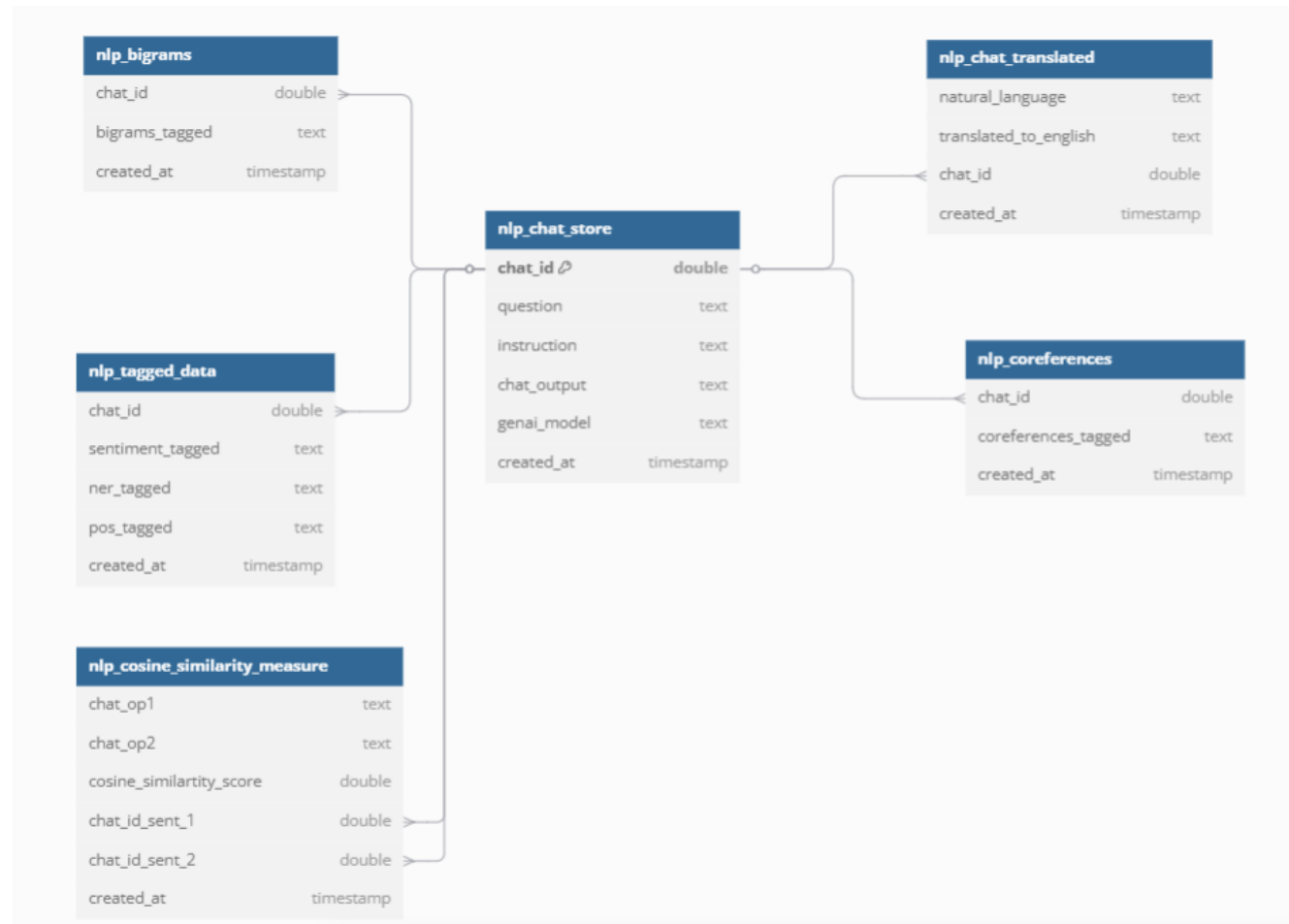This document implements few NLP tasks using GenAI plugins step using Pentaho Data Integration (PDI) which calls OpenAI based LLM. The tests are done via "gpt-4o-mini" model.

