

EXPLORATION

What's prodigy?

prodigy is a modern annotation tool for creating training data for machine learning models. It runs entirely on the user's own machine and is **fully scriptable** in Python.

prodigy includes a CLI and Python library that lets you run an annotation server, and serve a web app with a collection of interfaces to create annotations for text, HTML, images and audio data.

UPDATE
CALLBACK, CAN
UPDATE MODEL-
IN-THE-LOOP

```
recipe.py
PSEUDOCODE 0

import prodigy

@prodigy.recipe(
    "my-custom-recipe",
    dataset=("Dataset to save answers to", "positional", None, str),
    view_id=("Annotation interface", "option", "v", str)
)

def my_custom_recipe(dataset, view_id="text"):
    # Load your own streams from anywhere you want
    stream = load_my_custom_stream() ← LOAD STREAM

    def update(examples):
        # This function is triggered when Prodigy receives annotations
        print(f"Received {len(examples)} annotations!")

    return {
        "dataset": dataset,
        "view_id": view_id,
        "stream": stream,
        "update": update
    } ← RETURN DICT OF COMPONENTS
```

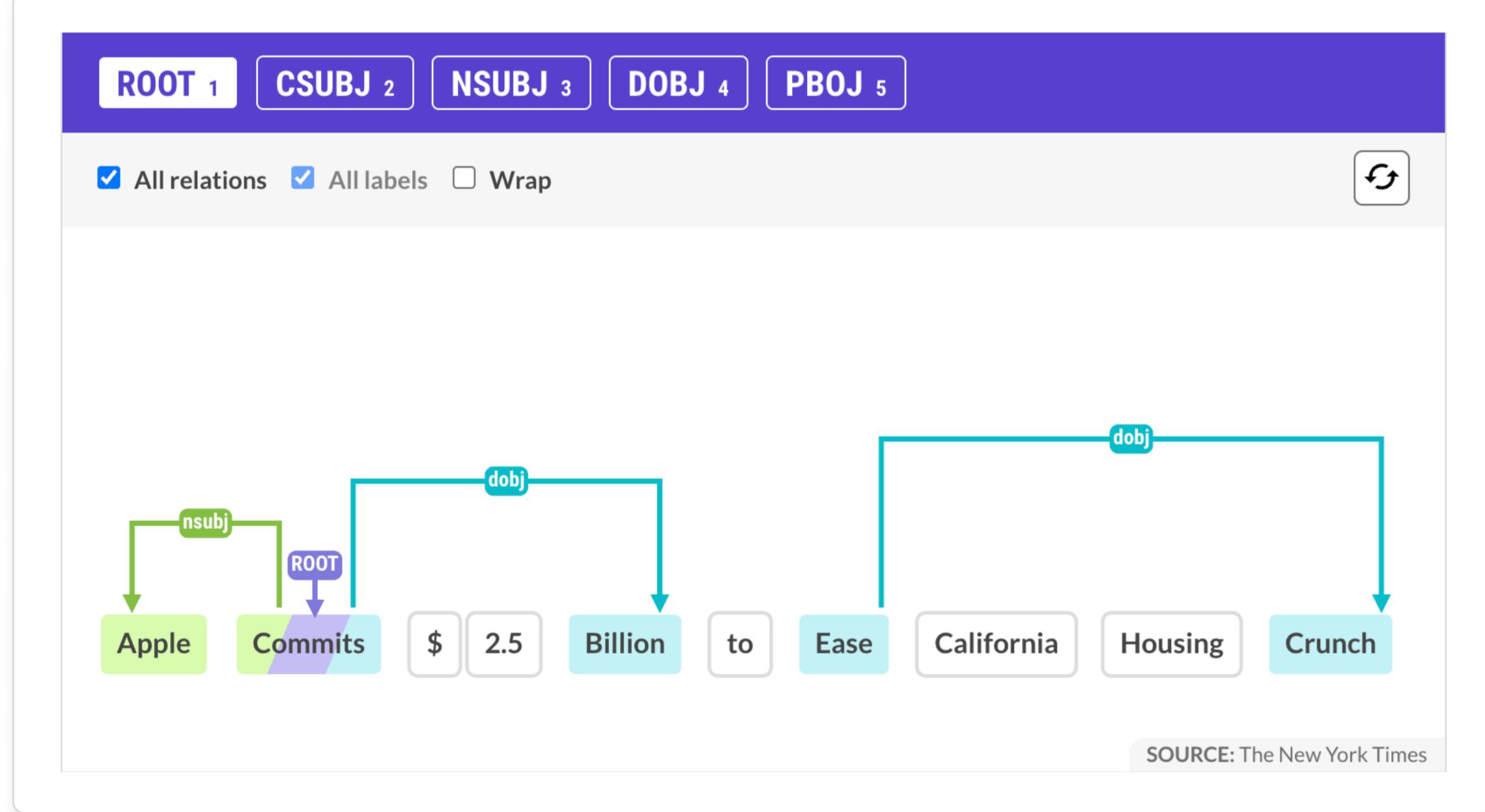
prodigy includes various pre-built workflows, called **recipes**, that load the input data, use a given interface to create the annotations and implement various logic for automation.

