

## CHAPTER 3: STEPPING FORWARD: TOWARDS A MORE SYSTEMATIC NPF WITH AUTOMATION

Laura P. Wolton\*, University of Colorado Denver, School of Public Affairs  
Deserai A. Crow, University of Colorado Denver, School of Public Affairs  
Tanya Heikkila, University of Colorado Denver, School of Public Affairs

### ABSTRACT

Advancements in automated text analysis have substantially increased our capacity to study large volumes of documents systematically in policy process research. The Narrative Policy Framework (NPF)—which promotes empirical analysis of narratives—has the potential to usher policy narrative research along the same path. Using the NPF and existing semi-automated analysis tools, we investigate the relationship between narrative components—namely, characters and proposed solutions—and the more “skeletal” frames that tie policy narrative elements to one another. To illustrate how these tools can advance policy narrative research, we auto-code 5,708 state and local news articles focusing on hydraulic fracturing of oil and gas. The findings suggest that the use and role of characters and policy solutions are portrayed in significantly different ways depending on the frame used. By using an autocoding approach, these findings increase our methodological and theoretical understanding of the relationship between narrative elements and frames in policy narratives. In discussing these findings, we also consider their implications for how issue frames matter theoretically in the NPF.

\*Corresponding author: [laura.wolton@ucdenver.edu](mailto:laura.wolton@ucdenver.edu)

To cite the chapter: Wolton, Laura P., Deserai A. Crow, and Tanya Heikkila. 2022. “Stepping Forward: Toward a More Systematic NPF with Automation”, in *Narratives and the Policy Process: Applications of the Narrative Policy Framework*, Michael D. Jones, Mark K. McBeth, and Elizabeth A. Shanahan (eds.), Montana State University Library, 40-90. [doi.org/10.15788/npf3](https://doi.org/10.15788/npf3)

### INTRODUCTION

For the last 60 years, text analysis methods have allowed advances in the research traditions of various fields, including behavioral sciences, consumer research, and media analysis

(Humphreys & Wang, 2017). The General Inquirer (Stone et al., 1962) was one of the first computer-assisted approaches designed for the analysis of text, employing a dictionary lookup method that “tagged” a sentence for the appearance of a word belonging to a specific category, and became a precursor to dictionary lookup methods available to content analysts today. Recent computer-assisted text analysis advances include semi- to fully-automated approaches, which emerge from multi and interdisciplinary research efforts from the computer sciences, social sciences, linguistics, and artificial intelligence.

While automated and semi-automated text analysis methods advanced for decades, we have only recently seen policy process scholars using these methods in published research (e.g., Lawlor, 2015; Heikkila & Weible, 2017; Olofsson et al., 2018; Berardo et al., 2020; Blair et al., 2016; Scott et al., 2020). Although some of these applications of automated text analyses have looked at themes and issues that arise in policy debates, less attention has focused on the structure of policy narratives. With widespread availability of policy narrative content in various media formats—from online news, downloadable publications, to social media—there is an ongoing need for automated tools to systematically evaluate larger datasets of narratives (Shanahan et al., 2018). This chapter will demonstrate that the incorporation of semi-automated and automated methods into the policy researcher’s analytical toolbox can advance policy narrative research. The chapter begins by discussing the differences and connections between narratives and frames, beginning with a detailed discussion of policy narratives. The example case used in this analysis is then described, followed by a step-by-step description of the automated method introduced in this chapter. The findings are presented followed by a discussion of the implications of this method for policy process scholarship in general and the study of narratives specifically.

The employment of automated methods in policy narrative research is congruent with the ambitions of the Narrative Policy Framework (NPF) as the framework itself promotes a systematic structure that allows for increased generalization and comparison amongst research findings. Through its operationalization of narrative elements, the NPF allows for a systematic and reliable investigation into the role of narratives in the policy process. As of yet, only two studies have used semi-automated techniques with the NPF—both with inductive methods of categorization for designing autocoding dictionaries (Merry, 2018; Crow & Wolton, 2020). The first objective of this chapter is to extend the previous work done by Crow and Wolton (2020) and, in doing so, propose a method to connect autocoded policy actors to NPF characters.

