



Renzo Alva Principe

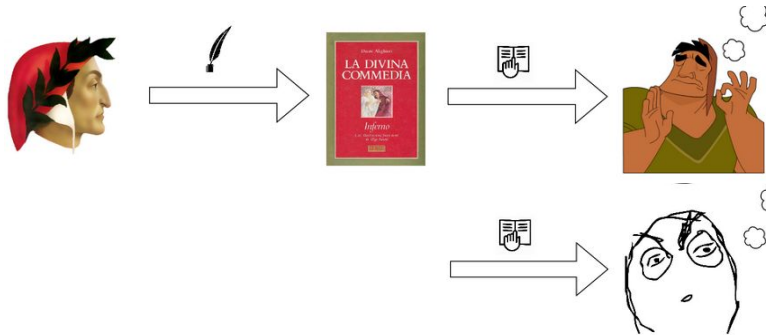
# Outline

- 1) Task
- 2) Distributional Hypothesis
- 3) Deep Learning
- 4) Pre-Trained Language Models
- 5) NLP Tasks
  - a) Sentiment Analysis
  - b) Semantic Text Similarity
- 6) Proposal

# The Task

## Human Communication:

- is mainly non verbal
- ... people have knowledge
- ... people have biases
- ... people have emotions
- ... people have a mood



## Task:

Demonstrate numerically that communication is inherently full of distortions and that information cannot be transmitted perfectly to others

# The Distributional Hypothesis

# What Drives Semantic Similarity?<sup>[1]</sup>

- **Accidental**

- Abominate
- Meander
- Inadvertent
- inhibit

- **FedEx**

- car
- UPS
- rotate
- Navy

- **Millennial**

- octopus
- fork
- water
- avocado

# What Drives Semantic Similarity?<sup>[1]</sup>

- **Accidental**

- Abominate
- Meander
- Inadvertent
- inhibit

- **FedEx**

- car
- UPS
- rotate
- Navy

- **Millennial**

- octopus
- fork
- water
- avocado

- **Meaning:** closeness in terms of meaning
- **Word Knowledge:** concepts have similar properties, often occur together, or occur in similar context
- **Psychology:** two concepts fit together within an over-arching psychological schema or framework

[\[1\] Sean Simpson, EMNLP Lecture 21, 2018](#)

# Bardiwac

Does anyone knows what it is?

Here some clues...

- 1) he handed her a glass of bardiwac
- 2) beef dishes are made to complement the bardiwac
- 3) Malbec, one of the lesser-known bardiwac grapes, responds well to Australia's sunshine
- 4) I dined off bread and cheese and this excellent bardiwac





Correct answer!

Bardiwac is a wine!



Just kidding ....



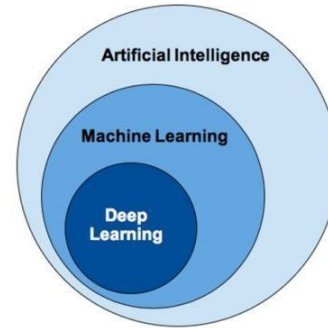
**Bardiwac is a fake word**

# The Distributional Hypothesis

*“Differences of meaning correlates with differences of distribution”*  
(Harris, 1970)

*“You shall know a word by the company it keeps!”*  
(Firth, 1957)

# Deep Learning



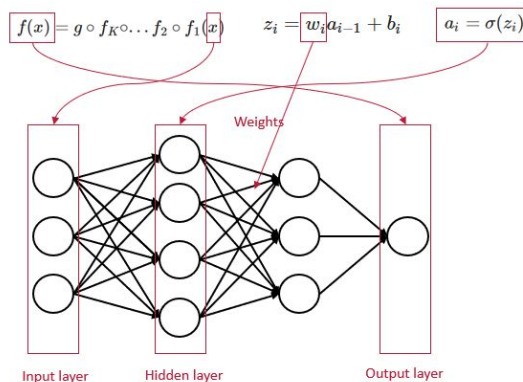
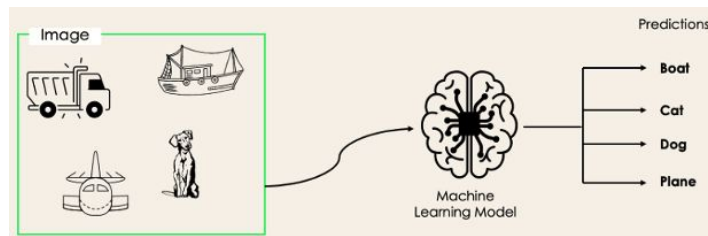
# Deep Neural Networks

