



Machine Learning Demo Wednesday

Prioritizing Inspections using ML

Alec Mishkin, Guillermo Vasquez, Stuti Polra, Casey Friedman, Scott Pringle, Theresa Smith

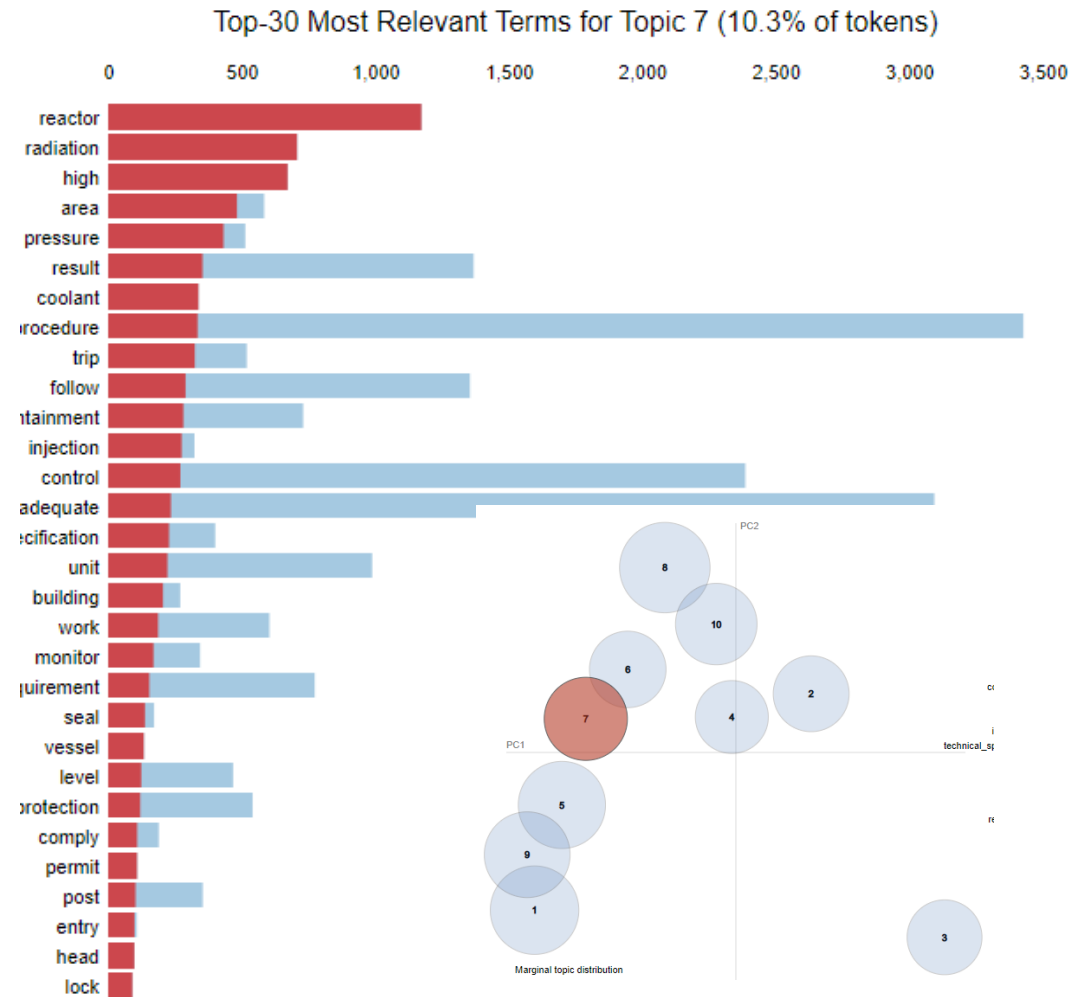
Wednesday, March 22, 2023

Agenda

- Clustering
- Data Analysis
 - Inspection Report Characterization
- Topic Modeling
- Tool Analysis
- Progress and Next Steps

Clustering

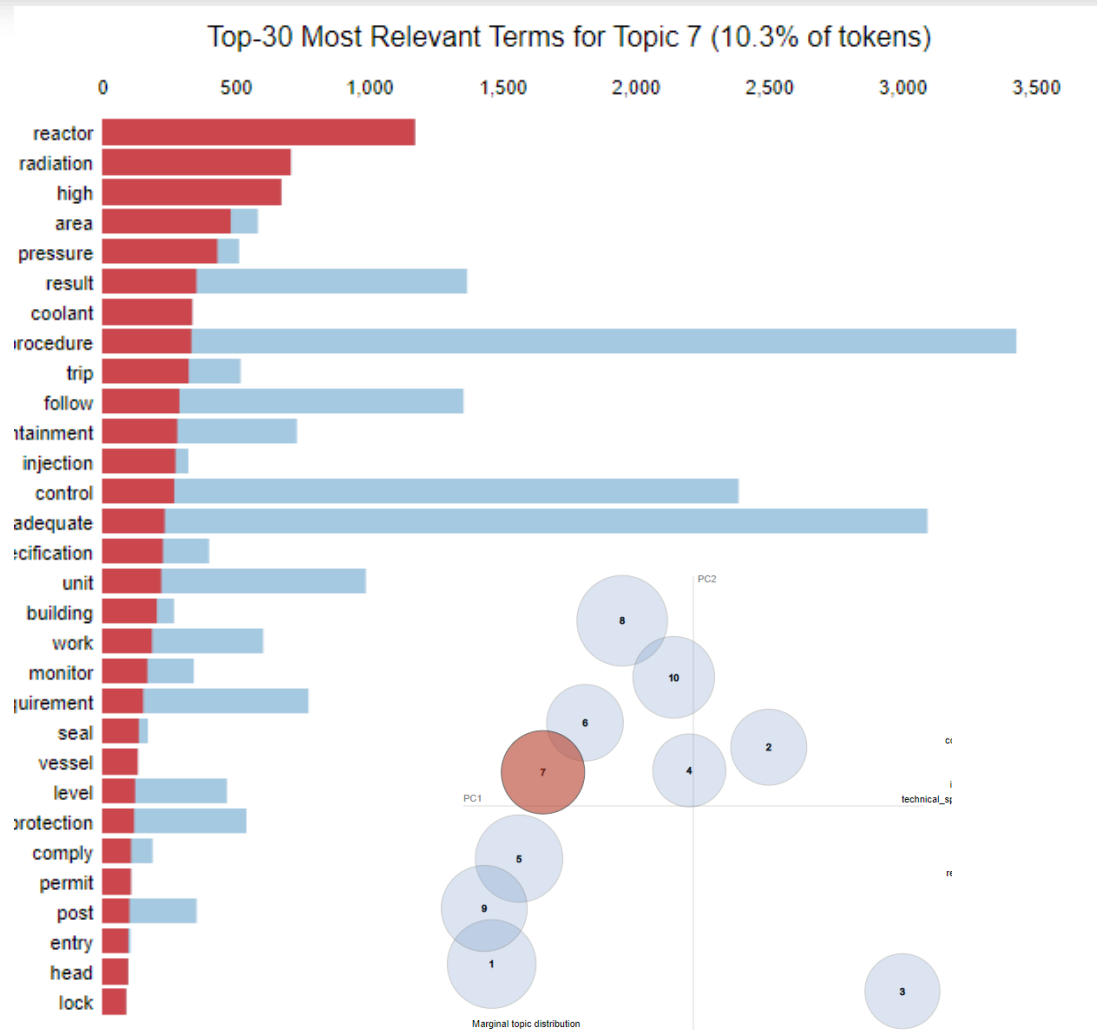
Clustering: Topic Modeling can lead to Safety Clusters



Topic Modeling

- Generic groupings based on words that frequently occur together in text
- The words at the left occurred together more often than they occurred with other words
- Once clusters of words (topics) are identified, documents can be aligned with those topics.
- A document may align with more than one topic – that is, more than one topic may have been discussed within the document.
- Topics can be characterized (given a name) by analyzing the words that define them.
- The topics discovered in our early experiments appear to align well with **Safety Clusters**.

