# Split and Gather Example: Analyzing Trump Immunity Case

This example demonstrates how to use the Split and Gather operations in DocETL to process and analyze a large legal document, specifically the government's motion for immunity determinations in the case against former President Donald Trump. You can download the dataset we'll be using here. This dataset contains a single document.

## Problem Statement

We want to analyze a lengthy legal document to identify all people involved in the Trump vs. United States case regarding presidential immunity. The document is too long to process in a single operation, so we need to split it into manageable chunks and then gather context to ensure each chunk can be analyzed effectively.

## Chunking Strategy

When dealing with long documents, it's often necessary to break them down into smaller, manageable pieces. This is where the Split and Gather operations come in handy:

1. **Split Operation**: This divides the document into smaller chunks based on token count or delimiters. For legal documents, using a token count method is often preferable to ensure consistent chunk sizes.

2. **Gather Operation**: After splitting, we use the Gather operation to add context to each chunk. This operation can include content from surrounding chunks, as well as document-level metadata and headers, ensuring that each piece maintains necessary context for accurate analysis.

> ✏️ **Pipeline Overview**
>
> Our pipeline will follow these steps:
>
> 1. Extract metadata from the full document
>
> 2. Split the document into chunks
>
> 3. Extract headers from each chunk
>
> 4. Gather context for each chunk
>
> 5. Analyze each chunk to identify people and their involvements in the case
>
> 6. Reduce the results to compile a comprehensive list of people and their roles

## Example Pipeline and Output

Here's a breakdown of the pipeline defined in trump-immunity_opt.yaml:

1. **Dataset Definition**: We define a dataset (json file) with a single document.

2. **Metadata Extraction**: We define a map operation to extract any document-level metadata that we want to pass to each chunk being processed.

3. **Split Operation**: The document is split into chunks using the following configuration:

```
- name: split_find_people_and_involvements
  type: split
  method: token_count
  method_kwargs:
    num_tokens: 3993
  split_key: extracted_text
```

This operation splits the document into chunks of approximately 3993 tokens each. This size is chosen to balance between maintaining context and staying within model token limits. `split_key` should be the field in the document that contains the text to split.

4. **Header Extraction**: We define a map operation to extract headers from each chunk. Then, when rendering each chunk, we can also render the headers in levels above the headers in the chunk--ensuring that we can maintain hierarchical context, even when the headers are in other chunks.

5. **Gather Operation**: Context is gathered for each chunk using the following configuration:

```
- name: gather_extracted_text_find_people_and_involvements
  type: gather
  content_key: extracted_text_chunk  ①
  doc_header_key: headers  ②
```

```
    doc_id_key: split_find_people_and_involvements_id  ③
    order_key: split_find_people_and_involvements_chunk_num  ④
    peripheral_chunks:
      next:
        head:
          count: 1
      previous:
        tail:
          count: 1
```

① The field containing the chunk content; the split_key with "_chunk" appended. Automatically exists as a result of the split operation. **This is required.**

② The field containing the extracted headers for each chunk. Only exists if you have a header extraction map operation. **This can be omitted if you don't have headers extracted for each chunk.**

③ The unique identifier for each document; the split operation name with "_id" appended. Automatically exists as a result of the split operation. **This is required.**

④ The field indicating the order of chunks; the split operation name with "_chunk_num" appended. Automatically exists as a result of the split operation. **This is required.**

This operation gathers context for each chunk, including the previous chunk, the current chunk, and the next chunk. We also render the headers populated by the previous operation.

6. **Chunk Analysis**: We define a map operation to analyze each chunk.

7. **Result Reduction**: We define a reduce operation to reduce the results of the map operation (applied to each chunk) to a single list of people and their involvements in the case.

