

# TOPIC MODELING WITH NATURAL LANGUAGE PROCESSING FOR IDENTIFICATION OF NUCLEAR PROLIFERATION-RELEVANT SCIENTIFIC AND TECHNICAL PUBLICATIONS



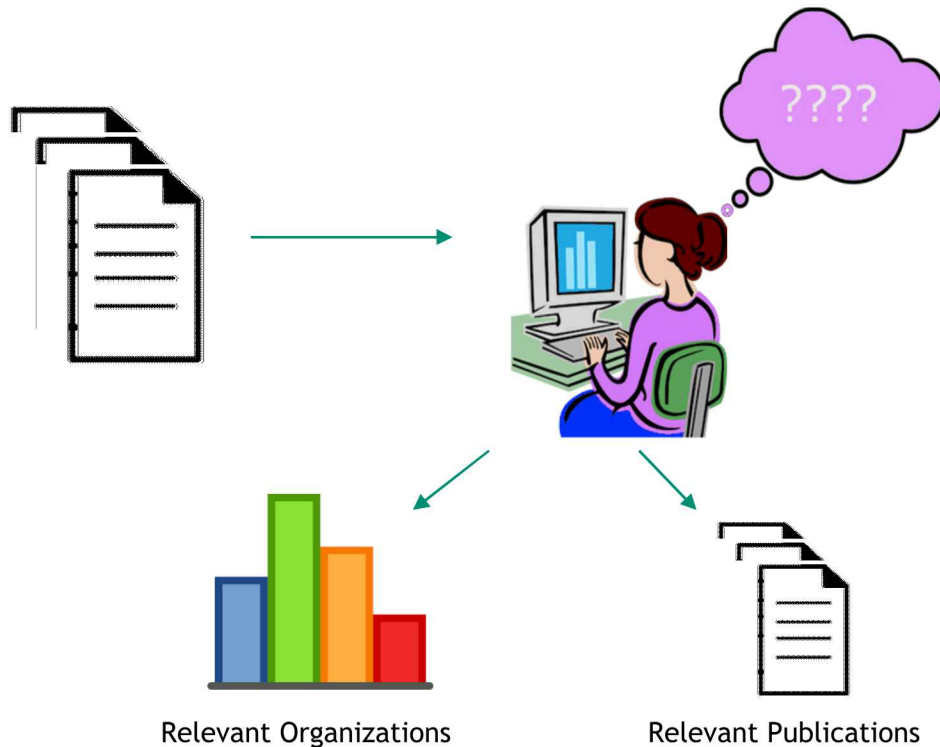
*INMM Annual Meeting 2020*

Presented by Jon Bisila

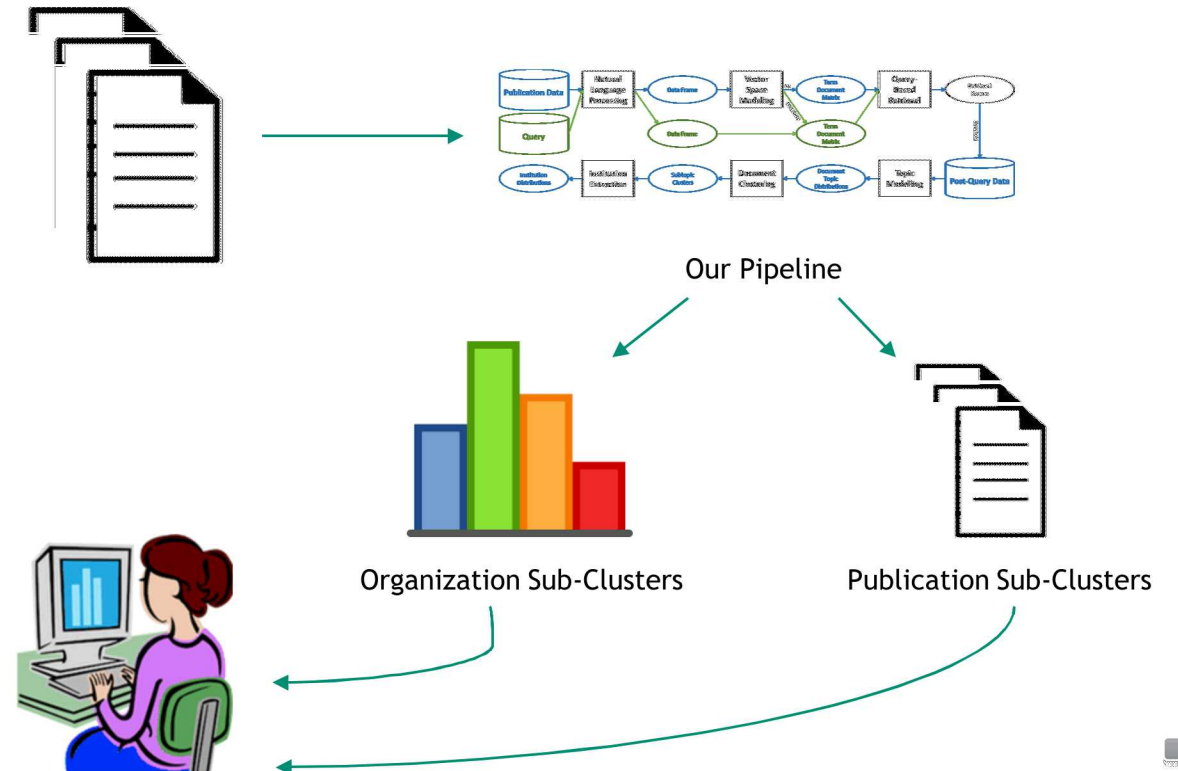
Team: Jon Bisila, Zoe Gastelum, Danny Dunlavy, Craig Ulmer

