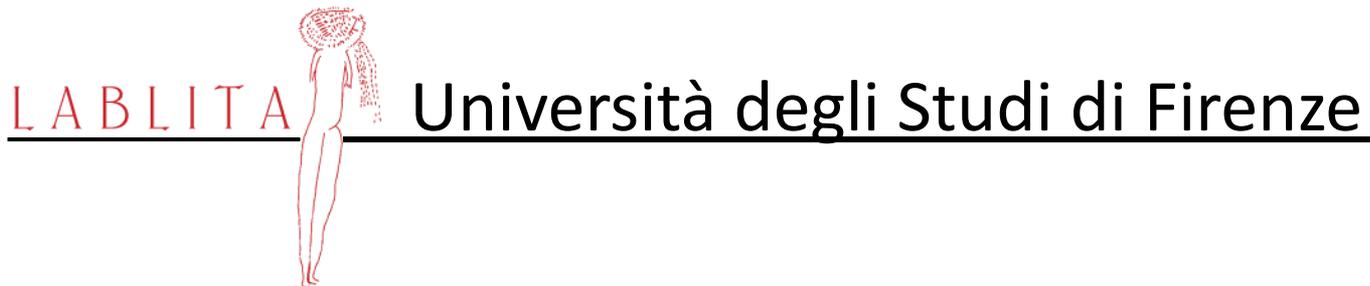


Risorse linguistiche e semantica computazionale

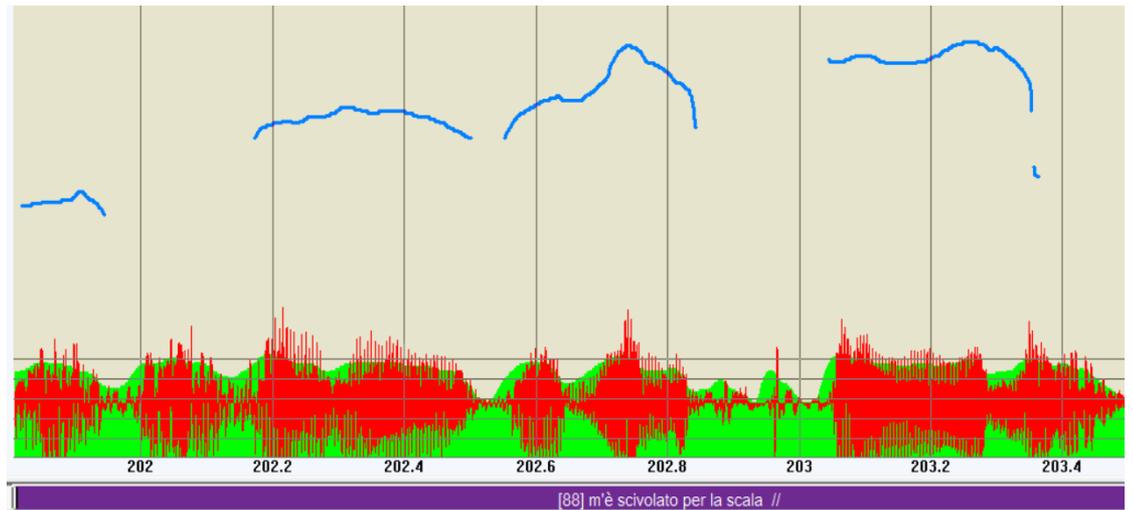
Andrea Amelio Ravelli | Lorenzo Gregori



Una lingua, varie manifestazioni

- Lingua scritta → [questa è una frase scritta]

- Lingua parlata →

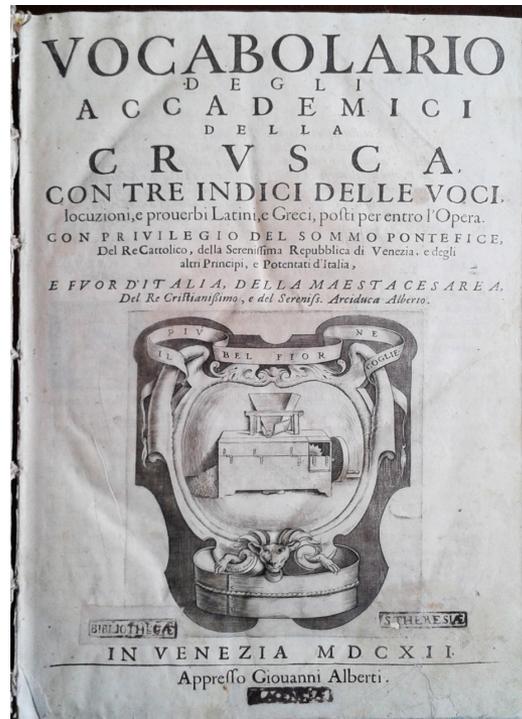


- Forme intermedie:

- parlato-scritto → chat, sms (mimesi del parlato nello scritto)
- scritto-parlato → discorso dei presidenti di camera e senato

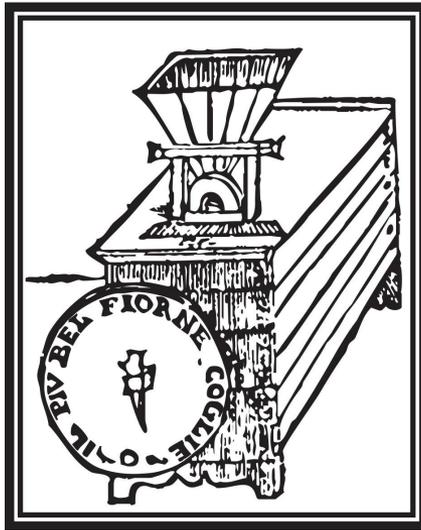
Processare la lingua

- 1612: 1^a ed. Vocabolario della Crusca
- 5 edizioni → spoglio di 1.684 opere letterarie
- Database di schede redatte a mano per ogni parola



