

2a within 3 words must have all criteria release we seek is in for c-141 this will be turned into a UI, so layperson does not need to know how to create Solr search strings

(3W("volume of release", (barrel* OR bbl*))) AND (3W("volume of release", "volume recovered")) AND C-141&MaxResults=1000

want release label and barrels to be close together release and recovered need to be close also there may be other forms with same label

The screenshot shows a multi-step process:

- 2**: A Solr query is entered into a search bar: `(3W("volume of release", (barrel* OR bbl*))) AND (3W("volume of release", "volume recovered")) AND C-141&MaxResults=1000`.
- 3**: A PDF document is displayed, showing a form with fields for release details. A red circle highlights a specific section of the form.
- 4**: The web scraper interface shows the XPath `//div[@class='snippet']` and a function to trim the content before and after the volume labels.
- 4a**: A function definition table is shown with columns for definition, xpath, and function.
- 4b**: A rule sheet is shown with columns for definition, xpath, and function.
- 5**: A CSV table with columns `file_name` and `release_volume_2` is shown, listing various release events and their volumes.
- 6**: A text box explains that a small volume (5 bbls) is included in the list for incident tracking.
- 7**: XML snippets are shown, illustrating how the data is structured in the final output.

1. If we want to capture release volume, we ask Solr to search between the two volume (release/recovered) labels.

2. We construct a search to grab the contents between release and recovered label. In this case we also require a form of term barrels to be between the labels. We will do this for MCF also.

If we return a number (e.g. 5), we don't know if it's oil or gas. Don't want to take a chance. A null value is better than a wrong value.

2a. We will make it easy for a layperson to search for structured content. What is the left value (volume of release)? What is the right (volume recovered)? What must the results contain? (bbl* OR barrel*). Any other must-have criteria? C-141.

2b. Any user will be able to come in and via a UI grab data, e.g. the paragraph between remedial action taken and cleanup action taken. These will form narratives for touch-points (periodic alerts for executives).

4. Using the existing Scraper, we isolate each search result "snippet" returned by Lucene.

4a. We trim this snippet so all that is left is the content between the two (release and recovered) labels.

4b. Next, we build out the rule sheet. We want document name (to tie back to API) and the trimmed released volume. We want to repeat this scrape for all results (over 10K) from Lucene. We will grab all C-141 fields at one time.

5. We get structured data back in a list that can be used to populate a database table.

6. PDFs can be multiple pages, even multiple C-141 forms. These PDFs will be broken into single page PDFs for scraping.

7. Each PDF is scraped and converted to an XML file. This XML file is used to transform data into the database(s). The root element is the PDF filename. There is one element for each page (processed separately). Executive alerts can also be driven off XML.

7a. A layperson can do all the scraping for C-141 forms with no programming knowledge. All 10K+ will be scraped at one time.

7b. If a <page> element does not have all required data, it will step higher in the XML until it finds relevant data. For example, the C-141 may not have API, but an earlier page might.

8. All company names can be parsed out of XML for "touchpoint" creation. Entire narratives can be sent to operator and/or third parties along with who at company handled.

5 bbls is so small, does not merit a "story" in the push notification (email) going to Chevron.

But it can be listed as part of an incident count and if reader wants to see it, single click to PDF

```
<?xml version="1.0"?>
<pjmw1302428740_1_ao.txt>
<page 1>
<type>email</type>
<to>jamos@blm.gov; Bratcher, Mike, EMNRD</to>
<from>Belvins, Bradley</from>
<address>bradley.belvins@chevron.com</address>
<subject>Pardue Farms initial C-141</subject>
<date>2013-01-06</date>
<api>30-015-26460</api>
</page 1>
<page 2>
<type>email</type>
<to>Blevins, Bradley G</to>
<from>Amos, James <jamos@blm.gov></from>
<address>bradley.belvins@chevron.com</address>
<subject>Re: Pardue Farms CTB</subject>
<date>2013-01-07</date>
</page 2>
<page 3>
<type>form</type>
<form>C-141</form>
<company>Chevron USA</company>
<surface_owner>Nicholas Farms</surface_owner>
<facility>tank battery</facility>
<type_release>produced water & oil</type_release>
<volume_recovered>0 bbls</volume_recovered>
<volume_release>5 bbls</volume_release>
<action>Release was caused by equipment failure. Approximately 4 bbls of produced water and 1 bbl of oil were released with no recovery. An Emergency Response Team arrived at the release area and began continuous abatement of the impacted area.</action>
etc...
</page 3>
<page 4>
etc...
```