Clustering samples – Potential Safety Clusters



Safety Clusters

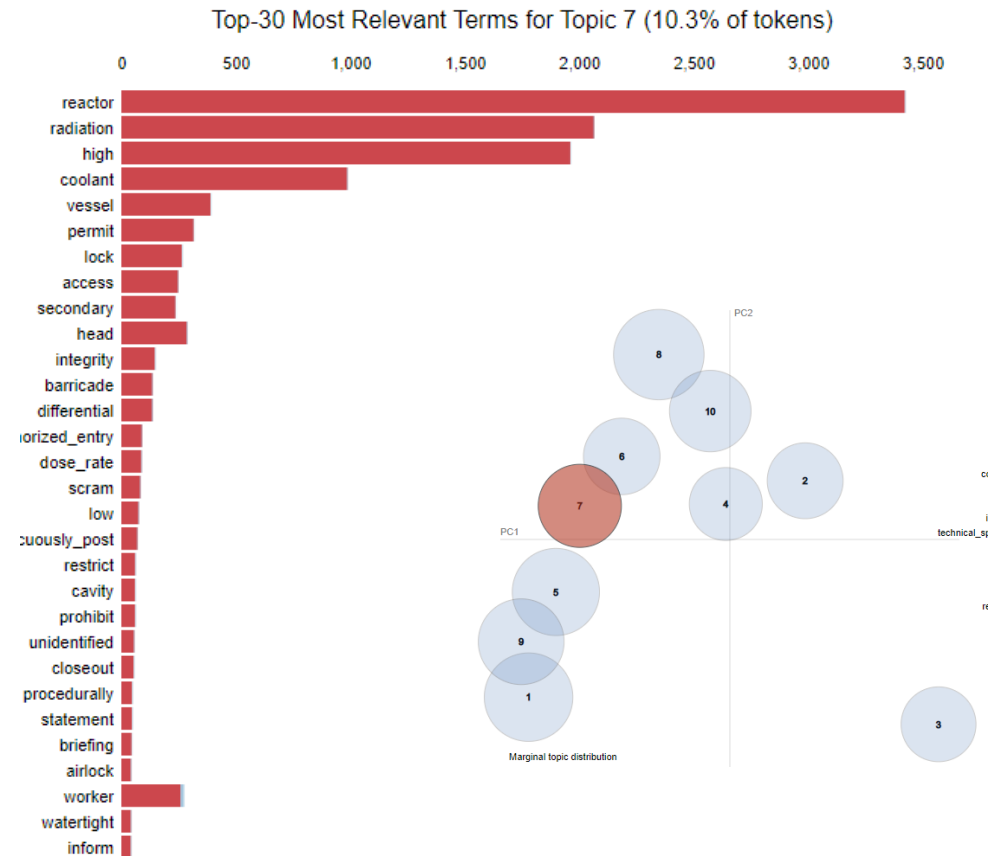
- Title - This clustering approach shows words that commonly occur together.
- This safety cluster has reactor and radiation as the most relevant unifying words.
- Additional experiments will be performed on larger text samples – potentially using the entire inspection report rather than just the title.
- Different approaches including bi-grams and tri-grams will be explored.
- Time periods will be introduced to ensure insights are current.
- These are very early results – several algorithms and parameter settings will be tried as we move into the experimentation phase.

Safety Cluster formation

Topics

- Clusters formed in Topic Modeling are characterized by the words that are unique to that collection
 - Reactor and Radiation appear in every title in this topic, and do not appear in any other titles
 - Same for High, Coolant, and Vessel
- Using the words most strongly connected to a given topic helps us characterize, or name that topic.
- That named topic can be considered a Safety Cluster
- A next step could be to characterize the reports in this Safety Cluster in terms of Cornerstones and Cross Cutting Areas