Below is a snapshot of various Natural Language Processing (NLP) tasks are categorized as below. Few of these NLP tasks we will address using PDI and LLM (OpenAI).
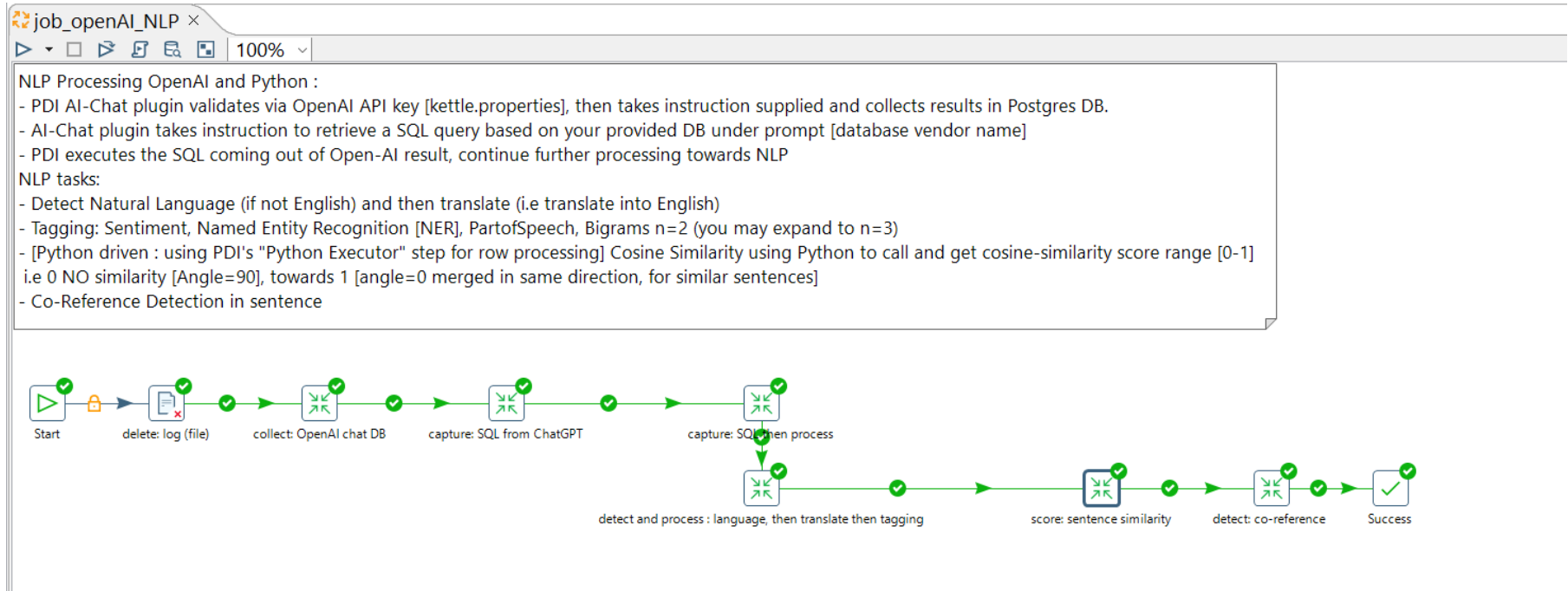


Image by NLPlanet

**Data Model:**

ERD: https://dbdiagram.io/d/ChatGpt_PDI-67c71518263d6cf9a030e12c

**Pentaho Data Integration jobs and transformtations:**

job_openAI_NLP ×

100%

NLP Processing OpenAI and Python :
- PDI AI-Chat plugin validates via OpenAI API key [kettle.properties], then takes instruction supplied and collects results in Postgres DB.
- AI-Chat plugin takes instruction to retrieve a SQL query based on your provided DB under prompt [database vendor name]
- PDI executes the SQL coming out of Open-AI result, continue further processing towards NLP
NLP tasks:
- Detect Natural Language (if not English) and then translate (i.e translate into English)
- Tagging: Sentiment, Named Entity Recognition [NER], PartofSpeech, Bigrams n=2 (you may expand to n=3)
- [Python driven : using PDI's "Python Executor" step for row processing] Cosine Similarity using Python to call and get cosine-similarity score range [0-1]
  i.e 0 NO similarity [Angle=90], towards 1 [angle=0 merged in same direction, for similar sentences]
- Co-Reference Detection in sentence

Start        delete: log (file)        collect: OpenAI chat DB        capture: SQL from ChatGPT        capture: SQL then process

detect and process : language, then translate then tagging        score: sentence similarity        detect: co-reference        Success

# #Generate Data (multiple natural languages) using OpenAI and store in a database for processing NLP

**#Testing to see SQL works**

# #Translation, Sentiment, NER, bigrams, PoS

# #Sentence Similarity (How similar are two sentences?)

**ChatGPT couldn't provide a score through prompt, so having this via Python while using data processing in PDI transformation (data in flight).**



## --Co-Reference