

OSINT

Ricerca e analisi da fonti aperte

Ricerca semantica e tecniche di confronto multilingue sui motori di ricerca



Dott. Ing. Manzoni Gabriel Esteban

Laurea Magistrale in Ingegneria Informatica - Politecnico di Milano
"Ambient Intelligence and Data Engineering"

Digital Solutions Architect presso SIGINT srl
Project Manager di DeePoint

CEH (Certified Ethical Hacker) Master - EcCouncil
CIFI (Certified Information Forensics Investigator) – IISFA
CPENT (Certified Penetration Testing Professional) – EcCouncil
LPT (Licensed Penetration Tester) Master - EcCouncil



DeePoint

Piattaforma di analisi **OSINT** (Open Source INTelligence) tramite applicazione di algoritmi di **Intelligenza Artificiale** legati al **NLP** (Natural Language Processing)

Principali funzioni:

- **Analisi** testi con estrazione di principali informazioni
- **Confronto** testi con comparazione semantica
- Analisi **Stilometrica**
- Ricerca **Social Media Network**
- Ricerca **Surface, Deep & Dark Web**
- E molto altro ...

INDICE

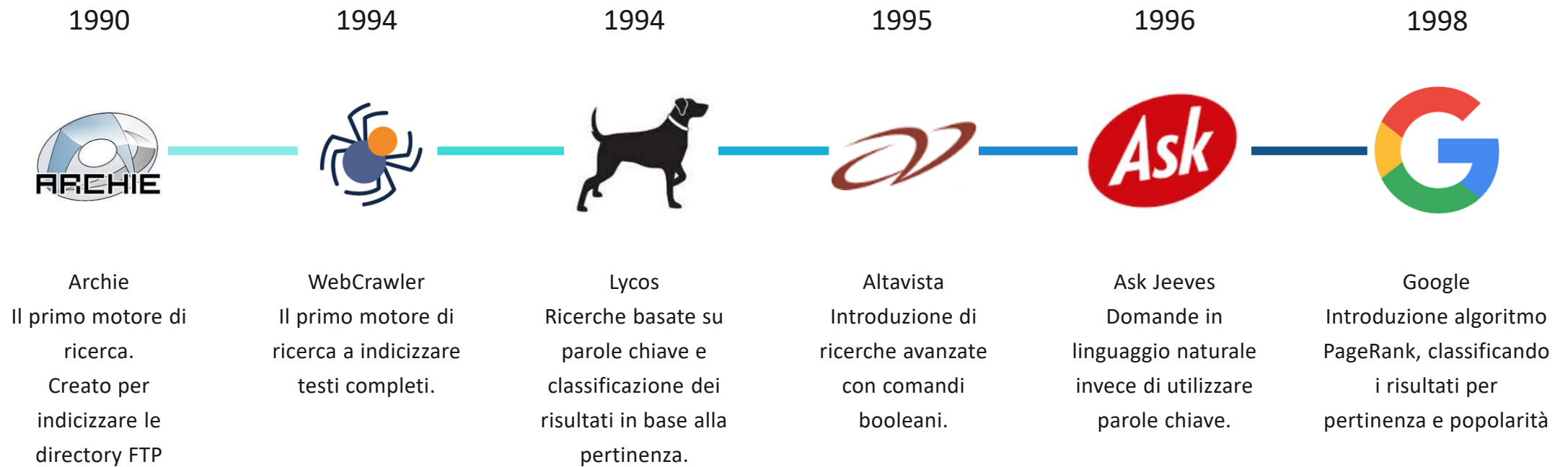
1. Introduzione ai motori di ricerca
2. Cos'è la ricerca semantica?
3. Modelli di Natural Language Processing
4. Tecniche e strumenti per il confronto multilingue
5. Sfide e soluzioni nell'integrazione della ricerca semantica
6. Case studies
7. Tendenze future e conclusioni



INTRODUZIONI AI MOTORI DI RICERCA



STORIA DEI MOTORI DI RICERCA E L'EVOLUZIONE VERSO LA RICERCA SEMANTICA



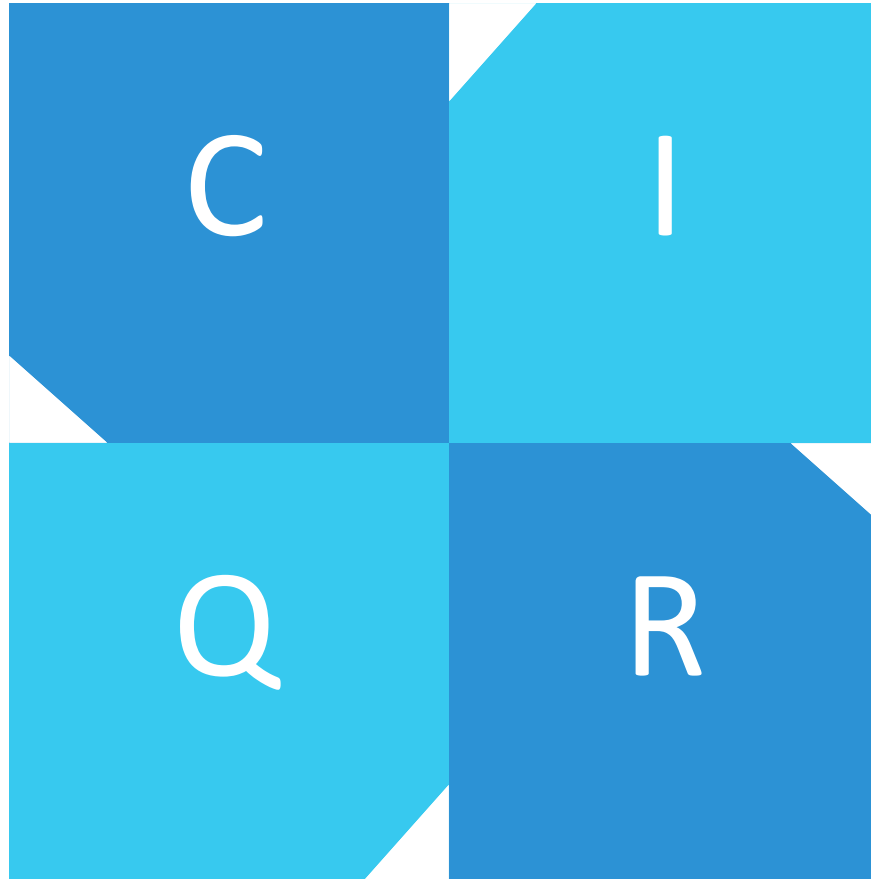
COME FUNZIONA UN MOTORE DI RICERCA?