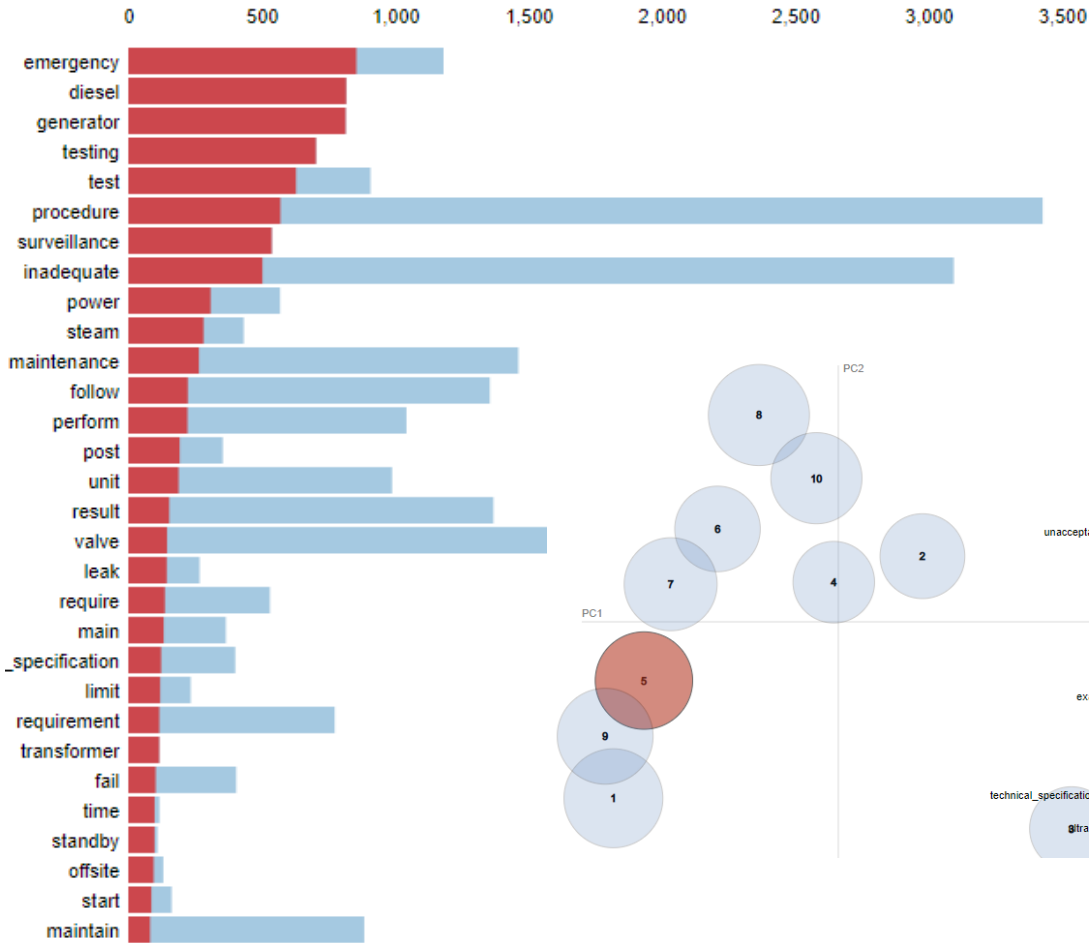
Topic Defining Terms => Safety Clusters



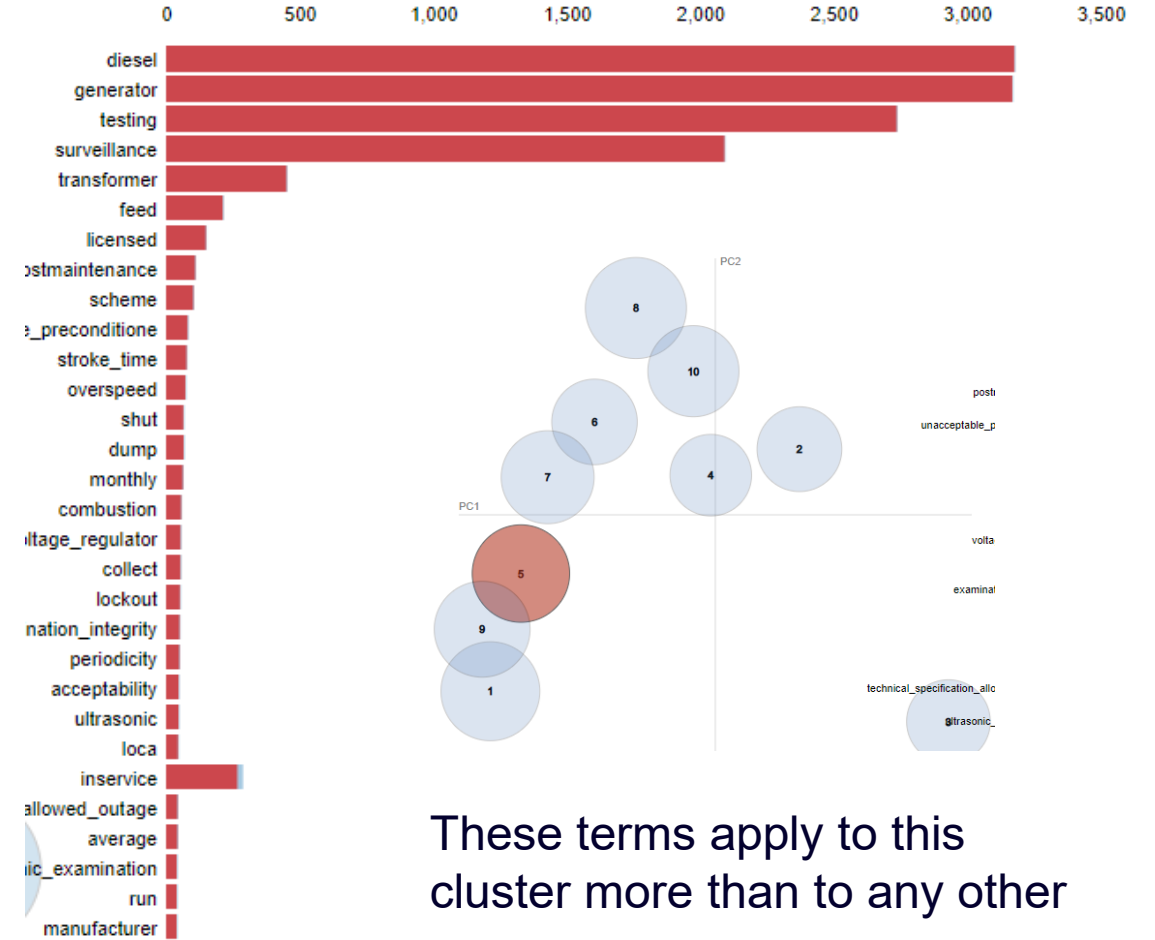
Safety Clustering Samples – Relative Importance

Terms on the right can be used to best characterize the cluster

Top-30 Most Relevant Terms for Topic 5 (11.3% of tokens)



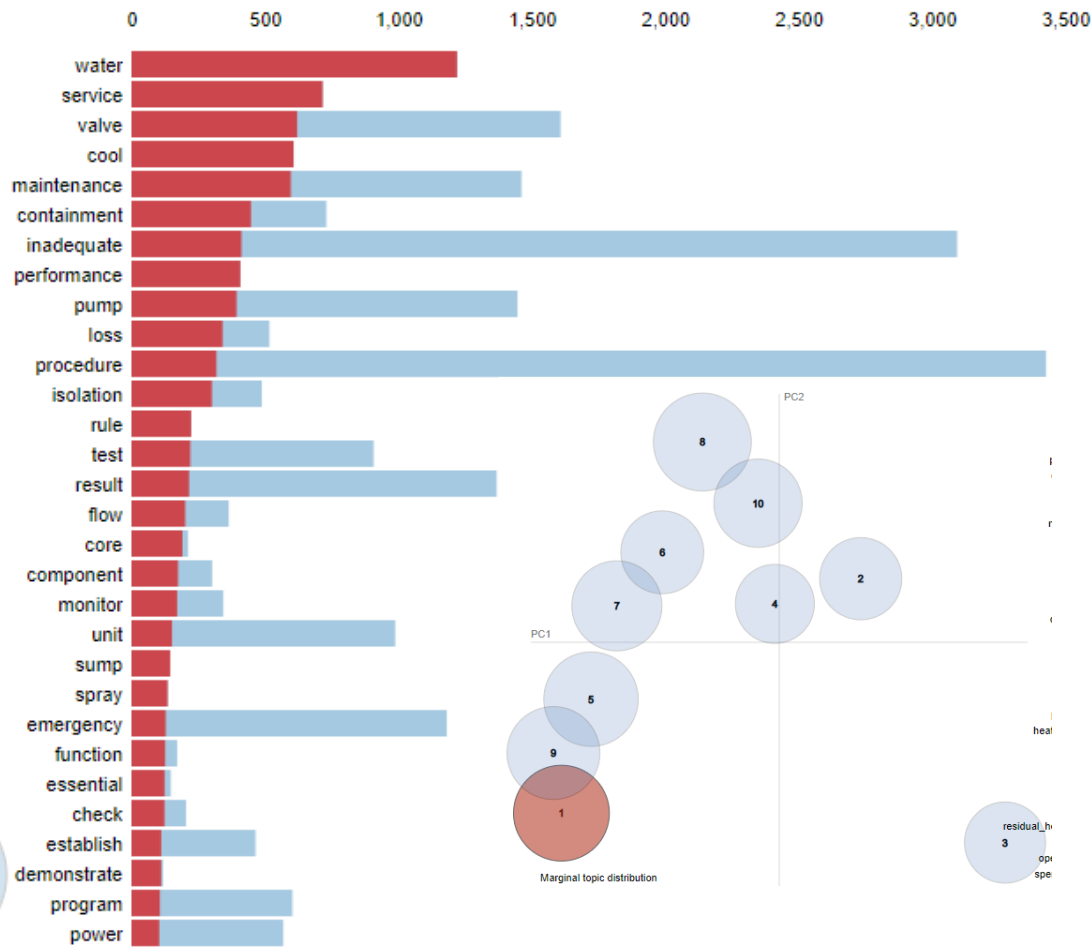
Top-30 Most Relevant Terms for Topic 5 (11.3% of tokens)