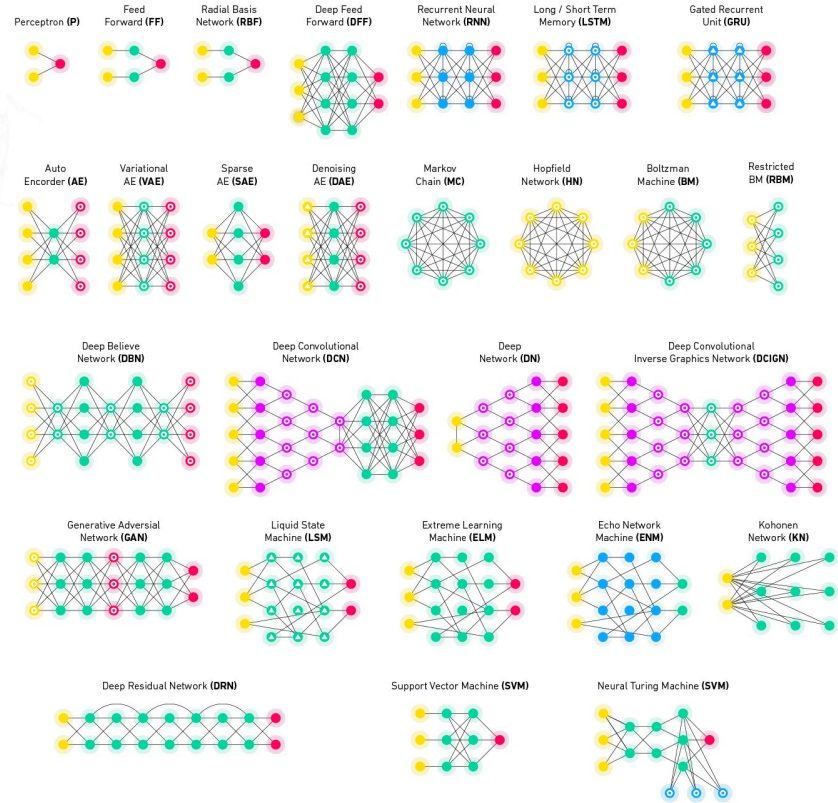
Machine Learning models are **function approximation algorithms**

Deep Learning:

- based on neural networks
- capable of learn **extremely complex functions**
- **automatic feature extraction**
- effort moved from feature engineering to architecture engineering
- requires large amounts of data

<https://playground.tensorflow.org>



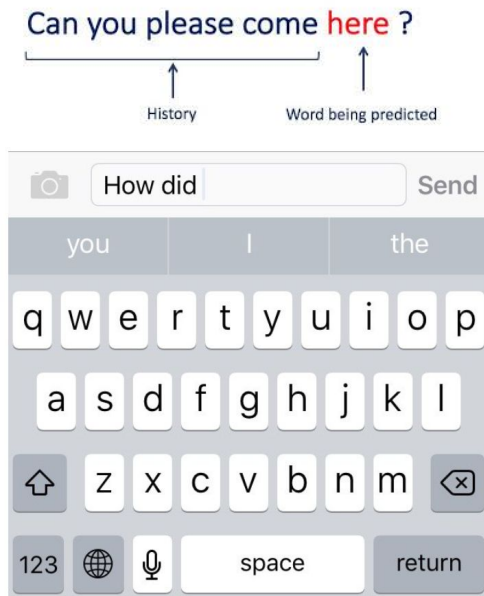


Tons of architectures!