## How can we help analysts more efficiently find what they are looking for?

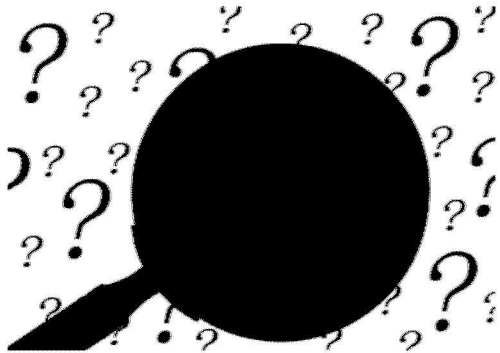
### Current: Manual Triage



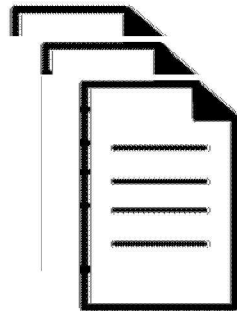
### Proposed: Assisted Triage



- Search strings used by analysts to find relevant information have limitations
- Sifting out relevant documents still requires manual triage
- Efficiently identifying organizations involved can be difficult



Effective Search Strings



Relevant Publications



Relevant Organizations

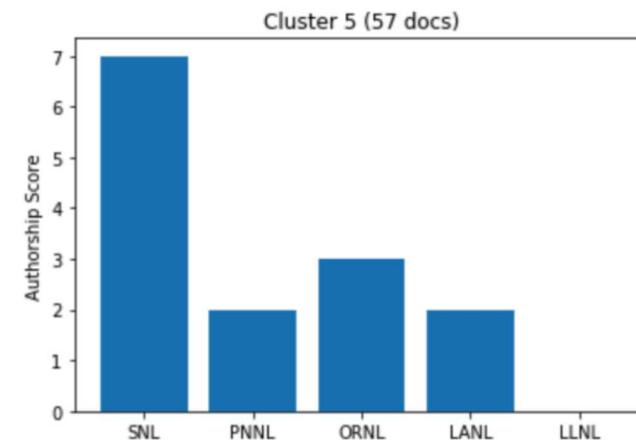


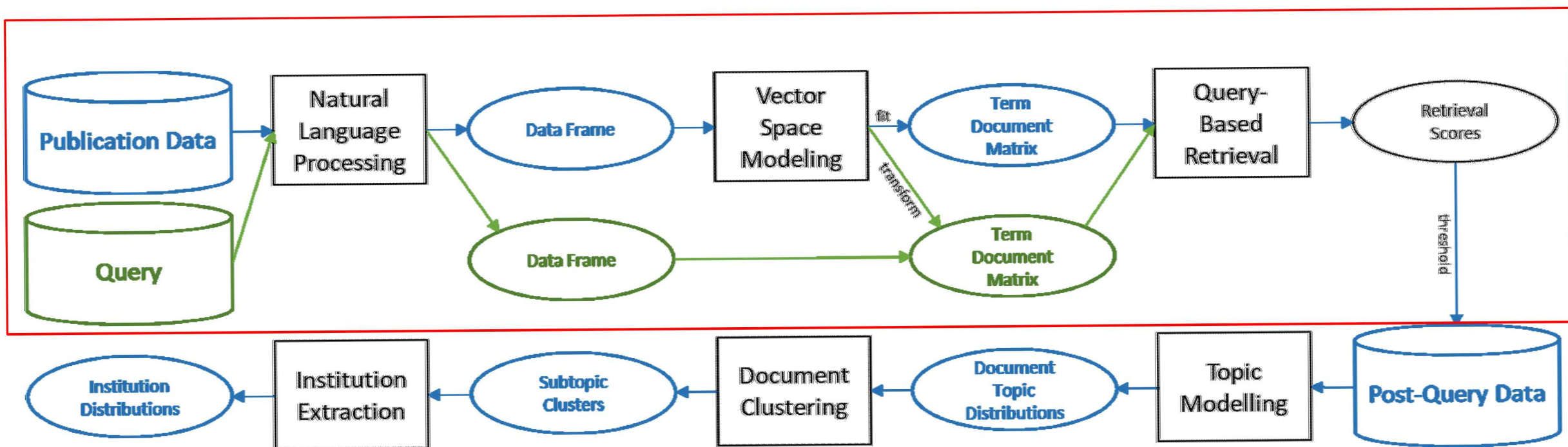
- Build from existing work to augment existing analyst search process
- Utilize topic modelling and document clustering to assist triage
- Present distributions of organizations in a more user-friendly manner

## Topic Modelling

`0.058*"compon" + 0.047*"contain" + 0.041*"polym" + 0.035*"solvent"`  
`0.042*"layer" + 0.039*"shell" + 0.031*"flow" + 0.028*"electrod"`  
`0.026*"modul" + 0.023*"optic" + 0.022*"system" + 0.018*"devic"`