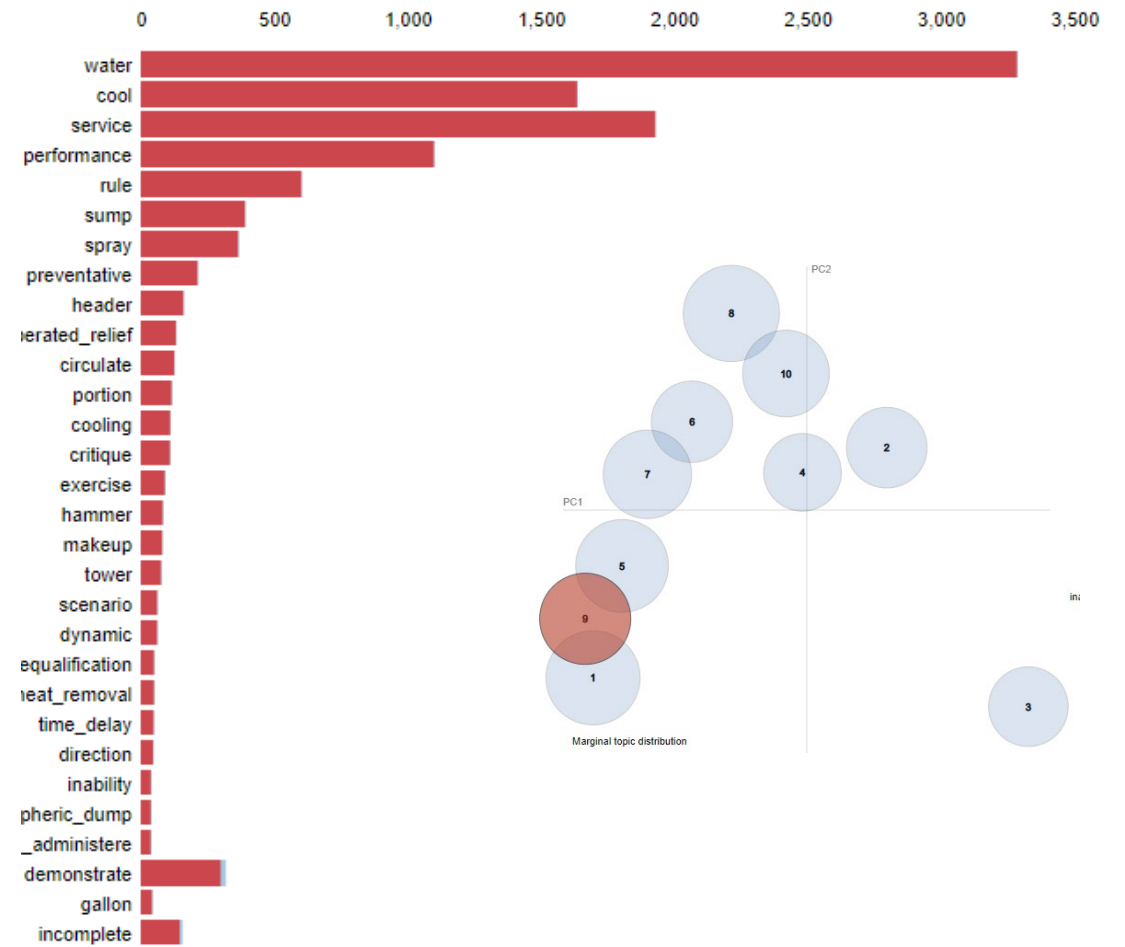
These terms apply to this cluster more than to any other

Clustering Samples – Overlap could indicate widespread issues

Top-30 Most Relevant Terms for Topic 1 (11.7% of tokens)



Top-30 Most Relevant Terms for Topic 1 (11.7% of tokens)



Data Analysis

Data Sources

- Reactor Locations
 - Site Name and Code, Unit, Docket, Reactor Type, Containment,
 - Region, State, City, Latitude, Longitude, Parent Company, Operating
 - Missing information on sites that have findings: FCS, CR, VOG3, CGS, PILG, KEWA, SANO, VY, TMI, OC
- Performance Indicators
 - Date (Year, Quarter), Docket, 17 indicators across 6 safety cornerstones
 - > 3 Initiating Events, 6 Mitigating Systems, 2 Barrier Integrity, 3 Emergency Preparedness, 1 Public Radiation Safety, 1 Occupational Radiation Safety
- Inspection Reports
 - 11,672 reports (4GB pdfs, 552MB extracted text)
 - > Text extracted from 10,671 pdfs; ~1000 pdfs are unreadable
 - Inspection Findings
 - > Year, Quarter, Issue Date, Report and Accession Numbers, Link
 - > Region, Site Code and Name, Docket Number
 - > Type, Item Severity Type Code, Significance
 - > Procedure, Cornerstone and its Attribute Type, Cross-cutting Area and its Aspect
 - > Identified By, Traditional Enforcement
 - > Title and Item Introduction

Data Sources

- **Event Notifications**
 - Event Notification Id and Number, Reactor Indicator, Site Name and Unit, Docket, Region, State, City, County, Time zone,
 - Rx Type; Event, Notification and Update Dates
 - Notified by, Operations Officer, Staff Names, Organizations, Text
 - Current (3) and Initial (3) Power; Critical Indicators (3); Scam Type (3)
 - Event and CRF Descriptions (4); I Mode (3) F Mode (3); Emergency Class
 - Agreement State, Licensee Name and Number, Containment Type (3); Document Type Code and CFR Code (4)
 - Security, Release Date, Note; Of Interest, Interest; Docket Number (2,3)
 - 38,864
- **Licensee Event Reports**
 - Plant Name, Event and Report Dates, LER Number, Accession Number, Title, Abstract
 - > 18872 with Title & Abstract in HTML, 18695 cleaned text with Title & Abstract separated
- **Part 21 Reports**
 - Log Number, Event/Accession Number, Report Date, Notifier, Description, Body, Links
 - 2,690
- **Action Matrix**
 - Date (Year, Quarter), Docket, total number of actions