Processare la lingua – ieri e oggi

XVII° secolo

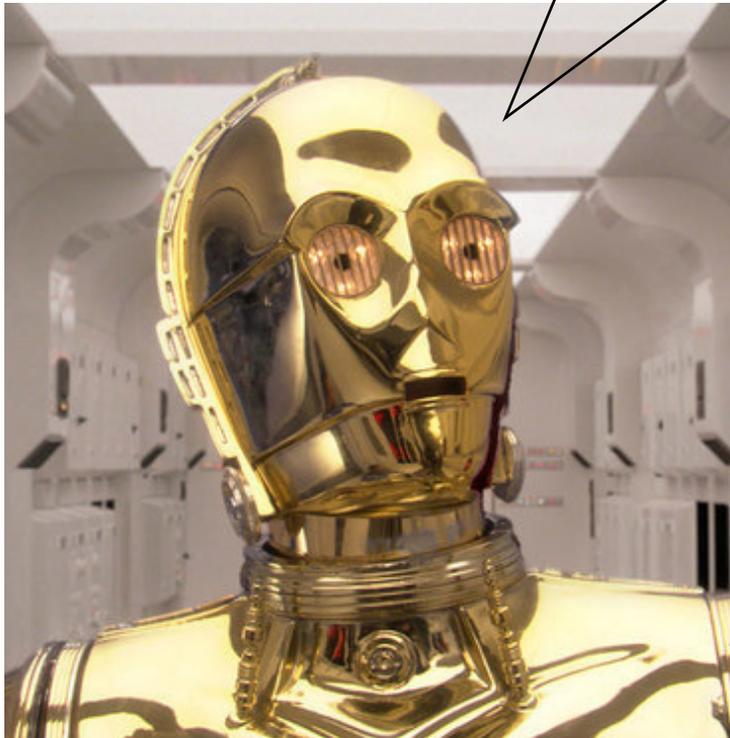


XX° secolo



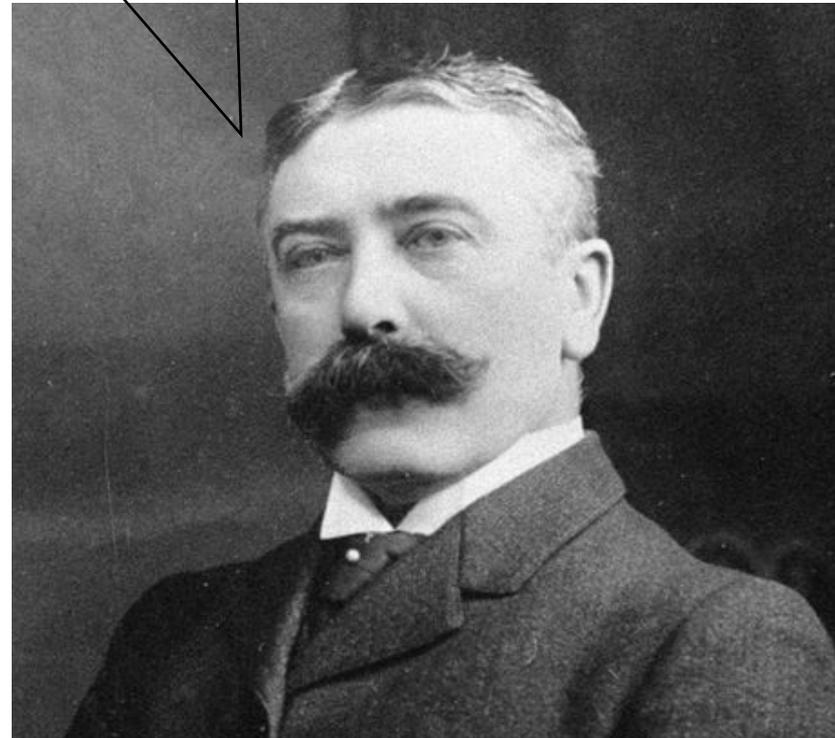
Processare la lingua – domani?

Salve, sono un droide
protocollare!



C-3PO – droide protocollare

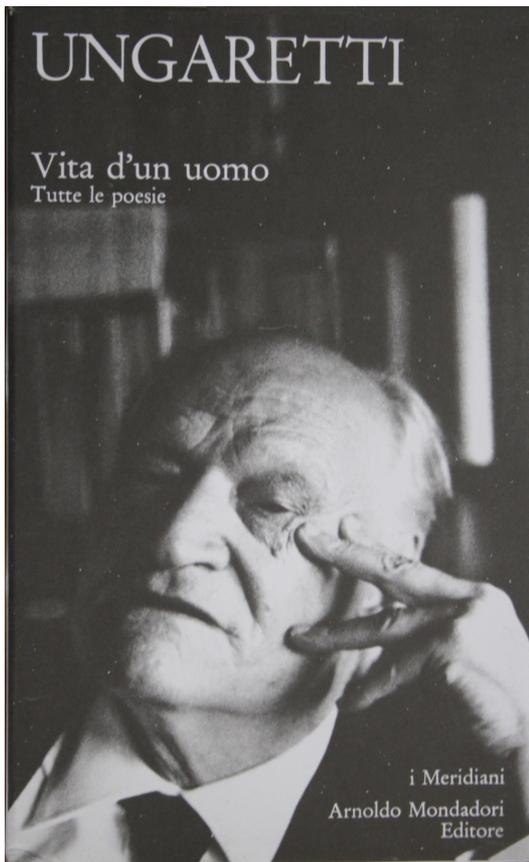
E io sono la Regina
d'Inghilterra!



Ferdinand de Saussure - linguista

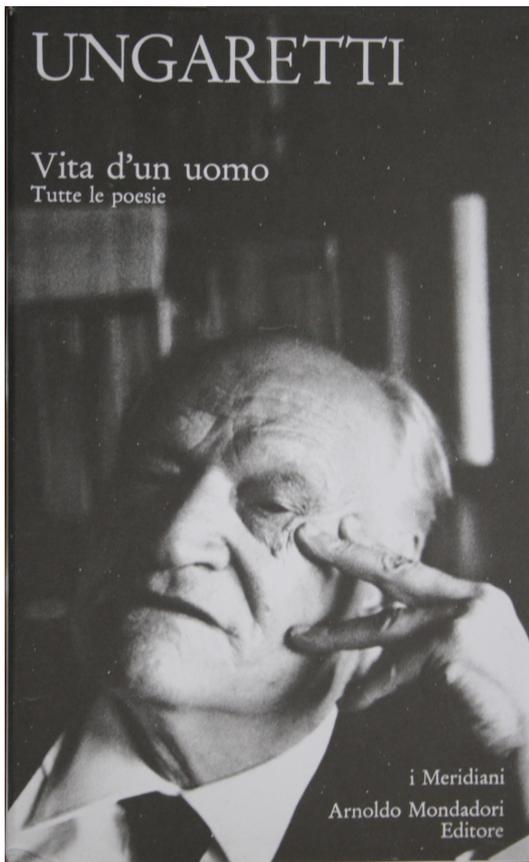
Che cos'è un corpus?

- Una raccolta di testi
 - Es. il corpus dell'opera ungarettiana



Che cos'è un corpus?

- Una raccolta di testi in formato elettronico
 - Es. il corpus dell'opera ungarettiana su un file .txt

