



criolho /
Drupal-Feeds-Using-LLMs



`<>` Code

Issues

Pull requests

Actions

Projects

Wiki

Security

L



master ▾

Drupal-Feeds-Using-LLMs / README.md



criolho Update README.md

afafafaf · 6 minutes ago



206 lines (121 loc) · 12.6 KB

Preview

Code

Blame



Raw

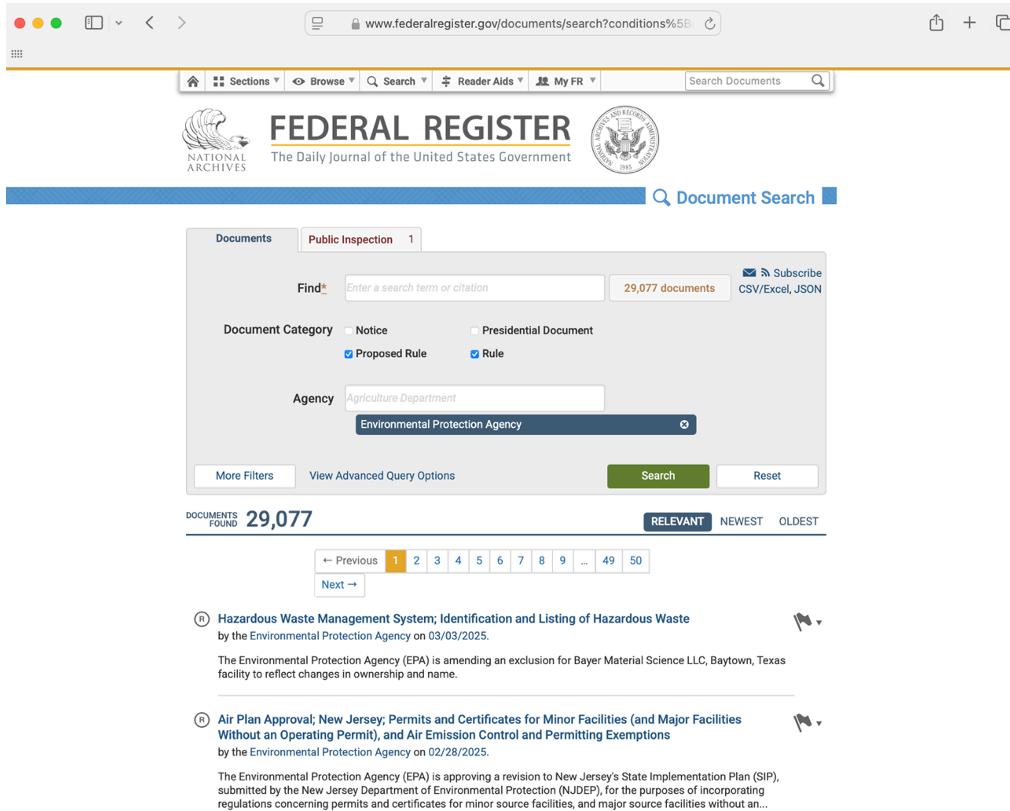


Drupal Feeds Using LLMs

If you regularly bring data into Drupal from external sources, you may be able to enrich the data using LLMs and a bit of custom programming. This repo demos a workflow utilizing Python scripts running on a Linux server, web scraping, LLMs and the Drupal Feeds module. Your tech stack may be different but the techniques discussed in this talk are easily translated to other environments.

The repo demonstrates two scrapers, one that interacts with the Federal Register API and the other which scrapes an EPA web page. The scrapers output JSON suitable for Drupal feeds import.

Federal Register



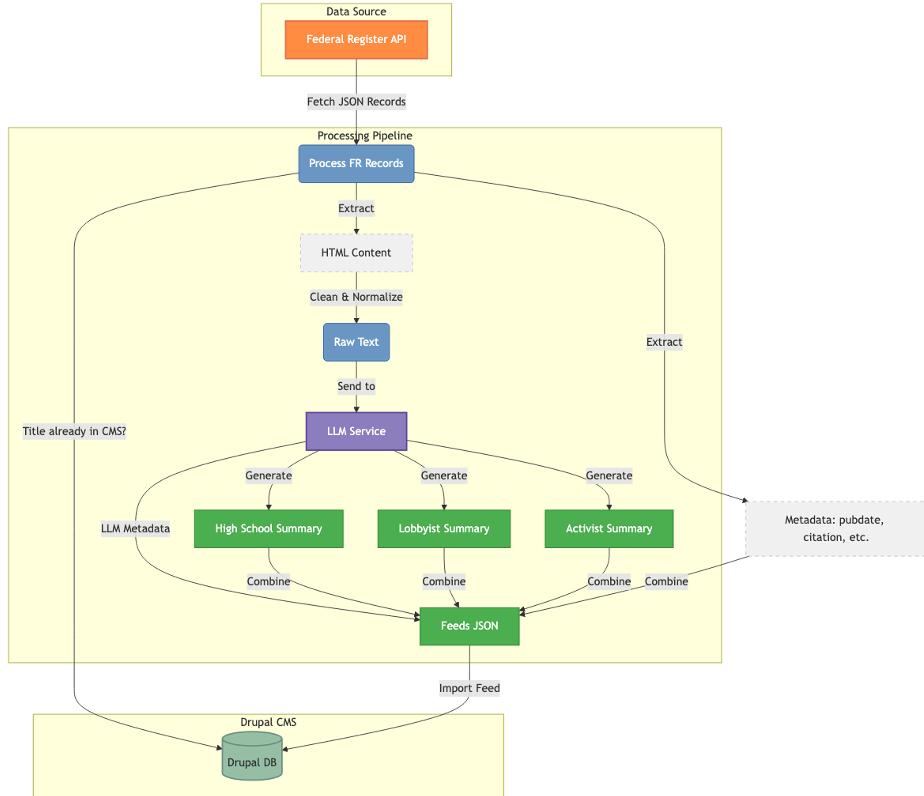
The screenshot shows the homepage of the Federal Register website. At the top, there's a navigation bar with links for 'Sections', 'Browse', 'Search', 'Reader Aids', and 'My FR'. Below the header, the 'FEDERAL REGISTER' logo is displayed, along with the text 'The Daily Journal of the United States Government' and the seal of the National Archives and Records Administration. A large search bar at the top right contains the placeholder 'Search Documents' with a magnifying glass icon. Below the search bar, there's a 'Document Search' section with tabs for 'Documents' and 'Public Inspection'. Under 'Documents', there's a search input field with the placeholder 'Enter a search term or citation' and a button showing '29,077 documents'. To the right of the search input are 'Subscribe' and 'CSV/Excel, JSON' buttons. Below the search input, there are filters for 'Document Category': 'Notice' (unchecked), 'Presidential Document' (unchecked), 'Proposed Rule' (checked), and 'Rule' (checked). There's also a dropdown for 'Agency' set to 'Agriculture Department', with 'Environmental Protection Agency' highlighted in blue. Below the agency dropdown are buttons for 'More Filters' and 'View Advanced Query Options'. At the bottom of the search section, there are buttons for 'Search' and 'Reset'. Below the search section, a summary shows 'DOCUMENTS FOUND 29,077' with buttons for 'RELEVANT', 'NEWEST', and 'OLDEST'. A navigation bar below shows page numbers from 1 to 50, with page 1 highlighted in yellow. Below the navigation bar are two news items:

- (R) Hazardous Waste Management System; Identification and Listing of Hazardous Waste by the Environmental Protection Agency on 03/03/2025. The Environmental Protection Agency (EPA) is amending an exclusion for Bayer Material Science LLC, Baytown, Texas facility to reflect changes in ownership and name.
- (R) Air Plan Approval; New Jersey; Permits and Certificates for Minor Facilities (and Major Facilities Without an Operating Permit), and Air Emission Control and Permitting Exemptions by the Environmental Protection Agency on 02/28/2025. The Environmental Protection Agency (EPA) is approving a revision to New Jersey's State Implementation Plan (SIP), submitted by the New Jersey Department of Environmental Protection (NJDEP), for the purposes of incorporating regulations concerning permits and certificates for minor source facilities, and major source facilities without an...

- Site: <https://federalregister.gov>
- Background: This is where the U.S. government publishes all proposed federal regulations or rules about to go into effect. The site has an excellent UI and is also known for its developer-friendly API and JSON output.
- Goals: connect to Federal Register API and get JSON metadata for recently-published rules. Then:
 - Make minor transformations to their field data for our matching fields in Drupal, e.g. prepending a citation to their title to match our Drupal node title best practices
 - Add 3 LLM-generated summaries individually focused on the needs of: 1) high school students, 2) corporate lobbyists, 3) environmental activists.

Federal Register processing pipeline

High level overview of the data processing game plan:



Inputs and Outputs

Federal Register developer interactive API:

https://www.federalregister.gov/developers/documentation/api/v1#/Federal%20Register%20Documents/get_documents_format

Federal Register JSON request for the EPA:

[https://www.federalregister.gov/api/v1/documents.json?
fields\[\]=%22abstract%22&fields\[\]=%22type%22&fields\[\]=%22pdf_url%22&fields\[\]=%22document_number%22&fields\[\]=%22publication_date%22&fields\[\]=%22agencies%22&fields\[\]=%22title%22&per_page=3&conditions%5Bpublication_date%5D%5Bgte%5D=2025-03-01&conditions%5Bagencies%5D%5B%5D=environmental-protection-agency&conditions%5Btype%5D%5B%5D=RULE&conditions%5Btype%5D%5B%5D=PRORULE](https://www.federalregister.gov/api/v1/documents.json?fields[]=%22abstract%22&fields[]=%22type%22&fields[]=%22pdf_url%22&fields[]=%22document_number%22&fields[]=%22publication_date%22&fields[]=%22agencies%22&fields[]=%22title%22&per_page=3&conditions%5Bpublication_date%5D%5Bgte%5D=2025-03-01&conditions%5Bagencies%5D%5B%5D=environmental-protection-agency&conditions%5Btype%5D%5B%5D=RULE&conditions%5Btype%5D%5B%5D=PRORULE)

```

{
  count: 5,
  description: "Documents published on or after 03/02/2025, from Environmental Protection Agency, and of type Rule or Proposed Rule",
  total_pages: 1,
  - results: [
    - {
      type: "Rule",
      publication_date: "2025-03-07",
      abstract: "The Environmental Protection Agency (EPA) is approving a State Implementation Plan (SIP) revision submitted through the South Carolina Department of Health and Environmental Control (SC DHEC) on September 26, 2023, regarding updates to the State's Cross-State Air Pollution Rule (CSAPR) emissions trading programs. The SIP revision incorporates by reference (IBRs) certain amendments EPA has made to the regulations for the Federal CSAPR trading programs for annual emissions of nitrogen oxides (NO<INF>x</INF>) and sulfur dioxide (SO<INF>x</INF>) for large electric generating units (EGUs). EPA is approving South Carolina's September 26, 2023, SIP revision because it is consistent with EPA's good neighbor CSAPR trading programs and the Clean Air Act (CAA or Act).",
      agency_names: [
        "Environmental Protection Agency"
      ],
      + cfr_references: [ ... ],
      citation: "90 FR 11478",
      effective_on: "2025-04-07",
      document_number: "2025-03609",
      pdf_url: "https://www.govinfo.gov/content/pkg/FR-2025-03-07/pdf/2025-03609.pdf",
      body_html_url: "https://www.federalregister.gov/documents/full_text/html/2025/03/07/2025-03609.html",
      title: "Air Plan Approval; SC; Updates to the Cross-State Air Pollution Rule",
      + topics: [ ... ]
    },
    ...
  ]
}

```

1. JSON from API

The JSON API response contains the raw document data, including the abstract, agency names, and various URLs.