The second objective of this chapter is to investigate the relationship between frames and narrative elements, expanding on an emerging body of NPF literature (Jones & Song, 2014; Lawlor & Crow, 2018; Crow & Lawlor, 2016; Merry, 2018). The framing literature suggests that communications using the same frame will contain the same narrative elements. For example, within one frame, similar categories of policy actors are likely to be portrayed in a similar character role (e.g., industry actors as villains, government actors as heroes). The solutions within one frame should similarly remain constant. Characters and solutions are likely transported in frames together because actors have a range of solutions that they have the ability or authority to perform (e.g., government can regulate, citizens can vote). Thus, to our second objective we ask: How does a narrative frame relate to the characters used and solutions presented? In answering this question, we not only offer insights on semi-automated narrative analysis tools, but also offer theoretical advancements for the NPF and framing literature.

### NPF, FRAMING, AND AUTOCODING

Arguably, improving automated techniques will amplify the potential for systematic and generalizable policy narrative research, particularly if the methods are transparent and accessible. Frameworks create a shared orientation and conceptual map that structure how researchers study a particular phenomenon (McGinnis & Ostrom, 2014). For policy communication researchers, the development of the NPF (Jones & McBeth, 2010) was a successful step towards promoting an increased degree of structure and consistency in policy narrative analysis. The interdisciplinary nature of policy narrative research and, thus, the potential for application of theories across policy domains highlights the need for increased replicability and validity in NPF studies (Shanahan et al., 2018).

Though NPF analysis has not commonly included frame analysis, several scholars have made efforts to connect the framework with framing theory (Jones, 2013; Crow & Lawlor, 2016; Olofsson et al., 2018). Frames are broader approaches to narrating a topic, focusing attention on particular themes, ideas, or issues. In doing so they can constrain narratives to the selection of a particular set of narrative elements that are commonly aligned with a given frame. Recent advances in the automated detection of frames in news articles suggest that, coupled with an increased understanding of elements, these methods can expand our understandings of the structures and elements of policy narratives.

### Narrative Policy Framework

The NPF is a means of investigating the structure and content of policy narratives (and the implications of that content to public opinion and policy) both qualitatively and quantitatively, promoting shared units of analysis and codebook design. The NPF asserts that *policy narratives* are structured stories that include a setting, characters, plot, and a moral of the story (the policy solution to the problem). However, to be considered a policy narrative, a communication or text must contain at least one character and a reference to a policy issue (Jones et al., 2014).

#### *Characters*

The term *characters*, according to the NPF definition, involves those entities “who act or are acted upon” and are categorized according to common roles that appear in narratives (Shanahan et al., 2018, p. 335). While only three narrative characters appear in the original framework of the NPF (heroes, villains, and victims), subsequent studies have added to the range of identifiable characters (Shanahan et al., 2017). Of those previous categories, this study limits character identifications to heroes, villains, victims, allies, opponents, and charismatic experts:

- *Heroes* fix or attempt to fix a problem and are praised in some way.
- *Villains* do the harm or are blamed for the policy problem or issue.
- *Victims* are harmed by the problem and are constructed to receive sympathy.
- *Allies* hold a policy position in agreement with the author. Interest groups—such as oil and gas industry associations and environmental groups—use messages intended to strengthen bonds among like-minded groups (Merry, 2016).
- *Opponents* hold a policy position in disagreement with the author (Merry, 2016). Though opponents support villains and may receive blame, the language that they are associated with is not expected to be as severe as that surrounding villains.

- *Charismatic experts* are strategically cited in policy narratives because they lend legitimacy and credibility to key evidence and reports. Lawton and Rudd (2014) argue the necessity of experts as an additional NPF character to improve empirical understanding of how evidence impacts policy decisions.

### *Moral of the Story or Solutions*

One of the four primary narrative structure elements (i.e., variables) in the NPF is the moral of the story, which elicits action or provides a policy solution (Stone, 2002; Ney & Thompson, 2000; Verweij et al., 2006). Solutions are most often related to characters in the NPF, as heroes “take action with purpose to achieve or oppose a policy solution” (Shanahan et al., 2018, p. 343). NPF research shows that though not necessary to the definition of *policy narrative*, solutions are present in a high proportion of narratives (Crow et al., 2017; Crow & Wolton, 2020). Research indicates that solutions are more likely to be accepted by readers who experience positive affect towards story characters (Jones & Song, 2014), pointing to the need for further research on both characters and solution. Thus, in exploring how frames relate to characters and solutions, we can help advance this line of research.

### **Frames**

Framing theory from the mass communication scholarship helps us connect the broad issue characteristics of focus (i.e., frames) to the more specific structural elements (e.g., characters, solutions) of policy narratives. As Crow and Lawlor (2016) discuss, frames are the overarching approach to narrating a story that constrains the specific choices made in construction of narratives. For example, when framing a discussion of the climate change issue, a narrator could use a religious frame and, therefore, cast characters such as a pastor or God when discussing the morality of addressing the climate crisis. While both policy narratives and frames may be used strategically to influence policy outcomes, narratives more overtly include story-like elements such as morals and contrasting characters.