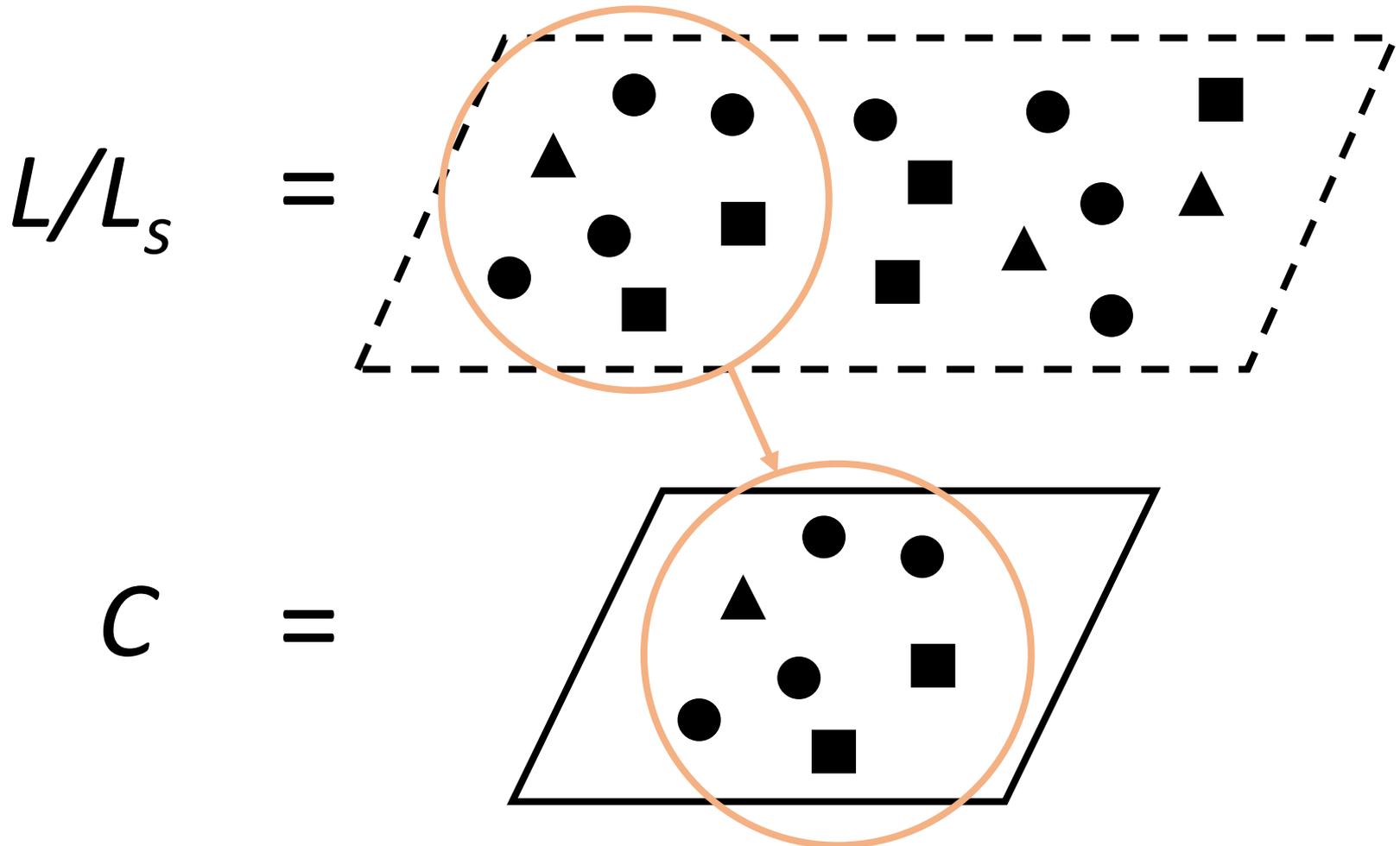


```
1 Vita d'un uomo - Tutte le Poesie, edizioni Mondadori, 1992
2
3 Poesie tratte da L'allegria 1914-1919
4
5 UN'ALTRA NOTTE
6 In quest'oscuro
7 colle mani
8 gelate
9 distinguo
10 il mio viso
11 Mi vedo
12 abbandonato nell'infinito
13
14 NOIA
15 Anche questa notte passerà
16 Questa solitudine in giro
17 titubante ombra dei fili tranviari
18 sull'umido asfalto
19 Guardo le teste dei brumisti
20 nel mezzo sonno
21 tentennare
22
23 AGONIA
24 Morire come le allodole assetate
25 sul miraggio
26 O come la quaglia
27 passato il mare
28 nei primi cespugli
29 perché di volare
30 non ha più voglia
31 Ma non vivere di lamento
32 come un cardellino accecato
33
34 [...]
```

Che cos'è un corpus?

- Una raccolta di testi in formato elettronico, opportunamente tokenizzati ed etichettati, gestibili e interrogabili informaticamente, selezionati allo scopo di essere rappresentativi della lingua presa in esame o di una sua specifica varietà.
- Un corpus è un campione rappresentativo e bilanciato di eventi comunicativi.

Campione rappresentativo e bilanciato di eventi comunicativi.



Caratteristiche dei corpora

- Generalità
 - Corpora di riferimento vs. corpora specialistici
- Cronologia
 - sincronia vs. diacronia
- Lingua
 - Monolingue vs. bi/plurilingue
- Medium
 - Scritto vs. parlato vs. multimediale

Corpora di riferimento vs. corpora specialistici

- Corpora di riferimento:
 - Campione trasversale e «omnicomprensivo» rispetto alle varietà di *L*
 - Es. British National Corpus (BNC), Perugia Corpus (PEC)
- Corpora specialistici:
 - Campione di una particolare varietà o dominio d'uso
 - Linguaggio giornalistico → Repubblica Corpus
 - Apprendenti di L2 → Lessico italiano Parlato da Stranieri (LIPS)
 - patologie linguistiche → Corpus di italiano Parlato Patologico Schizofrenico (CIPPS)

Sincronia vs. diacronia

- **Corpora sincronici:**
 - Rappresentano la lingua L in un dato momento → campione di eventi comunicativi prodotti nello stesso arco temporale
 - Es. Lessico Italiano Televisivo (LIT); Corpus di Riferimento dell'Italiano Scritto (CORIS)
- **Corpora diacronici:**
 - Rappresentano la lingua L in diversi momenti → campioni di eventi comunicativi in più finestre temporali
 - Es. DIA-LIT, DiaCORIS

Monolingue vs. bi/multilingue

- Corpora monolingue
 - Eventi comunicativi di una sola lingua
- Corpora multilingue
 - Paralleli: lo stesso evento comunicativo è rappresentato (in traduzione) in più di una lingua → Europarl
 - Comparabili: eventi comunicativi in più lingue collezionati secondo le stesse specifiche → C-ORAL-ROM

Scritto vs. parlato vs. multimediale

- Scritto: campioni esclusivi della varietà scritta di una lingua
- Parlato: campioni esclusivi della varietà parlata di una lingua
- Multimediale: campioni di eventi comunicativi completi di informazione acustica e visiva

Web Corpora

- Collezionare miliardi di parole, dai testi del web, è possibile con poche righe di codice
- Il corpus design è fondamentale per assicurare un dato quanto più rappresentativo possibile

A cosa serve un corpus?

- Frequenze
- Concordanze
- Collocazioni
- Word Sketches

Sketch

Lemma: Scegli sketch: **Calcola sketch**

postV_N	135092	
porta	9801	10,34
udienza	6438	9,79
strada	5554	8,62
occhio	3906	8,61
battente	1627	8,57
inchiesta	1947	8,48
finestra	1879	8,43
fuoco	2433	8,33
bocca	1524	8,1
dibattito	1514	7,99
varco	1046	7,92
danza	1165	7,78
sipario	927	7,75
fascicolo	973	7,61
serata	1425	7,61
concerto	1981	7,59
spiraglio	812	7,58
orizzonte	913	7,51
scenario	1004	7,44
prospettiva	1204	7,42

Corpus: Informazione

postN_V				
cucinare	<u>19154</u>	<u>0</u>	13,55	0
subire	<u>83</u>	<u>0</u>	5,81	0
raccontare	<u>108</u>	<u>0</u>	5,5	0
potere	<u>582</u>	<u>582</u>	5,18	5,23
stare	<u>203</u>	<u>254</u>	4,78	5,23
vivere	<u>103</u>	<u>112</u>	5,28	5,76
trovare	<u>146</u>	<u>182</u>	4,74	5,22
provenire	<u>47</u>	<u>40</u>	5,03	5,56
soffrire	<u>42</u>	<u>34</u>	5	5,55
appartenere	<u>33</u>	<u>30</u>	4,64	5,35
indossare	<u>32</u>	<u>29</u>	4,75	5,59
tendere	<u>32</u>	<u>31</u>	4,73	5,63
perdere	<u>54</u>	<u>79</u>	4,38	5,3
chiamare	<u>55</u>	<u>86</u>	4,47	5,5
piare	<u>31</u>	<u>38</u>	4,59	5,75
tentare	<u>39</u>	<u>89</u>	4,62	6,48
minacciare	<u>0</u>	<u>31</u>	0	5,4
picchiare	<u>0</u>	<u>20</u>	0	5,42
estrarre	<u>0</u>	<u>20</u>	0	5,44
preferire	<u>0</u>	<u>45</u>	0	5,5
aprire	<u>0</u>	<u>136</u>	0	5,9
attaccare	<u>0</u>	<u>60</u>	0	5,98
confessare	<u>0</u>	<u>38</u>	0	6,15
uccidere	<u>0</u>	<u>218</u>	0	6,7
sparare	<u>0</u>	<u>91</u>	0	6,9

