagrafici, clinici e linguistici sono stati trattati nel rispetto della normativa vigente (Regolamento UE 2016/679 – Regolamento generale sulla protezione dei dati personali; D.Lgs. 196/2003 – Codice in materia di protezione dei dati personali): le informazioni richieste dal progetto e i risultati delle indagini sono stati annotati nella Scheda Raccolta Dati (CRF) e utilizzati in forma pseudo-anonimizzata attraverso l'attribuzione di un codice progressivo univoco ad ogni partecipante.

### 3. *Materiali e metodo*

#### 3.1 Il corpus

L'eloquio dei partecipanti è stato acquisito mediante la somministrazione di tre *task* linguistici: due compiti narrativi, per i quali è stato richiesto – rispettivamente – il racconto di una gita o di viaggio compiuto nella propria vita (frase stimolo: “vuoi parlarmi di un viaggio o di una gita che hai fatto?”) e quello delle tradizioni legate alla giornata di Natale (frase stimolo: “come passi di solito la giornata di Natale?”); un compito descrittivo, per il quale è stato impiegato come stimolo visivo una figura inclusa nel test neuropsicologico “Esame del Linguaggio II” (Ciurli *et al.*, 1996) (frase stimolo: “potresti descrivermi questo disegno?”).

Per il gruppo di pazienti, le sessioni di registrazione sono state svolte nella RSA ospitante, quando possibile in stanze vuote/isolate al fine di limitare al minimo i rumori di sottofondo; le sessioni di registrazione dei soggetti del gruppo di controllo, invece, sono state svolte all'interno delle loro abitazioni. La somministrazione dei tre *task* ha richiesto in media circa 30 minuti.

Il *corpus* raccolto consta in totale di 8 h e 50 min di sonoro, di cui 7 h e 50 min acquisite per i gruppi reclutati e 1 h attribuibile alle produzioni dell'intervistatrice e di coloro che – in rari casi – hanno partecipato alle conversazioni anche se non direttamente intervistati (es. compagni di stanza dei pazienti sottoposti ad intervista nella RSA; familiari presenti durante le interviste dei soggetti di controllo).

### 3.2 Annotazione

Tutte le registrazioni sono state trascritte ortograficamente mediante il *software* ELAN (2021), secondo le convenzioni dell'Italiano Standard, impiegando il formato di trascrizione L-AcT (Cresti *et al.*, 2018).

Contestualmente alla trascrizione ortografica, è stata portata a termine la segmentazione dell'eloquio in enunciati, selezionati come unità di riferimento dell'analisi discorsiva in quanto concepibili – seguendo Austin (1962) – come minime entità linguistiche pragmaticamente interpretabili. I confini di enunciato sono stati annotati mediante l'identificazione dei *break* prosodici, indici percettivamente salienti poiché contraddistinti da una serie di caratteristiche acustiche, tra cui *reset* della F0, allungamenti finali e aumenti dell'intensità (cfr. Cresti & Moneglia, 2018).

La segmentazione prosodica del parlato acquisito e l'annotazione di alcuni fenomeni tipici della lingua parlata sono state dunque realizzate mediante una serie di diacritici proposti dal suddetto quadro teorico e brevemente specificati nella tabella 2.

Tabella 2 - *Diacritici – Framework teorico L-AcT*

<i>Diacritico</i>	<i>Spiegazione</i>
//	<i>break</i> prosodico terminale (conclusione di una sequenza discorsiva)
?	<i>break</i> prosodico terminale (conclusione di una sequenza discorsiva con profilo prosodico interrogativo)
/	<i>break</i> prosodico non terminale (conclusione di un'unità prosodica, ma non di un enunciato)
[/]	falsa partenza o un <i>retracting</i> con ripetizione
+	<i>break</i> prosodico coincidente con la terminazione non intenzionale dell'unità discorsiva progettata dal parlante, la cui comprensione risulta dunque potenzialmente compressa
#	esitazione temporanea silente o interruzione del flusso discorsivo (limite minimo: 250 ms)
&	parole incomplete e/o frammenti fonetici
< >	<i>overlapping</i> tra gli interattanti
xxx	parole incomprensibili
hhh	elementi paralinguistici / extralinguistici (es. risata, pianto)
yyy	parole non trascrivibili per <i>privacy</i>

La *pipeline* computazionale progettata da Gagliardi e Tamburini (2022) è stata successivamente impiegata al fine di estrarre i Biomarker Linguistici Digitali (*Digital Linguistic Biomarkers*, DLBs), *feature* linguistiche direttamente derivabili dalla produzione linguistica dei parlanti, la cui analisi consente di quantificare in maniera oggettiva la degradazione della facoltà di linguaggio in una data popolazione clinica (*infra*).

La procedura di estrazione di tali marcatori linguistici prevede alcune imprescindibili fasi di preprocessazione automatica dei dati, ovvero: segmentazione dell'eloquio acquisito in porzioni di silenzio e porzioni di parlato; tokenizzazione, lemmatizzazione, *PoS-tagging* e *parsing* sintattico a dipendenze, annotazioni realizzate sfruttando un *toolkit* di NLP *open source* (STANZA, v. 1.3.0; Qi *et al.*, 2020).

Prima di procedere con l'impiego della *pipeline*, gli *input* sono stati manualmente ripuliti: mediante il software *Audacity*, dai file audio sono state eliminate tutte le sezioni non riconducibili all'eloquio dei soggetti reclutati (interventi del ricercatore, rumori di sottofondo, sovrapposizioni tra interattanti) e la frequenza di campionamento di tutti i file è stata ridotta da 96 KHz / 24 bit a 44KHz / 16 bit (compatibile con la *pipeline*); dalle trascrizioni in formato .txt sono stati eliminati i diacritici precedentemente inseriti, ad eccezione dei *break* prosodici terminali (“/”, sostituiti con il punto “.”) e delle interruzioni improvvise degli enunciati (“+”, sostituiti con i tre puntini sospensivi “...”).

Inoltre, le espressioni inizialmente oscurate per *privacy* (= “yyy”) sono state sostituite con un nome proprio (sempre lo stesso, scelto casualmente), al fine di continuare a tutelare la *privacy* dei soggetti senza intaccare la processazione delle produzioni linguistiche.

### 3.3 Analisi delle *feature* linguistiche

In prima istanza, la durata degli enunciati prodotti dai soggetti reclutati è stata valutata alla luce di una serie di variabili cliniche, ovvero sesso ed età dei partecipanti, presenza/ assenza di declino cognitivo, tipologia e severità della demenza diagnosticata.