CRAWLING

Spiders o robots esplorano il web per trovare nuove pagine.

QUERY

L'utente inserisce una query e il motore di ricerca recupera le informazioni più pertinenti dal suo indice.



INDEXING

Le pagine trovate vengono analizzate e archiviate in un database.

RANKING

Algoritmi, come l'algoritmo PageRank di Google, determinano la rilevanza delle pagine indicizzate in base a specifiche query.

LA NECESSITÀ DELLA RICERCA SEMANTICA

Mancanza di precisione: le parole chiavi non sono sufficienti

Mancanza di contesto: necessario identificare il significato della query

Molteplicità di linguaggi: necessario ricevere come input query in lingue diverse

Ambiguità linguistiche: una stessa parola può avere più significati in lingue diverse

Interrogazioni complesse: le ricerche richieste sono sempre più complesse

COS'È LA RICERCA SEMANTICA?

«La **semantica** è quella parte della linguistica che studia il significato delle parole (*semantica lessicale*), degli insiemi delle singole lettere (negli e degli alfabeti antichi) e delle frasi (**semantica frasale**) e dei testi.»

[Wikipedia - Semantica](#)

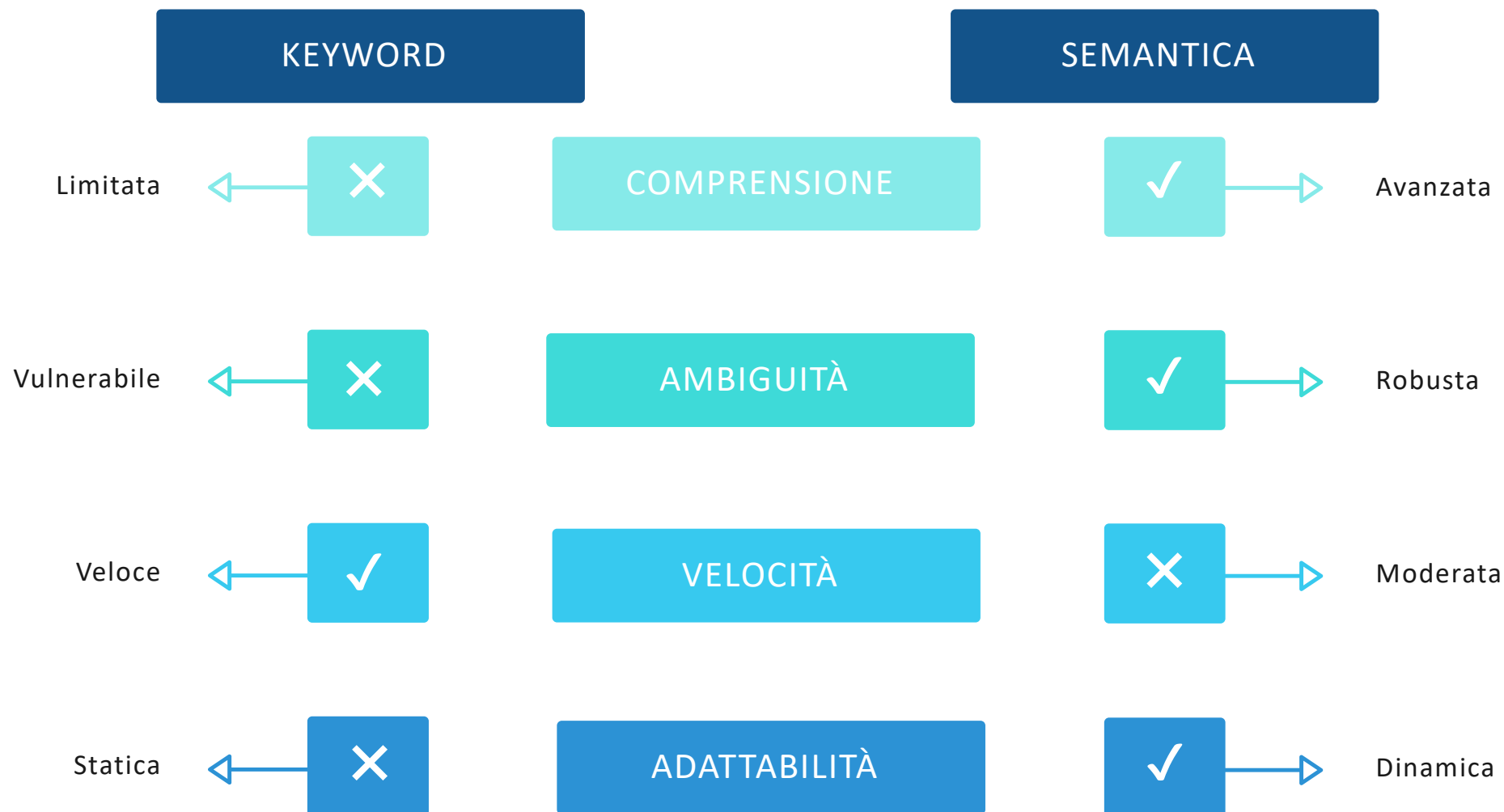
Gli algoritmi di Intelligenza Artificiale utilizzano la semantica computazionale per potenziare strumenti e sistemi capaci di comprendere, interpretare e reagire al linguaggio umano in contesti reali. Questa integrazione permette alle macchine di svolgere compiti come la ricerca semantica, la traduzione automatica, l'analisi dei sentimenti e la risposta alle domande, mimando la capacità umana di comprendere il significato e il contesto delle parole e delle frasi.

«La **semantica computazionale** è lo studio di come automatizzare il processo di costruzione e ragionamento con l'ausilio di rappresentazioni del significato di espressioni di una lingua naturale.»