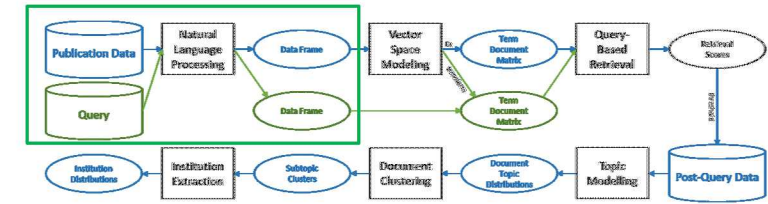
## Distribution of Organizations



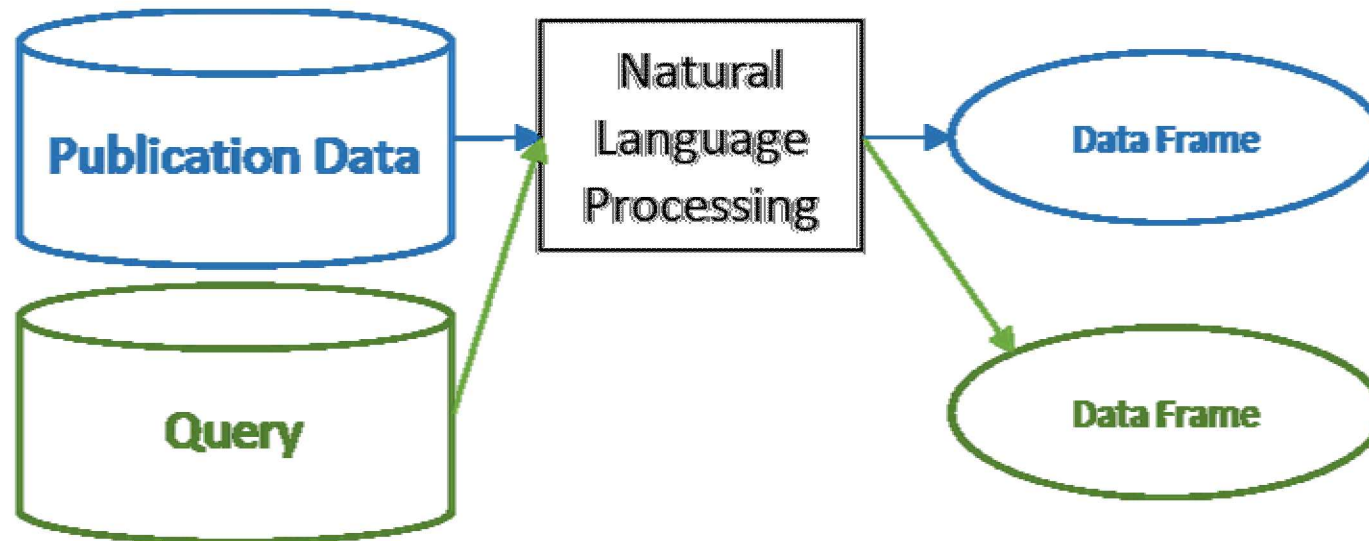




# Methods: Natural Language Processing



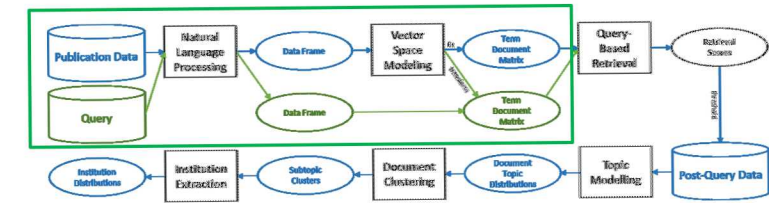
- Example transformations\* of Publication Data:
  - Using word roots instead of words: *proliferate, proliferation, proliferates...* --> *proliferat*
  - Remove words that do not help discriminate topics: *a, an, “the”, “in”, etc.*



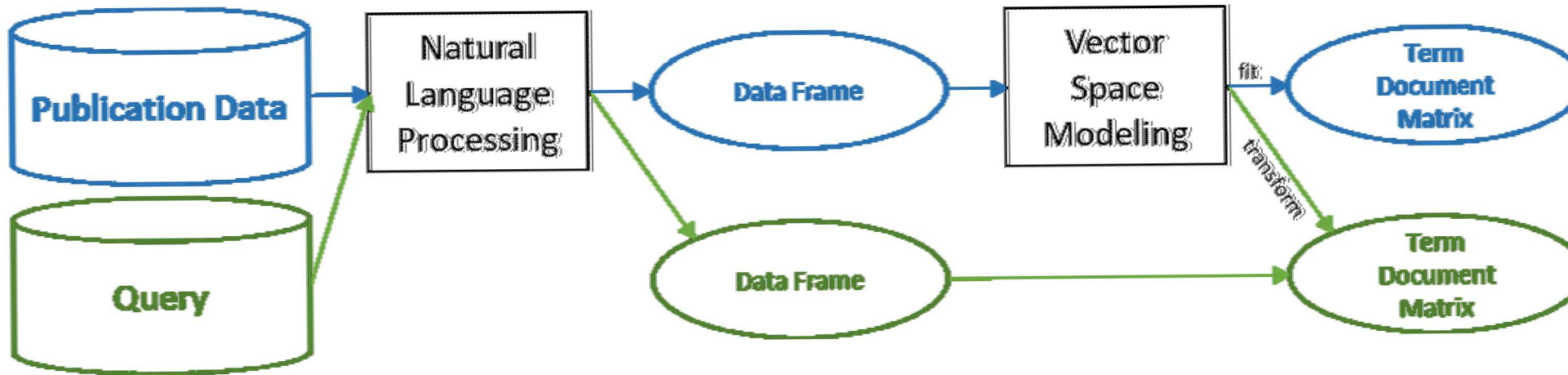
\* Note: Details regarding all natural language processing (NLP) tasks performed are in the full paper