Successivamente, si è operato un confronto sulla ricca serie di DLBs estratti (tabella 3)<sup>1</sup>, sia in relazione alla produzione linguistica complessiva, sia in relazione ai singoli *task*. Tale metodologia ha permesso dunque di individuare:

1. i correlati linguistici più “solidi” del deterioramento cognitivo, poiché significativamente differenti tra le due coorti nonostante l'eterogeneità dei *task* somministrati e la maggiore quantità di eloquio analizzato;
2. i “potenziali” correlati linguistici del deterioramento cognitivo, ovvero *feature* linguistiche meritevoli di futuri approfondimenti perché caratterizzanti l'eloquio dementigeno solamente in uno o due *task* ma non nella produzione linguistica complessiva.

<sup>1</sup> Per una definizione esaustiva dei molteplici DLBs estratti si rimanda a Calzà *et al.*, 2021.

La significatività delle divergenze riscontrate rispetto all'eloquio normofasico è stata analizzata mediante l'applicazione di *test* di inferenza statistica parametrici e non parametrici (es. T-test di Student, Test di Wilcoxon-Mann-Whitney) in R (R Core Team, 2014).

Tabella 3 - *Feature indagate*

<i>DIGITAL LINGUISTIC BIOMARKERS (n°151)</i>	
<i>Livello fonetico-acustico (20 DLBs)</i>	<ul style="list-style-type: none"> <li>– Silence Segments Duration (Satt <i>et al.</i>, 2013)</li> <li>– Speech Segments Duration (Satt <i>et al.</i>, 2013)</li> <li>– Temporal Regularity of Voiced Segments (Satt <i>et al.</i>, 2013)</li> <li>– Verbal Rate (Singh <i>et al.</i>, 2001; Roark <i>et al.</i>, 2011)</li> <li>– Transformed Phonation Rate (Singh <i>et al.</i>, 2001; Roark <i>et al.</i>, 2011)</li> <li>– Standardized Phonation Time (Singh <i>et al.</i>, 2001; Roark <i>et al.</i>, 2011)</li> <li>– Standardized Pause Rate (Singh <i>et al.</i>, 2001; Roark <i>et al.</i>, 2011)</li> <li>– Root Mean Square Energy (López-de-Ipiña <i>et al.</i>, 2013)</li> <li>– Pitch (López-de-Ipiña <i>et al.</i>, 2013)</li> <li>– Spectral centroid (López-de-Ipiña <i>et al.</i>, 2013)</li> <li>– Higuchi Fractal Dimension (López-de-Ipiña <i>et al.</i>, 2013)</li> </ul>
<i>Livello lessicale (34 DLBs)</i>	<ul style="list-style-type: none"> <li>– Content Density (Roark <i>et al.</i>, 2011)</li> <li>– Part-of-speech Rate (Bucks <i>et al.</i>, 2000)</li> <li>– Reference Rate to Reality (Vigorelli, 2004)</li> <li>– Personal, Spatial and Temporal Deixis Rate (March <i>et al.</i>, 2006; Cantos-Gómez, 2009)</li> <li>– Relative pronouns and negative adverbs Rate</li> <li>– Lexical Richness: Type-Token Ratio (Holmes &amp; Singh, 1996), Brunet's Index (Brunet, 1978) e R – Honor's Statistic (Honoré, 1979)</li> <li>– Action Verbs Rate (Gagliardi, 2014)</li> <li>– Frequency of use Tagging (De Mauro, 2000)</li> <li>– Propositional Idea Density (Snowdon <i>et al.</i>, 1996; Roark <i>et al.</i>, 2011)</li> <li>– Mean number of words for utterance</li> </ul>
<i>Livello semantico<sup>2</sup> (84 DLBs)</i>	<ul style="list-style-type: none"> <li>– Language</li> <li>– Emotiveness</li> <li>– Cognition</li> <li>– Sociality</li> <li>– Physicality</li> <li>– Lifestyle</li> <li>– Spatio-temporal references</li> <li>– Conversation style</li> </ul>
<i>Livello sintattico (9 DLBs)</i>	<ul style="list-style-type: none"> <li>– Number of dependent elements linked to the noun</li> <li>– Global dependency distance (Roark <i>et al.</i>, 2007; 2011)</li> <li>– Syntactic complexity (Szmrecsanyi, 2004)</li> <li>– Syntactic Embeddedness</li> <li>– Maximum Depth of The Structure</li> <li>– Utterance Length</li> </ul>

<sup>2</sup> Tale gruppo di feature viene computato sulla base dell'analisi LIWC (*Linguistic Inquiry and Word Count*; cfr. Chung & Pennebaker, 2007). Nello specifico, le indagini sono basate sulla versione italiana del dizionario LIWC (Agosti & Rellini, 2007).

<i>DIGITAL LINGUISTIC BIOMARKERS (n°151)</i>	
<i>Leggibilità</i> <sup>3</sup> (4 DLBs)	– Basic readability – Lexical readability – Syntactic readability – Total readability

#### 4. *Analisi quantitative: risultati*

##### 4.1 Il profilo linguistico dei pazienti con deterioramento cognitivo

Le analisi condotte sui campioni di linguaggio raccolti hanno evidenziato un numero considerevole di divergenze: le successive tabelle espongono i dati statistici ottenuti analizzando la durata degli enunciati alla luce delle variabili cliniche oggetto di interesse (tabella 4), nonché indagando i valori dei molteplici tratti linguistici estratti tra i due gruppi reclutati (tabella 5).

Tabella 4 - *Durata degli enunciati: variabili cliniche*

<i>Variabile clinica</i>	<i>Test di inferenza statistica</i>		
	<i>Mann-Whitney per campioni indipendenti</i>	<i>Kruskall – Wallis</i>	<i>Correlazione di Kendall</i>
<i>Diagnosi di demenza</i>	< 0.001		
<i>Tipo di demenza</i>		< 0.001	
<i>MMSE</i>			< 0.001
<i>Severità della demenza</i>	<i>MOCA</i>		< 0.001
	<i>F_SE</i>		< 0.001
	<i>F_FO</i>		< 0.001
<i>Sesso</i>	<i>GC</i>	0.08361	
	<i>GP</i>	< 0.001	
<i>Età</i>	<i>GC</i>		< 0.001
	<i>GP</i>		0.00333

<sup>3</sup> Tale *set* di *feature* viene direttamente computato grazie allo strumento di valutazione della leggibilità che, per la lingua italiana, attualmente è considerato il migliore sotto il profilo dell'affidabilità: READ-IT, sviluppato presso il CNR-ILC di Pisa (Dell'Orletta *et al.*, 2011).