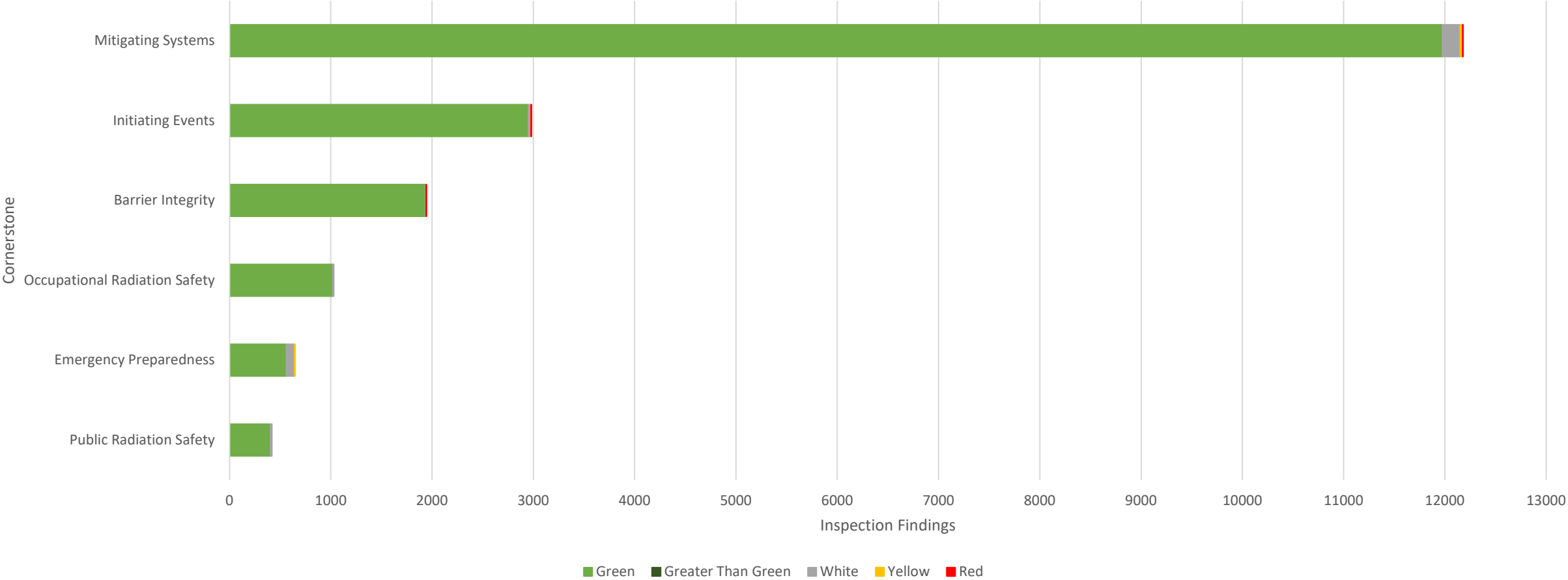
Data Sources

- Reactor Locations
 - Site Name and Code, Unit, Docket, Reactor Type, Containment,
 - Region, State, City, Latitude, Longitude, Parent Company, Operating
- Missing information on sites that have findings:
 - Inspection reports for these sites are not available on the main webpage
 - > FCS: Fort Calhoun
 - > CR: Crystal River
 - > PILG: Pilgrim
 - > KEWA: Kewaunee
 - > SANO: San Onofre
 - > VY: Vermont Yankee
 - > TMI: Three Mile Island
 - > OC: Oyster Creek
 - Will attempt to download from the links in inspection findings file

Inspection Report Characterization

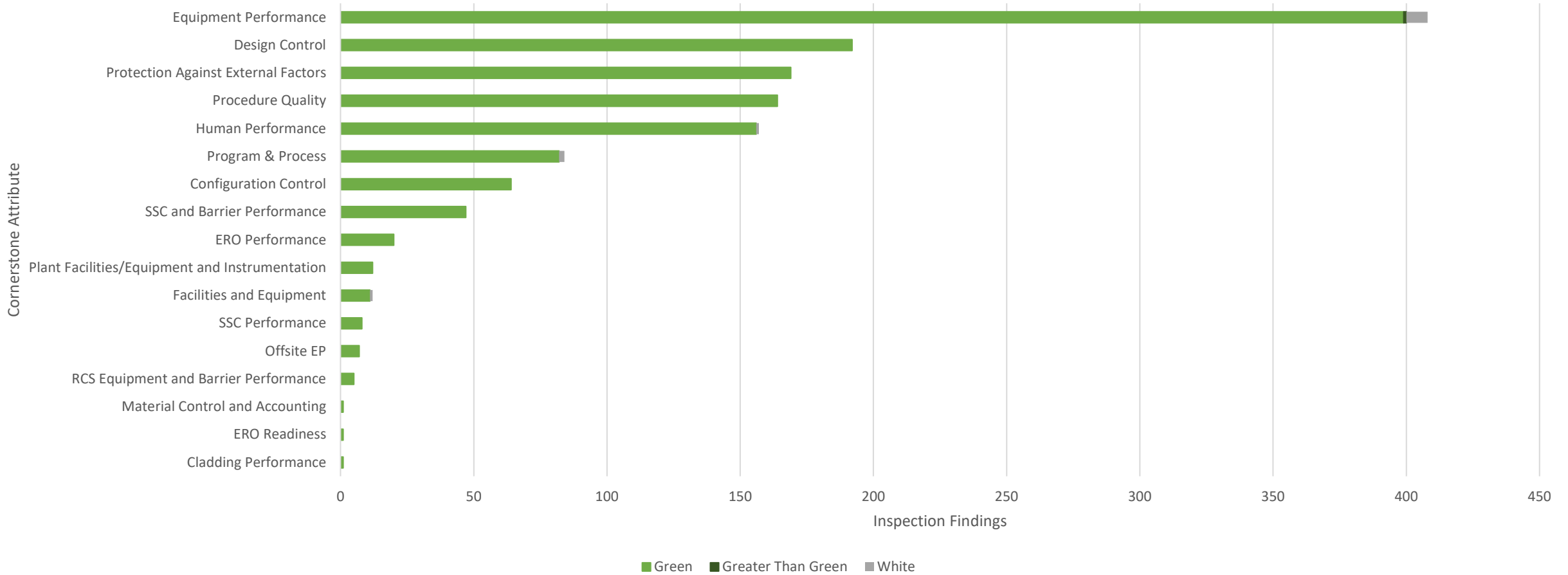
Inspection Findings: Cornerstone

Inspection Findings by Cornerstones



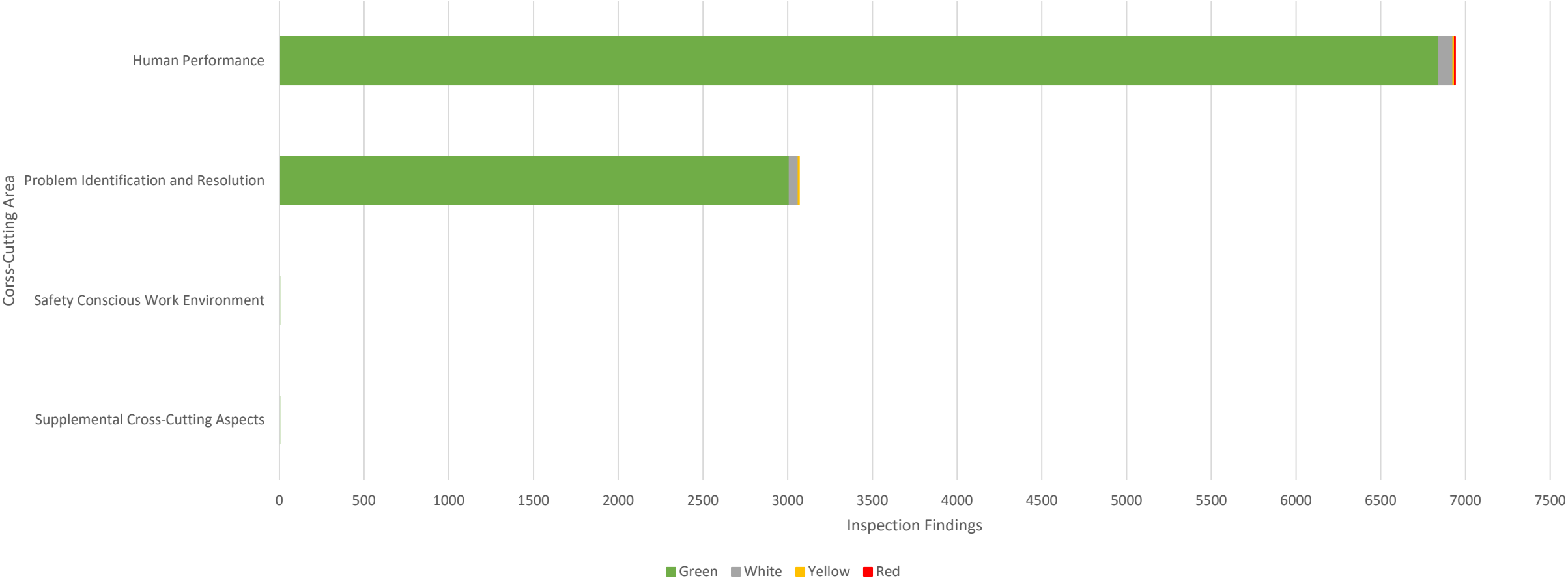
Inspection Findings: Cornerstone Attribute