Munged title (no AI)

The munged title is the original title of the document, "Air Plan Approval; SC; Updates to the Cross-State Air Pollution Rule".

Text fetch (no AI)

The text fetch is a summary generated by AI, stating "...This plan is called the State Implementation Plan and it helps control pollution...".

Summaries (AI-generated)

The summaries are AI-generated descriptions of the document's content, such as "These pollutants contribute to environmental contamination, affecting ecosystems...".

2. JSON for Drupal

The JSON for Drupal is a transformed version of the API response, likely containing fields like 'body' and 'summary'.

Data transformations

The diagram shows the flow from the API response to the Drupal transformation step, where the data is processed and summarized.

Drupal

On the Drupal side of things we have a content type Federal Register that has matching fields:

Field	Widget
>Title	Textfield
Short Summary	Text area with a summary
High School Summary	Text area (multiple rows)
Activist Summary	Text area (multiple rows)
Lobbyist Summary	Text area (multiple rows)
Document Number	Textfield
Agency	Check boxes/radio buttons
Federal Register Type	Check boxes/radio buttons
Effective On	Date and time

And our Federal Register Feed that will consume our JSON output looks like this:

Context*		Summary	Configure	Unique	Remove
Source	Target				
title	Title (title)		<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
abstract	Short Summary (body): Text	Format: Full HTML	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
- Select a source -	Short Summary (body): Summary				
document_number	Document Number (field_document_number)		<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
type	Federal Register Type (field_federal_register_type)	Reference by: Name Autocreate entities: No	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
agency_names	Agency (field_agency)	Reference by: Name Autocreate entities: No	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
citation	Citation (field_citation)		<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
effective_on	Effective On (field_effective_on)	Default timezone: UTC	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
article_text	Raw Text (field_raw_text): Raw text of Federal Register doc		<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
high_school_summary	High School Summary (field_high_school_summary)	Format: Full HTML	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

News Meta-summary

We can also have our script create a News summary that iterates over all the Federal Registers we are importing and ask the LLM to summarize all the other summaries, a kind of meta-summary. This should not be published as is – that's "cheating" – but it could certainly form the skeleton on which a human editor could build. Imagining we are using a simple Drupal "Article" content type we might see:

Environmental Protection Agency Regulatory Review from February 28, 2025 to March 07, 2025

LLM-generated meta-summary - a human should edit!

14 March 2025

The EPA is actively engaged in various air quality and hazardous waste management actions, as evidenced by recent Federal Register notices. Several actions involve State Implementation Plan (SIP) revisions, including approvals for South Carolina (90 FR 11478, 2025-03-07) concerning the Cross-State Air Pollution Rule (CSAPR) and New Jersey (90 FR 10872, 2025-02-28) regarding permits for minor and major facilities. These revisions aim to reduce emissions of nitrogen oxides (NOX) and sulfur dioxide (SO2) and regulate stationary sources to meet National Ambient Air Quality Standards (NAAQS). Activists should **scrutinize** whether these regulations go far enough to protect **public health** and **ecosystems**, especially for **vulnerable populations**. It's crucial to **advocate** for stronger emission reduction targets and increased **monitoring** to ensure **accountability**.

Other actions include proposed approvals for Alabama's SIP revision (90 FR 11500, 2025-03-07), which modifies the definition of "volatile organic compounds" (VOC), and a reopened comment period for Ohio's regional haze plan (90 FR 10876, 2025-02-28) focusing on permit conditions for power plants. The Alabama revision warrants careful examination of the **long-term effects** of added non-VOCs on **public health** and the **environment**. The Ohio plan provides an opportunity to **advocate** for stronger **environmental protections** and **challenge** permit conditions that may not adequately protect **public health** and the **environment** from **emissions** contributing to **air pollution** and **acid rain**.

The EPA is also addressing hazardous waste management, including the deletion of sites from the Superfund National Priorities List (NPL) (90 FR 11218, 2025-03-05) and an amendment reflecting a name change for a facility in Baytown, TX (90 FR 11025, 2025-03-03). While these actions may seem like progress or administrative updates, it's essential to **scrutinize** the details and ensure continued **monitoring** to prevent future risks to **public health** or **ecosystems**. Activists should **demand** that facilities adhere to the strictest standards to prevent **environmental contamination** and protect **soil health**, paying close attention to potential impacts on **vulnerable populations** and **environmental justice** concerns.

Finally, the EPA has reopened the comment period for the proposed regulation addressing the risks associated with C.I. Pigment Violet 29 (PV29) (90 FR 11142, 2025-03-04). This presents an opportunity for **environmental** activists to **scrutinize** the proposed measures and **voice concerns** regarding potential **environmental** and **public health** impacts. It is important to **advocate** for stronger regulations to protect **vulnerable populations** and the **environment**, focusing on potential **environmental contamination**, impacts on **ecosystems**, and risks to **public health**. **Demand** greater **transparency** and **accountability** in the regulation of PV29 and push for stricter **emissions controls** and **monitoring** to minimize potential harm.

Sample Federal Register

After importing here's what a sample FR looks like:

90 FR 10876 - Air Plan Approval; Ohio; Regional Haze Plan for the Second Implementation Period

28 February 2025

Direct from Federal Register JSON

Short Summary

The Environmental Protection Agency (EPA) is reopening the public comment period for 15 days for a proposed rule published August 30, 2024. Reopening the public comment period is focused exclusively on the references to various permit conditions within three Director's Final Findings and Orders (DFFOs) included in the docket that were issued by the Ohio Environmental Protection Agency (Ohio EPA) as those conditions apply to the limits in the DFFOs for Cardinal Power Plant, Ohio Valley Electric Corp.--Kyger Creek, and General James M. Gavin Power Plant.