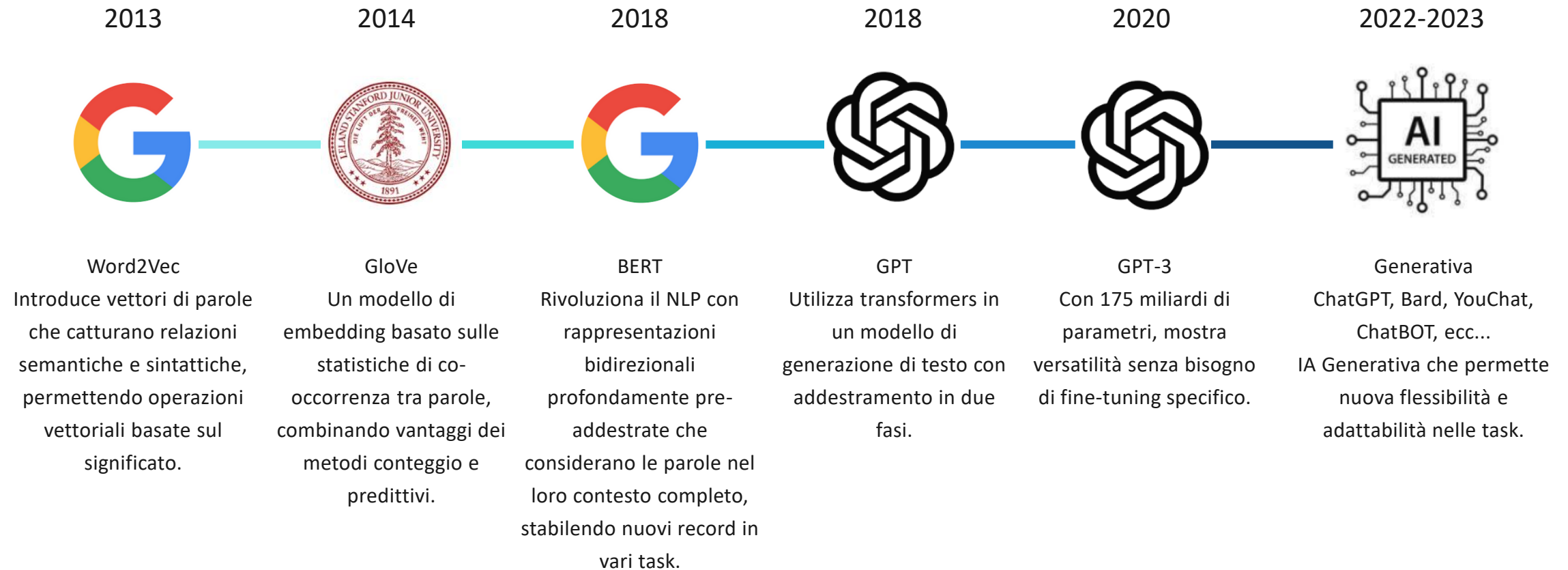
[Wikipedia – Semantica computazionale](#)