Inspection Findings by Cornerstone Attribute



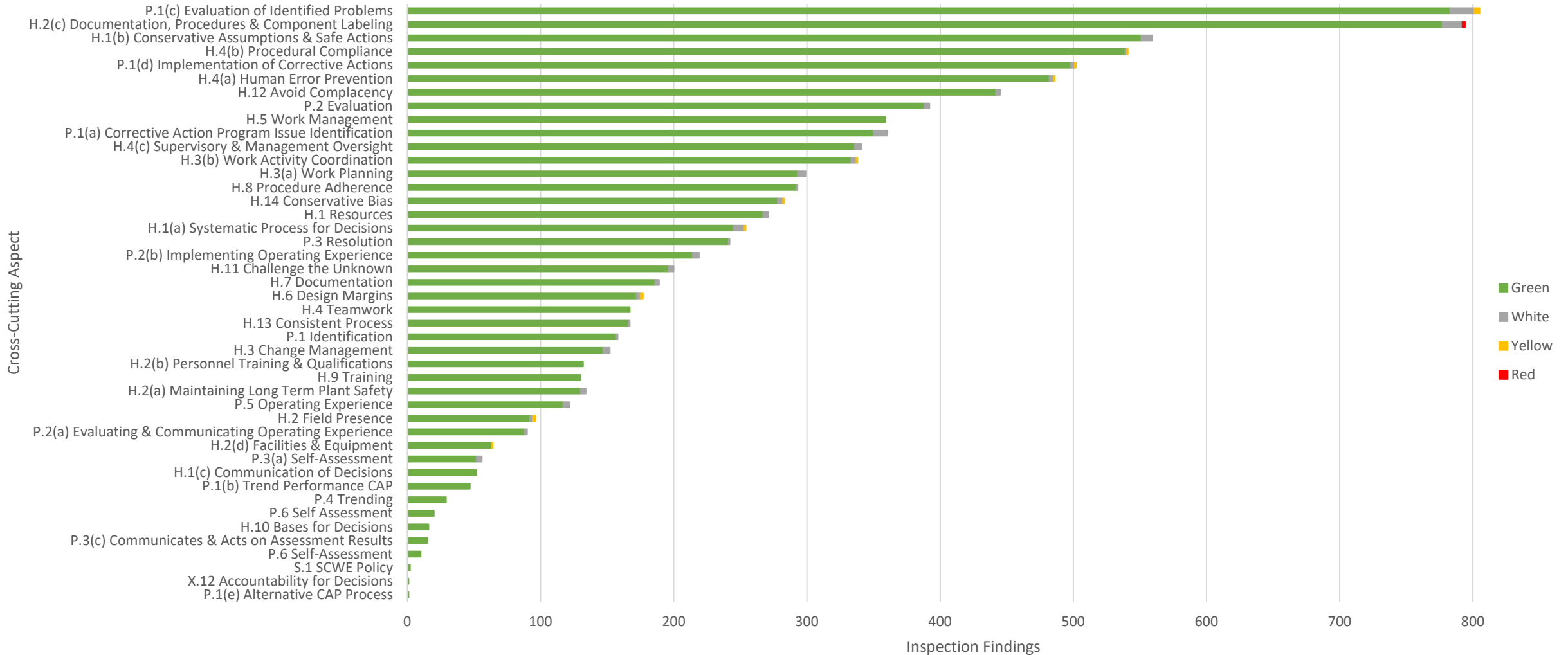
Inspection Findings: Cross-cutting Area

Inspection Findings by Cross-Cutting Area

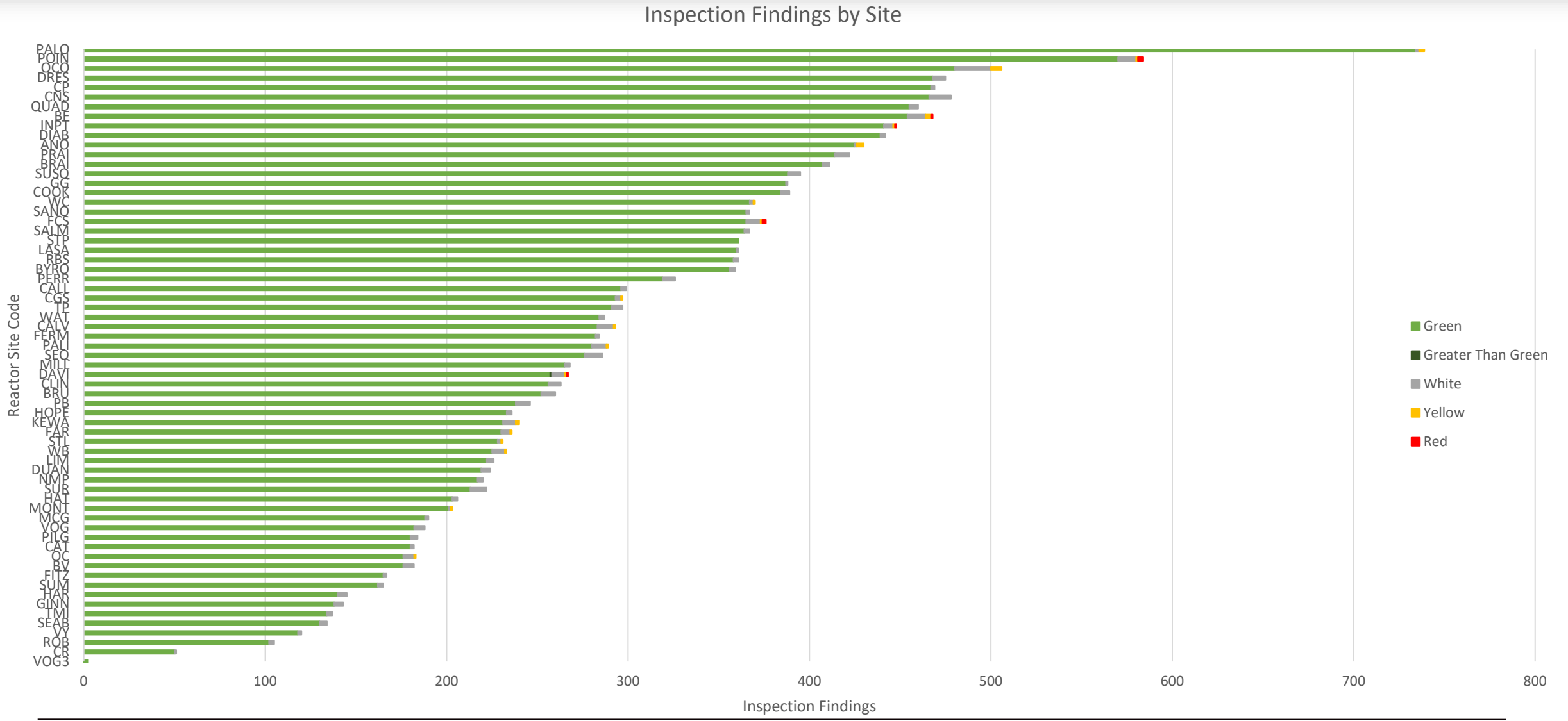


Inspection Findings: Cross-Cutting Aspect

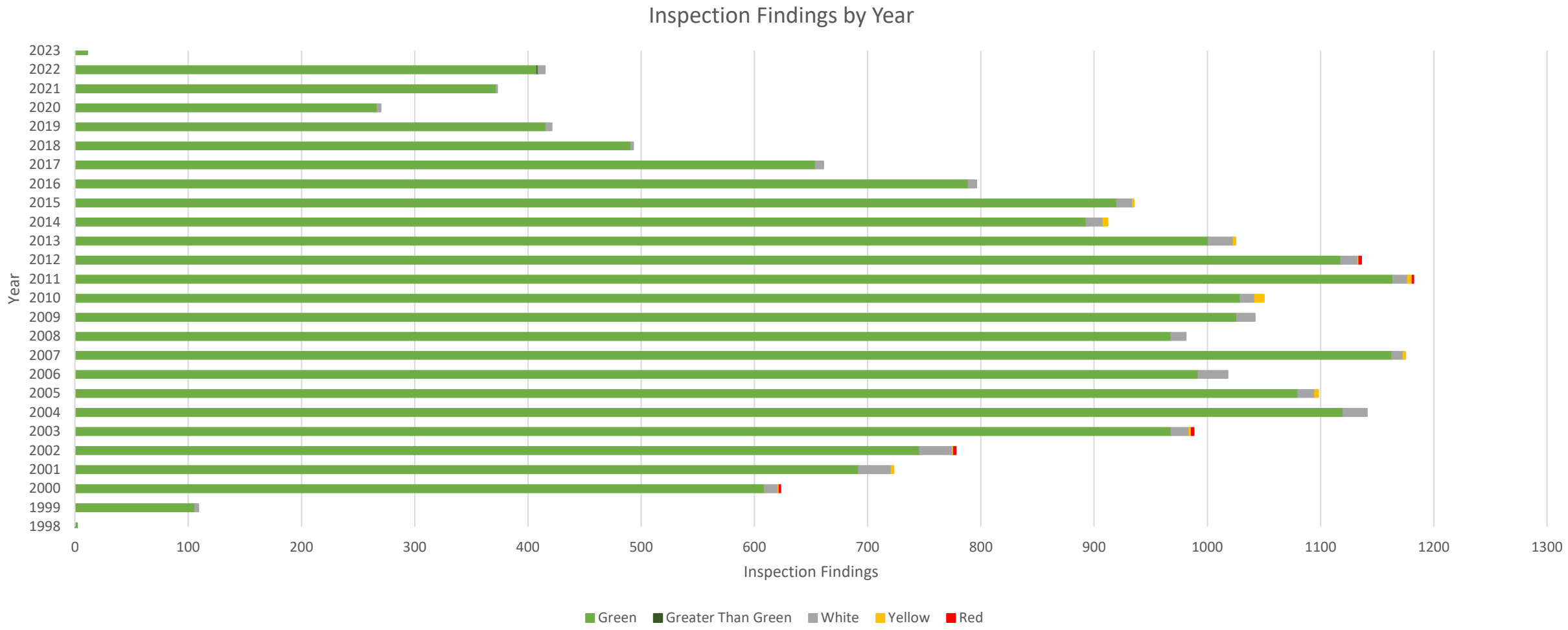
Inspection Findings by Cross-Cutting Aspect



Inspection Findings: Site



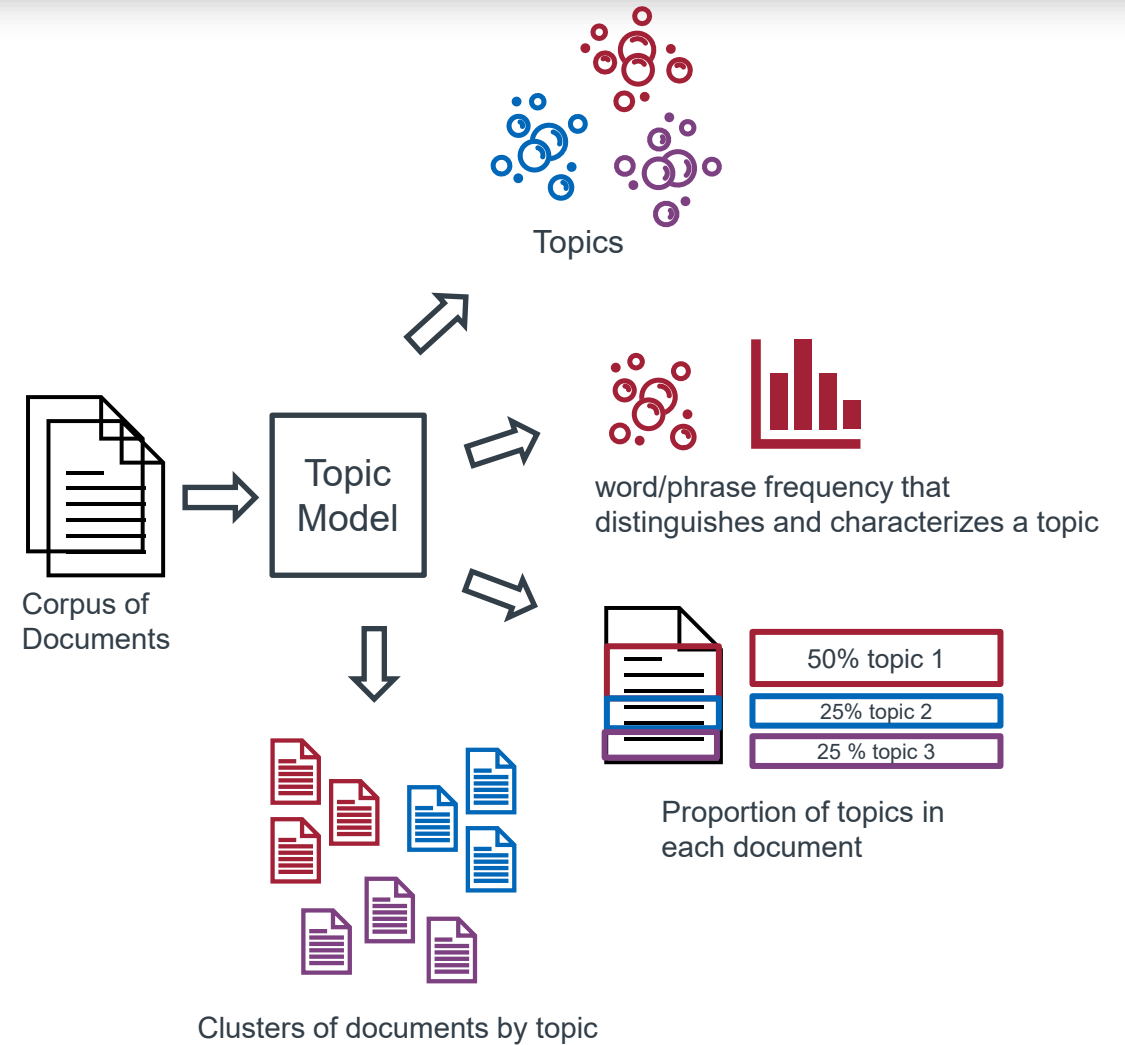
Inspection Findings: Year



Topic Modeling

Topic Modeling

- Unsupervised discovery of topics from a collection of text documents
- Latent Dirichlet Allocation (LDA)
 - Describe a document as a bag-of-words
 - Model each document as a mixture of latent topics
 - Topic is represented as a distribution over the words in the vocabulary
- Variants of Topic Modeling can be explored
- Text-embeddings from Language Models and Neural Topic Modeling can be used to improve the quality of results