Tabella 5 - *Digital Linguistic Biomarkers*

<i>Livello</i>	<i>DLB</i>	<i>T-test / U-Test</i>	
		<i>p-value</i>	<i>Significatività<sup>4</sup></i>
<i>Fonetico-acustico</i>	SPE_SPEMEDIAN	0.001435	**
	SPE_TPR	0.001824	**
	SPE_SILMEDIAN	0.00304	**
	SPE_RMSEM	0.003717	**
	SPE_VR	0.004864	**
	SPE_SPEMEAN	0.006644	**
	SPE_SPCENTRM	0.009047	**
	SPE_SILSD	0.01121	*
	SPE_SILMEAN	0.02108	*
	<i>Lessicale</i>	LEX_PoS_DET	0.0006802
LEX_ContDens		0.0003114	***
LEX_PoS_PUNCT		0.0002644	***
LEX_PoS_ADJ		0.00006	***
LEX_PoS_ADP		0.00006	***
LEX_PoS_		1.1e-05	****
LEX_RPRO		0.01136	***
LEX_OCW		0.0002117	***
LEX_CCW		0.0002117	***
LEX_PoS_INTJ		0.00129	**
LEX_TDEIXIS		0.001291	**
LEX_PoS_AUX		0.002091	**
LEX_PoS_X		0.002227	**
LEX_ACTVRB		0.007714	**
LEX_PoS_ADV		0.02447	*

<sup>4</sup> Livelli di significatività: \* p-value < 0.05, \*\* p-value < 0.01, \*\*\* p-value < 0.001, \*\*\*\* p-value < 0.0001.

<i>Livello</i>	<i>DLB</i>	<i>T-test / U-Test</i>	
		<i>p-value</i>	<i>Significatività<sup>4</sup></i>
<i>Semantico</i>	LWC_Articol	0.03364	***
	LWC_Passato	0.0001192	***
	LWC_Prepos	0.0001554	***
	LWC_Tempo	0.0001854	***
	LWC_Ottimis	0.0004921	***
	LWC_Causa	0.0004983	***
	LWC_Sen_Pos	0.0007168	***
	LWC_Affett	0.001085	**
	LWC_Fisico	0.001609	**
	LWC_No_i_Verb	0.002149	**
	LWC_Sopra	0.004391	**
	LWC_Corpo	0.004105	**
	LWC_Scuola	0.004782	**
	LWC_Proc_Sen	0.007572	**
	LWC_P_pass	0.009773	**
	LWC_Consen	0.01018	*
	LWC_Tu_Verbo	0.01069	*
	LWC_Svago	0.01142	*
	LWC_Occupaz	0.02282	*
	LWC_No_i	0.02474	*
	LWC_Mangiare	0.02506	*
	LWC_Condizio	0.02594	*
	LWC_Movimen	0.02633	*
	LWC_riempiti	0.0263	*
	LWC_Sport	0.02854	*
	LWC_Raggiun	0.03484	*
	LWC_Rif_gen	0.03596	*
	LWC_Casa	0.03848	*
	LWC_Social	0.04381	*
	<i>Sintattico</i>	SYN_MAXDEPTHM	0.0002984
SYN_GRAPHDISTM		0.0004101	***
SYN_MAXDEPTHSD		0.007158	**
SYN_SLENM		0.01094	*
<i>Leggibilità</i>	REA_BASE	3,73E-06	****
	REA_LEXICAL	0.002643	**

Sintetizzando le evidenze riscontrate, si può affermare che l'eloquio dementigeno è caratterizzato da un numero maggiore di enunciati (tuttavia, significativamente più brevi) e da un numero maggiore di interruzioni improvvise della pianificazione discorsiva.

La severità della demenza, concepita in termini di punteggio ottenuto alle prove psicometriche somministrate, correla significativamente con la durata degli enunciati.

Per i pazienti, fattori anagrafici quali sesso ed età risultano significativamente influenti: nello specifico, la durata degli enunciati prodotti rivela una correlazione negativa con l'età e risulta divergere in base al sesso del soggetto. Poiché la condizione normotipica è rappresentata dal gruppo di controllo e, all'interno di tale popolazione, la focalizzazione sulle variabili anagrafiche non consente di evidenziare differenze statisticamente significative, si può ipotizzare che i risultati ottenuti per i pazienti siano in realtà specchio dell'azione di altri fattori: l'ineguale distribuzione dei tipi di demenza (e severità delle stesse) tra i pazienti dei due sessi, così come tra pazienti ascrivibili a una medesima fascia d'età, nonché possibili comorbidità di cui soffrono i soggetti anche quando ascrivibili a uno stesso quadro dementigeno.

Relativamente ai *Digital Linguistic Biomarker*, le indagini statistiche condotte sulla produzione linguistica complessiva dei soggetti permettono di affermare che l'eloquio patologico si differenzia da quello normotipico per 59/151 tratti, così suddivisibili: 9 DLBs acustici, 15 DLBs lessicali, 29 DLBs LWC, 4 DLBs sintattici, 2 DLBs di leggibilità.

Sintetizzando i risultati ottenuti, il parlato dementigeno risulta caratterizzato da una maggiore durata media e mediana dei segmenti di silenzio, nonché da una loro maggiore variabilità, e da una minore durata media e mediana dei segmenti di *speech*. La velocità di fonazione, riconducibile a un maggiore numero di pause, l'energia dei segmenti di *speech* e la *Verbal Rate* sono risultate significativamente minori.

Dal punto di vista lessicale, è osservabile – innanzitutto – una minore percentuale di parole di classe aperta: in particolare, emerge un ridotto impiego di aggettivi e avverbi, ma non di nomi e verbi, classi per le quali si osservano valori comparabili con quelli estratti per i soggetti controllo. La maggiore percentuale di parole di classe chiusa registrata è riconducibile prevalentemente a un maggiore ricorso alle interiezioni, poiché si osserva una minore occorrenza di preposizioni semplici e articolate, di determinanti, ma anche di avverbi interrogativi/congiunzioni subordinanti come “perché” (indice di un minore riferimento ai processi causativi). Senza dubbio, gli aspetti appena delineati concorrono a spiegare la minore densità di contenuto osservabile per l'eloquio dei pazienti.

Si evidenzia, ancora, una significativa riduzione delle espressioni deittiche di natura temporale, confermata da una generale riduzione delle espressioni temporali, soprattutto riferite/formulate al passato.

