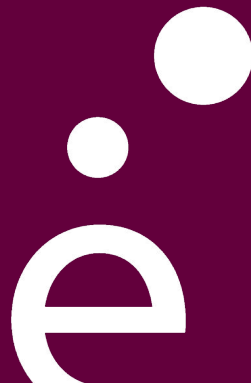


Information Retrieval Project

C12N - Microorganisms, enzymes and more...

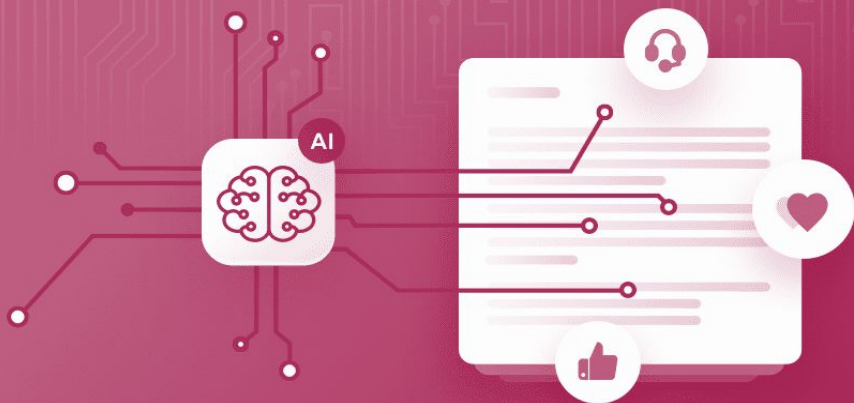
Jacobo Ruiz Ocampo

Diego Andres Torres Guarin



TERM DETECTION

- Domain specific NER
TT
- Pipeline
- Annotations
- Example and limits



Domain specific NER TT



- Domain:

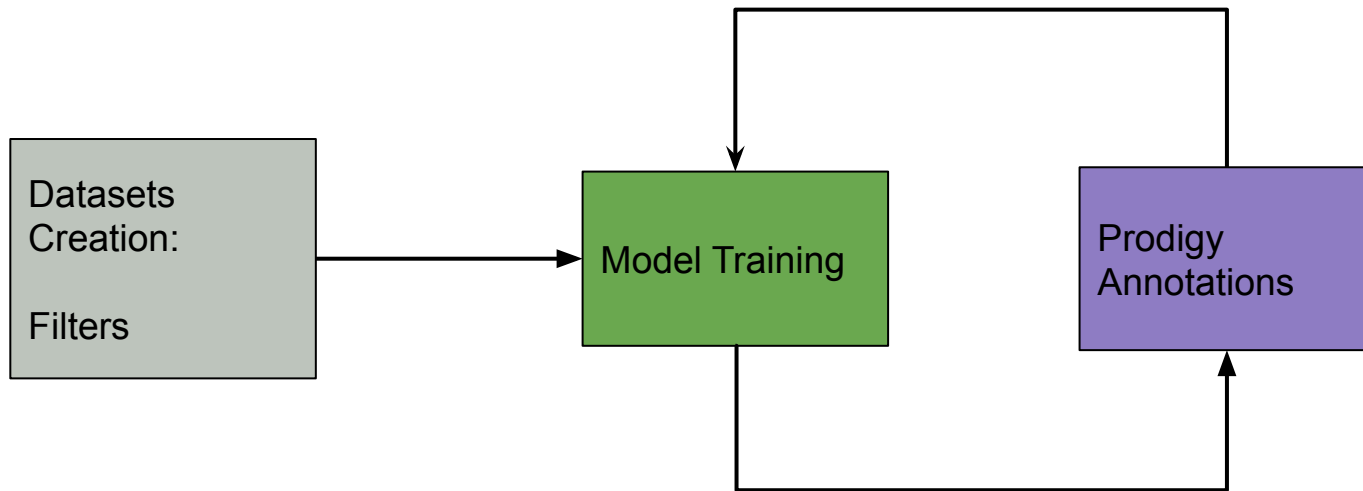
MICROORGANISMS OR ENZYMES; COMPOSITIONS THEREOF; PROPAGATING, PRESERVING, OR MAINTAINING MICROORGANISMS; MUTATION OR GENETIC ENGINEERING; CULTURE MEDIA