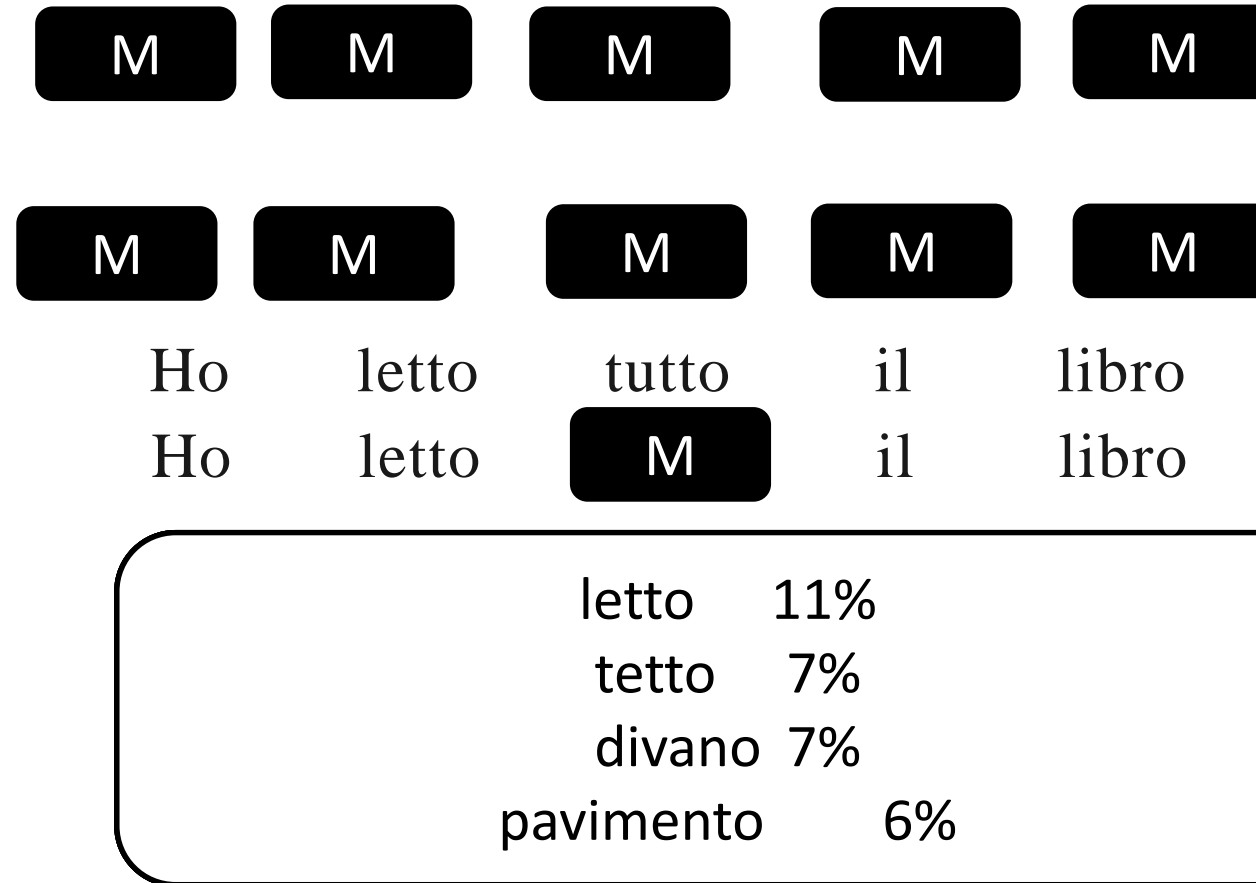
DIFFERENZE RICERCA KEYWORD E SEMANTICA



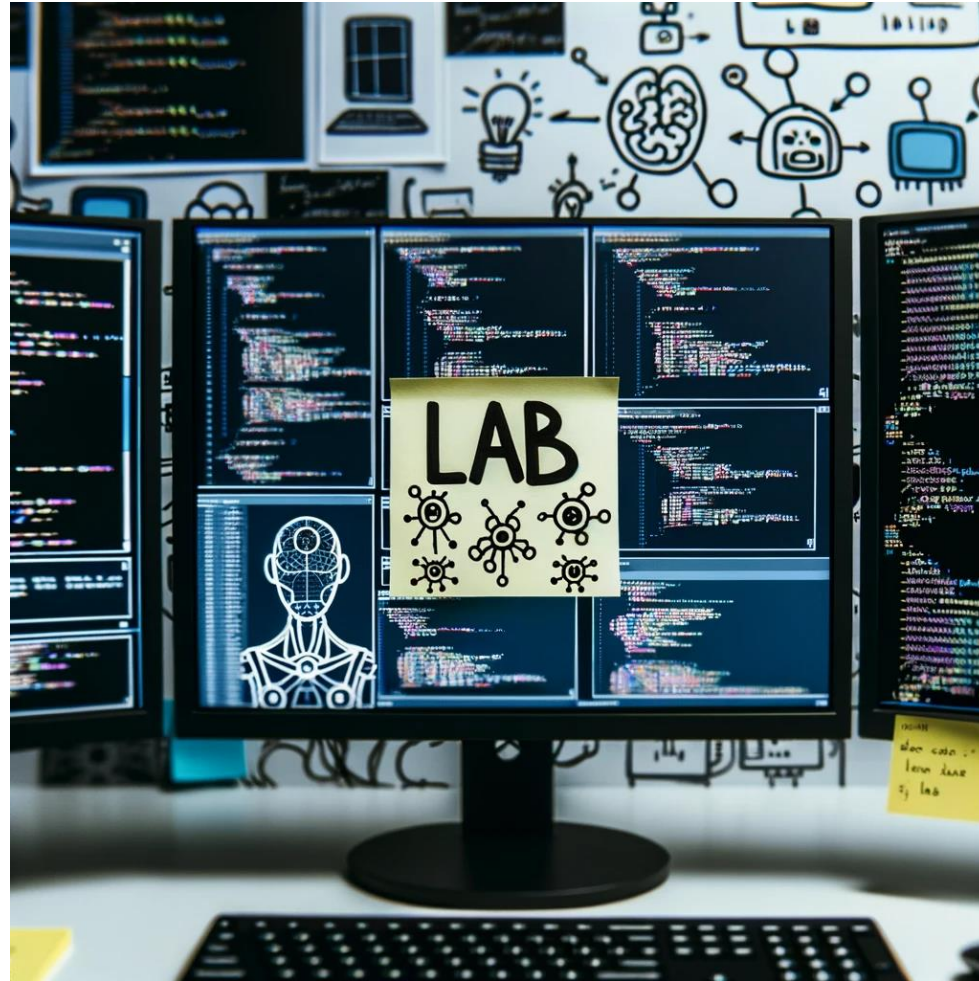
EVOLUZIONE DELLA RAPPRESENTAZIONE SEMANTICA



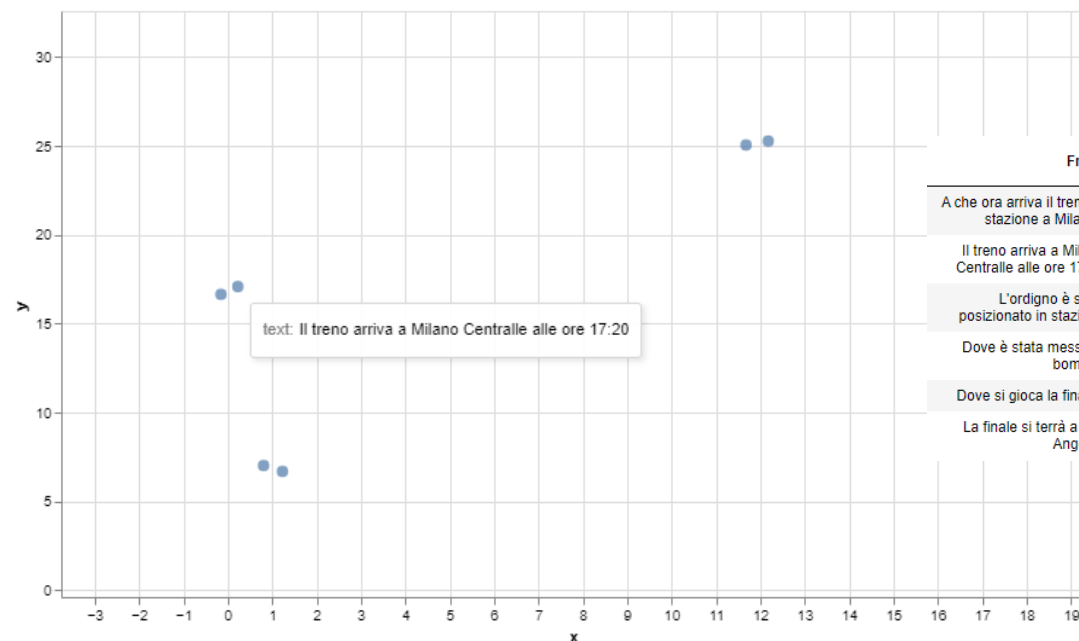
ADDESTRAMENTO BERT



FILL_MASK_EXAMPLE



EMBEDDING PER LA RICERCA SEMANTICA

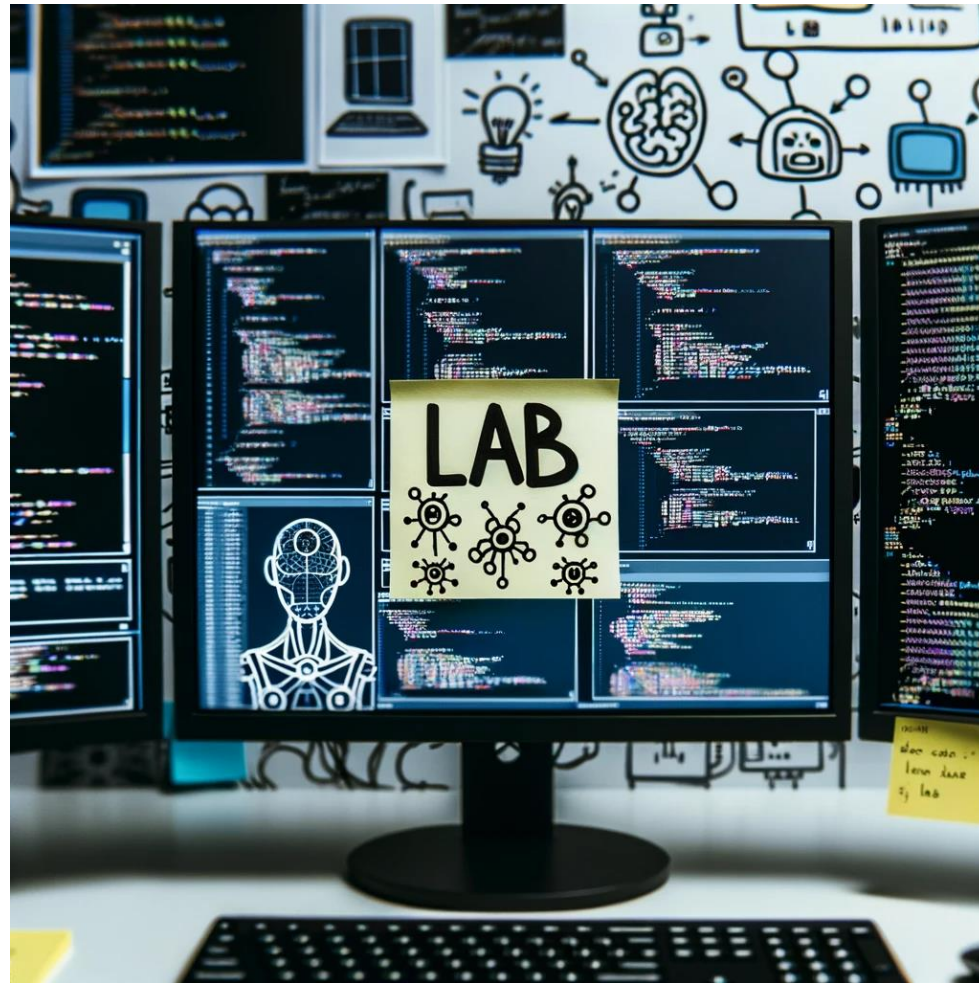


Frase	A che ora arriva il treno in stazione a Milano?	Il treno arriva a Milano Centrale alle ore 17:20	L'ordigno è stato posizionato in stazione	Dove è stata messa la bomba?	Dove si gioca la finale?	La finale si terrà a Los Angeles
A che ora arriva il treno in stazione a Milano?	1.00	0.62	0.40	0.29	0.27	0.22
Il treno arriva a Milano Centrale alle ore 17:20	0.62	1.00	0.29	0.14	0.19	0.27
L'ordigno è stato posizionato in stazione	0.40	0.29	1.00	0.57	0.11	0.15
Dove è stata messa la bomba?	0.29	0.14	0.57	1.00	0.32	0.13
Dove si gioca la finale?	0.27	0.19	0.11	0.32	1.00	0.58
La finale si terrà a Los Angeles	0.22	0.27	0.15	0.13	0.58	1.00

embedding.

o ranking per similitudine

SEMANTIC_COMPARISON