La maggiore percentuale di punteggiatura evidenziata dalle analisi, invece, è da interpretare come maggiore numero di enunciati e più frequenti interruzioni improvvise delle pianificazioni discorsive (*supra*). Questa *feature*, infatti, non va intesa come indicazione dell'uso dei segni di interpunzione (anche perché si tratta di testi orali), ma è da ricondursi alle modifiche testuali avvenute negli *step* di preprocessazione, durante i quali i *break* prosodici terminali (“/”) e le interruzioni improvvise degli enunciati (“+”) sono stati rispettivamente sostituiti dai segni di punteggiatura “.” e “...”.

Sebbene nel complesso non emergano differenze significative circa l'uso di verbi, l'eloquio dementigeno risulta essere comunque contraddistinto da una percentuale

minore di verbi d'azione e di verbi coniugati alla prima persona plurale, al netto di una maggiore percentuale di verbi coniugati alla seconda persona singolare. Minore è anche l'uso del modo condizionale, delle forme ausiliarie e del participio passato.

Lo stesso dicasi per l'uso dei pronomi: nel complesso le percentuali di occorrenza non differiscono, anche se si registra un minore impiego di forme pronominali relative e di forme pronominali di prima persona plurale.

L'analisi semantica degli enunciati prodotti dai pazienti ha evidenziato aspetti psicopatologici classicamente associati alle sindromi demenziali. Nello specifico, si osserva un generale decremento del riferimento alla sfera affettiva, dell'espressione di sensazioni positive e di atteggiamenti propositivi e consensuali. Gli atti referenziali risultano complessivamente ridotti, soprattutto se concernenti la sfera della socialità, i processi cognitivi di natura sensoriale (e soprattutto visiva), le attività di moto e lo spazio circostante.

Buona parte delle tematiche connesse alle attività quotidiane (occupazione, scuola, svago, ambiente domestico e sport) viene verbalizzata in minore misura; ridotta è anche la tendenza a discutere relativamente alla fisicità in generale (e, nello specifico, al corpo) nonché ad attività strettamente connesse con la cura di sé, come il nutrirsi.

Lo stile conversazionale, infine, si caratterizza per un minore numero di riempitivi.

Si conferma la complessiva preservazione delle competenze sintattiche, seppure sia evidenziabile una minore lunghezza media dell'enunciato e, in generale, una minore complessità delle strutture sintagmatiche, al netto di una maggiore variabilità circa la massima profondità dell'*embedding*.

Proprio relativamente alla lunghezza media dell'*embedding* si registra un dato interessante, poiché per i pazienti è risultato significativamente maggiore. Sebbene in letteratura sia ampiamente riconosciuta la generale preservazione delle competenze sintattiche nei soggetti affetti da demenza (ovvero la produzione di strutture sintattiche corrette ma semplificate), anche solo intuitivamente appare poco verosimile che i pazienti siano in grado di formulare costruzioni con un maggiore grado di *embedding* rispetto ai soggetti cognitivamente integri. A meno che non si tratti semplicemente di un errore a livello del *parsing* sintattico automatico, l'informazione ottenuta potrebbe trovare chiarimento mediante un'analisi qualitativa di livello macrolinguistico: il fatto che si registri un maggiore livello di incassatura sintattica non implica che, anche quando formalmente ben costruite, le costruzioni esibiscano lo stesso grado di correttezza a livello pragmatico-discorsivo. *Deficit* tutt'altro che rari, infatti, investono la capacità di produrre testi coesi e coerenti (a livello locale e/o globale), da cui consegue l'inefficacia comunicativa delle produzioni dei pazienti. Come si vedrà più avanti (§ 5), tali compromissioni si riflettono anche sulla struttura sintattica delle enunciazioni.

I dati relativi alla leggibilità confermano sostanzialmente quanto emerso analizzando i singoli livelli linguistici: l'eloquio dementigeno si caratterizza per una maggiore leggibilità di base, da ricondurre fondamentalmente alla minore lunghezza media degli enunciati e alla produzione di parole più brevi, così come

per una minore leggibilità lessicale, imputabile al fatto che i testi prodotti siano maggiormente imperniati su lemmi appartenenti al lessico fondamentale, ad alta frequenza e ad alta disponibilità (i.e., maggiore ambiguità lessicale dovuta alla bassa specificità del lessico impiegato), ma anche alla presenza di false partenze, frammenti fonetici e parafasie di diversa tipologia.

Il profilo linguistico-comunicativo finora tracciato può plausibilmente arricchirsi di ulteriori caratteristiche: *feature* linguistiche emerse come significativamente divergenti solo in uno dei *task* (eventualmente anche in due), ma non nella produzione linguistica complessiva. In questa sede vengono dunque presentati come “potenziali” correlati linguistici della demenza, poiché certamente meritevoli di futuri approfondimenti.

Da un punto di vista fonetico-acustico, si osserva una minore *Pause Rate* (indice di un maggiore numero di pause rispetto al numero di parole prodotte), una maggiore variabilità del Centroide Spettrale, così come un valore maggiore della *Higuchi Fractal Dimension* (indicativa del grado di somiglianza/identità tra strutture contigue nelle bande spettrali computate su sequenze temporali discrete) e una minore variabilità di tale indice.

Dal punto di vista lessicale, la lunghezza media dell'enunciato e le occorrenze di congiunzioni subordinanti, del verbo essere, dei nomi, dei verbi transitivi e del tempo presente, nonché il peso medio della frequenza d'uso delle parole sono risultati significativamente minori. L'uso dei nomi propri, invece, è certamente da ridefinirsi in futuro, poiché i dati emersi sono contrastanti.

Sul piano semantico, l'eloquio dementigeno talvolta pare caratterizzarsi per una minore espressione di processi cognitivi in generale, di riferimenti alla possibilità, di emozioni negative (specialmente della rabbia), di riferimenti agli affetti familiari e alle situazioni inclusive. Significativamente maggiori risultano, invece, sia il riferimento a processi introspettivi che l'espressione di inibizione.

Relativamente al piano sintattico, si dimostra significativamente minore per le produzioni linguistiche dei pazienti un'ulteriore misura della complessità sintattica, computata in termini di elementi linguistici suscettibili di incrementare la profondità delle strutture (es. congiunzioni subordinanti, pronomi *wh*).