Here is the full pipeline configuration, with the split and gather operations highlighted. Assuming the sample dataset looks like this:

```
[
  {
    "pdf_url":
"https://storage.courtlistener.com/recap/gov.uscourts.dcd.258148/gov.uscourts.
dcd.258148.252.0.pdf"
  }
]
```

## Full Pipeline Configuration

```yaml
1    datasets:
2      legal_doc:
3        type: file
4        path: /path/to/your/dataset.json
5        parsing: ❶
6          - function: azure_di_read
7            input_key: pdf_url
8            output_key: extracted_text
9            function_kwargs:
10               use_url: true
11               include_line_numbers: true
12
13   default_model: gpt-4o-mini
14
15   system_prompt:
16     dataset_description: the Trump vs. United States case
17     persona: a legal analyst
18
19   operations:
20     - name: extract_metadata_find_people_and_involvements
21       type: map
22       model: gpt-4o-mini
23       prompt: |
24         Given the document excerpt: {{ input.extracted_text }}
25         Extract all the people mentioned and summarize their involvements
26   in the case described.
27       output:
28         schema:
29           metadata: str
30
31     - name: split_find_people_and_involvements
32       type: split
33       method: token_count
34       method_kwargs:
35         num_tokens: 3993
36       split_key: extracted_text
37
38     - name: header_extraction_extracted_text_find_people_and_involvements
39       type: map
40       model: gpt-4o-mini
41       output:
42         schema:
43           headers: "list[{header: string, level: integer}]"
44       prompt: |
45         Analyze the following chunk of a document and extract any headers
46   you see.
47
48         { input.extracted_text_chunk }
49
50         Examples of headers and their levels based on the document
51   structure:
52         - "GOVERNMENT'S MOTION FOR IMMUNITY DETERMINATIONS" (level 1)
53         - "Legal Framework" (level 1)
54         - "Section I" (level 2)
55         - "Section II" (level 2)
56         - "Section III" (level 2)
57         - "A. Formation of the Conspiracies" (level 3)
```

```yaml
58          - "B. The Defendant Knew that His Claims of Outcome-Determinative
59    Fraud Were False" (level 3)
60            - "1. Arizona" (level 4)
61            - "2. Georgia" (level 4)
62
63      - name: gather_extracted_text_find_people_and_involvements
64        type: gather
65        content_key: extracted_text_chunk
66        doc_header_key: headers
67        doc_id_key: split_find_people_and_involvements_id
68        order_key: split_find_people_and_involvements_chunk_num
69        peripheral_chunks:
70          next:
71            head:
72              count: 1
73          previous:
74            tail:
75              count: 1
76
77      - name: submap_find_people_and_involvements
78        type: map
79        model: gpt-4o-mini
80        output:
81          schema:
82            people_and_involvements: list[str]
83        prompt: |
84          Given the document excerpt: {{ input.extracted_text_chunk_rendered
85    }}
86          Extract all the people mentioned and summarize their involvements
87    in the case described. Only process the main chunk.
88
89      - name: subreduce_find_people_and_involvements
90        type: reduce
91        model: gpt-4o-mini
92        associative: true
93        pass_through: true
94        synthesize_resolve: false
95        output:
96          schema:
97            people_and_involvements: list[str]
98        reduce_key:
99          - split_find_people_and_involvements_id
100       prompt: |
101         Given the following extracted information about individuals
102   involved in the case, compile a comprehensive list of people and their
103   specific involvements in the case:
104
105         {% for chunk in inputs %}
106         {% for involvement in chunk.people_and_involvements %}
107         - {{ involvement }}
108         {% endfor %}
109         {% endfor %}
110
111         Make sure to include all the people and their involvements. If a
112   person has multiple involvements, group them together.
113
114   pipeline:
115     steps:
116       - name: analyze_document
117         input: legal_doc
118         operations:
```

```
119              - extract_metadata_find_people_and_involvements
120              - split_find_people_and_involvements
                 - header_extraction_extracted_text_find_people_and_involvements
                 - gather_extracted_text_find_people_and_involvements
                 - submap_find_people_and_involvements
                 - subreduce_find_people_and_involvements

          output:
             type: file
             path: /path/to/your/output/people_and_involvements.json
             intermediate_dir: /path/to/your/intermediates
```

> **1** This is an example parsing function, as explained in the Parsing docs. You can define your own parsing function to extract the text you want to split, or just have the text be directly in the json file.