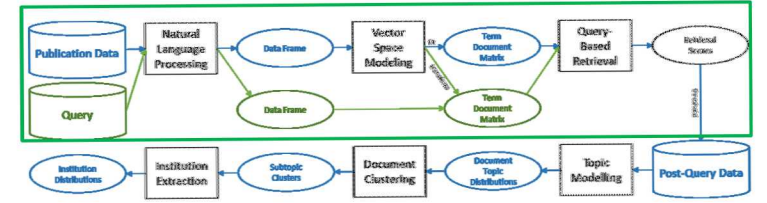
# Methods: Vector Space Modeling



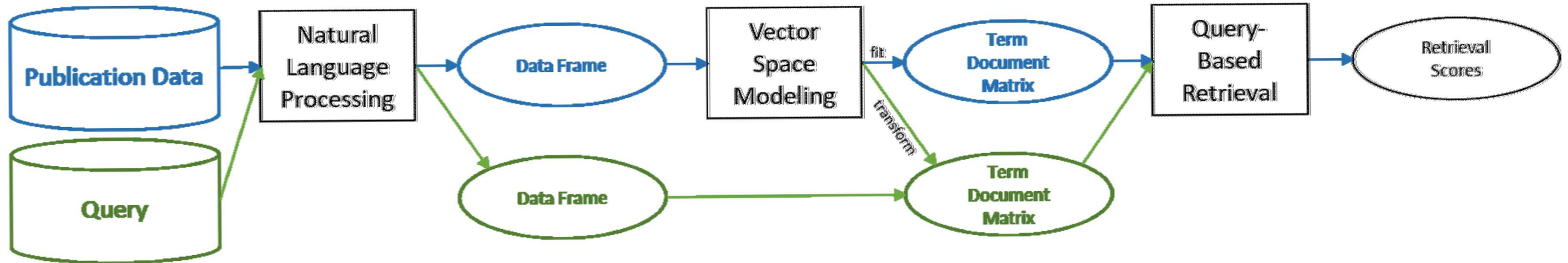
- Vector Space Models (VSMs) provide another interpretation of the data
- *Term-document matrix*: columns represent terms, rows represent documents



# Methods: Query-Based Information Retrieval

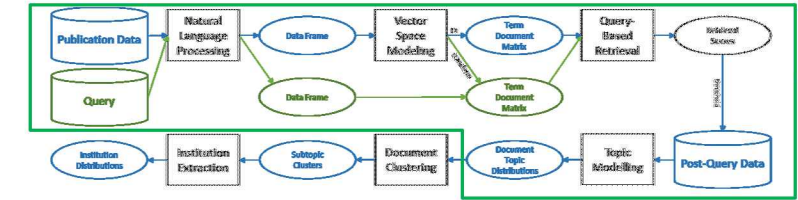


- Simulate Boolean search efficiently using single matrix-vector product
- Provide a score for each document in relation to the query

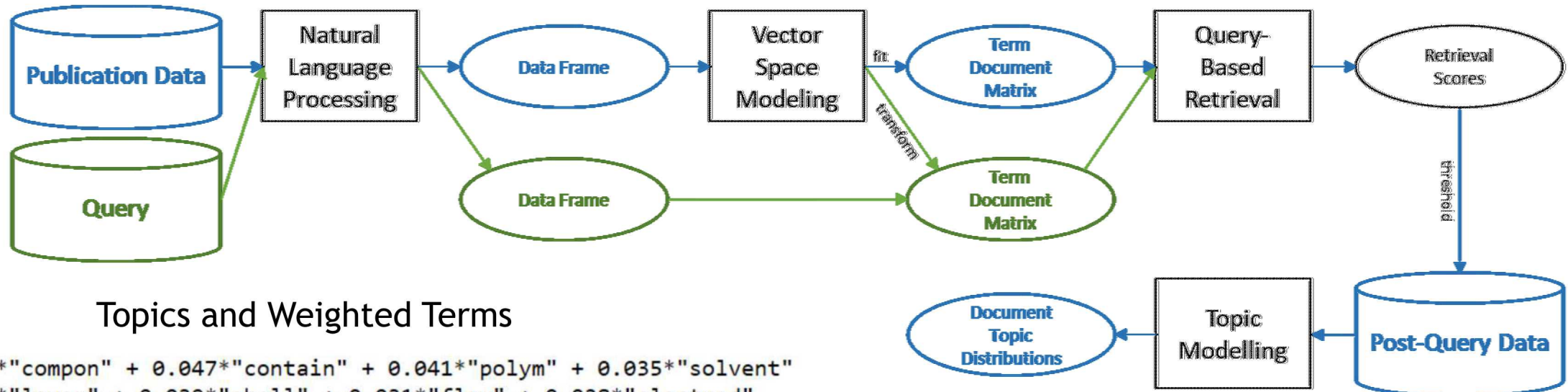




# Methods: Topic Modeling



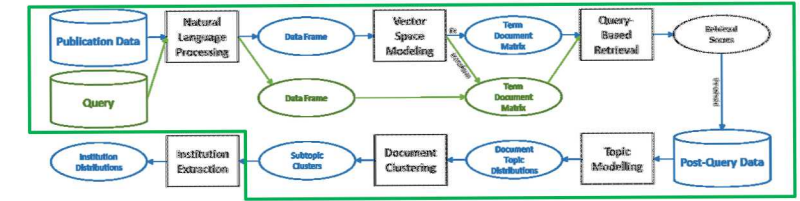
- Build topic models using Latent Dirichlet Allocation (LDA)
- Documents represented as distributions of topics
- Topics composed of distributions of terms