Quest'ultimo dato viene confermato dalla maggiore leggibilità sintattica dei testi orali prodotti dai pazienti, emersa per il *task 2* e riconducibile alla minore complessità delle strutture realizzate (in termini di articolazione interna dei periodi e delle proposizioni) e alla minore lunghezza delle relazioni di dipendenza (lunghezza media delle dipendenze e media delle dipendenze massime).

#### 4.2 Quadri dementigeni e correlati linguistici

Nel complesso, i *Digital Linguistic Biomarker* risultati significativamente differenti tra i pazienti e i coetanei cognitivamente integri possono essere interpretati come correlati linguistici “più forti” della demenza: sono emersi nonostante l'eterogeneità del gruppo dei pazienti e dunque, presumibilmente, rappresentano aspetti linguistici più comunemente oggetto di degradazione, indipendentemente

dal processo neurodegenerativo sottostante. Allo stesso tempo, è bene sottolineare che proprio tale composizione mista può aver oscurato le peculiarità linguistiche di ogni quadro dementigeno: nonostante i sottogruppi di pazienti (ad eccezione dei soggetti con morbo di Alzheimer) non siano numerosi, la breve analisi in merito alla significatività delle *feature* estratte ha permesso di evidenziare per ognuno di essi un determinato corollario di correlati linguistici. Tuttavia i risultati devono essere interpretati con grande cautela: non sono infatti generalizzabili e certamente necessitano di future conferme.

Innanzitutto, la durata degli enunciati varia sensibilmente a seconda del quadro dementigeno al quale il paziente è ascrivibile. Da questo punto di vista, un risultato particolarmente rilevante concerne l'esito del confronto tra i dati estratti per i soggetti con demenza non ulteriormente specificata e quelli estratti per tutti gli altri pazienti: l'assenza di significatività registrata unicamente nel raffronto con i casi di demenza mista potrebbe suggerire l'identificazione dei casi di demenza non specificata con forme di demenza mista non ufficialmente diagnosticate. Se validata da indagini future, tale procedura potrebbe rappresentare un interessante supporto alla diagnosi differenziale in campo clinico.

In secondo luogo, la produzione linguistica dei soggetti affetti da Alzheimer risulta caratterizzata dalla maggiore percentuale di *feature* significativamente differenti tra eloquio dementigeno ed eloquio normotipico (72%); il 33% distingue il gruppo di soggetti con demenza vascolare, il 21% il gruppo di soggetti con demenza non specificata, il 7% il gruppo di soggetti con demenza mista e il 5% il soggetto con demenza fronto-temporale.

Sebbene la maggior percentuale di tratti linguistici significativi registrata per i soggetti con Morbo di Alzheimer possa riflettere il fatto che numerose *feature* proposte in letteratura per l'identificazione delle patologie dementigene siano state individuate proprio mediante l'analisi del linguaggio alzheimeriano, i risultati ottenuti evidenziano piuttosto la necessità di acquisire nuovi dati: non è probabilmente un caso che la percentuale di tratti significativi (e la significatività stessa del *p-value*) decresca parallelamente al decrescere del numero di informanti per ogni sottogruppo.

Gli aspetti forse più interessanti nell'ottica di una discriminazione diagnostica *language-based* riguardano, da un lato, l'evidenziazione di caratteristiche significativamente differenti tra le singole coorti di pazienti (es. la percentuale maggiore di interiezioni dei soggetti con demenza vascolare rispetto ai soggetti con demenza non specificata, così come la minore complessità sintattica delle produzioni linguistiche di quest'ultima coorte rispetto a quanto si osserva per il soggetto con demenza fronto-temporale) e, dall'altro, l'identificazione di tratti significativamente divergenti solo tra alcuni sottogruppi di pazienti (es. la percentuale di disfluenze e di congiunzioni coordinanti, in entrambi i casi maggiori per i soggetti con demenza vascolare rispetto ai soggetti con malattia di Alzheimer) ma non tra il gruppo di pazienti e quello di controllo.

### 5. *Analisi qualitative: risultati*

L'analisi linguistica qualitativa ha consentito l'osservazione di ulteriori peculiarità dell'eloquio dementigeno.

Particolarmente evidente è la compromissione pragmatico-discorsiva: agli usi frequentemente agrammaticali di pronomi e congiunzioni, ad esempio, possono ricondursi la coesione referenziale e la coesione locale/globale deficitarie dei testi orali prodotti dai pazienti.

Per quanto riguarda i pronomi, è ben nota in letteratura la tendenza dei soggetti con demenza a impiegarne in alta percentuale: la ragione va sicuramente ricercata nella maggiore accessibilità e nel minore dispendio articolatorio che caratterizzano le forme pronominali, aspetti che permettono al paziente di portare a termine un atto referenziale cognitivamente più "economico" e altrimenti di difficile attuazione in virtù dei forti *deficit* di denominazione. La criticità più evidente, però, è relata al fatto che la pronominalizzazione nell'eloquio dementigeno spesso non riveste la sua funzione canonica: nonostante le occorrenze possano essere numericamente simili a quelle dei soggetti sani, i pazienti spesso ricorrono ai pronomi senza che sia effettivamente introdotto un referente nel cotesto linguistico antecedente/successivo e/o senza che ne venga indicato uno nel contesto extralinguistico. La comprensione dei riferimenti pronominali, anche di natura deittica, risulta perciò spesso impossibile, come esemplificato dal seguente estratto:

- (1) SOGG\_02 – Task 1<sup>5</sup>  
 8. si mettono sempre con # tutte quelle schifezze / # &po mi vengono a prendere / # e che se no # quello che &li devono fare //

Lo stesso dicasi per l'uso delle congiunzioni, sebbene sia necessaria una puntualizzazione: è più verosimile che i pazienti riescano a costruire costruzioni paratattiche grammaticalmente ben formate piuttosto che strutture ipotattiche ugualmente ben costruite, in virtù della maggiore dispendiosità cognitiva di queste ultime<sup>6</sup>. Detto ciò, non è scontato che la funzione delle congiunzioni subordinanti, anche se inserite nel discorso e computate come indici di complessità sintattica, sia poi adeguatamente concretizzata, similmente a quanto si è osservato per l'uso dei pronomi: anche se ne manifesta la progettazione, infatti, è possibile che il paziente non completi la realizzazione della subordinazione. Sono diversi gli esempi di questo tipo, ovvero enunciati nei quali il soggetto enuclea una costruzione ipotattica, salvo poi lasciarla in sospeso:

<sup>5</sup> Il SOGG\_02, una paziente affetta da morbo di Alzheimer, impiega sì forme pronominali (sia esplicite che sottintese), ma per nessuna di esse è stato specificato/indicato un referente che possa rendere pienamente comprensibile il messaggio formulato.