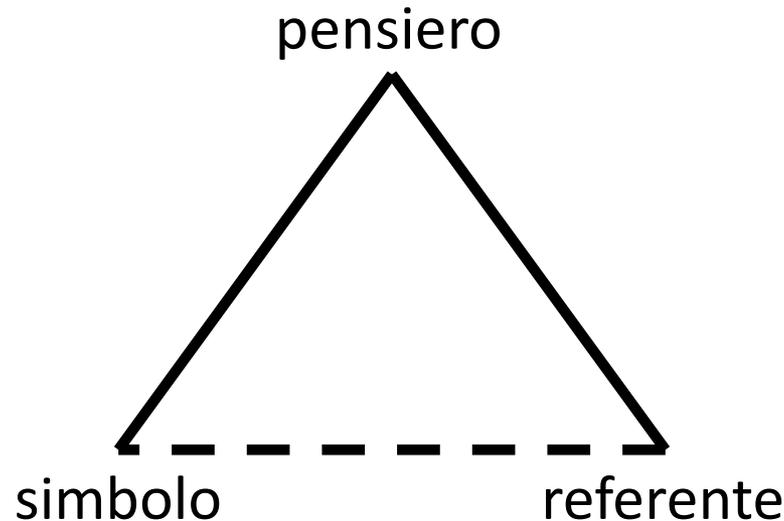
donna
-6.0
-4.0
-2.0
0.0
2.0
4.0
6.0
uomo

Corpus: Cinema

AofN				
bello	<u>619</u>	<u>0</u>	9,17	0
amato	<u>144</u>	<u>0</u>	8,77	0
occidentale	<u>137</u>	<u>0</u>	8,73	0
incinto	<u>65</u>	<u>0</u>	7,8	0
forte	<u>136</u>	<u>0</u>	7,8	0
affascinante	<u>73</u>	<u>0</u>	7,66	0
unico	<u>131</u>	<u>0</u>	7,36	0
giovane	<u>1174</u>	<u>226</u>	9,84	7,4
misterioso	<u>81</u>	<u>79</u>	7,64	7,43
maturo	<u>46</u>	<u>47</u>	7,26	7,06
anziano	<u>95</u>	<u>96</u>	8,16	7,96
moderno	<u>50</u>	<u>54</u>	6,98	6,9
sposato	<u>76</u>	<u>87</u>	8,08	8,02
perfetto	<u>90</u>	<u>127</u>	7,47	7,82
ricco	<u>50</u>	<u>115</u>	6,62	7,67
terzo	<u>0</u>	<u>107</u>	0	7,36
libero	<u>0</u>	<u>81</u>	0	7,37
medio	<u>0</u>	<u>72</u>	0	7,53
spietato	<u>0</u>	<u>77</u>	0	7,7
normale	<u>0</u>	<u>85</u>	0	7,73
invisibile	<u>0</u>	<u>87</u>	0	7,88
nero	<u>0</u>	<u>166</u>	0	7,99
qualunque	<u>0</u>	<u>94</u>	0	8,03
potente	<u>0</u>	<u>120</u>	0	8,17
comune	<u>0</u>	<u>176</u>	0	8,44

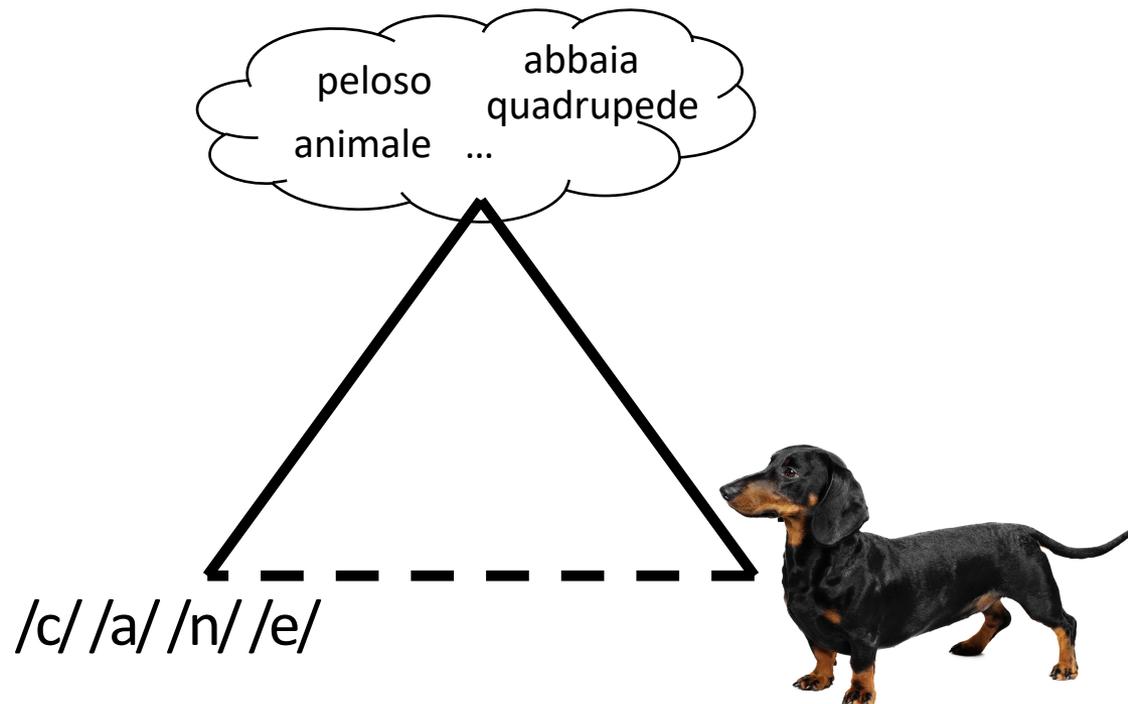
Semantica

- Studio del significato degli eventi comunicativi
- Rapporto tra simbolo, pensiero e referente



Semantica

- Studio del significato degli eventi comunicativi
- Rapporto tra simbolo, pensiero e referente



