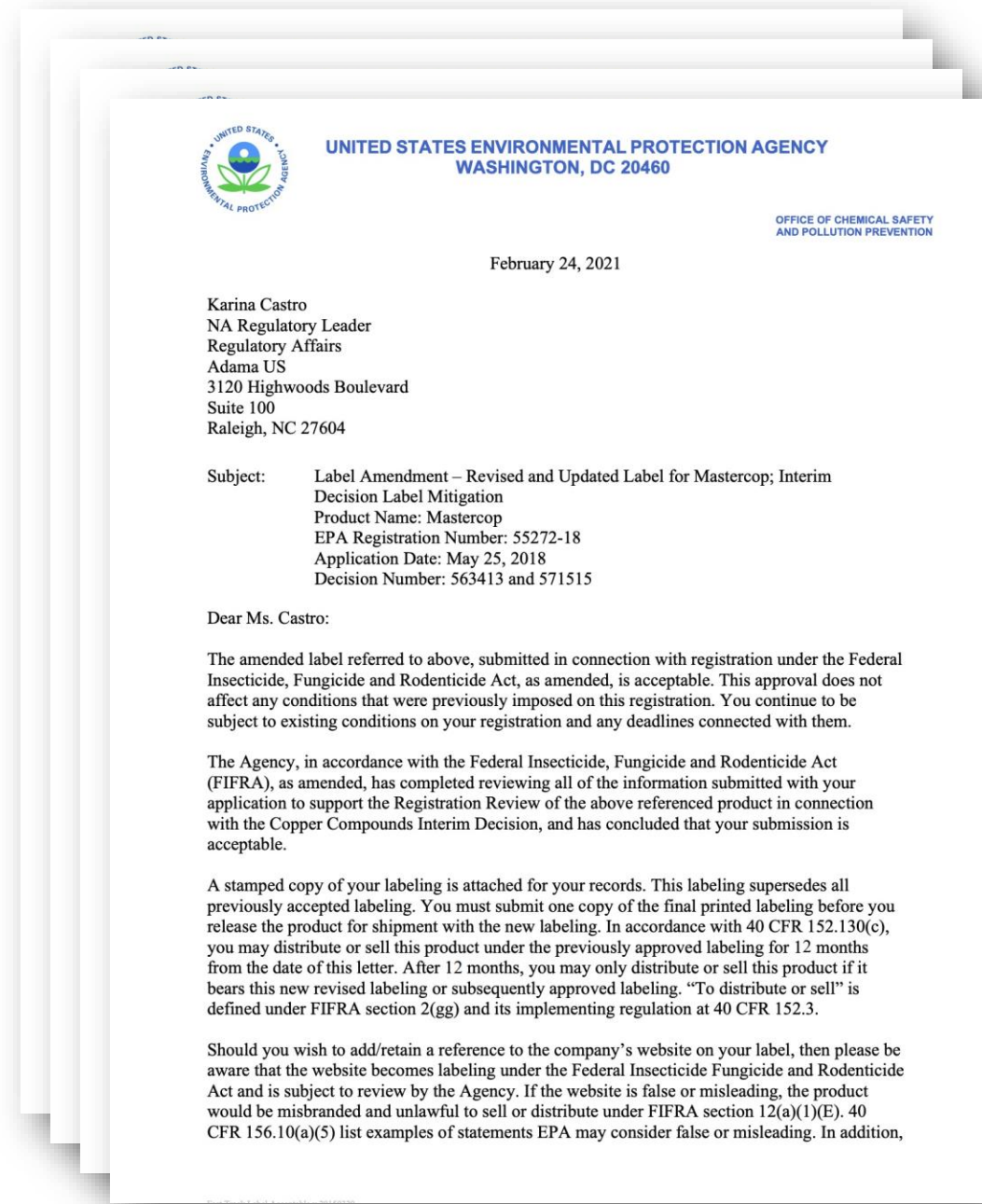


## Problem Definition



The U.S. EPA contains over **360K** pesticide-label documents—each packed with critical application and safety instructions for growers—but they exist only as inconsistent PDFs that are cumbersome to search, forcing farmers to comb through lengthy documents for dosage, timing, or safety guidance.