This article contains AI-generated summaries.

Disclaimer

Agency

Environmental Protection Agency

AI classification

High School Summary

The Environmental Protection Agency (EPA) is giving the public another chance to share their thoughts on a proposed rule about air quality in Ohio. This rule is part of a larger plan to reduce haze and improve visibility in national parks and wilderness areas. The EPA wants feedback specifically on some permit conditions for three power plants: Cardinal Power Plant, Ohio Valley Electric Corp.--Kyger Creek, and General James M. Gavin Power Plant. These permits set limits on how much pollution these plants can release. The EPA initially asked for comments on this rule, but some important information about the permits wasn't available at the time. Now that this information is available, the EPA is reopening the comment period so people can review it and share their opinions. This is important because it affects the air we breathe and the health of our environment.

Lobbyist Summary

AI summarization

The EPA has reopened the comment period for a proposed rule concerning Ohio's Regional Haze State Implementation Plan (SIP). This reopening focuses exclusively on permit conditions within Director's Final Findings and Orders (DFFOs) for Cardinal Power Plant, Ohio Valley Electric Corp.--Kyger Creek, and General James M. Gavin Power Plant. The permit conditions, initially unavailable, pertain to emission limits outlined in the DFFOs. **Of interest to lobbyists:** The EPA seeks comments on these specific permit conditions as they relate to the limits in the DFFOs. **Actionable items:** Review the newly available permit conditions for the Cardinal, Kyger Creek, and Gavin power plants. **Prepare and submit comments** addressing the permit conditions' impact on the proposed rule. Understand the potential **compliance** implications for the affected power plants. Assess potential **costs** associated with meeting the permit conditions. Evaluate the **economic implications** of the proposed rule and permit conditions on the power sector in Ohio. **Opportunities to engage** with the EPA during the reopened comment period to **influence the final rule**.

Activist Summary

The EPA is providing a second opportunity to comment on a proposed rule regarding Ohio's plan to reduce regional haze. This renewed comment period is centered on specific permit conditions for three power plants: Cardinal Power Plant, Ohio Valley Electric Corp.--Kyger Creek, and General James M. Gavin Power Plant. These permits dictate the allowable pollution levels for these facilities. The EPA initially sought public input, but crucial permit details were missing. Now that these details are accessible, the EPA is reopening the comment period.

Civil Enforcement Actions by the EPA

The screenshot shows the official website of the United States Environmental Protection Agency (EPA) at www.epa.gov/enforcement/civil-and-cleanup-enforcement-cases-and-settlements. The page title is "Civil and Cleanup Enforcement Cases and Settlements". The left sidebar contains a navigation menu under "Enforcement" with sections like "Enforcement Basics", "National Enforcement and Compliance Initiatives", "Enforcement and Compliance Assurance Annual Results for Fiscal Year 2024", "Air Enforcement", "Water Enforcement", "Waste, Chemical and Cleanup Enforcement", "Criminal Enforcement", "Enforcement at Federal Facilities", "Data and Results", and "Policy, Guidance and Publications". Below this is a link to "Enforcement: Contact Us". The main content area features a heading "Civil and Cleanup Enforcement Cases and Settlements" and a sub-section stating "Currently available civil cases are listed below. Each case has a brief description and a link to detailed information about the case. You can filter the cases by statute or use the provided search box to search the table. You may also sort the data table by selecting a header." A search bar and a "Search" button are located at the top right of the table. The table itself has columns for "Respondent", "Description", "Order Type", and "Date". It lists two cases: one for Oasis Mobile Home Park Safe Drinking Water Act and another for Stericycle, Inc. RCRA Settlement Summary.

Respondent	Description	Order Type	Date
Oasis Mobile Home Park Safe Drinking Water Act	On January 16, 2025, the U.S. Environmental Protection Agency (EPA) and the Department of Justice (DOJ) reached a settlement agreement with the operators of the Oasis Mobile Home Park (Park) for violations of the Safe Drinking Water at the Park (SDWA).	Consent Decree	2025-01-16
Stericycle, Inc. RCRA Settlement Summary	On January 16, 2025, the Environmental Protection Agency (EPA) and the Department of Justice announced a settlement agreement with Stericycle, Inc., a national provider of hazardous waste transportation, storage, and disposal services, for violations of the Resource Conservation and Recovery Act's (RCRA) regulations regarding hazardous waste transport, storage, and recordkeeping.	Stipulation and Order of Settlement and Complaint	2025-01-16
	On January 16, 2025, the Environmental Protection		

- Site: <https://www.epa.gov/enforcement/civil-and-cleanup-enforcement-cases-and-settlements>
- Background: Centralized source for EPA civil actions for violations. The EPA [historically] enforces a wide range of environmental laws, including the Clean Air Act (CAA) and Clean Water Act (CWA).
- Goals:
 - Better summaries
 - Extract penalty info
 - What specific laws were violated?
 - Categorize environmental issues

Inputs and Outputs

Turn 14 Clean Air Act Settlement Summary

On January 17, 2025, the Environmental Protection Agency (EPA) and Department of Justice reached a settlement agreement with Turn 14 Distribution, Inc. ("Turn 14"), for violations of the Clean Air Act from the sale of devices that "defeat" emissions control systems in cars and trucks.

Defeating vehicle emission controls causes the release of excess air pollution including nitrogen oxides (NOx), nonmethane hydrocarbons (NMHC), carbon monoxide (CO), and particulate matter (PM) above legal limits. Excess pollution from vehicles with defeat devices harms public health and impedes efforts by EPA, tribes, states, and local agencies to plan for and attain air quality standards. Under the settlement, the defendant will pay \$3.6 million.