How to use prodigy?

START THE ANNOTATION SERVER

```
RECIPE USED      SPACY MODEL      SOURCE DATA
$ prodigy dep.correct deps_news en_core_web_sm ./news_headlines.jsonl
--label ROOT,csubj,nsbj,dobj,pobj --update
UPDATING MODEL
```

LABEL EXAMPLES



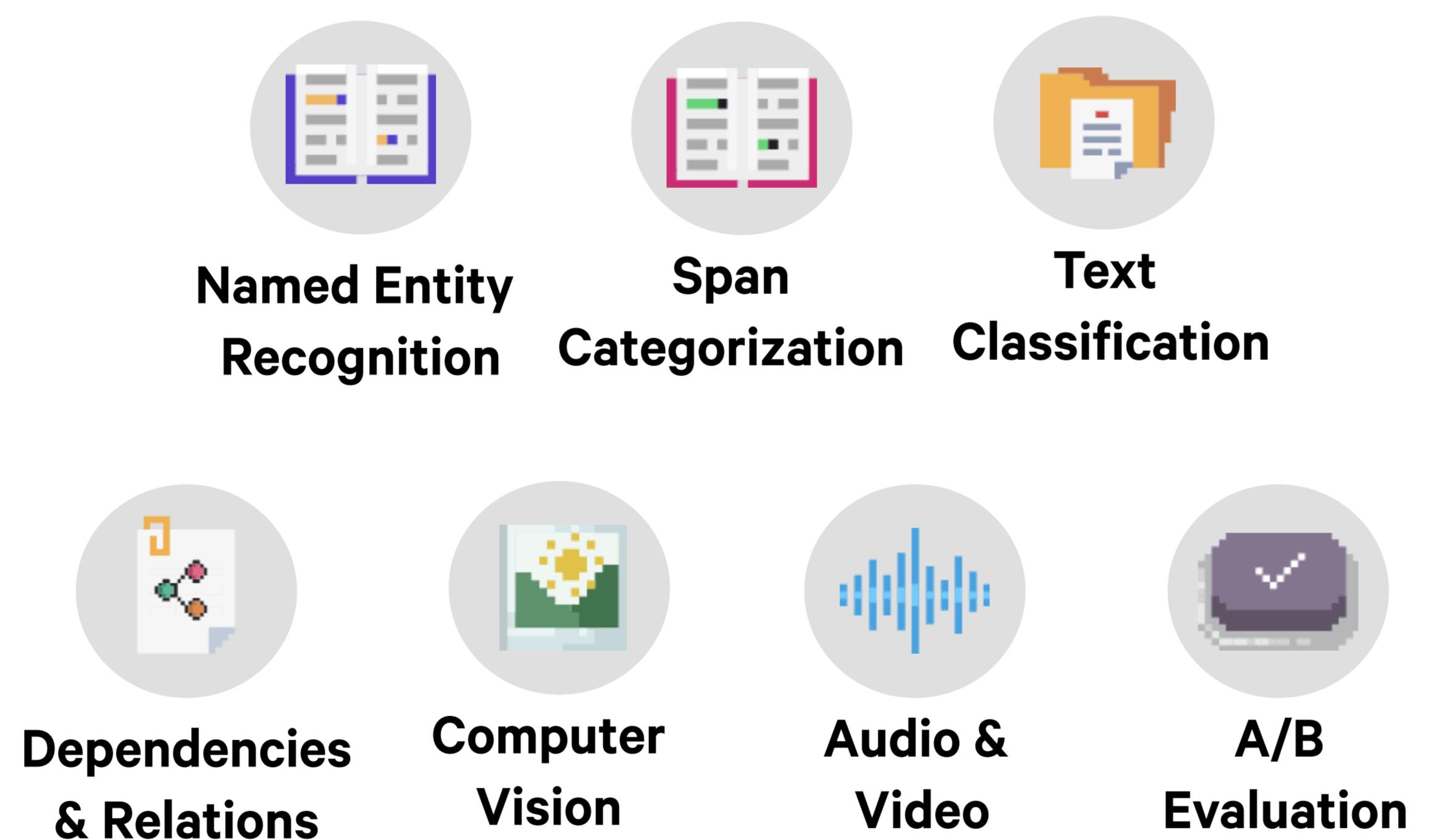
DATA-TO-SPACY AND SPACY TRAIN

```
$ prodigy data-to-spacy ./corpus
--ner news_ner_person,news_ner_org,news_ner_product
--textcat news_cats2018,news_cats2019 --eval-split 0.3
Using language 'en'

===== Generating data =====
Components: ner, textcat
Merging: training and evaluation data for 2 components
- [ner] Training: 685 | Evaluation: 300 (from datasets)
- [textcat] Training: 538 | Evaluation: 230 (30% split)
Training: 1223 | Evaluation: 530
Labels: ner (3) | textcat (5)
Saved 1223 training examples
./corpus/train.spacy
```