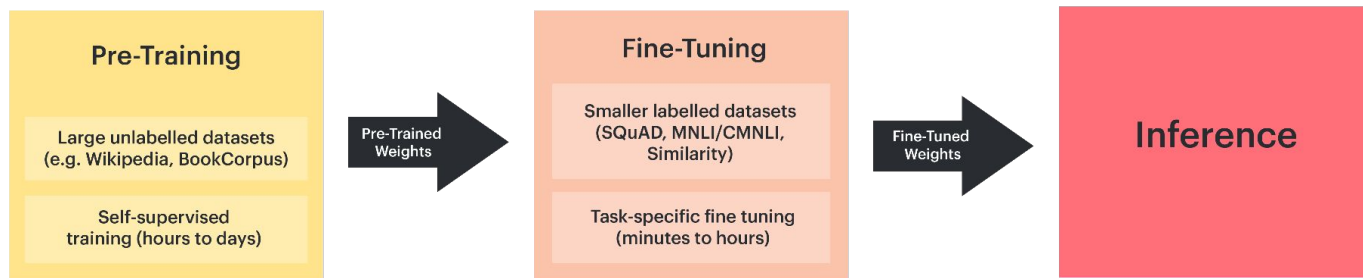
# Pre-Trained Language Models

# Language Models

- is a probabilistic statistical model that determines the probability of a given sequence of words occurring in a sentence based on the previous words.
- language modelling is usually a non-supervised process
- Language Models are the backbone of NLP:  
→ they are a way of transforming qualitative information about text into quantitative information that machines can understand



# Pre-Training and Fine-Tuning



Pre-Training:

- allows the model to learn general knowledge
- sometimes specific tasks are chosen to this aim
- extremely expensive
- requires large amount of data

Fine-Tuning:

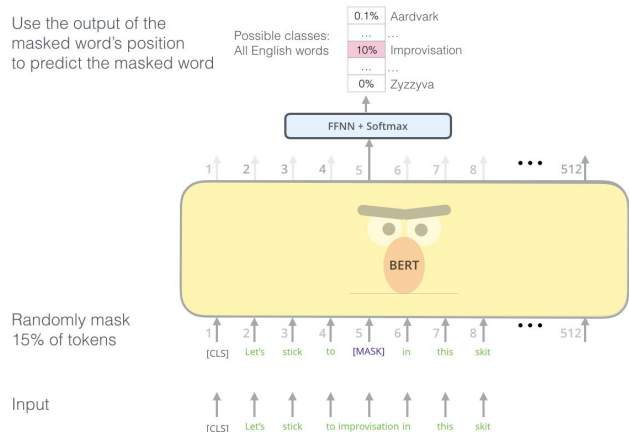
- task specialization
- domain adaptation
- cheap and requires much less data

This paradigm is not limited to NLP, it is also well-known in Computer Vision (see [ImageNet challenge](#))

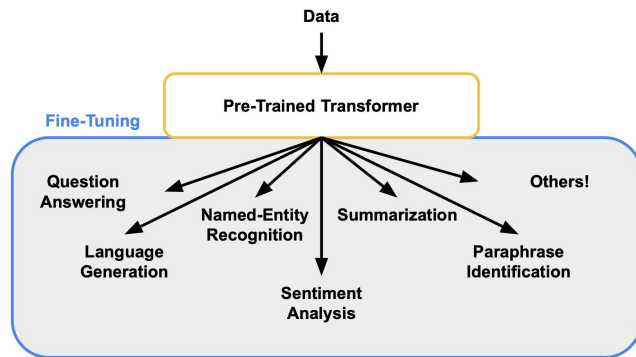


# BERT<sub>[2]</sub>: A Pre-Trained Language Model

Use the output of the masked word's position to predict the masked word



... it is a function approximation!



[\[2\] Devlin et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. 2018, Arxiv](#)

# Pre-Trained Language Models

=

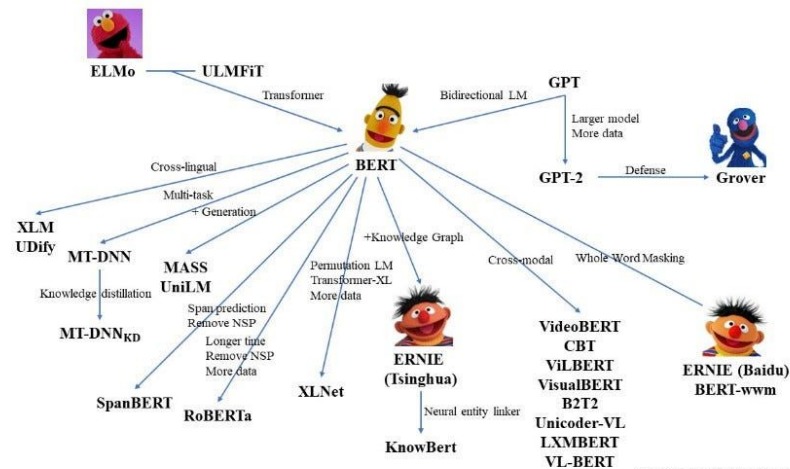
Distributional Hypothesis

+

Deep Learning

+

Pre-Training



By Xiaozhi Wang & Zhengyan Zhang @THUNLP

# NLP Tasks



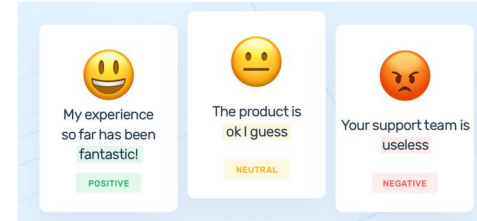
# Sentiment Analysis

Is used to determine whether data is positive, negative or neutral.