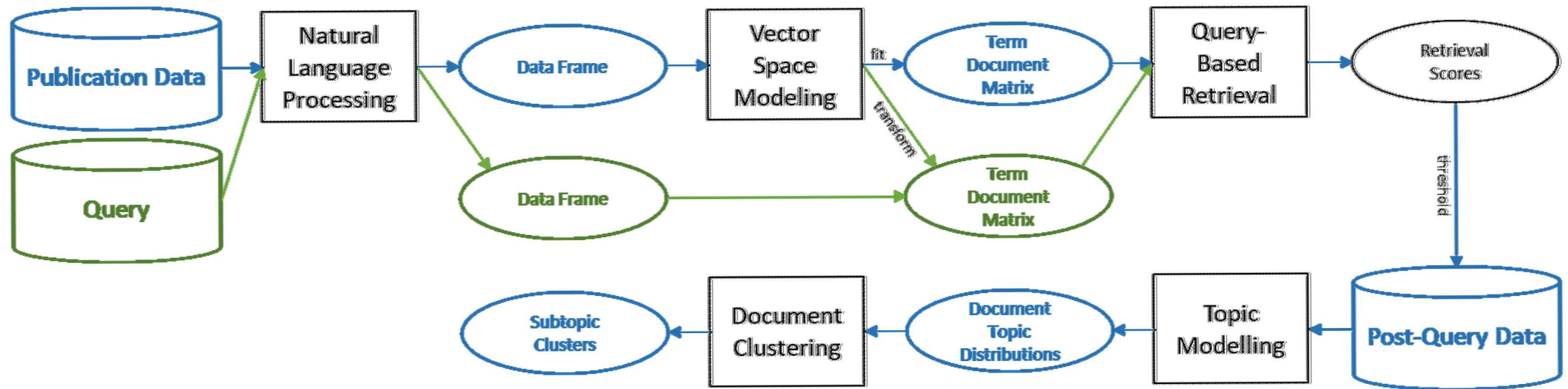
0.058\*"compon" + 0.047\*"contain" + 0.041\*"polym" + 0.035\*"solvent"  
 0.042\*"layer" + 0.039\*"shell" + 0.031\*"flow" + 0.028\*"electrod"  
 0.026\*"modul" + 0.023\*"optic" + 0.022\*"system" + 0.018\*"devic"