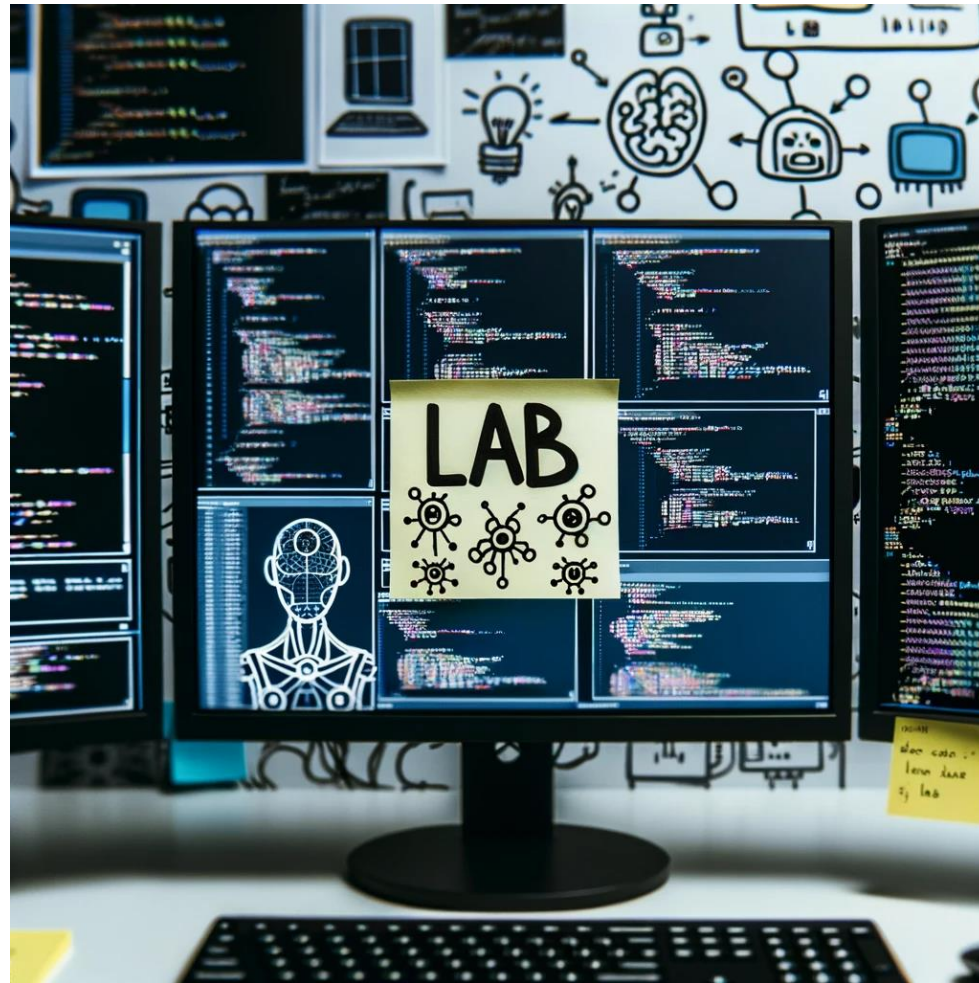
HUB DI MODELLI DI NLP

- | | |
|-----------------------|--|
| 1. HuggingFace | <u>https://huggingface.co/</u> |
| 2. TensorFlow Hub | <u>https://www.tensorflow.org/hub</u> |
| 3. Spacy | <u>https://spacy.io/</u> |
| 4. Stanford NLP Group | <u>https://nlp.stanford.edu/</u> |
| 5. Facebook's Fairseq | <u>https://ai.meta.com/tools/fairseq/</u> |
| 6. Google Research | <u>https://research.google/</u> |
| 7. OpenAI | <u>https://openai.com/</u> |

DATI ESTRATTI IN NLP

1. Named Entity Recognition (NER): entità suddivise in categorie predefinite come persone, luoghi , organizzazione, ecc...
2. Keyword: parole chiave le quali possono essere parole o frasi più rilevanti
3. Part-Of-Speech (POS): classifica delle parole in categorie grammaticali
4. Lemmatizzazione: estrazione della parola nella forma base (bellissima->bello)
5. Estrazione di relazioni: relazioni quali per esempio soggetto->azione->oggetto
6. Coerenza: estrazione di frasi che riprendono lo stesso concetto o l'opposto