Running the pipeline with `docetl run pipeline.yaml` will execute the pipeline and save the output to the path specified in the output section. It cost $0.05 and took 23.8 seconds with gpt-4o-mini.

Here's a table with one column listing all the people mentioned in the case and their involvements:

**Final Output**                                                                                    ⌄

**People Involved in the Case and Their Involvements**

DONALD J. TRUMP: Defendant accused of orchestrating a criminal scheme to overturn the 2020 presidential election results through deceit and collaboration with private co-conspirators; charged with leading conspiracies to overturn the 2020 presidential election; made numerous claims of election fraud and pressured officials to find votes to overturn the election results; incited a crowd to march to the Capitol; communicated with various officials regarding election outcomes; exerted political pressure on Vice President Pence; publicly attacked fellow party members for not supporting his claims; involved in spreading false claims about the election, including through Twitter; pressured state legislatures to take unlawful actions regarding electors; influenced campaign decisions and narrative regarding the election results; called for action to overturn the certified results and demanded compliance from officials; worked with co-conspirators on efforts to promote fraudulent elector plans and led actions that culminated in the Capitol riot.

MICHAEL R. PENCE: Vice President at the time, pressured by Trump to obstruct Congress's certification of the election; informed Trump there was no evidence of significant fraud; encouraged Trump to accept election results; involved in discussions with Trump regarding election challenges and strategies; publicly asserted his constitutional limitations in the face of Trump's pressure; became the target of attacks from Trump and the Capitol rioters; sought to distance himself from Trump's efforts to overturn the election.

CC1: Private attorney who Trump enlisted to falsely claim victory and perpetuate fraud allegations; participated in efforts to influence political actions in targeted states; suggested the defendant declare victory despite ongoing counting; actively involved in making false fraud claims regarding the election; pressured state officials; spread false claims about election irregularities and raised threats against election workers; coordinated fraudulent elector meetings and misrepresented legal bases.

CC2: Mentioned as a private co-conspirator involved in the efforts to invalidate election results; proposed illegal strategies to influence the election certification; urged others to decertify legitimate electors; involved in discussions influencing state officials; pressured Mike Pence to act against certification; experienced disappointment with Pence's rejection of proposed strategies; presented unlawful plans to key figures.

CC3: Another private co-conspirator involved in scheming to undermine legitimate vote counts; promoted false claims during public hearings and made remarks

**People Involved in the Case and Their Involvements**

inciting fraud allegations; encouraged fraudulent election lawsuits and made claims about voting machines; pressured other officials regarding claims of election fraud.

CC5: Private political operative who collaborated in the conspiracy; worked on coordinating actions related to the fraudulent elector plan; engaged in text discussions regarding the electors and strategized about the fraud claims.

CC6: Private political advisor providing strategic guidance to Trump's re-election efforts; involved in communications with campaign staff regarding the electoral vote processes.

P1: Private political advisor who assisted with Trump's re-election campaign; advocated declaring victory before final counts; maintained a podcast spreading false claims about the election.

P2: Trump's Campaign Manager, providing campaign direction during the election aftermath; informed the defendant regarding false claims related to state actions.

P3: Deputy Campaign Manager, involved in assessing election outcomes; coordinated with team members discussing legal strategies post-election; marked by frequent contact with Trump regarding campaign operations.

P4: Senior Campaign Advisor, part of the team advising Trump on election outcome communication; expressed skepticism about allegations of fraud; contradicted Trump's claims about deceased voters in Georgia.

P5: Campaign operative and co-conspirator, instructed to create chaos during vote counting and incited unrest at polling places; engaged in discussions about the elector plan.

P6: Private citizen campaign advisor who provided early warnings regarding the election outcome; engaged in discussions about the validity of allegations.

P7: White House staffer and campaign volunteer who advised Trump on potential election challenges and outcomes; acted as a conduit between Trump and various officials; communicated political advice relevant to the election.

P8: Staff member of Pence, who communicated about the electoral process and advised against Trump's unlawful plans; was involved in discussions of political strategy surrounding election results.