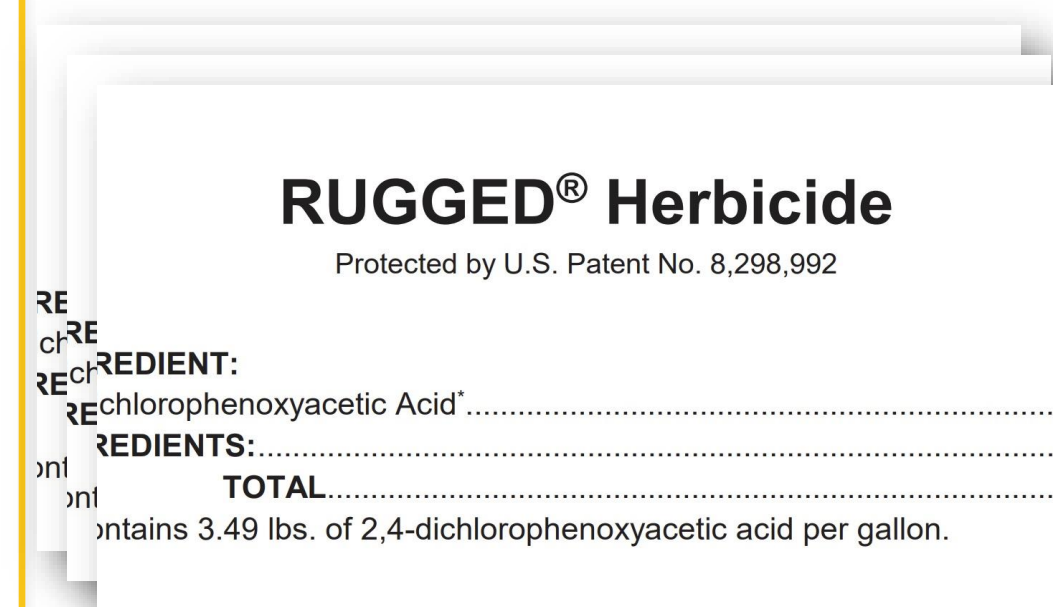
Our solution harnesses **state-of-the-art artificial intelligence** to automatically parse and extract important data—crop, target pest or disease, application rate, safety information ... —and loads them into an **open-source database**. This AI-driven pipeline turns text into actionable insights, so growers can focus on farming rather than paperwork.

## Methodology



### STEP 1: Label Filtering

Filter down the EPA database entries into a subset of 5K pesticide labels. We narrow down to the scope defined by our client.



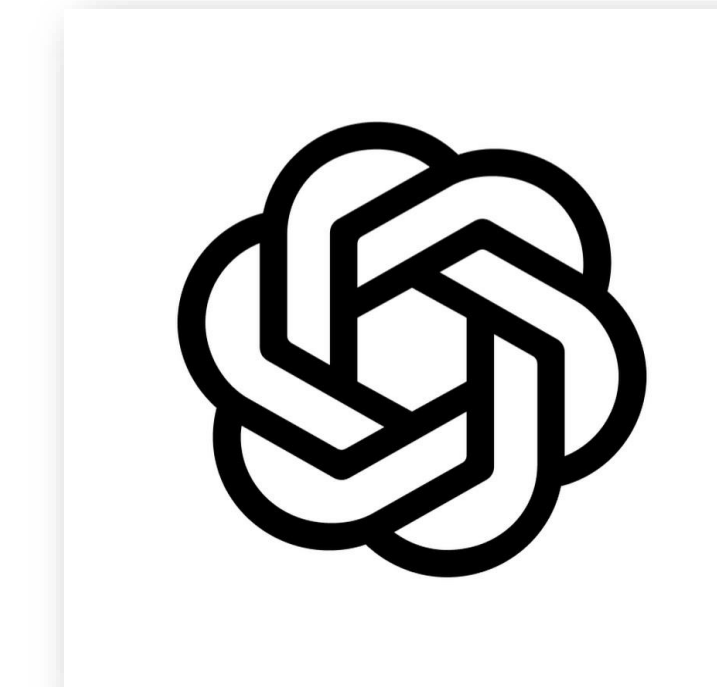
### STEP 2: PDF Parsing

Transform the 5K PDF pesticide labels, including tables, into text files using PDF Parsers and OCR.