What recipes are built-in?



prodigy integrates with our NLP library **spaCy** out-of-the-box and provides workflows for creating training data for spaCy models. But since it's extensible, you can develop your own workflows with another framework or export out your data.

spaCy

```
•••
[nlp]
lang = "en"
pipeline = ["llm"]

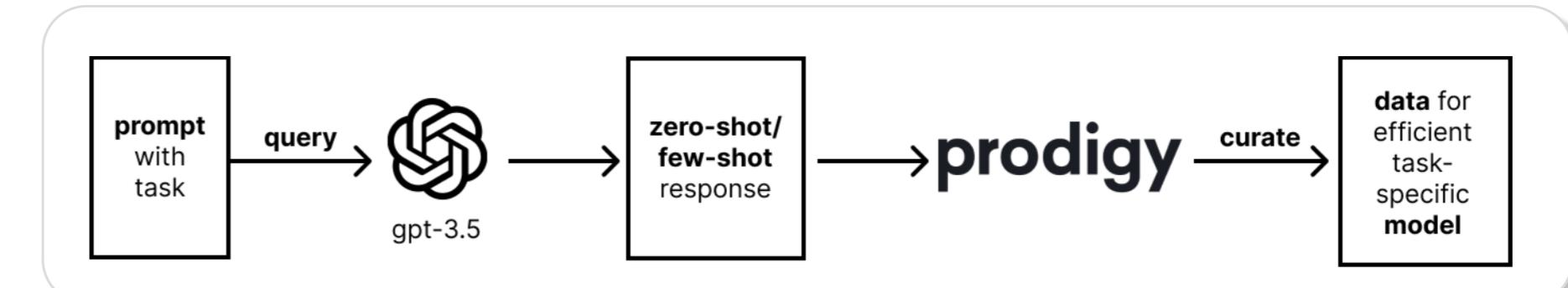
[components]
[components.llm]
factory = "llm"

[components.llm.task]
@llm_tasks = "spacy.TextCat.v1"
labels = COMPLIMENT,INSULT

[components.llm.backend]
@llm_backends = "spacy.REST.v1"
api = "OpenAI"
config = {"model": "text-davinci-003", "temperature": 0.3}
```