**People Involved in the Case and Their Involvements**

P9: White House staffer who became a link between Trump and campaign efforts regarding fraud claims; provided truthful assessments of the situation; facilitated communications during post-election fraud discussions.

P12: Attended non-official legislative hearings; involved in spreading disinformation about election irregularities.

P15: Assistant to the President who overheard Trump's private comments about fighting to remain in power after the 2020 election; involved in discussions about various election-related strategies.

P16: Governor of Arizona; received calls from Trump regarding election fraud claims and the count in Arizona.

P18: Speaker of the Arizona State House contacted as part of efforts to challenge election outcomes; also expressed reservations about Trump's strategies.

P21: Chief of Staff who exchanged communications about the fraudulent allegations; facilitated discussions and logistics during meetings.

P22: Campaign attorney who verified that claims about deceased voters were false; participated in discussions around the integrity of the election results.

P26: Georgia Attorney General contacted regarding fraud claims; openly stated there was no substantive evidence to support fraud allegations; discussed Texas v. Pennsylvania lawsuit with Trump.

P33: Georgia Secretary of State; defended election integrity publicly; stated rumors of election fraud were false; involved in discussions about the impact of fraudulent elector claims in Georgia.

P39: RNC Chairwoman; advised against lobbying with state legislators; coordinated with Trump on fraudulent elector efforts; refused to promote inaccurate reports regarding election fraud.

P47: Philadelphia City Commissioner; stated there was no evidence of widespread fraud; targeted by Trump for criticism after his public statements.

**People Involved in the Case and Their Involvements**

P52: Attorney General who publicly stated that there was no evidence of fraud that would affect election results; faced pressure from Trump's narrative.

P50: CISA Director; publicly declared the election secure; faced backlash after contradicting Trump's claims about election fraud.

P53: Various Republican U.S. Senators participated in rallies organized by Trump; linked to his campaign efforts regarding the election process.

P54: Campaign staff member involved in strategizing about elector votes; discussed procedures and expectations surrounding election tasks and claims.

P57: Former U.S. Representative who opted out of the fraudulent elector plan in Pennsylvania; cited legal concerns about the actions being proposed.

P58: A staff member of Pence involved in communications directing Pence regarding official duties, managing conversations surrounding election processes.

P59: Community organizers who were engaged in discussions relating to Trump's electoral undertakings.

P60: Individual responses to Trump's directives aimed at influencing ongoing election outcomes and legislative actions.

## Optional: Compiling a Pipeline into a Split-Gather Pipeline

You can also compile a pipeline into a split-gather pipeline using the `docetl build` command. Say we had a much simpler pipeline for the same document analysis task as above, with just one map operation to extract people and their involvements.

```
default_model: gpt-4o-mini

datasets:
  legal_doc:  ①
    path: /path/to/dataset.json
    type: file
    parsing:  ②
      - input_key: pdf_url
        function: azure_di_read
        output_key: extracted_text
        function_kwargs:
```

```yaml
          use_url: true
          include_line_numbers: true

operations:
  - name: find_people_and_involvements
    type: map
    optimize: true
    prompt: |
      Given this document, extract all the people and their involvements in
the case described by the document.

      {{ input.extracted_text }}

      Return a list of people and their involvements in the case.
    output:
      schema:
        people_and_involvements: list[str]

pipeline:
  steps:
    - name: analyze_document
      input: legal_doc
      operations:
        - find_people_and_involvements

  output:
    type: file
    path: "/path/to/output/people_and_involvements.json"
```

( 1 )  This is an example parsing function, as explained in the Parsing docs. You can define
your own parsing function to extract the text you want to split, or just have the text
be directly in the json file. If you want the text directly in the json file, you can have
your json be a list of objects with a single field "extracted_text".

( 2 )  You can remove this parsing section if you don't need to parse the document (i.e., if
the text is already in the json file in the "extracted_text" field in the object).

In the pipeline above, we don't have any split or gather operations. Running `docetl build pipeline.yaml [--model=gpt-4o-mini]` will output a new pipeline_opt.yaml file with the split and gather operations highlighted--like we had defined in the previous example. Note that this cost us $20 to compile, since we tried a bunch of different plans...