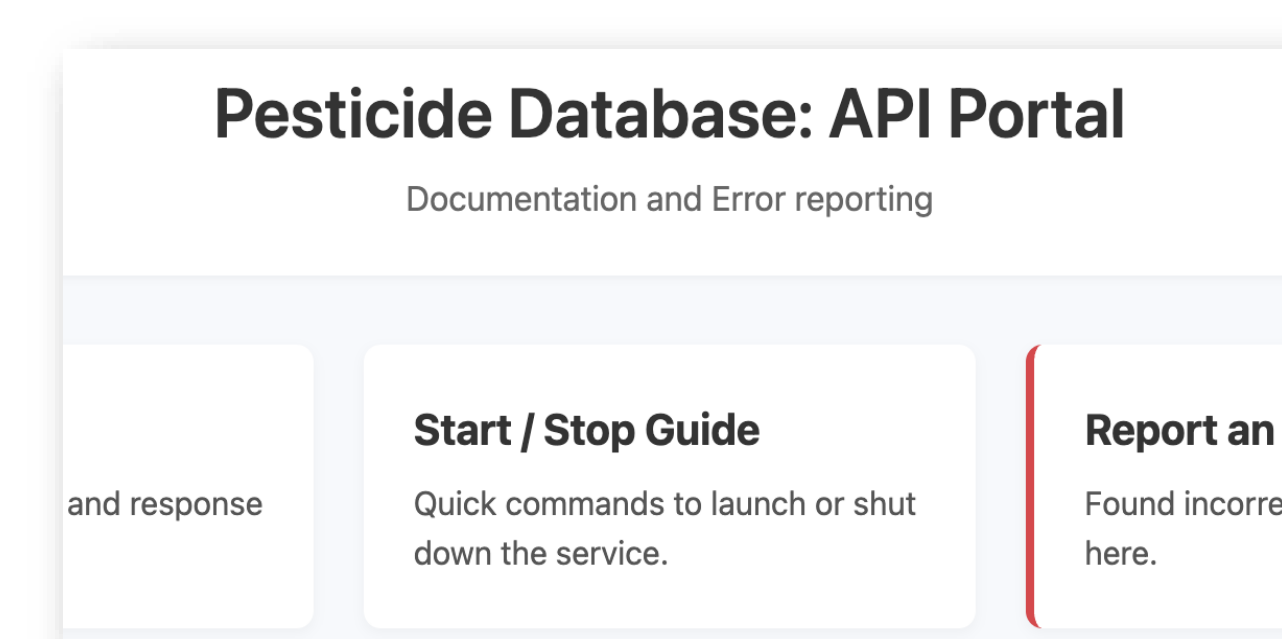
### STEP 3: AI Inference

Take in the text files and extract out the pesticide information by inferencing a LLM (reasoning model).



### STEP 4: Database

Parse the JSON and upload the data into the database, where it can be accessed via a public API.



## Our Goals

Empower farmers, researchers, and agricultural professionals with reliable insights on pesticide application, safety guidelines, and compliance.



90% Accuracy on automated extraction pipeline

## AI Model Research

At first, we considered a variety of possible models to use, including models from OpenAI, Google, Anthropic, and others. Additionally, we experimented with open-source models from Hugging Face.

AI models aren't perfect, and many of the models we tested were subject to hallucinations and accuracy struggles, including missing or incorrect information. In the pursuit of an accurate and cost-effective solution, we developed a variety of parsing techniques, tested over a dozen models, and explored several prompt engineering strategies.

After many weeks of experimentation, we found OpenAI's newest reasoning model o4-mini was the most powerful in our price range.



Cost per label on average: 5.4 cents  
Overall accuracy: >90% (based on internal experiments)

o4-mini