DATA_EXTRACTION



Pausa di 10 minuti

La lezione riprenderà alle 17:30



TIPI DI MODELLI DI NLP

1. Text Classification:
2. Question Answering:
3. Translation: traduttore
4. Summarization: riassunti
5. Text Generation: generazione di testo
6. Sentence Similarity: similarità
7. Conversational: chatbot
8. Fill-Mask: indovina la parola mancante

Funzione base di tutte le altre, utile direttamente per correzioni grammaticali, suggerimenti, identificazione di messaggi in codice

Contesto

Embedding

MODALITÀ DI APPRENDIMENTO DEI MODELLI DI NLP

1. Supervised Learning: vengono forniti un numero elevato di esempi «etichettati» per addestrare la rete neurale da zero
2. Fine-tuning Learning: vengono forniti un ristretto numero di esempi «etichettati» per affinare una rete neurale pre-addestrata
3. Unsupervised Learning: viene fornito testo «non etichettato» con l'obiettivo che l'algoritmo apprenda automaticamente strutture o modelli senza supervisione
4. Semi-Supervised Learning: si adotta un approccio ibrido tra Supervised e Unsupervised
5. Reinforcement learning: dopo l'addestramento di una rete neurale vengono proposti ulteriori esempi con risposta al fine di aumentarne l'accuratezza
6. Transfer Learning: le capacità apprese da una rete neurale vengono trasferite verso un'altra rete neurale per facilitarne l'apprendimento
7. Zero-Shot Learning: viene proposta una nuova task su cui la rete neurale non è stata addestrata
8. One or Few Shot Learning: vengono proposti solo uno o pochi esempi da cui la rete neurale deve imparare

COME FUNZIONA UN MOTORE DI RICERCA?

CRAWLING

Spiders o robots esplorano il web per trovare nuove pagine.

QUERY

L'utente inserisce una query e il motore di ricerca recupera le informazioni più pertinenti dal suo indice.

Query per estrarre il testo più simile
o estrarre la risposta direttamente
dal testo

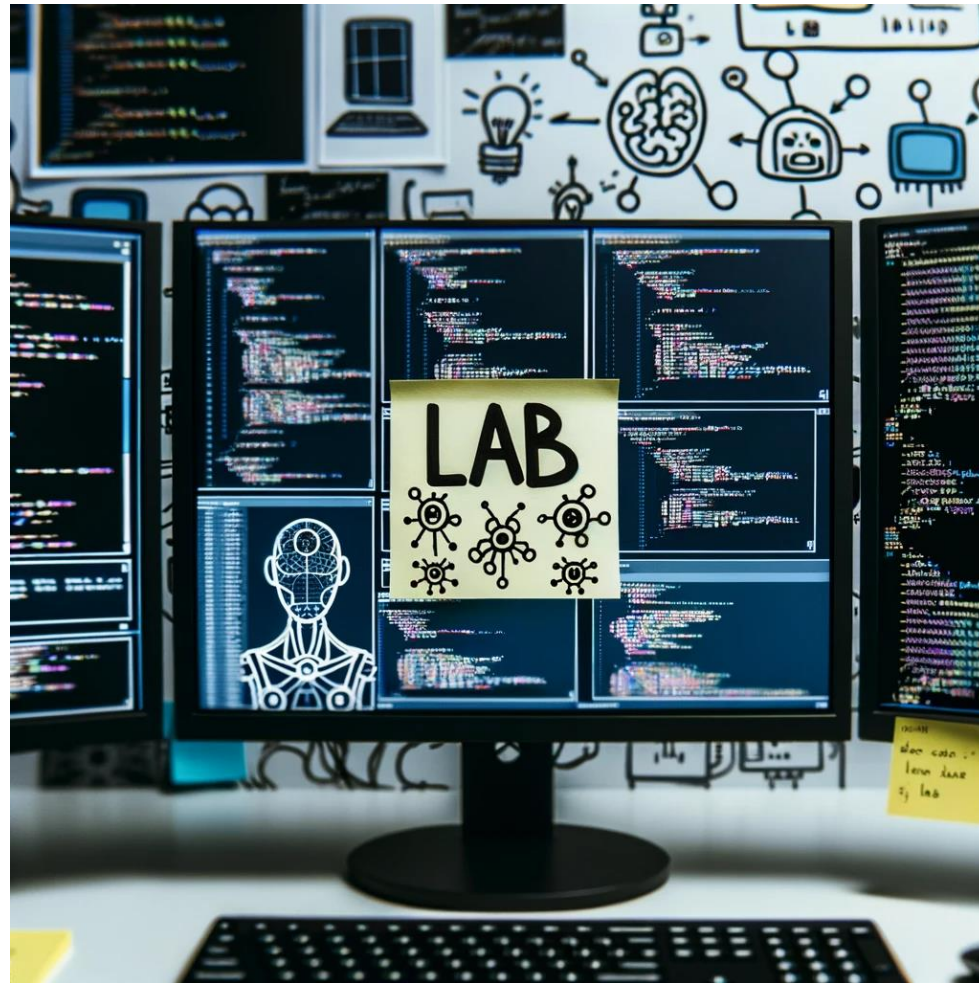
INDEXING

Le pagine trovate vengono analizzate e archiviate in un database.