Framing entails the “selection of certain aspects of a perceived reality and makes them more salient in a communicating text, in such a way as to promote a particular problem definition, causal interpretation, moral evaluation, and/or treatment recommendation” (Entman, 1993, p. 52). That is, there is a larger structure to any story that organizes facts into a storyline and emphasizes selected facets to communicate a particular viewpoint (Crow & Lawlor, 2016). Framing theory considers frames to be the structural “bones” of stories—a major theme that bounds the story within which the narrative elements are assembled.

An underlying assumption of framing theory is that there are options for how to present information (Scheufele & Iyengar, 2014). Thus, the effects of framing on audience perceptions are due to choices made in the way information is presented—not necessarily the facts that are revealed. Journalists are not the only communicators who frame stories—individuals do the same in how we convey information and focus on selective parts of a story at the expense of other potential angles. For example, journalists in the United States strive for “unbiased” presentation of facts but are also taught to shape stories in an audience-focused manner; to stir emotion and empathy they select from a list of issue characteristics choosing the most compelling (Crow & Lawlor, 2016). These choices are consequential because focusing on one aspect of a policy issue

may limit the conversation of the public and policymakers around that issue (Crow & Lawlor, 2016; Lawlor, 2015).

There are a number of approaches to studying frames and competing or conflicting definitions. Using framing analysis, we can study sub-topics of a broad issue (issue frames), different conceptual aspects of an issue from moral to political and others (news frames), broad themes versus specific incidences (thematic vs. episodic frames), among others. In this analysis, we will employ issue frames, but others could be used depending on the focus of a given study. A researcher needs to specify the approach and follow the method presented in this chapter with a clear definition of frames in mind.

### *Frames and Elements Theoretically Travel Together*

The definition of framing presented by Entman (1993) suggests that elements of the NPF, such as solutions and characters, should be correlated to the use of frames because framing involves choosing (perhaps not always deliberately) a certain way to present information (e.g., religious vs. economic). The co-occurrence of narrative characters and solutions within a specific frame is a somewhat obvious deduction. For example, if the Environmental Protection Agency is cast as the hero character in an environmental frame, the solution attached to this hero will likely be one related to regulation—the agency has only a range of policy solution activities that it can perform. If the agency is cast as a villain in an economic frame, it may be associated with a similar solution—regulating—while the hero of the story is most likely opposing this government intervention. Finding empirical evidence that narrative elements and frames co-occur would be a meaningful contribution to framing literature as well as NPF so that we can better understand the relationship between these two communication tools and eventually the corollary effects of those tools in tandem and separately.

### **Applying Automation to Policy Narratives**

While automation has been increasingly applied to other areas of research, automated analysis with the NPF has been held back by some aspects of the framework itself, as well as limitations of automated text analysis in general. We address a few perceived roadblocks and propose ways to move through them, focusing on the NPF elements of *characters* and *solutions*. Although the automated detection of frames is relatively established, we also stress the importance of conceptual clarity *prior* to semi-automated or automated coding.

### *Character/Actor Association Fluidity*

The NPF defines characters in such a way to allow for fluidity, so that various policy actors or groups may play the role of a single character and that a single policy actor may hold different character roles in different narratives. Because this definitional approach to characters means that one character is never necessarily tied to one policy actor or group, NPF researchers have found it difficult to apply autocoding and capture this fluidity (Shanahan et al., 2018). Our methods, which take the context surrounding the policy actor into consideration, address this issue by allowing the actor to be associated with an NPF character in each frame. Additionally, particularly in cases of large volumes of narratives, the number of policy actors and groups is potentially very large. We use categorization strategies to reduce the complexity of the association between policy actors and

characters. Our method takes advantage of information beyond the name of a policy actor to identify and characterize them (further explained in the *General Method Description* section of this chapter).

### *Solution Variability*

While the NPF is definitionally clear on the moral of the story—a proposed specific or general policy fix (e.g., “go vote”)—prescriptive coding schema for *moral of the story* is less developed. Currently, a basic codebook on solutions may include a range of policy solutions that may vary according to policy domain, level of government, or other factors (see Appendix A in Shanahan et al., 2018). For automated coding, this presents a problem with the potentially high variability of solutions, particularly for large volumes of policy narratives. Similar to the work of Crow and Wolton (2020), which created code categories by policy domain, we designed a category schema based on the primary policy tools the solutions are proposed to (or do) implement. Additionally, we allow for solutions that are less concrete than specific actions or policies to be placed in a general category. Our method of autocoding uses parts-of-speech analysis and term-frequency generation as a primary method of identifying solutions (further explained in the *General Method Description* section of this chapter).