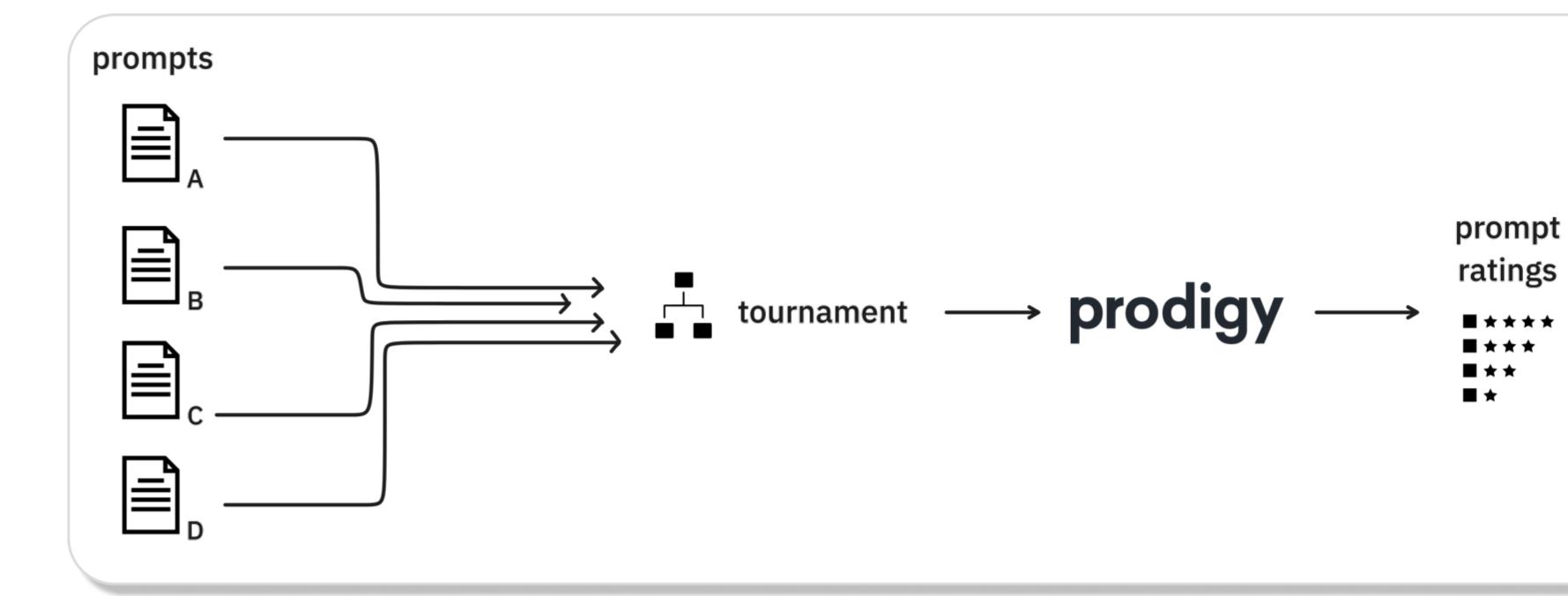
SPACY-LLM CONFIG FILE

What's new in v1.12?



LLM & Prompt Engineering

prodigy supports recipes and tools for Large Language Models and prompt engineering via prompt tournaments. In these tournaments, your prompts will compete as you annotate the one that gives the best results.