<sup>6</sup> Le occorrenze di congiunzioni coordinanti risultano essere significativamente maggiori rispetto a quelle delle congiunzioni subordinanti, sia nei testi orali prodotti dai pazienti (Test di Kruskal-Wallis; *p-value* < 0.001) che in quelli prodotti dai soggetti di controllo (Test di Kruskal-Wallis; *p-value* < 0.001).

- (2) SOGG\_03 – Task 2a  
1. sono loro che # stanno # &tr [/1] trovando delle [/1] delle avvolte / nelle cose così / perché +
- (3) SOGG\_07 – Task 2  
3. se devo andare a # tanto per dire / # io vengo da te # prima / # dopo dieci giorni nemmeno +

Nei casi in cui, invece, le strutture sintattiche siano formalmente corrette, tangenzialità, prolissità e/o deragliamento ideativo plausibilmente spiegano perché le enunciazioni dei pazienti appaiono sia inefficaci dal punto di vista comunicativo, sia caratterizzate da un maggior grado di *embedding* da un punto di vista quantitativo (cfr. § 4): come esemplificato dai seguenti passaggi, spesso i contenuti proposizionali veicolati risultano totalmente tangenziali rispetto al *topic* della frase principale (e, in generale, del *task*), e sono ricchi di informazioni irrilevanti, implausibili e/o ridondanti.

- (4) SOGG\_10 – Task 3  
10. questo è un attrezzo # del medioevo / # che veniva usato per stirare [/1] # tirare la luna / # una volta ca [/1] fatto il cardinamento della lana / # a per si [/1] si fanno [/2] delle [/1] delle a forma di palla / # dei rotoli a forma di palla / # e allora xxx più tesa è la [/1] corda e più vengono bene le palle a forma di [/3] a forma di palla //
- (5) SOGG\_10 – Task 2  
14. eh / una volta lì / non c'erano # medici # non c'erano i cosiddetti medici di famiglia / # allora la &buonanim di mio padre che / da questo io ho compreso che poi # era molto affezionato ai figli / # e difatti ne ha avuto nove / due sono morti piccolini / # e all'epoca non c'erano le cure che ci sono oggi / # e allora si moriva facilmente / quando si era &ba [/1] bimbi //

Il SOGG\_10, durante la somministrazione del task 3, si pronuncia nel modo sopra riportato durante la descrizione dei ferri per il lavoro a maglia, rappresentati nella figura usata per l'elicitazione linguistica. Per quanto non si ravvisino agrammaticalità sintattiche, l'enunciazione non può dirsi corretta dal punto di vista semantico-lessicale e pragmatico-discorsivo. Sebbene sia stata espressa la funzione dell'oggetto da identificare (i ferri), nel discorso figura un'espressione certamente non pertinente (probabilmente una forma di parafasia verbale, ovvero “cardinamento”), così come appare non del tutto appropriata la definizione dell'oggetto come “attrezzo del medioevo”. Gli enunciati prodotti durante il *task* 2, in aggiunta, ben esemplificano il deragliamento semantico e la verbalizzazione di informazioni fondamentalmente non pertinenti rispetto al *topic* del *task* (le tradizioni legate alla giornata del Natale), nonostante la correttezza grammaticale e la complessità delle costruzioni sintattiche realizzate.

Terminando questa sezione, si sottolinea che ulteriori caratteristiche ben evidenziabili dalla lettura del *corpus* riguardano:

1. l'uso della ripetizione come “riempitivo defaticante”, al quale i soggetti ricorrono per sopperire alle difficoltà di prosecuzione dello scambio dialogico e che, a ben vedere, rappresenta la manifestazione “minima” di un più generale fenomeno di iterazione argomentativa, più evidente nei casi di minore severità della malattia

- e testimonianza dell'esacerbazione dell'andamento epicicloidale che caratterizza i testi orali anche in condizioni normotipiche;
2. i riflessi linguistici dei severi *deficit* di denominazione da cui i pazienti sono affetti, a volte esplicitamente dichiarati, alle volte compensati da strategie come la circonlocuzione o l'impiego di espressioni deittiche, così come le frequenti parafasie fonemiche/fonetiche, semantiche e neologistiche;
  3. la degradazione della memoria autobiografica, che si esplica nella trasposizione diacronica del passato nel presente e nella misidentificazione di affetti familiari in oggetti/persone estranee/non animate. Se l'uso di pronomi e aggettivi personali e possessivi sembra rispecchiare una percezione più o meno salda della propria identità, per i pazienti la collocazione spazio-temporale delle proprie esperienze di vita e dei propri affetti è invece evidentemente deficitaria.

## 6. Conclusioni

Complessivamente, i risultati sono stati coerenti con quanto riportato da studi condotti su altre lingue, mentre alcune delle *feature* proposte non sono risultate statisticamente significative (es. la maggiore percentuale di pronomi, di espressioni deittiche spaziali, la minore ricchezza lessicale, cfr. Boschi *et al.*, 2017; Millington & Luz, 2021): ciò a testimonianza di quanto ogni profilazione dei *deficit* linguistico-comunicativi debba essere tracciata alla luce della specifica lingua impiegata dai parlanti.

L'indagine qualitativa del *corpus* concorre a dettagliare il profilo delle disabilità linguistico-comunicative associabili alla demenza, permettendo di evidenziare *deficit* a livello pragmatico-discorsivo non desumibili esclusivamente dai dati quantitativi computati (es. errori di coesione referenziale e di coerenza locale/globale), nonché ulteriori peculiarità, tra cui notevole presenza di parafasie e false partenze, attuazione di strategie in compensazione della severa anomia e contestualizzazione spazio-temporale degli eventi deficitaria.

Tutte le *feature* risultate già significative nel confronto tra le condizioni di *Mild Cognitive Impairment* (MCI; cfr. Petersen *et al.*, 1999) e di integrità cognitiva per soggetti italofofoni (cfr. Calzà *et al.*, 2021) sono risultate altrettanto discriminanti in questo studio: fanno eccezione solo due indici di ricchezza lessicale (TTR e BrunetW). La motivazione potrebbe risiedere in errori di annotazione (e successiva computazione) automatica, o più semplicemente nelle differenti fasce d'età coinvolte nei due studi: con l'avanzare dell'età, la ricchezza lessicale decresce anche in condizioni di normalità cognitiva e, dunque, per i soggetti più anziani tali aspetti semantico-lessicali potrebbero non essere più significativamente differenzianti.