Approaches:

- **lexicon-based:** counts the number of positive, neutral, and negative words and assigns a score based on intensity and frequency
- **rule-based:** depends on handcrafted rules and lexicons that might not apply for all texts
- \* **machine-learning based:** patterns identification on hand-crafted features
- \* **deep-learning based:** uses DL to automatically extract features
- \* **PTLMs based:** like DL but exploits the power of transfer learning

\*supervised approaches: need training data



Challenges:

- tone:
- sarcasm
- negations
- emojis

# Semantic Text Similarity

Text similarity calculates how two documents are close to each other. Closeness may be **lexical** or in **meaning**

- The dog bites the man
- The man bites the dog



lexical similarity very high... almost identical  
semantic similarity is very low ... totally different

Approaches:

- **Knowledge base methods:** leverages lexical structured representation of concepts connected by semantic relations, further offering an ambiguity free semantic measure
- **Corpus based:** use the distributional hypothesis to get rid of ambiguity
- **DL methods and PT LMs:** exploit transfer learning power along with DL architectures

Usages:

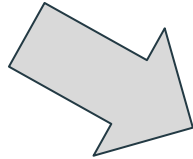
- Information Retrieval
- Document Similarity and Clustering

# My Proposal

# Sentiment Analysis and Semantics Analysis

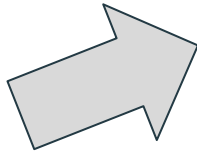
Problem characteristics:

- italian language
- no labelled data
- extremely small data points
- words definitions as input



proposal:

- sentiment analysis
- semantic analysis



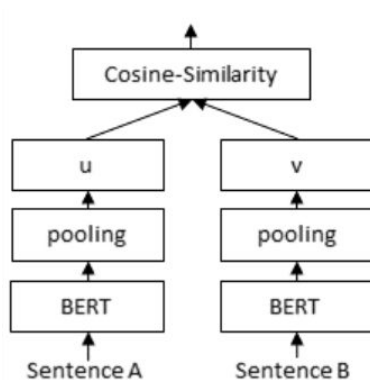
- already trained models for SA and STS
- Italian words definitions as a compass



# Semantic Text Similarity: SBERT<sup>[3]</sup>

Sentence BERT uses Bi-Encoders.

They pass to a BERT **independently** the sentences A and B, which result in the sentence embeddings  $u$  and  $v$ . These sentence embedding can then be compared using cosine similarity.



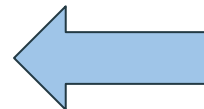
Training data:

multilingual translation of the [STS benchmark dataset](#)

score	2 example sentences	explanation
5	<i>The bird is bathing in the sink.</i> <i>Birdie is washing itself in the water basin.</i>	The two sentences are completely equivalent, as they mean the same thing.
4	<i>Two boys on a couch are playing video games.</i> <i>Two boys are playing a video game.</i>	The two sentences are mostly equivalent, but some unimportant details differ.
3	<i>John said he is considered a witness but not a suspect.</i> <i>"He is not a suspect anymore." John said.</i>	The two sentences are roughly equivalent, but some important information differs/missing.
2	<i>They flew out of the nest in groups.</i> <i>They flew into the nest together.</i>	The two sentences are not equivalent, but share some details.
1	<i>The woman is playing the violin.</i> <i>The young lady enjoys listening to the guitar.</i>	The two sentences are not equivalent, but are on the same topic.
0	<i>The black dog is running through the snow.</i> <i>A race car driver is driving his car through the mud.</i>	The two sentences are completely dissimilar.