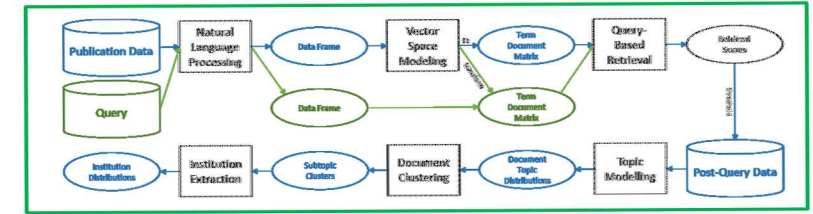
# Methods: Document Clustering



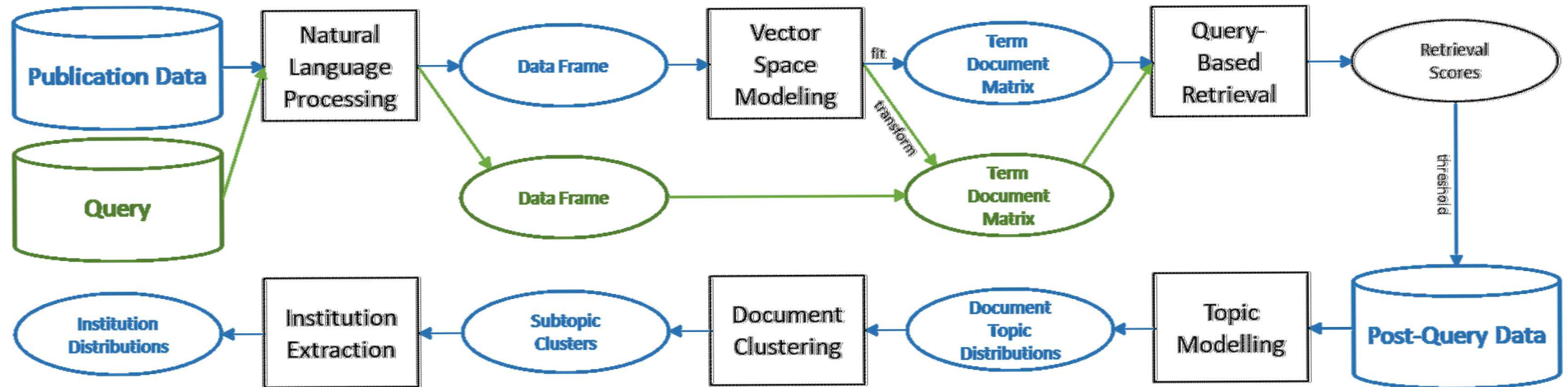
- Cluster documents using Louvain community detection algorithm
- Output: rank-ordered list of documents that best summarize a cluster
- Also provide rank-ordered list of summary topics and terms



# Methods: Institution Extraction



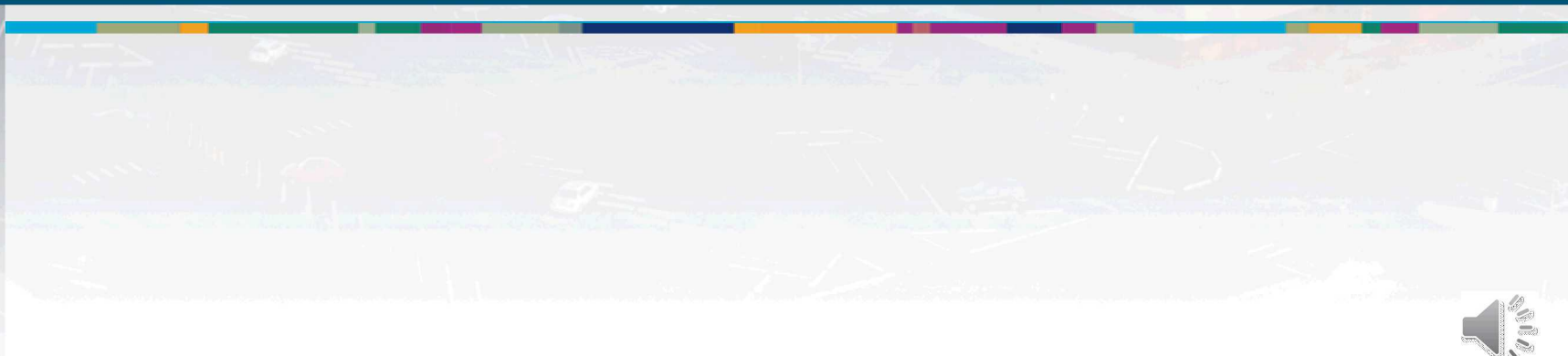