On this page:

- [Overview of Turn 14](#)
- [Summary of the Violations](#)
- [Summary of Environmental and Health Impacts](#)
- [Overview of the Consent Decree](#)
- [Comment Period](#)
- [Contact Information](#)

Grab raw text of details page, any PDFs and send to LLM

Settlement Resources

- [Turn 14 Consent Decree \(pdf\) \(4.94 MB\)](#)
- [Turn 14 Complaint \(pdf\) \(304.24 KB\)](#)

Overview of Turn 14

Turn 14 is one of the country's largest wholesale distributors of automotive parts and suppliers in the United States, with warehouses located in Hatfield, PA, Reno, NV, Arlington, TX, and Indianapolis, IN. Turn 14 purchases products from manufacturers and sells them to other distributors or retail outlets through its private, online, business-to-business website.

Summary of the Violations

Between January 2016 and August 2021, Turn 14 sold over 140,000 of these tampering devices throughout the United States.

Fetch details page

Respondent	Description	Order Type	Date
Turn 14 Clean Air Act Settlement Summary	On January 17, 2025, the Environmental Protection Agency (EPA) and Department of Justice reached a settlement agreement with Turn 14 Distribution, Inc. ("Turn 14"), for violations of the Clean Air Act from the sale of devices that "defeat" emissions control systems in cars and trucks.	Consent Decree and Complaint	2025-01-17

JSON for Drupal Feed

1. From web page

```
{ "documents": [ { "order_type": "Consent Decree and Complaint", "date": "2025-01-17", "link": "https://www.epa.gov/enforcement/turn-14-clean-air-act-settlement-summary", "pdf_links": "https://www.epa.gov/system/files/documents/2025-03/turn-14-dn-2-notice-of-lodging-complaint.pdf, https://www.epa.gov/system/files/documents/2025-03/turn-14-dn-1-complaint.pdf", "title": "EPA Enforcement Turn 14 Clean Air Act Settlement Summary", "raw_text": "Turn 14 Clean Air Act Settlement Summary On January 17, 2025, the Environmental Protection Agency (EPA) and Department of Justice reached a settlement agreement with Turn 14 Distribution, Inc. ("Turn 14"), for violations of the Clean Air Act from the sale of devices that \"defeat\" emissions control systems in cars and trucks. Defeating vehicle emission controls causes the release of excess air pollution including nitrogen oxides (NOx), nonmethane hydrocarbons (NMHC), carbon monoxide (CO), and particulate matter (PM) above legal limits. Excess pollution from vehicles with defeat devices harms public health and impedes efforts by EPA, tribes, states, and local agencies to plan for and attain air quality standards. Under the settlement, the defendant will pay $3.6 million.", "summary": "<p>On January 17, 2025, the <b>Environmental Protection Agency (EPA)</b> and <b>Department of Justice</b> reached a settlement agreement with <b>Turn 14 Distribution, Inc.</b>, for violations of the <b>Clean Air Act</b>. The violations stemmed from the sale of devices that \"defeat\" emissions control systems in cars and trucks. These \"defeat devices\" cause the release of excess air pollution, including nitrogen oxides (NOx), nonmethane hydrocarbons (NMHC), carbon monoxide (CO), and particulate matter (PM) above legal limits. Between January 2016 and August 2021, <b>Turn 14</b> sold over 140,000 of these tampering devices throughout the United States, violating Section 203(a)(3)(B) of the <b>Clean Air Act</b>. Car, truck, and engine manufacturers install emission controls on vehicles and engines to comply with <b>Clean Air Act</b> emission standards. <b>Turn 14</b> sold devices designed to remove or defeat emission controls from numerous vehicle models.....</p>", "penalty": "3600000.0", "environmental_issues": [ "Automobiles and Trucks" ], "laws": { "federal_law": [ { "type": "Statute", "citation": "42 U.S.C. § 7522(a)(3)(B)" }, { "type": "Rule", "citation": "28 C.F.R. § 50.7" }, { "type": "Rule", "citation": "40 C.F.R. § 86.1803-01" } ] }, "flattened_federal_laws": "Statute - 42 U.S.C. § 7522(a)(3)(B), Rule - 28 C.F.R. § 50.7, Rule - 40 C.F.R. § 86.1803-01" }, "llm": "models/gemini-2-0-flash", "time": "Sat, 08 Mar 2025 12:08:26", } ] }
```

2. From LLM: Summary, Penalty, Issues, Laws

Document content type we'll be using for the EPA Civil Actions and its associated Feed

Home > Administration > Structure > Content types > Document

Manage fields

Edit Manage fields Manage form display Manage display Manage permissions

[+ Create a new field](#) [+ Re-use an existing field](#)

“Document” content type

Label	Machine name	Field type
AI Tags	field_ai_tags	Entity reference Reference type: Taxonomy term Vocabulary: AI Tags
Body	body	Text (formatted, long, with summary)
Environmental Issues	field_environmental_issues	Entity reference Reference type: Taxonomy term Vocabulary: Environmental Issues
Feeds item	feeds_item	Feed Reference type: Feed
Laws	field_laws	Text (plain)
LLM	field_llm	Entity reference Reference type: Taxonomy term Vocabulary: LLM
Penalty	field_penalty	Number (float)
Raw Text	field_raw_text	Text (plain, long)
Source URL	field_source_url	Link

Mappings Document

Edit **Mapping** Tamper Custom sources Manage fields Manage form display Manage display

Define which elements of a single item of a feed (= Sources) map to which content pieces in Drupal (= Targets). To avoid importing duplicates, make sure that at least one definition has an *Unique target*. A unique target means that a value can only occur once. For example, only one item with the URL <http://example.com/content/1> can exist.