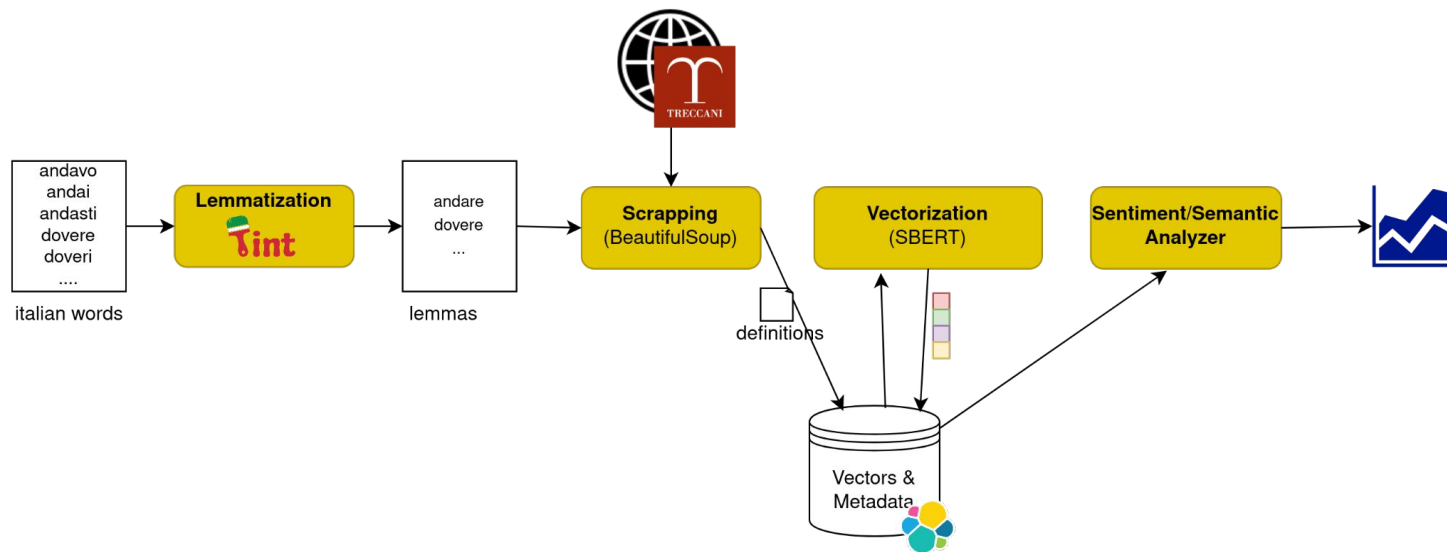
# Sentiment Analysis

- [neuralt/bert-base-italian-cased-sentiment:](#)
  - base model: BERT based italian cased
  - fine-tuning data: [Sentipolc EVALITA 2016](#). (45K pre-processed tweets)
- [citizenlab/twitter-xlm-roberta-base-sentiment-finetuned:](#)
  - base model: Multilingual XLM-Roberta
  - fine-tuning data: [toxicity prediction dataset](#)
- [cardiffnlp/twitter-xlm-roberta-base-sentiment:](#)
  - base model: multilingual XLM-roBERTa-base
  - fine-tuning data: 198M multilingual tweets
- [lxyuan/distilbert-base-multilingual-cased-sentiments-student:](#)
  - Teacher model: [MoritzLaurer/mDeBERTa-v3-base-mnli-xnli](#)
  - Teacher hypothesis template: "The sentiment of this text is {}."
  - Student model: distilbert-base-multilingual-cased



the chosen one

# Logical Architecture



Demo: <http://2.44.8.105:8501/>

# Conclusions

- Sentiment Analysis
  - differences are evident but motivations are not always clear
- Semantic Analysis
  - SBERT is a good model to be used for Semantic Text Similarity
  - differences and similarities between participants are evident
  - .. however there are some cases where is hard to explain similarity
  - coherence and incoherence is visible

Understanding and communication phenomena result in different conceptualisations

This is a nice way to show differences about people understanding... but IMHO going further requires a psychological analysis

# Limitations & Future Work

## Human limitations:

- **exhaustivity**: is what has been written really what the participants know?
- **expressivity**: did the participants express themselves well?
- **commitment**: did the participants make a serious effort?

## Models limitations:

- Sentiment analysis domain differences: tweets vs definitions
- Semantics analysis domain differences: Treccani definitions vs custom definitions
- a SBERT model trained on definitions would work better

## Future Work:

- Exploit knowledge from LLMs to get embeddings
- use LLMs to rephrase Treccani definitions and fine-tune a SBERT model
- Try to use prompting to get sentiment scores

# Links

- <https://www.pinecone.io/learn/series/nlp/sentence-embeddings/>
- <https://www.sbert.net/examples/applications/cross-encoder/README.html>
- <https://www.searchcandy.uk/nlp/sentence-bert/>