Un aspetto cruciale del disegno sperimentale presentato in questa sede concerne indubbiamente le modalità di annotazione del *corpus* acquisito e, conseguentemente, di estrazione dei DLBs: l'annotazione a tutti i livelli linguistici è avvenuta in maniera automatica, senza che – per questioni di tempo – si procedesse con una revisione manuale degli *output* prodotti dalla *pipeline*. Trattandosi di procedure automatiche, la percentuale di errore è certamente più alta rispetto a un'annotazione manuale,

ma la questione più spinosa è di natura sociolinguistica: sebbene i soggetti reclutati siano italofofoni, sono tutt'altro che rare le inflessioni dialettali (a livello soprattutto lessicale); non è da escludere, dunque, che vi possano essere stati errori di *PoS-Tagging*, lemmatizzazione e, conseguentemente, di *parsing* sintattico, dovuti al mancato riconoscimento delle forme dialettali usate dai soggetti.

Nonostante queste problematiche, i risultati ottenuti sono molto promettenti, poiché la semplice annotazione automatica ha permesso di evidenziare un considerevole numero di tratti significativamente correlati alla demenza. Sarebbe certamente interessante ripetere le analisi condotte dopo aver revisionato manualmente l'annotazione automatica del *corpus*: da un lato, per individuare eventuali divergenze rispetto ai DLBs risultati significativi in questo studio e, dunque, per incrementare/diminuire i correlati linguistici della demenza per la lingua italiana; dall'altro, per avere una stima della bontà dei risultati ottenibili esclusivamente con una totale automatizzazione delle procedure di annotazione ed estrazione delle *feature* linguistiche.

In futuro si dimostrerà sicuramente fruttuoso ampliare il *corpus* già raccolto, reclutando un maggiore numero di pazienti: 1) ascrivibili a quadri dementigeni differenti, in modo da approfondirne le singole profilazioni linguistico-comunicative; 2) affetti da demenza di severità differente, così da osservare al meglio variazioni delle competenze linguistiche legate alla gravità della compromissione cognitiva; 3) con provenienza regionale differente, affinché si possano individuare anche eventuali divergenze relate alla ricca variazione diatopica della lingua italiana; 4) con vari livelli di istruzione, al fine di valutare la possibile persistenza di tratti linguistici primariamente imputabili all'utilizzo di varietà *substandard* dell'italiano e ad una considerevole esposizione a sistemi linguistici dialettali.

L'acquisizione di nuovi dati sarà naturalmente preziosa per la generalizzabilità dei risultati e per il miglioramento dei processi di annotazione linguistica automatica, poiché i sistemi di NLP implicati (es. ASR, *Pos-Tagging*, ecc.) esigono grandi quantità di dati relativi alle produzioni linguistiche non normotipo per essere addestrati e ottenere migliori prestazioni.

Infine, incrementare i *Digital Linguistic Biomarker* da indagare (per le più recenti rassegne al riguardo: de la Fuente Garcia *et al.*, 2020; Petti *et al.*, 2020; Vigo *et al.*, 2022) e integrare i risultati quantitativi con riflessioni qualitative consentirà indubbiamente di delineare una caratterizzazione linguistico-comunicativa del deterioramento cognitivo più precisa e completa.

### *CRedit Author Statement*

EM: Investigation, Data Curation (i.e., linguistic samples collection), Writing – Original Draft; VG, FM, IN, FC, MM, GDB: Resources (i.e., patient enrollment); GG: Conceptualization, Methodology, Supervision, Writing – Review & Editing.

*Riferimenti bibliografici*

ALZHEIMER'S DISEASE INTERNATIONAL (2021). *World Alzheimer Report 2021: Journey through the diagnosis of dementia*. London: Alzheimer's Disease International.

AMERICAN PSYCHIATRIC ASSOCIATION (2013). *Diagnostic and statistical manual of mental disorders (5th ed.)*. Washington (DC): American Psychiatric Publishing. DOI: 10.1176/appi.books.9780890425596

AUSTIN, J.L. (1962). *How to do things with words*. Oxford: Clarendon Press.

BECKER, J.T., BOILER, F., LOPEZ, O.L., SAXTON, J. & MCGONIGLE, K.L. (1994). The natural history of Alzheimer's disease: description of study cohort and accuracy of diagnosis. In *Archives of neurology*, 51(6), 585-594. DOI: 10.1001/archneur.1994.00540180063015

BENVENUTI, N., BOLIOLI, A., MAZZEI, A., VIGORELLI, P. & BOSCA, A. (2020). The "Corpus Anchise 320" and the analysis of conversations between healthcare workers and people with dementia. In DELL'ORLETTA, F., MONTI, J. & TAMBURINI, F. (a cura di), *Proceedings of the Seventh Italian Conference on Computational Linguistics CLiC-it 2020*. Torino: Accademia University Press.

BERRUTO, G. (2021). *Sociolinguistica dell'italiano contemporaneo*. Roma: Carocci Editore.

BOSCHI, V., CATRICALÀ, E., CONSONNI, M., CHESI, C., MORO, A. & CAPPA, S.F. (2017). Connected speech in neurodegenerative language disorders: a review. In *Frontiers in Psychology*, 8, 269. DOI: 10.3389/fpsyg.2017.00269

CALTAGIRONE, C., GAINOTTI, G., CARLESIMO, G.A. & PARNETTI, L. (1995). Batteria per la valutazione del deterioramento mentale: I. Descrizione di uno strumento di diagnosi neuropsicologica. In *Archivio di Psicologia, Neurologia e Psichiatria*, 56(4), 461-470.

CALZÀ, L., GAGLIARDI, G., ROSSINI FAVRETTI, R. & TAMBURINI, F. (2021). Linguistic features and automatic classifiers for identifying Mild Cognitive Impairment and dementia. In *Computer Speech & Language*, 65, 101-113. DOI: 10.1016/j.csl.2020.101113

CARLESIMO, G., CALTAGIRONE, C., GAINOTTI, G. & NOCENTINI, U. (1995). Batteria per la valutazione del deterioramento mentale: standardizzazione ed affidabilità diagnostica nell'identificazione di pazienti affetti da sindrome demenziale. In *Archivio di Psicologia Neurologia e Psichiatria*. 56, 471-488.

CIURLI, P., MARANGOLO, P. & BASSO, A. (1996). *Esame del Linguaggio II. Manuale e materiale d'esame*. Firenze: Giunti.