On *Read only* targets a value can only be set the first time. Blank sources can be used for tampering: see the [documentation](#).

See the [Mapping documentation](#) for more information.

Define which elements of a single item of a feed (= Sources) map to which content pieces in Drupal (= Targets). To avoid importing duplicates, make sure that at least one definition has an *Unique target*. A unique target means that a value can only occur once. For example, only one item with the URL <http://example.com/content/1> can exist.

On *Read only* targets a value can only be set the first time. Blank sources can be used for tampering: see the [documentation](#).

See the [Mapping documentation](#) for more information.

Context*

\$documents[*]

The base query to run. See the [Context query documentation](#) for more information.

Document Feed

Source	Target	Summary	Configure	Unique
title	Title (title)		*	<input checked="" type="checkbox"/>
summary (body)	Body (body): Text	Format: Full HTML	*	
- Select a source -	Body (body): Summary		*	
environmental_issues[*] (environmental_issues)	AI Tags (field_ai_tags)	Reference by: Name Autocreate entities: No	*	
environmental_issues[*] (environmental_issues)	Environmental Issues (field_environmental_issues)	Reference by: Name Autocreate entities: No	*	
flattened_federal_laws	Laws (field_laws)		*	<input type="checkbox"/>
llm	LLM (field_llm)	Reference by: Name Autocreate entities: No	*	

Drupal Taxonomies

Here are the sample vocabularies and terms used in the project. Note that we dynamically import the Environmental Issues tags to our Python code at runtime, add them to our Pydantic model and use them as part of our prompt to the LLM to shape its behavior ("categorize documents according to these tags, and only these tags").

Taxonomy

Taxonomy is for categorizing content. Terms are grouped into vocabularies. For example, a vocabulary called "Fruit" "Banana".

Taxonomy is for categorizing content. Terms are grouped into vocabularies. For example, a vocabulary called "Fruit" "Banana".

Sample taxonomies used for the project

Vocabulary name	Description
Environmental Issues	Topics related to the environment
AI Tags	Tags helpful for keeping track of how AI is applied to our content
LLM	Large Language Models used for generation
Agency	Federal and State government agencies
Federal Register Type	Rule, Proposed Rule or Notice

Environmental Issues

You can reorganize the terms in *Environmental Issues* using their drag-and-drop position.

Name
Automobiles and Trucks
Boats and Ships
Chemicals
Construction Equipment
Drinking Water

AI Tags

You can reorganize the terms in *AI Tags* using their drag-and-drop position.

Name
AI-Generated Categories
AI-Generated Entity Extraction
AI-Generated Text

LLM

You can reorganize the terms in *LLM* using their drag-and-drop position.

Name
claude-3-7-sonnet-latest
gpt-4o
gpt-4o-mini
models/gemini-2.0-flash

These tags are pulled in dynamically and fed to our LLM as part of our prompt

Example civil action Document import:

EPA Enforcement - Fayat Clean Air Act Settlement Summary

This article contains AI-generated summaries.

Disclaimer

This document outlines a proposed administrative settlement agreement and order on consent (Settlement) between the **U.S. Environmental Protection Agency (EPA)** and **Ashland LLC**. The Settlement addresses the release or threatened release of hazardous substances at the **Chemical Commodities, Inc. (CCI)** Superfund Site in Olathe, Kansas. **Ashland** will perform a Remedial Investigation/Feasibility Study (RI/FS) at the Site, as required by **42 U.S.C. § 9604(a)(1)**, and will reimburse the **EPA** for certain response costs. The estimated cost for **Ashland's** performance of the RI/FS is approximately \$3,000,000.

AI Summary

The Settlement includes provisions for stipulated penalties for failure to comply with the terms of the agreement. Specifically, failure to submit timely and adequate reports or plans, or failure to meet other requirements of the Settlement, may result in penalties ranging from \$500 to \$2,500 per day for the first week of violation, \$2,500 to \$5,000 per day for the second week, and \$5,000 to \$7,500 per day for each day thereafter. The Settlement also includes a covenant not to sue **Ashland** concerning the Site, subject to certain reservations and exceptions.

Environmental Issues

Hazardous Waste

AI classification

Laws

Statute - 42 U.S.C. § 9604(a)(1)

AI entity extraction

Source URL

<https://www.epa.gov/enforcement/fayat-clean-air-act-settlement-summary>

LLM Considerations

Closed models, open source – they're constantly leapfrogging one another. At present, closed source foundation models produce more **reliable structured output**. This may change in the future.

The code here has options for OpenAI's "gpt-4o", Anthropic's "claude-3-7-sonnet-latest" and Google's "gemini-2.0-flash". Currently we default to **Gemini Flash** because it is fast, reliable, inexpensive (as of March 2025: input \$0.10 / 1,000,000 input tokens, \$0.40 / 1,000,000 output tokens) and has a long context window (1,000,000 tokens) that can readily accommodate multi-hundred page PDFs.

Optimization – a custom `title_exists()` function checks to see if we already have a node in Drupal; if so, we avoid making unnecessary LLM calls. This is essential when performing web scraping at scale.