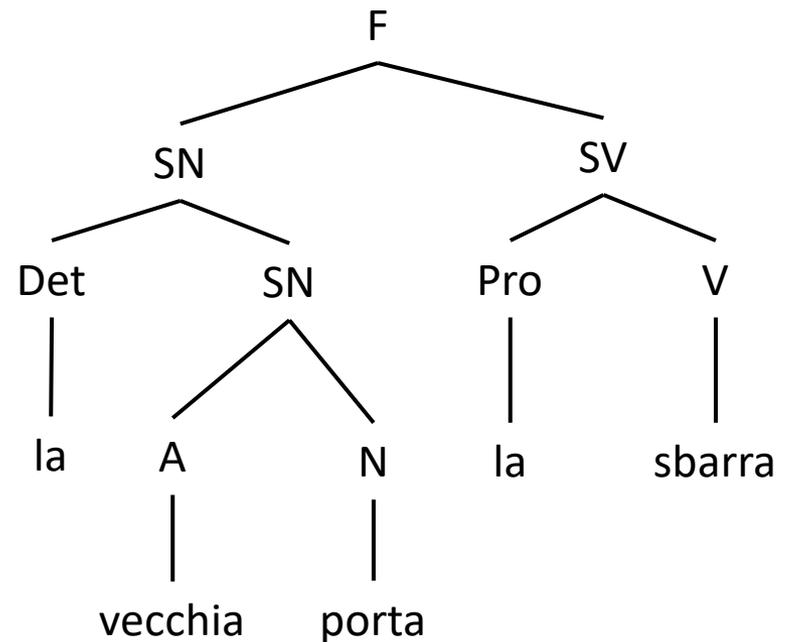
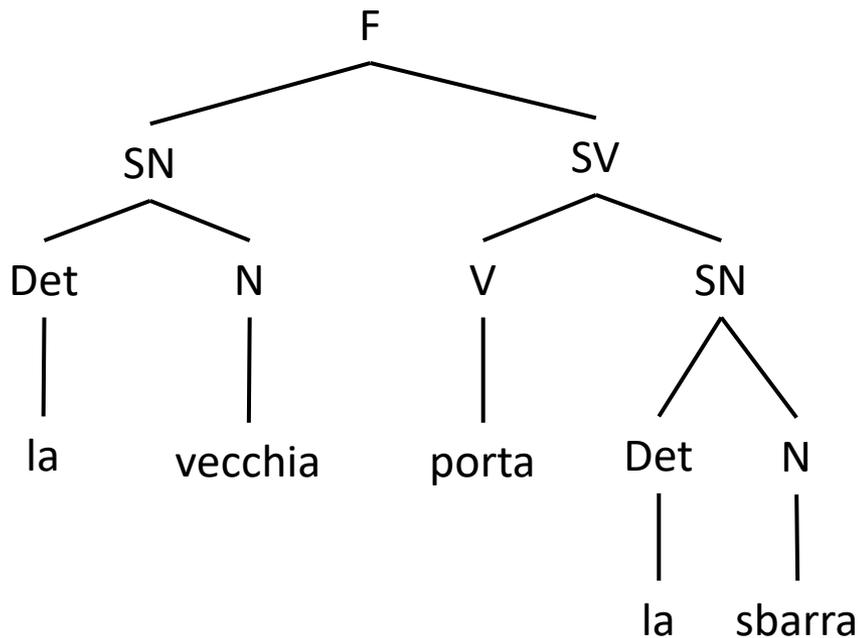
Ambiguità del senso e polisemia

La vecchia porta la sbarra



Ambiguità del senso e polisemia

La vecchia porta la sbarra



Ambiguità del senso e polisemia

- Come si risolve? → Word Sense Disambiguation
- Reti semantiche e ontologie

WordNet

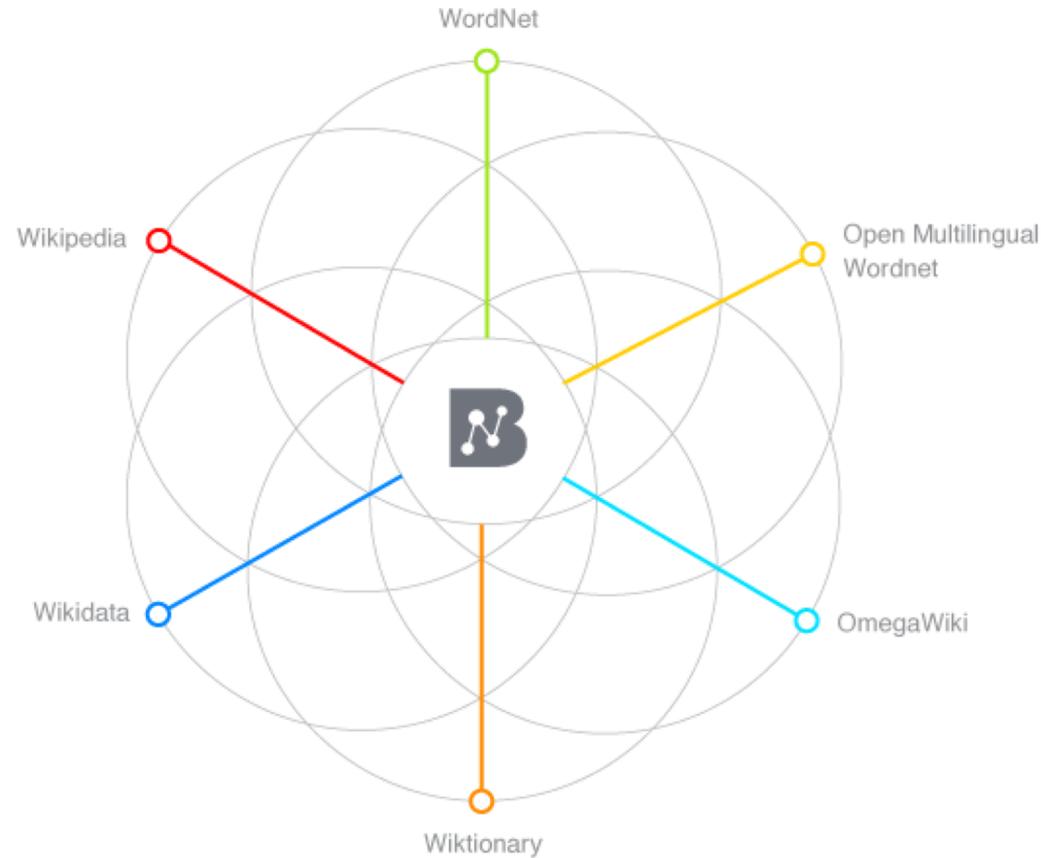
- Più importante rete semantica della lingua inglese → 117.000 synsets
- Synset: insieme di lemmi che condividono lo stesso significato
- I synset sono collegati attraverso relazioni ontologiche e semantiche

S: (v) put, set, place, pose, position, lay (put into a certain place or abstract location) *"Put your things here"; "Set the tray down"; "Set the dogs on the scent of the missing children"; "Place emphasis on a certain point"*

S: (v) put (cause to be in a certain state; cause to be in a certain relation) *"That song put me in awful good humor"; "put your ideas in writing"*

S: (v) frame, redact, cast, put, couch (formulate in a particular style or language) *"I wouldn't put it that way"; "She cast her request in very polite language"*

BabelNet



BabelNet 4.0: General statistics

Number of languages:	284
Total number of Babel synsets:	15,788,626
Total number of Babel senses:	832,469,391
Total number of concepts:	6,117,108
Total number of Named Entities:	9,671,518
Total number of lexico-semantic relations:	1,307,706,673
Total number of glosses (textual definitions):	72,542,300
Total number of images:	53,879,884
Total number of Babel synsets with at least one domain:	2,637,414
Total number of Babel synsets with at least one picture:	10,524,280
Total number of sources:	47

babelnet.org

BabelNet

bn:00090224v • VERB • Concept

EN **place**   • **lay**   • **put**   • **set**   • **pose**  

Put into a certain **place** or **abstract location**  *WordNet*  *More definitions*

Put your things here  *WordNet*  *More examples*

IS A [displace](#)

HAS KIND [appose](#) • [arrange](#) • [barrel](#) 