PRODIGY'S OPENAI NER CORRECT RECIPE

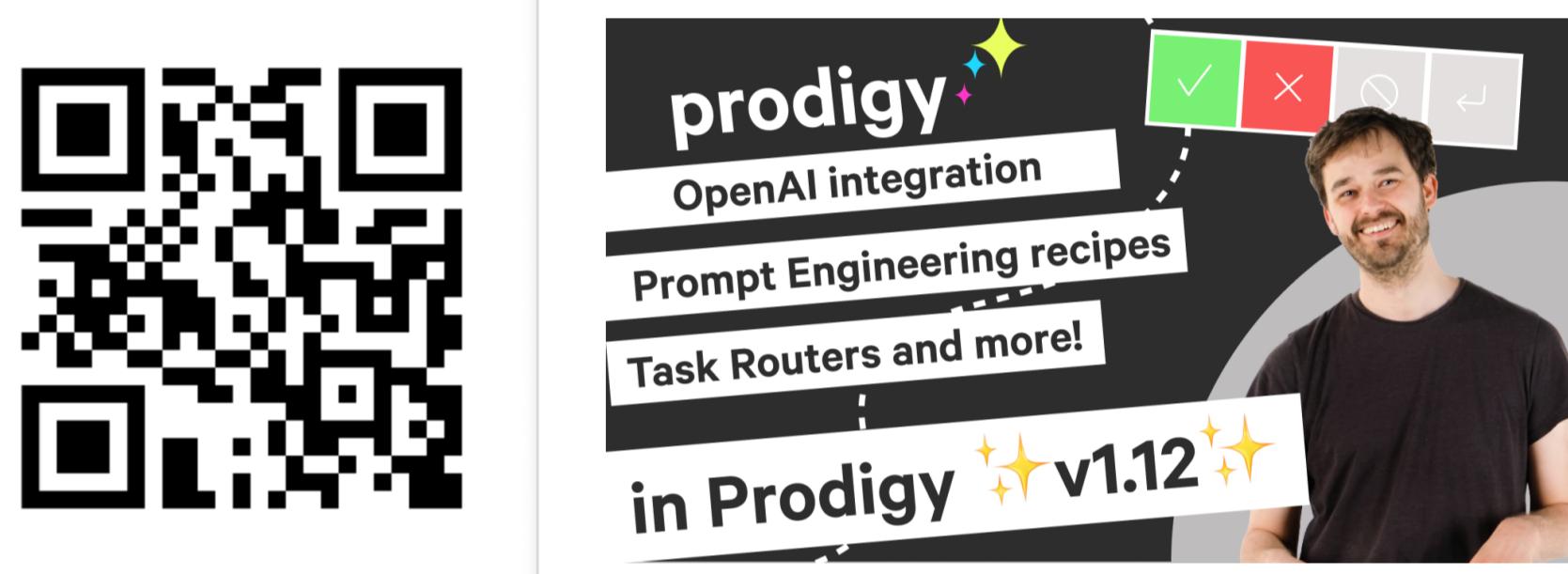
```
CELL 1 ORGANISM_SUBSTANCE 2 PATHOLOGICAL_FORMATION 3
MULTITISSUE_STRUCTURE 4 ORGANISM_SUBDIVISION 5 ORGAN 6
CELLULAR_COMPONENT 7 ANATOMICAL_SYSTEM 8 TISSUE 9
DEVELOPING_ANATOMICAL_STRUCTURE 10
IMMATERIAL_ANATOMICAL_ENTITY 11
```

Our growing understanding of the role that unfavorable patterns of gene expression play in the etiology of **neurodegenerative disease**. **PATHOLOGICAL_FORMATION** emphasizes the need for strategies to selectively block the biosynthesis of harmful **proteins** **ORGANISM_SUBSTANCE** in the **brain** **ORGANISM_SUBDIVISION**.

Show the prompt for OpenAI
Show the response from OpenAI

Cell:
Organism_substance: proteins
Pathology: neurodegenerative disease
Multi-tissue_structure:
Organism_subdivision: brain
Organ:
Cellular_component:
Anatomical_system:
Tissue:
Developing_anatomical_structure:
Immaterial_anatomical_entity:

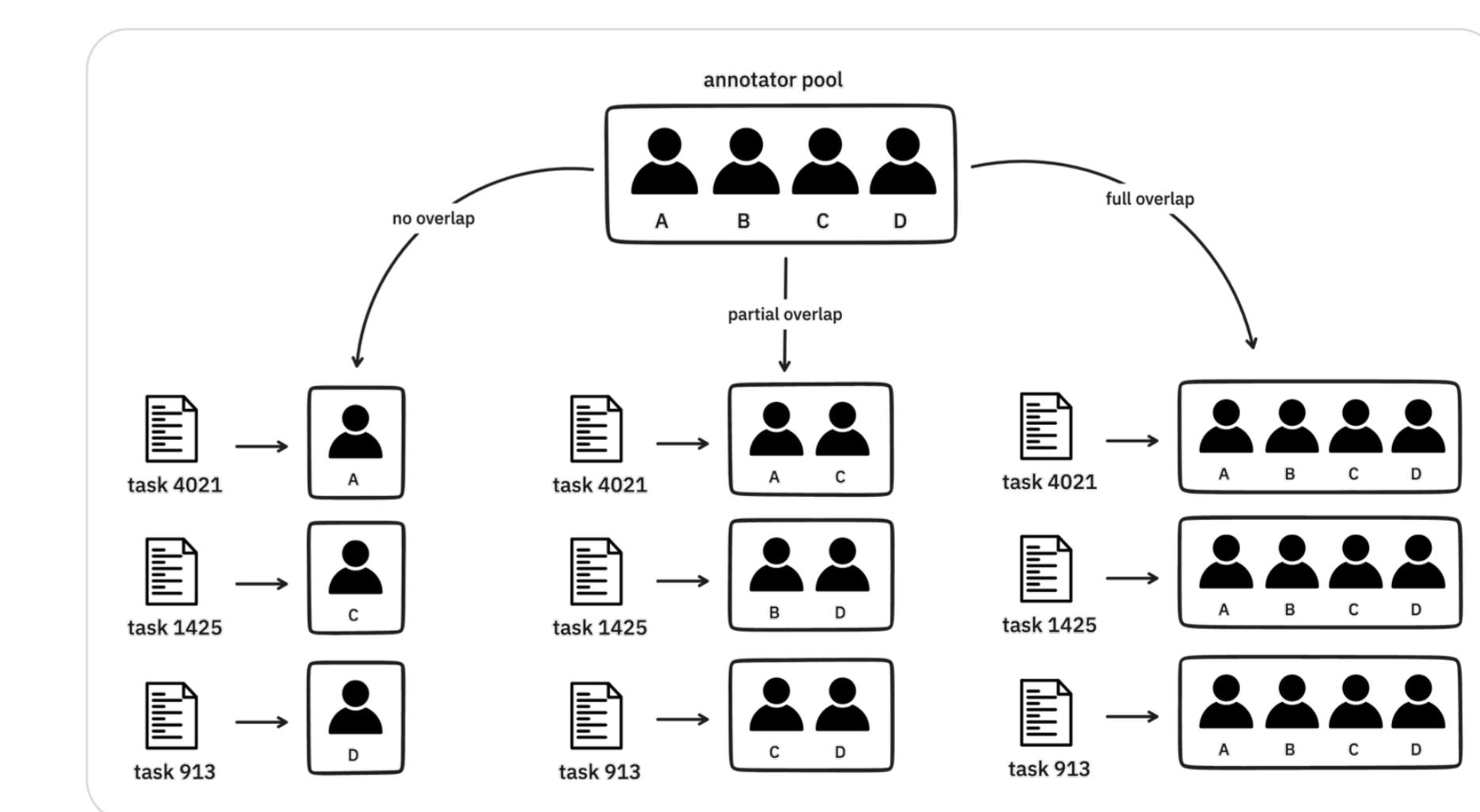
RAW RESPONSE FROM LLM



SCAN HERE TO WATCH
VINCENT'S V1.12 VIDEO

Task Routing

prodigy allows you to distribute the annotation workload among multiple annotators or workers. You can make sure examples are seen by multiple annotators or even set up **custom routing** so that specific examples get seen by specific annotators.



What's future work?

The team is working hard on several large projects. For spaCy, we just released **spaCy v3.6** and **spacy-llm v0.4**, an extension for integrating LLMs into structured NLP pipelines. For Prodigy, we're working on **Prodigy v2.0**, which includes enhanced data validation, spaCy-llm integration, built-in IAA metrics, and Structured examples, and much more. We're also restarting Beta testing for **Prodigy Teams**, our long-awaited SaaS version of Prodigy. We're excited for a public release soon!

```
Task router based on model confidence
PSEUDOCODE 0

def task_router_conf(ctrl: Controller, session_id: str, item: Dict[str]):
    """Route tasks based on the confidence of a model."""
    # Get all sessions known to the Controller now
    all_annotators = ctrl.all_session_ids

    # Calculate a confidence score from a custom model
    confidence_score = model(item['text'])

    # If the confidence is low, the example might be hard
    # and then everyone needs to check
    if confidence_score < 0.3:
        return all_annotators
    # Otherwise just one person needs to check.
    # We re-use the task_hash to ensure consistent routing of the task.
    idx = item['task_hash'] % len(all_annotators)
    # Return list with a single annotator reference
    return [all_annotators[idx]]
```

EXAMPLE OF CUSTOM TASK ROUTER

SCAN HERE TO
REQUEST A FREE
PRODIGY
RESEARCH
LICENSE