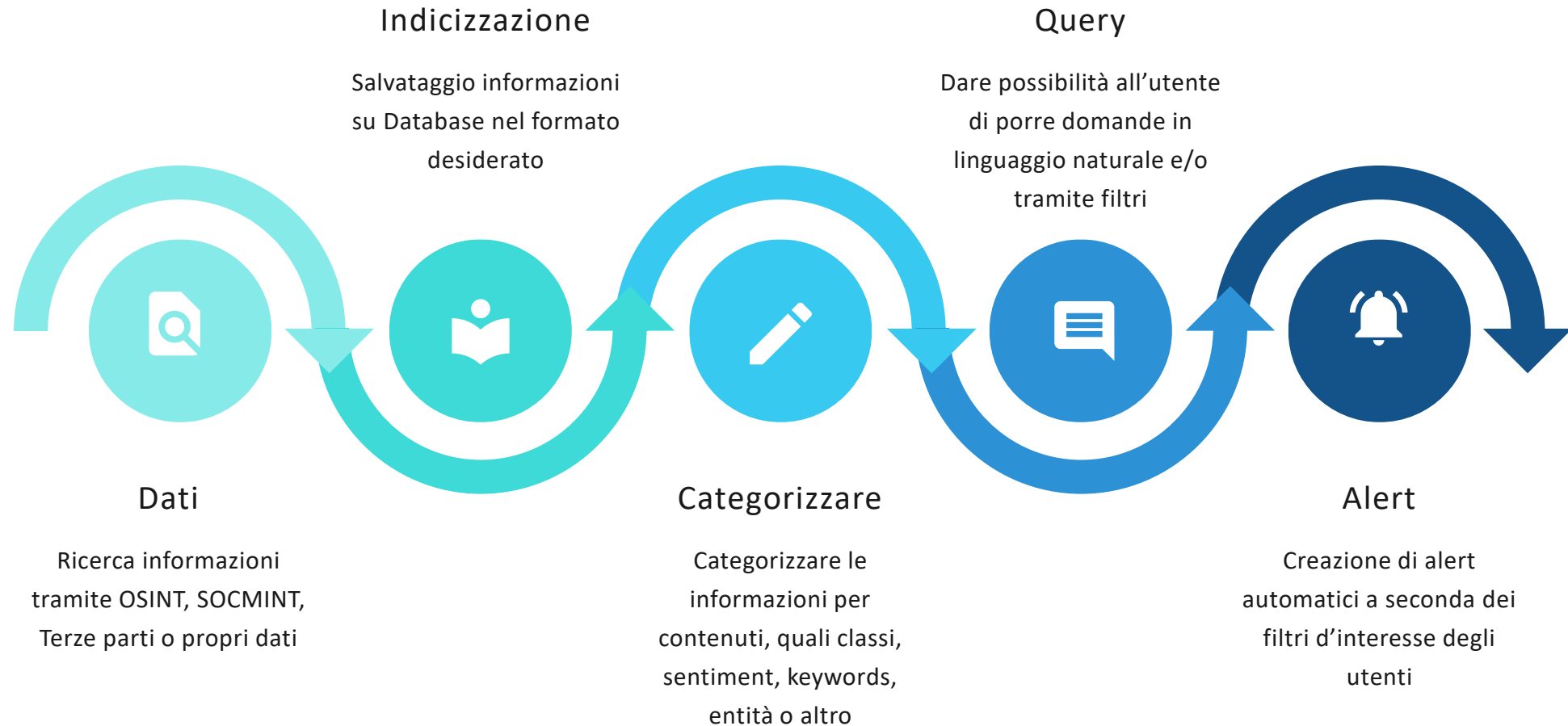
RANKING

Algoritmi, come l'algoritmo PageRank di Google, determinano la rilevanza delle pagine indicizzate in base a specifiche query.

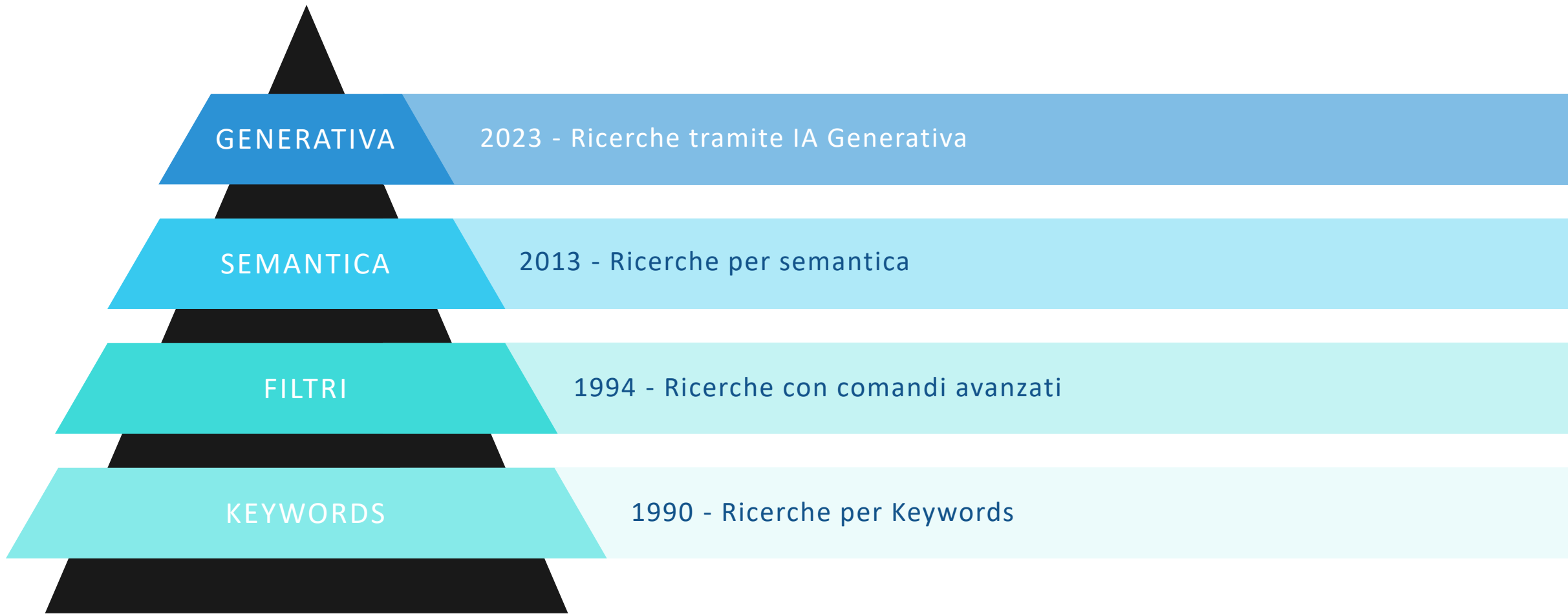
EXAMPLE_SEARCH_ENGINE



COSTRUIRE MOTORE DI RICERCA ON-PREMISE



MODELLI DI RICERCA



DIFFERENZE RICERCA SEMANTICA E IA GENERATIVA



PRODOTTI DI IA GENERATIVA

1. OpenAI
2. Google
3. Microsoft
4. HuggingFace
5. GitHub

Dedicato a scrivere codice
Utilizza GPT-4

openai.com/

ai.google.com/

www.bing.com/

huggingface.co/chat/

github.com/features/copilot

MODELLI DI IA GENERATIVA

1. GPT-3.5 & GPT-4	OpenAI (ChatGPT)	Proprietario
2. Palm	Google (BARD)	Open-Source(*)
3. Llama	Meta	Open-Source
4. Falcon	HuggingFace	Open-Source
5. StableLM	Stability AI	Open-Source
6. Alpaca	Stanford	Open-Source
7. MPT	MosaicML	Open-Source