Translations

ZH 放, 放置, 置, 安置, 搁, 摆, 把, 设置

EN place, lay, put, set, pose, position

FR mettre, placer, poser

DE stellen, einordnen, platzieren, plazieren

EL ακουμπώ, βάζω, τοποθετώ, *που*

IT mettere, appoggiare, collocare, piazzare, porre, posare, posizionare

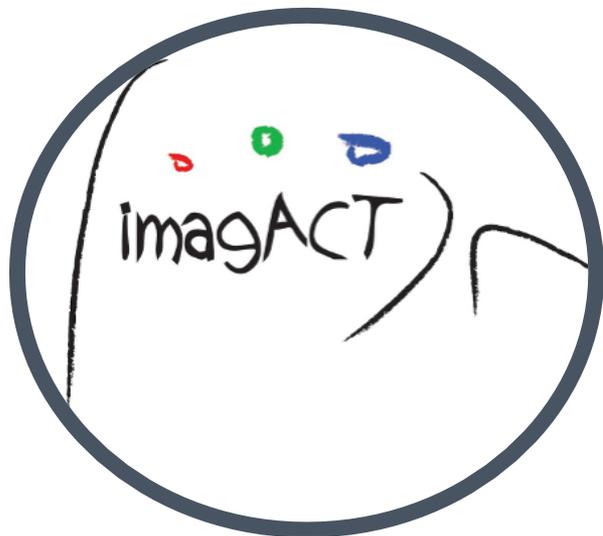
JA 置く, 付ける, 入れる, 割り当てる, 収める, 定める, 配置する, 置く, 設定

RU класть, положить, поместить, помещать, поставить, ставить

ES poner, colocar, situar, dejar, depositar, posar, posicionar, ubicar

babelnet.org

IMAGACT



www.imagact.it

- **Ontologia Multimediale e Multilingue dell'azione**
- **Video-based Translation**
- **Framework di disambiguazione dei verbi d'azione**



UNIVERSITÀ
DI SIENA
1240

Finanziato da:



Costruzione della risorsa

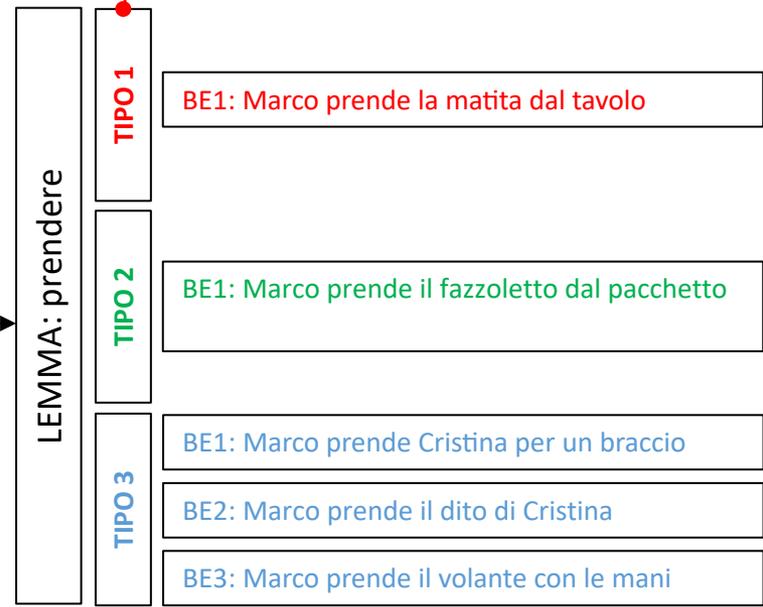


1 - Corpora di parlato ITA-EN



2 - Analisi delle occorrenze

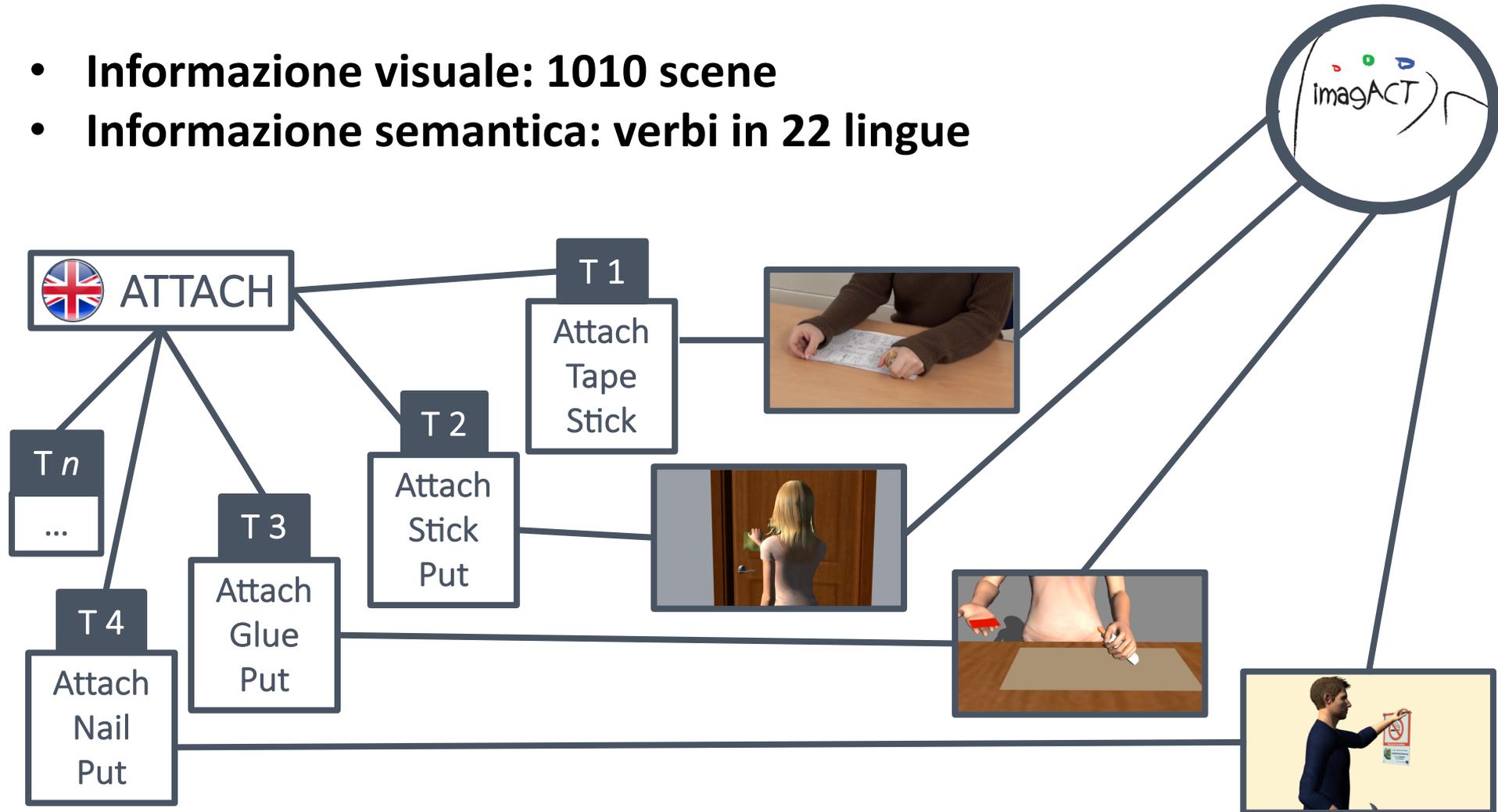
spargere	spread
attaccare	stick
indossare	wear
prendere	take
portare	bring
mettere	put