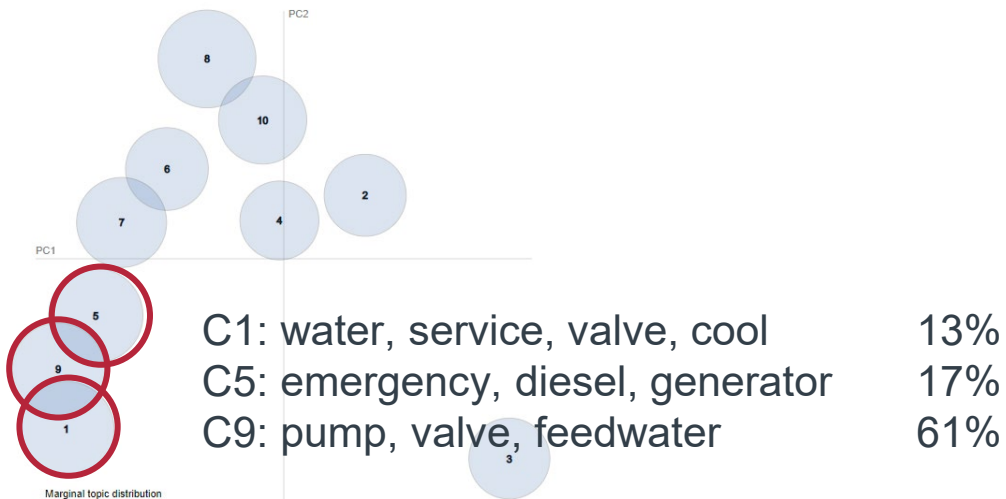
Demo

Actionable Insights

Actionable Insights from Cluster Analysis

Inspection Report ML15037A011

Indian Point



- The algorithm reveals the safety topics from findings
- These are the safety clusters
- One report may mention multiple safety topics
- The algorithm links the report to those safety clusters
- Each safety topic has a probability of match

Actionable Insights within a facility

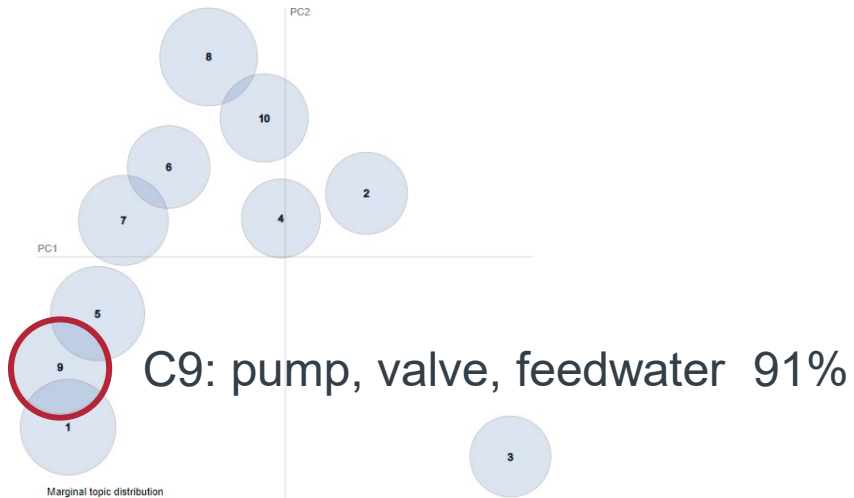
- Indicators for inspectors based on all of the Safety Clusters that a facility is linked to.
- All issues that are part of a topic cluster tend to occur together – ensure inspections cover each of these areas.
 - It appears that each inspection has a focus, these safety clusters could help validate that focus.

Actionable Insights from Cluster Analysis

Inspection Reports:

ML042190340 Grand Gulf

ML101330214 Indian Point



- Both reactors are strongly tied to safety Cluster #9 – Pump, Valve, and feedwater

Actionable Insights across facilities

- Indicators that tie reactors together based on safety findings
- Analyze other factors to seek similarities in designs, procedures, etc. to reveal potential hazards.

Actionable Insights from classification analysis

Classification

Input from:

- Performance Indicators
- Event Notifications
- Event Reports
- Part 21 Reports

Classify on:

- Significance
- Cornerstone Areas
- Crosscutting Areas

Significance	Probability
No Finding	47%
Green	28%
White	14%
Yellow	9%
Red	2%

Actionable Insights for Classification

- Use alternate sources of input.
 - Exclude data on the target feature (significance)
 - Use the remaining features to assign probabilities to the significance categories
 - Could also classify on cornerstone and crosscutting areas
- Can use historical inspection reports to verify accuracy

Actionable Insights from predictive analysis

Predictive

Input from:

- Performance Indicators
- Event Notifications
- Event Reports
- Part 21 Reports

Predictions on:

- Cornerstone Areas / Attributes
- Crosscutting Areas / Aspects

Actionable Insights for Future

- Use alternate sources of input.
 - Exclude data on the target feature (cornerstones)
 - Use the remaining features to discover cornerstone trends and patterns
- Assigns a probability to findings in inspections
- Can use historical inspection reports to verify accuracy

Are these types of actionable insights what you would like to see from the study?

What other goals do you have?

Tool Analysis

Survey of Platforms and Tools

If you want simple, but limited:

Platform Supplied Algorithms

Pros

- Integrated into user experience
- Can be combined into standard pipelines
- Work well in common cases (social media)
- Well known algorithms

Cons

- Do not handle technical domains well
- Difficult to tailor pipelines
- Limited selection of algorithms
- Limited selection of pre-training datasets
- Advanced algorithms not available
- Strengths in one area (Topic Modeling) may not translate to other areas (Classification, clustering, regression, anomaly detection)

If you want flexible and timely:

Python Library Algorithms

Pros

- Flexibility to select from multiple algorithms
- Flexibility to select from multiple pre-trainings
- Advanced models available
- Flexibility to address highly technical domains
- Flexibility to leverage internal parameters
- Work well in Machine Learning Notebooks
- Can be deployed in any cloud or on premises

Cons

- Requires some knowledge of Python
- Algorithms must be researched and selected

Survey of Platforms and Tools

- Platforms and Tools:
 - Amazon's AWS SageMaker
 - Microsoft Azure's AI/ML Services
 - Google's Cloud AI Products
 - MatLab
- Initial survey indicates that all platforms provide the ability to use
 - Python libraries needed for pre-processing textual data
 - Python libraries that provide algorithms for unsupervised learning with textual and numerical data
 - Various pre-trained models and easily fine-tune them with our data
- Python notebooks that are currently used locally can be launched and scaled on all platforms
- Refine evaluation factors and their weights & ranges with continued exploration of unsupervised techniques

What do you see as the primary role of the people executing this type of analysis?

Are there longer-term aspirations for ML initiatives?

Progress

SOW Task Status

Phase I: March 6, 2023 - April 9, 2023

	Status
Describe the Problem	In progress
Search the Literature	In progress
Select Candidates	In progress
Select Evaluation Factors	In progress
Develop evaluation factor weights	In progress
Define evaluation factor ranges	In progress
Perform assessment	Not started
Report Results	Not started
Deliver Trade study report	Not started

Phase II: March 20, 2023 - May 7, 2023

	Status
Platform/system selection and installation	Not started
Data acquisition and preparation	In progress
Feature pipeline engineering	In progress
Clustering method experimentation & selection	In progress
Cluster pipeline engineering	In progress
Anomaly detection (as needed)	Not started
Model Development, Training, Evaluation	Not started
Test harness development	Not started
PoC integration and demonstration	Not started
Trial runs and evaluation	Not started
Demonstrate PoC capability	Not started

Phase III: April 19, 2023 - June 16, 2023

	Status
Live data ingestion	Not started
Model execution	Not started
Cluster evaluation	Not started
Critical Method documentation	Not started
Technical Report Document	Not started
Deliver final report with findings	Not started

Next Steps

Next Steps

Determine algorithms and approaches needed to perform the study:

- Exploring LDA with other text data: full inspection reports, event notifications, LERs, part 21's
- Explore neural topic modeling with BERTopic and Contextualized Topic Model
- Explore other ML areas as directed

Define the selection criteria for tools / environments based on study needs:

- Compare the services and methods available for topic modeling across our 4 platforms