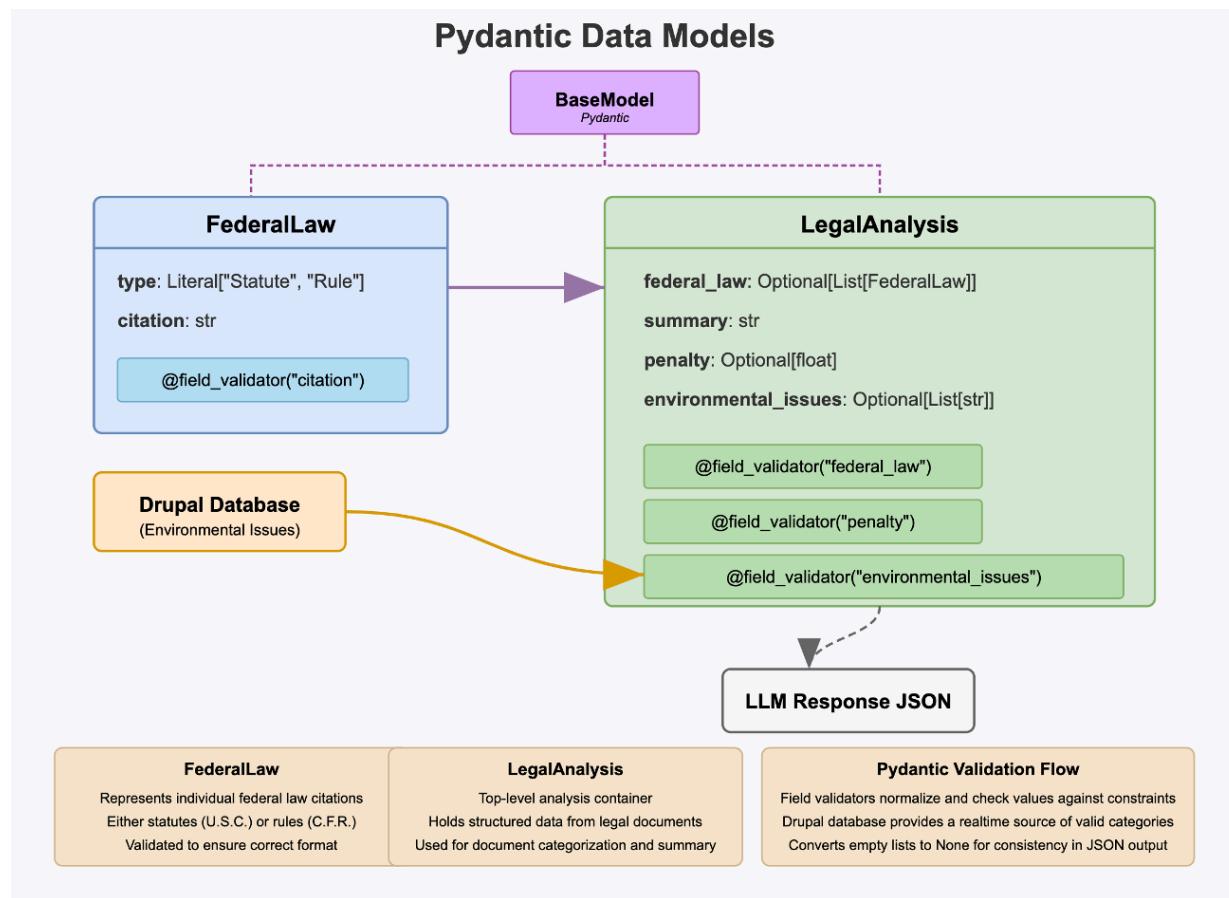
Pydantic and Instructor

Two Python libraries aid greatly in steering and validating the output of LLMs:

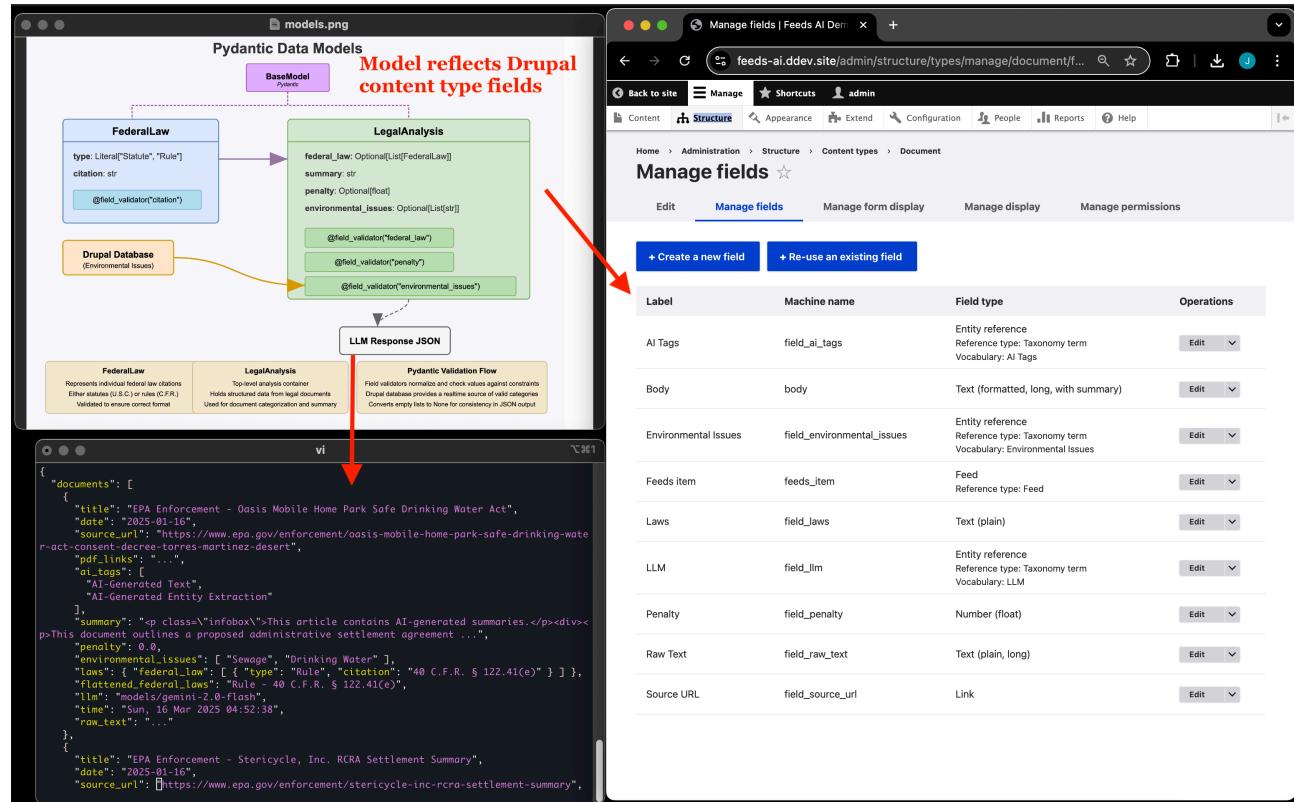
- [Pydantic](#) ensures that data meets certain criteria before being sent downstream. It's used in many contexts, not just with LLMs. For example, in validating inputs to an API or outputs to a database. A key feature for LLM usage is that Pydantic allows you to define a "data model" with rich annotations suitable for use with the Instructor library.
- [Instructor](#) helps steer LLMs toward a desired behavior. In particular, you can take a deeply annotated Pydantic model and use it to provide detailed prompt instructions to an LLM, then verify that the return data from the LLM conforms to your expectations.