- Extract the institutions associated with the documents in each cluster
- Provide histogram of the distribution of document counts per institution







# Case Study: U.S. Department of Energy Publications



## Case Study: U.S. Department of Energy Publications

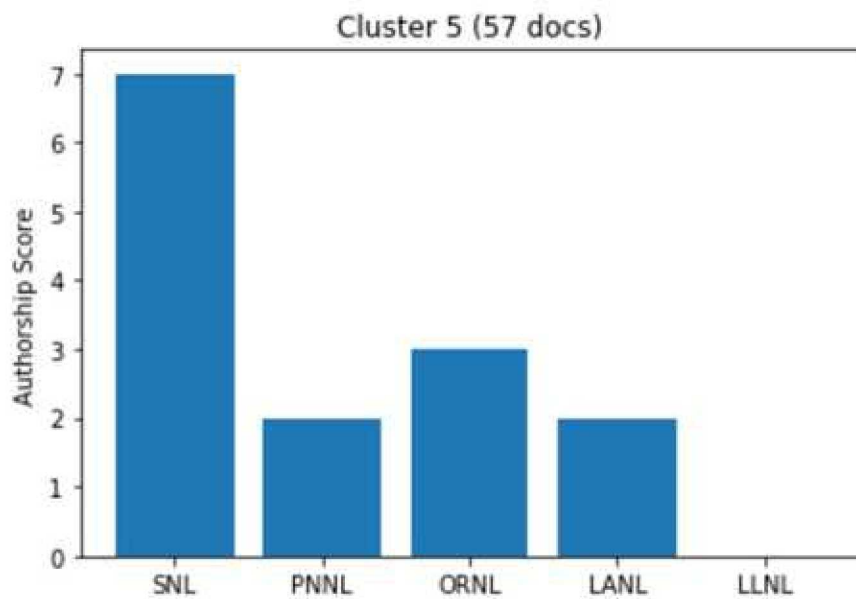
- Generated an expert-derived search string related to nuclear fuel reprocessing