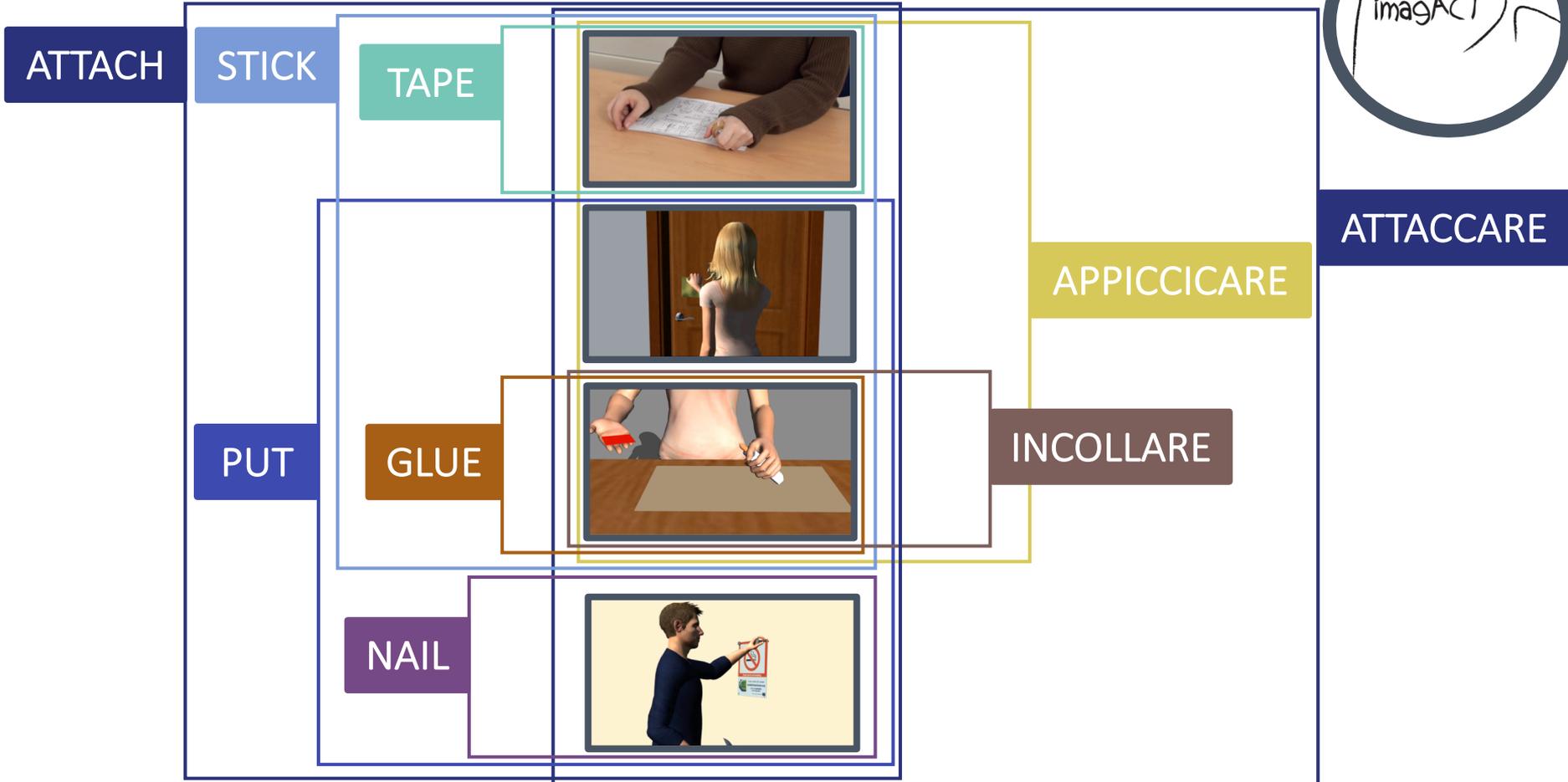
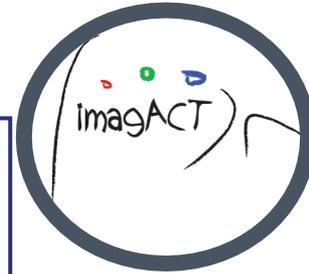
3 - Creazione dei tipi azionali

IMAGACT

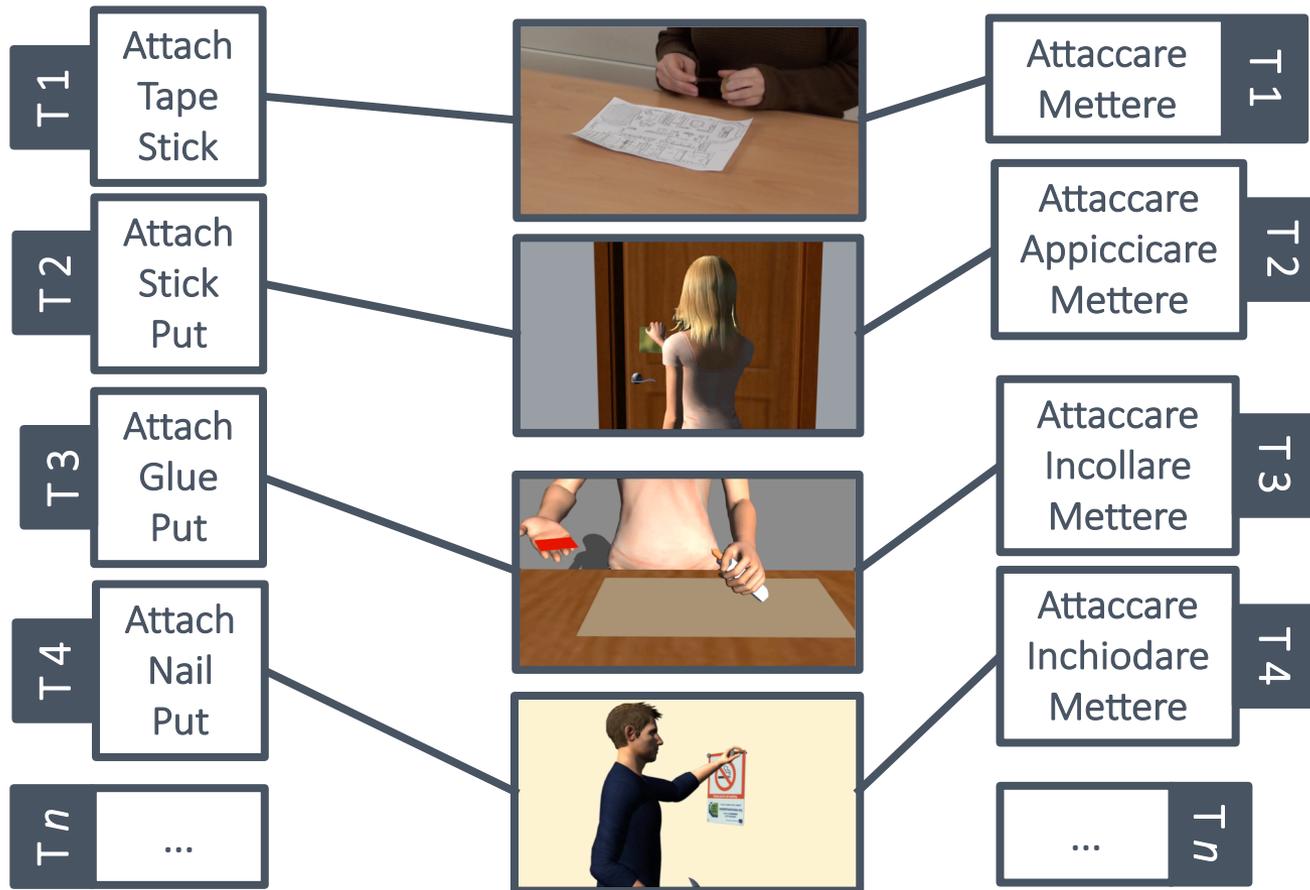
- **Informazione visuale: 1010 scene**
- **Informazione semantica: verbi in 22 lingue**