These libraries help us to tell LLMs precisely how we want them to structure their reply, e.g. JSON with particular fields and types of data; and validate the data before passing it on for downstream use. In particular, we want anything GenAI-related to pass muster before importing it to Drupal. Drupal, which is all about structured content, does well by partnering with a robust framework for getting structured output from LLMs.

Here are the models used for the EPA scraper:



And an overview of how the Pydantic models, Drupal content type and LLM output JSON work together:



Pydantic Field validators

These are optional functions you can write that automatically run against fields before Pydantic signs off on the data. You can use them to simply confirm that data matches certain criteria, or you can modify the data to match your needs. The point is that this is a structured way of enforcing standards for your data.

What we're choosing to validate:

- **citation** – make sure legal citations are in a standardized format, e.g. 40 C.F.R. §§ Part 1039" should be transformed to "40 C.F.R. § 1039
- **penalty** – make sure it's a numeric float with at most 2 decimal places
- **environmental_issues** – suppose we have a Drupal taxonomy we're already using. We can fetch the terms dynamically from the Drupal DB and use them both to tell the LLM what terms we want it to look for and then to make sure that's what it did.

Miscellaneous “Best Practices”

There are few “standards” yet for how websites should manage AI-generated content. You may want to consider:

- Creating a vocabulary “AI Tags” to help you keep track of nodes to which you’ve applied GenAI. For example: AI-Generated Text, AI-Generated Categories, or AI-Generated Entity Extraction
- It is quite likely that you’ll use different LLMs over time. You may want to have a vocabulary for these as well with terms such as: claude-3-7-sonnet-latest, gpt-4o, gpt-4o-mini, gemini-2.0-flash
- If you’re going to the trouble of extracting lots of raw text for use in GenAI, even though in its raw state it may not be suitable for end users it might be good for a) fulltext search; b) future GenAI passes over the same nodes
- People are justifiably wary of what’s being pushed on them – consider including preamble / disclaimer text such as “This article contains an AI-generated summary” that fully informs people what they’re getting.

An important additional consideration is tracking Drupal LLM usage. Drupal is noteworthy precisely because of its structured approach to data, including fields, and we may want to respect (and take advantage of) that affordance of granularity when considering an LLM usage audit strategy. It's possible that Drupal isn't the right place to perform auditing and bookkeeping and that it'll be easier to use an external database. Companies in the future are likely to have external LLMOps systems that are used for more assets than just Drupal so that might be the way to go.

- In the examples presented here note that we're using LLMs at the Drupal field level - whereas the "AI Tags" as applied in these examples are assigned at the node level. Having field-level tags for tracking AI may be a bit much to do within Drupal itself and an external auditing database may be a better bet.
- MLOps is an emerging discipline. Keeping track of what content in Drupal has been "touched" by LLMs raises many questions. For any given Drupal field there are multiple things you may want to track. Given that the behavior of LLMs is complicated, varies over time and model version, the prompt, etc. here is a minimum you might consider tracking:
 - which LLM model and version did you use, e.g. OpenAI gpt-4o, Anthropic Claude 3.7, etc.
 - the actual text of the prompt you used - rather than repeat this text over and over you might use a taxonomy
 - token usage
 - date / time you ran the model
- You might also consider storing the raw text blobs from PDFs externally as well so the Drupal DB doesn't bloat unnecessarily).

Drupal demo

You can create the simple demo site that corresponds to the scraping demo per <https://github.com/crioelho/Drupal-Feeds-Using-LLMs/blob/master/feeds-ai/README.md>.

To Run the Scrapers

Note that these are barebones scrapers that only pull content from the first page of results. A fully-fledged scraping process would involve paging. That isn't the focus of this repo but the underlying code can readily be extended to accommodate paging.

1. You need Python 3.12 or higher
2. Create a new virtual env
3. Install requirements
4. Set env variables, e.g. use a .env file. You'll also need an API key for whatever LLMs you choose to you. Then set an env variable such as OPENAI_KEY, GEMINI_KEY, ANTHROPIC_KEY, etc.
5. Run scrapers

```
# Example .env file:  
DB_HOST=127.0.0.1
```



```
DB_PORT=49677
DB_USER=db
DB_PASSWORD=db
DB_NAME=db

# Set up virtual environment, activate it and install prerequisites:
python -m venv feeds_scraper
source feeds_scraper/bin/activate
pip install -r requirements.txt

# Run fr scraper for agency = EPA (see config/agencies.json file)
retrieve Federal Registers dating 3/1/25 or later and generate a news
summary.

python fr.py -a epa -d 2025-03-01 -n

# Grab first 5 rows of EPA enforcement actions
python epa.py --numrecs 5
```