alloy, anion, arms, boxes, bromine, cadmium, calcination, cell, ceric, chloride, chopping, concrete, cuts, decladding, density, electrolysis, exchange, extraction, fluoride, fluoridized, fuel, fuels, glass, glove, handling, hardened, hexone, high, inconel, irradiated, isobutyl, ketone, lead, lithium, manipulator, master, metal, methyl, MIBK, molten, nitrate, nuclear, oxidation, plutonium, precipitation, purification, radiation, redox, reduction, remote, reprocess, rods, salts, separation, shield, slave, solvent, spent, trifluoride, uranium, volatility, window

- Data:
  - 2,000 Office of Science and Technical Information (OSTI) documents from 2017
  - Seven relevant “tracer” documents, identified by subject matter experts
- Main question to answer using our approach:
  - Do these documents cluster together?







#### Cluster 5 topics and terms

0.058\*"compon" + 0.047\*"contain" + 0.041\*"polym" + 0.035\*"solvent"  
 0.042\*"layer" + 0.039\*"shell" + 0.031\*"flow" + 0.028\*"electrod"  
 0.026\*"modul" + 0.023\*"optic" + 0.022\*"system" + 0.018\*"devic"

#### Cluster 5 titles

Production of alpha-hydroxy carboxylic acids and esters from higher sugars using tandem catalyst systems  
 Electrochemical Nucleation and Growth of Uranium and Plutonium from Molten Salts  
 Variants of polypeptides having cellulolytic enhancing activity and polynucleotides encoding same

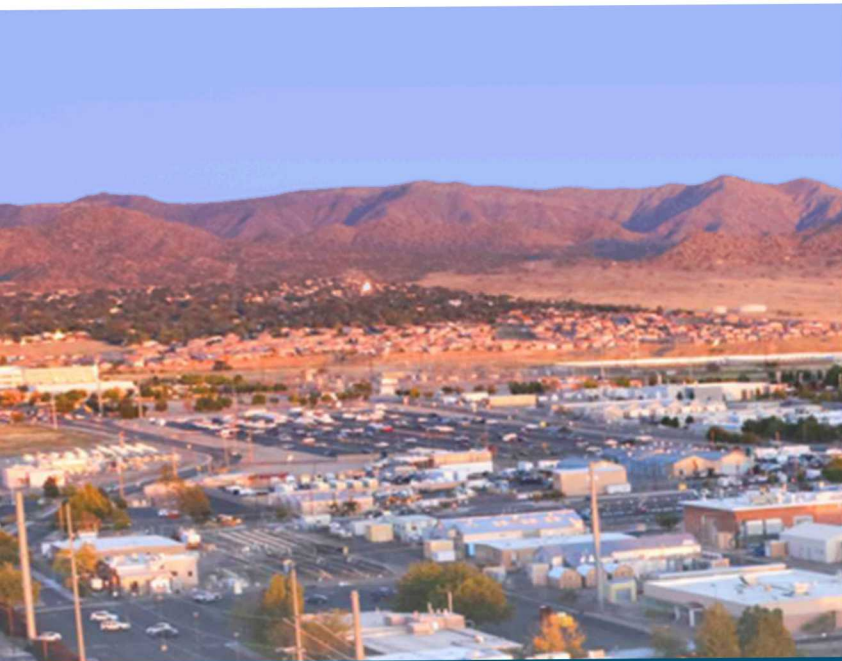
- Of the seven tracer documents, five of them appear in the same cluster
- Benefit of our approach:
  - Analysts can quickly scan the topics, terms, and titles to determine relevancy of each cluster



# Future Work

- Uncertainty Quantification: how certain are we about the stability and reliability of our choices of algorithms?
- Different vector space models
- Alternative document retrieval methods





Questions?

Jon Bisila  
jbisila@sandia.gov