IMAGACT



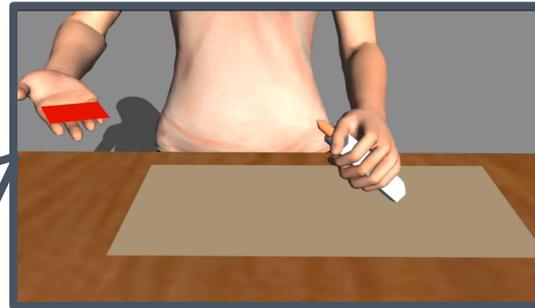
IMAGACT



IMAGACT



 ATTACH
STICK
PUT
GLUE



 ATTACCARE
APPICCIARE
INCOLLARE

 UNIR
PEGAR
FIJAR

 贴 (tiē)
粘 (zhān)

...



Dizionario multilingue

Verb

prendere

Input language

Italian

Output language

English

Show types

IT

prendere ↻

pigliare ↻

EN

take ↻



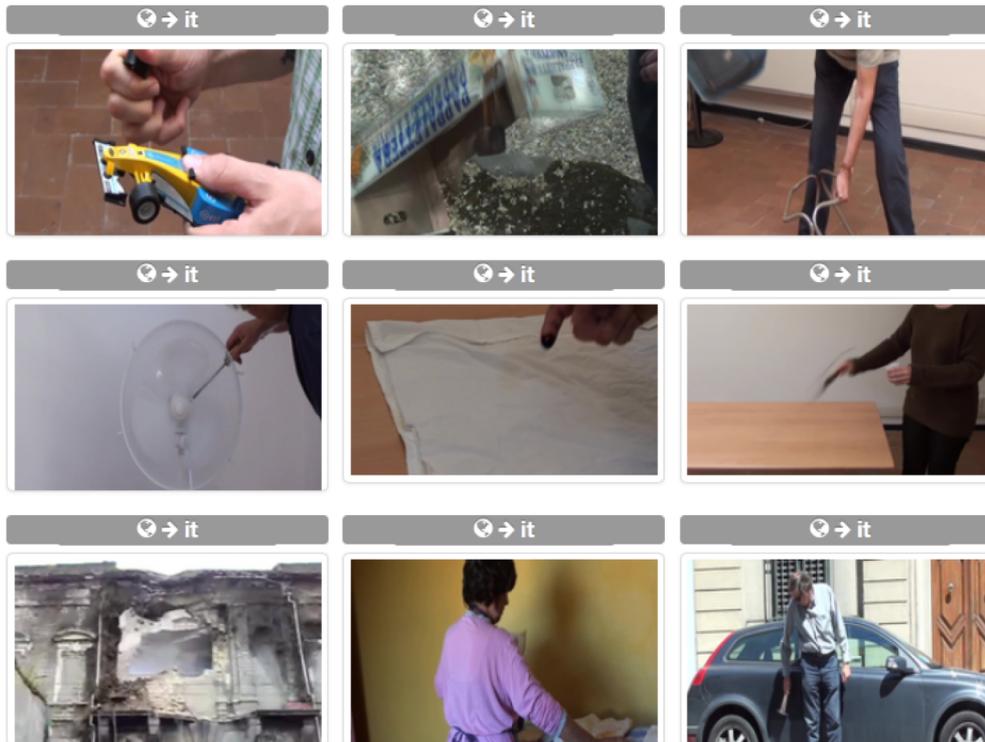


Galleria Multimediale

Output language

Italian

Deterioration of an object: 28





Strumento di comparazione

prendere

→ en



Fabio prende Marta per un braccio

→ en



Fabio prende il coniglio con la trappola

→ en



Marta prende la farfalla con il retino

prendere ⇔ **take**



Fabio prende la tazza dal ripiano



Fabio prende la pallina cadute a terra



Marta prende il libro dal tavolo

take

it ←



John takes the chair over to the table

it ←



The police take the prisoner with them

it ←



John takes the blind man across the street

Estensione di IMAGACT

Come lo dici
nella tua lingua?



Estensione di IMAGACT- Statistiche

Num Verbs (Total)	Italian 	English 	Chinese 	Spanish 	Portuguese 	Danish 
	522	554	414	736	793	646
>30 scenes	6 1,1%	7 1,3%	0 0,0%	8 1,1%	9 1,1%	1 0,2%
11-30	29 5,6%	30 5,4%	1 0,2%	43 5,8%	46 5,8%	21 3,3%
5-10	94 18,0%	98 17,7%	8 1,9%	80 10,9%	90 11,3%	57 8,8%
2-4	195 37,4%	222 40,1%	80 19,3%	243 33,0%	241 30,4%	179 27,7%
1 scene	198 37,9%	197 35,6%	325 78,5%	362 49,2%	407 51,3%	388 60,1%

Num Verbs (Total)	German 	Serbian 	Polish 	Hindi 	Arab 
	992	1124	1193	466	574
>30 scenes	0 0,0%	2 0,2%	0 0,0%	2 0,4%	4 0,7%
11-30	22 2,2%	31 2,8%	19 1,6%	17 3,6%	32 5,6%
5-10	61 6,1%	101 9,0%	119 10,0%	37 7,9%	90 15,7%
2-4	294 29,6%	363 32,3%	382 32,0%	128 27,5%	221 38,5%
1 scene	615 62,0%	627 55,8%	673 56,4%	282 60,5%	227 39,5%

Distribuzione spaziale delle azioni

Points: 999 | Dimension: 500 | Selected 101 points



Marta sbatte le uova

ID 5d19f1f9
Caption Marta sbatte le uova
URL <http://www.imagact.it/imagact/v?id=workspace://SpacesStore/538dc352-2086-4099-85cd-e0e7ca512349>

Show All Data Isolate 101 points Clear selection

Search by **action**

neighbors 100
distance COSINE EUCLIDEAN

Nearest points in the original space:

- Marta gira la zuppa velocemente 0.147
- Marta sbatte l'uovo 0.208
- Marta gira le zucchine con un cucchiaino 0.291
- Marta gira la zuppa 0.336
- Marta monta l'albume a neve 0.490
- Marta monta l'albume a neve 0.492
- Marta mischia le carte 0.640
- Marta mischia i cereali al cacao ed i cereal... 0.657
- Marta mischia uova e farina in un piatto 0.706
- Il fantino frusta il cavallo 0.809
- Marta mischia i due colori con il pennello 0.836
- Marta mischia il risotto 0.860
- Marta frulla la frutta 0.886
- I due gruppi si mischiano 0.892
- Fabio chiude forte la porta 0.919
- Marta batte forte il libro sul tavolo 0.926
- L'operaio del salumificio affumica lo spe... 0.947
- Fabio batte i piedi a terra 0.973
- L'educatrice picchia Fabio con la bacche... 0.974
- Marta batte l'ombrello contro il muro 0.975
- Marta colpisce Fabio con uno schiaffo 0.979
- Fabio tamburella le dita sul tavolo 0.980
- L'operaio attacca i due pezzi di metallo c... 0.980
- Fabio fotografa Marta e Marco 0.982