- NER TT:

```
1 # Load best model
2 nlp_ner = spacy.load("./results/spacy_output/model-best")
3
4 # Just text snippet
5 doc = nlp_ner(sample)
6
7 # Show NER results
8 spacy.displacy.render(doc, style="ent", jupyter=True)
```

Table TT 5. The retinal cell markers TT and dilutions TT used in the studyCell TypeCell MarkerDilutionsMüller cellAnti-CRALBP1:1000Anti-Vimentin1:100Anti-GFAP1:1000PhotoreceptorAnti-Op sin Red/Green1:250Anti-Op sin Blue1:250Neuron TT in INLAnti-PKCa1:200AstrocytesAnti-GFAP1:1000 TT

Pipeline



Annotating with prodigy



The 40 g/L glucose bottles were sparged with **N2/CO2** **TT** prior to incubation, whereas the 110 g/L bottles were not.

SCORE: 1.00

accept (a)



Correct annotations: 488
Incorrect annotations: 376
Correct percentage: 56.48%
Number of sentences annotated: 125

Examples

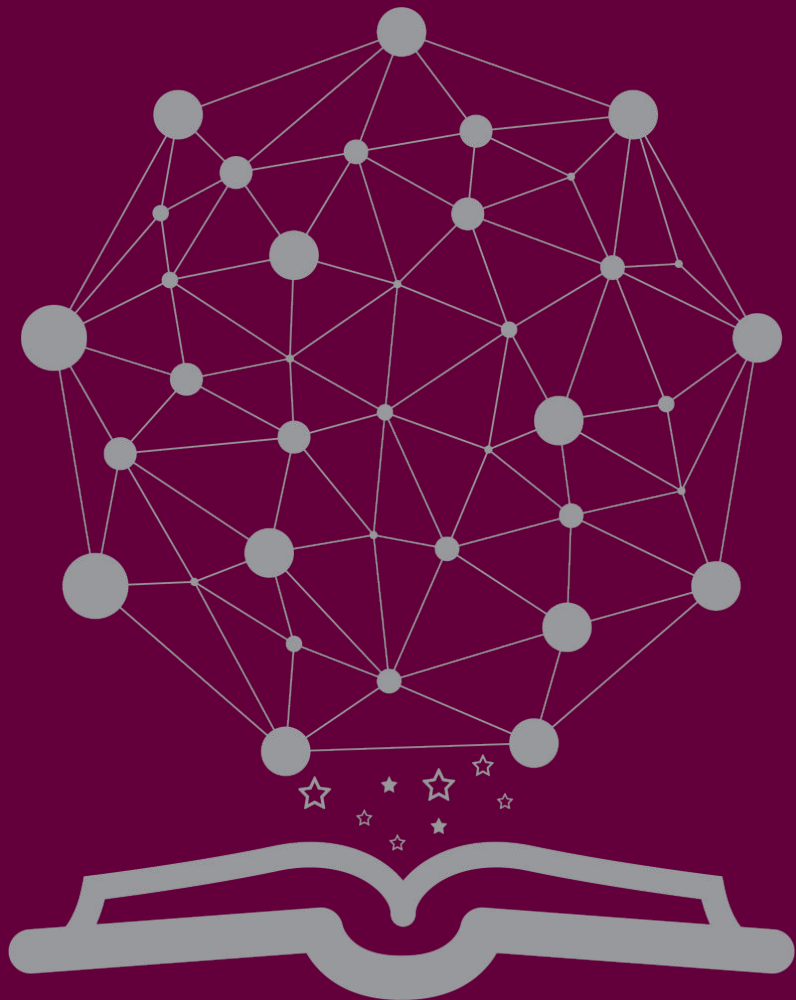


Fine tuned model on one iteration of annotations:

The CpfI based CRISPR systems provided herein can be used to introduce targeted double-strand or single-strand breaks and/or to introduce gene activator and/or repressor systems and without being limitative, can be used for gene targeting, gene replacement, targeted mutagenesis, targeted deletions or insertions, targeted inversions and/or targeted translocations. By co-expression of multiple targeting RNAs directed to achieve multiple modifications in a single cell, multiplexed genome modification can be ensured. This technology can be used to high-precision engineering of plants with improved characteristics, including enhanced nutritional quality, increased resistance to diseases and resistance to biotic and abiotic stress, and increased production of commercially valuable plant products or heterologous compounds.