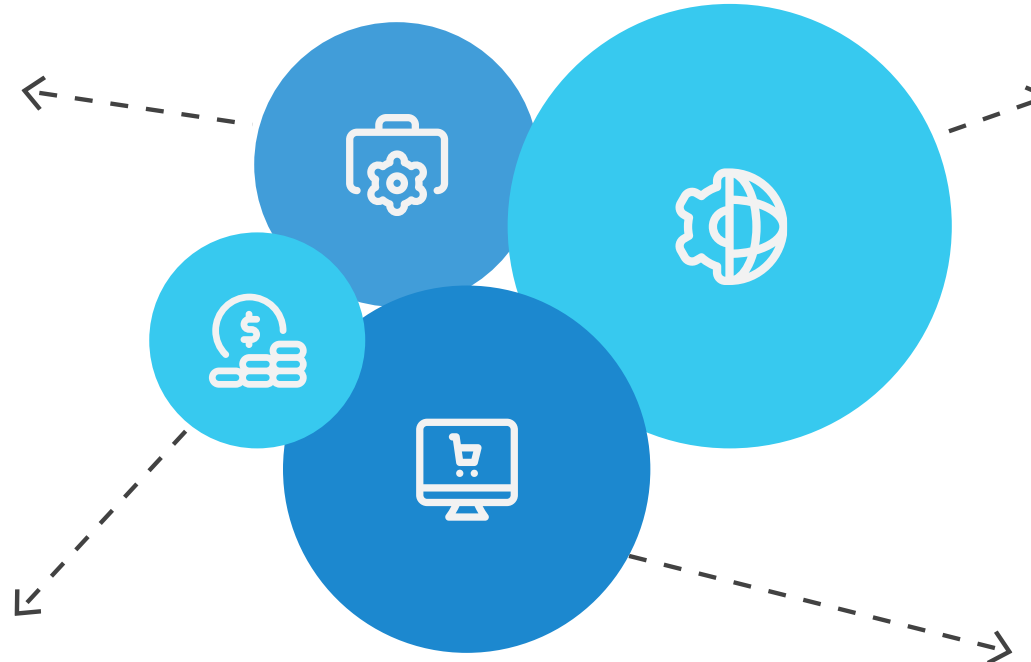
SFIDE FUTURE

Ethic AI

Maggiore attenzione su questioni etiche e morali dell'AI tramite applicazione e rispetto dell'EU AI Act

Scalabilità

Necessario cambiare il trend in corso, ovvero incrementare dimensione del dataset necessario e della potenza di calcolo quanto piuttosto migliorare l'apprendimento su dataset più piccoli e ridurre la potenza di calcolo necessaria



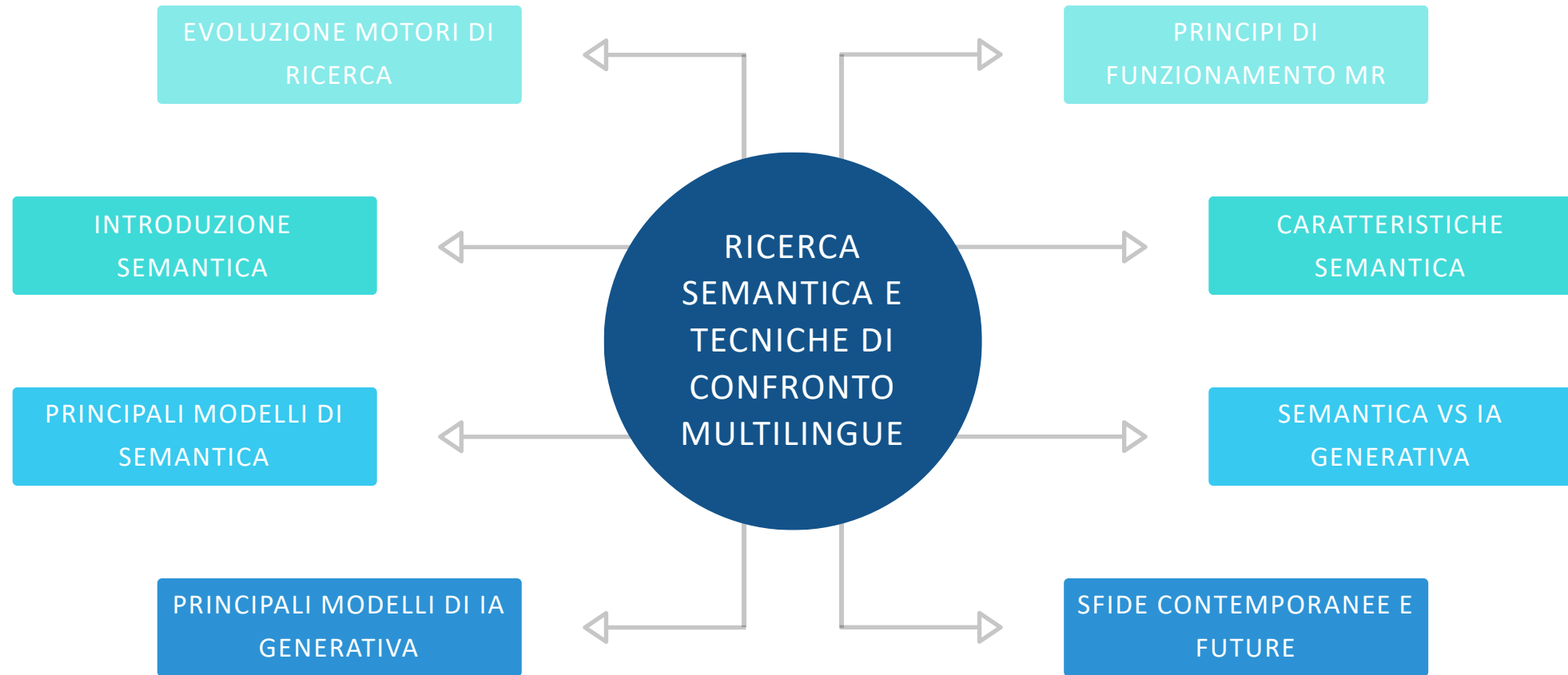
Riduzione Allucinazioni

Le AI Generative sono ancora vulnerabili alle allucinazioni. La priorità per il futuro è arrivare a ridurle drasticamente

Slow-Thinking

Le attuali AI Generative sono forti per la parte di Fast-Thinking ma deboli in quella di Slow-Thinking. Per ottenere maggiore qualità in task complesse occorre cambiare questo aspetto

RIEPILOGO







Gabriel Esteban Manzoni



g.manzoni@sigint.srl.it



https://github.com/gmanzoni/POL_OSINT_Semantica