CONTI, S., BONAZZI, S., LAIACONA, M., MASINA, M. & VANELLI CORALLI, M. (2015). Montreal Cognitive Assessment (MoCA) – Italian version: regression-based norms and equivalent scores. In *Neurological Sciences*, 26, 209-214. DOI: 10.1007/s10072-014-1921-3

COSTA, A., SINFORIANI, E. (2020). *Le demenze. Manuale di diagnosi e trattamento*. Milano: NEMS.

CRESTI, E., MONEGLIA, M. (2018). Chapter 13. The illocutionary basis of Information Structure: The Language into Act Theory (L-Act). In ADAMO, E. ET AL. (a cura di), *Information Structure in Lesser-described Languages: Studies in prosody and syntax*. Amsterdam: John Benjamins Publishing Company, 360-402.

CRYSTAL, D. (1981). *Clinical Linguistics*. Wien: Springer Verlag.

DE LA FUENTE GARCIA, S., RITCHIE, C.W. & LUZ, S. (2020). Artificial Intelligence, Speech, and Language Processing Approaches to Monitoring Alzheimer's Disease: A Systematic Review. In *Journal of Alzheimer's disease*, 78(4), 1547-1574. DOI: 10.3233/JAD-200888

- DE MAURO, T. (2000). *Il dizionario della lingua italiana*. Torino: Paravia.
- DELL'ORLETTA, F., MONTEMAGNI, S. & VENTURI, G. (2011). READ-IT: Assessing readability of Italian texts with a view to text simplification. In NORMAN, A. (a cura di), *Proceedings of the Second Workshop on Speech and Language Processing for Assistive Technologies*. Stroudsburg (PA): ACL, 73-83.
- DOVETTO, F.M., GUIDA, A., PAGLIARO, A.C., GUARASCI, R., RAGGIO, L., SORRENTINO, A. & TRILLOCCO, S. (2022). Corpora di Italiano Parlato Patologico dell'età adulta e senile. In CRESTI, E., MONEGLIA, M. (a cura di), *Corpora e Studi Linguistici. Atti del LIV Congresso della Società di Linguistica Italiana*. Milano: Officinaventuno, 165-177.
- MAX PLANCK INSTITUTE FOR PSYCHOLINGUISTICS (2021). ELAN. Nijmegen: Max Planck Institute for Psycholinguistics, The Language Archive. [Computer Software] Version 6.2. <https://archive.mpi.nl/tla/elan>
- GAGLIARDI, G., TAMBURINI, F. (2021). Linguistic biomarkers for the detection of Mild Cognitive Impairment. In *Lingue e linguaggio*, 1/2021, 3-31.
- GAGLIARDI, G. (2019). *Linguistica per le professioni sanitarie*. Bologna: Pàtron.
- HERNÁNDEZ-DOMÍNGUEZ, L., RATTÉ, S., SIERRA-MARTÍNEZ, G. & ROCHE-BERGUA, A. (2018). Computer-based evaluation of Alzheimer's disease and mild cognitive impairment patients during a picture description task. In *Alzheimer's & Dementia: Diagnosis, Assessment & Disease Monitoring*, 10(1), 260-268. DOI: 10.1016/j.dadm.2018.02.004
- HONORE, A. (1979). Some simple measures of richness of vocabulary. In *Association of Literary and Linguistic Computing Bulletin*, 7, 172-177.
- MAGNI, E., BINETTI, G., BIANCHETTI, A., ROZZINI, R. & TRABUCCHI, M. (1996). Mini-Mental State Examination: a normative study in Italian elderly population. In *European Journal of Neurology*, 3(3), 198-202. DOI: 10.1111/j.1468-1331.1996.tb00423
- PAGANO, C., BOLOGNINI, N. (a cura di) (2020). *Neuropsicologia delle demenze*. Bologna: Il Mulino.
- PETERSEN, R.C., SMITH, G.E., WARING, S.C., IVNIK, R.J., TANGALOS, E.G. & KOKMEN, E. (1999). Mild cognitive impairment: clinical characterization and outcome. In *Archives of neurology*, 56(3), 303-308. DOI: 10.1001/archneur.56.3.303
- PETTI, U., BAKER, S. & KORHONEN, A. (2020). A systematic literature review of automatic Alzheimer's disease detection from speech and language. In *Journal of the American Medical Informatics Association – JAMIA*, 27(11), 1784-1797. DOI: 10.1093/jamia/ocaa174
- QI, P., DOZAT, T., ZHANG, Y. & MANNING, C.D. (2018). Universal Dependency Parsing from Scratch. In *Proceedings of the CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*. Stroudsburg (PA): ACL, 160-170. DOI: 10.18653/v1/K18-2016
- R CORE TEAM (2014). *R: A language and environment for statistical computing*. Wien: R Foundation for Statistical.
- ROARK, B., MITCHELL, M. & HOLLINGSHEAD, K. (2007). Syntactic complexity measures for detecting Mild Cognitive Impairment. In BRETONNEL COHEN, K., DEMNER-FUSHMAN, D., FRIEDMAN, C., HIRSCHMAN, L. & PESTIAN, J. (a cura di), *Proceedings of the Workshop BioNLP 2007: Biological, translational, and clinical language processing*. Stroudsburg (PA): ACL, 1-8.

ROARK, B., MITCHELL, M., HOSOM, J.-P., HOLLINGSHEAD, K. & KAYE, J.A. (2011). Spoken Language Derived Measures for Detecting Mild Cognitive Impairment. In *IEEE Transactions on Audio Speech, and Language Processing*, 19(7), 2081-2090. DOI: 10.1109/TASL.2011.2112351

SATT, A., SORIN, A., TOLEDO-RONEN, O., BARKAN, O., KOMPATSIARIS, I., KOKONOZI, A. & TSOLAKI, M. (2013). Evaluation of Speech-Based Protocol for Detection of Early-Stage Dementia. In *Proceedings of Interspeech 2013*. Grenoble: ISCA, 1692–1696.

SPINNLER, H., TOGNONI, G. (1987). Standardizzazione e taratura italiana di test neuropsicologici. Milano: Masson Italia periodici (*Supplementum 8 – Italian journal of neurological sciences*).

VIGO, I., COELHO, L. & REIS, S. (2022). Speech- and Language-Based Classification of Alzheimer's Disease: A Systematic Review. In *Bioengineering*, 9(1):27. DOI: 10.3390/bioengineering9010027

WORLD HEALTH ORGANIZATION (2012). *Dementia: a public health priority*. Geneva: World Health Organization. <https://apps.who.int/iris/handle/10665/75263>