Model trained once:

The CpfI based CRISPR systems provided herein can be used to introduce targeted double-strand or single-strand breaks and/or to introduce gene activator and/or repressor systems and without being limitative, can be used for gene targeting, gene replacement, targeted mutagenesis, targeted deletions or insertions, targeted inversions and/or targeted translocations. By co-expression of multiple targeting RNAs directed to achieve multiple modifications in a single cell, multiplexed genome modification can be ensured. This technology can be used to high-precision engineering of plants with improved characteristics, including enhanced nutritional quality, increased resistance to diseases and resistance to biotic and abiotic stress, and increased production of commercially valuable plant products or heterologous compounds.



RELATION DETECTION

- Failed attempt
- Types of relations
- Implementation
- Example and limits
- Outlook

Failed attempt



Use NB5 with our entities =



source		target	edge
hours	hours or 24 hours	within	
hours or 24 hours		at least 10	forms
at least 10		at least 100	times
at least 100		at least 1000	times
at least 1000		at least 10000	ester
one		HLAME	produce



Types of Relations (0.94 per line)



Verb:

Subj + Verb + Obj

- Composition includes molecule
- Gene encodes protein
- Cells express antigen

To be:

Subj + To be + Attr

- Oxidation is a pretreatment
- Protein may be a polypeptide
- Carrier be a liquid

Preposition:

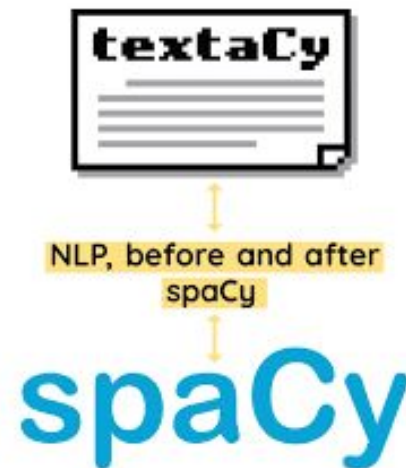
Subj + Prep + Pobj

- Fragment consists of acid
- Cells used for screening
- HBeAg function as toleragen

Part of Speech tagging and Dependency parsing

For each consecutive pair of entities:

1. Check if they correspond to (subj, obj) (subj, attr) (subj, pobj)
2. Verify that they are in the same sentence
3. Check if the text in between has any (verb, aux, prep)
4. If match, clean the relation and add it



Beyond microorganisms



Term extractor (TE):

- Trained on specific data
- Will only generalize to similar text

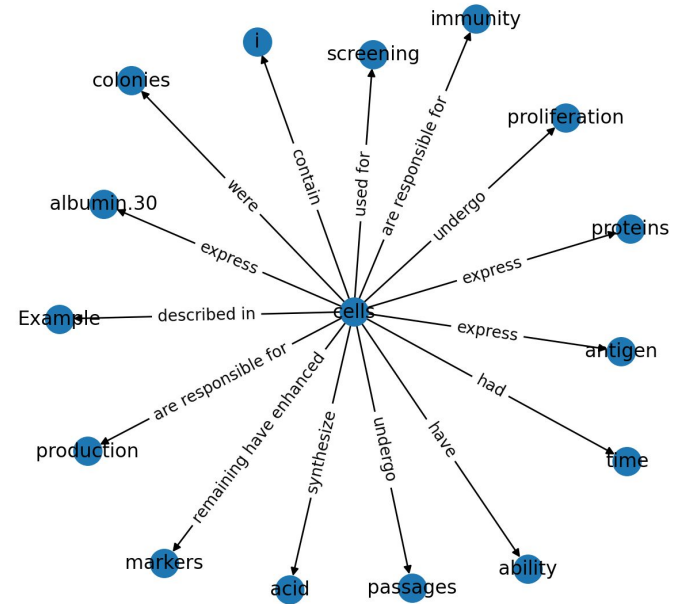
Relation extractor (RE):

- No domain-specific rules
- Generalizes to other data if TE does

Search engine?:

1. Create the graph
2. Extract the entities in the query
3. Retrieve the subgraph containing entities
4. Maybe allow for n-n search

QUERY: Source=cells



THANK YOU!