### *Finding Conceptual Clarity on Frames*

Frames have been effectively identified with semi-automated methods by numerous scholars in various policy domains including public health, immigration, and air pollution (Lawlor, 2015; Olofsson et al., 2018; Poirier et al., 2020; Yu et al., 2020; Berardo et al., 2020). Two general approaches are used in autocoding frames: 1) inductive dictionary development and 2) Latent Dirichlet Allocation (LDA) modeling (these are further explained in the *General Method Description* section in this chapter). As with any research involving frames, it is essential to be clear which type of frame is under investigation *prior* to semi-automated or automated coding. Understanding the differences in framing definitions since there are numerous definitions that often conflict or overlap with one another, as well as taking steps to appropriately bound investigation by the objectives of the research, are especially beneficial when approaching large datasets.

## **THE POLICY DOMAIN: SHALE OIL AND GAS DEVELOPMENT IN THE UNITED STATES**

Energy development, particularly of shale oil and gas, is a contentious issue in the United States. The United States is the world’s top producer of petroleum and natural gas hydrocarbons and has held this position since 2009 when hydraulic fracturing and horizontal drilling allowed for rapid expansion of shale development (US Energy Information Administration, 2018). Proponents of oil and gas development argue that expansion of shale gas provides energy independence, jobs, and economic benefits. For all its beneficial economic contributions, those who oppose shale oil and gas development argue that it can lead to property rights conflicts, environmental harm, and risks to public health and worker safety, among other issues. For instance, some of the recorded impacts of hydraulic fracturing include increased local air pollution from dust, odor, ozone, and volatile organic carbon (VOC) emissions; methane emissions; water contamination; high water use; structural damage from earthquakes triggered by wastewater reinjection; damage to wildlife habitat; abandoned orphan wells; oil, gas, and wastewater spills; noise pollution; and worker

fatalities (Joyce & Wirfs-Brock, 2015; Konkol, 2016; Hand, 2015; Witze, 2015; Adgate et al., 2014; Gallegos et al., 2015; Bamberger & Oswald 2015; Moskowitz, 2015; Mason et al., 2015).

Analyzing news articles on a salient, controversial issue, such as shale oil and gas, is a useful way to explore a wide range of policy narrative elements, such as character portrayals, solutions, and variable frames. Additionally, shale oil and gas development draws significant attention at local as well as state levels in the United States —yielding potentially different narratives across locations and additional variation in the narrative elements used. Lastly, the variance in regulatory stringency of shale oil and gas development across states allows for further insight into a range of policy perspectives.

### GENERAL METHODS DESCRIPTION

To highlight the potential for semi- or fully-automated techniques and investigate the relationship between NPF elements and frames, we detail a multi-step method. In addition to describing how to detect characters, solutions, and frames, we briefly describe how to check a large data sample of written media to ensure that they are policy narratives. The method presented here is less a procedure than it is a flexible guideline. The researcher's theoretical lens should ultimately shape the selection of the methodological approach, variables, and relationships to be studied.

We describe six general phases shown in Figure 1: 1) data selection, collection, and cleaning data; 2) designing and refining the dictionaries; 3) subsetting to only policy narratives; 4) final autocoding; 5) subsetting to a single frame; and 6) mapping of the policy actors into NPF character sentiment ranges. Novice- to intermediate-level programmers will likely do well to complete Phases 1 through 4 in MAXQDA,<sup>1</sup> export all segments as a spreadsheet to *R*, and complete Phases 5 and 6 in *R* or similar programming language.

#### Phase 1: Data Selection, Collection, and Cleaning

Textual data selected for automated content analysis must be in computer-readable formats, the dataset filtered to include only relevant material, and each item cleaned to contain only the text to be analyzed. For text documents, common computer-readable online formats that are also readable by humans include hNews, HTML, XML, and PDF. Some media archive formats, image formats, and hand-written documents are not computer readable and, therefore, will not work with this method of narrative analysis. The process of filtering media, such as removing items that do not fit a study's selection criteria, can be completed by hand or automatically (as in this study). Options exist for the collection of news articles through news archive services such as NexisUni, application programming interfaces (APIs), or webscraping (automatically downloading webpages). These collection methods may yield relatively consistent retrieved files that can be batch cleaned with a programming language such as *R* or *Python*. Cleaning may include the removal of items that will not be analyzed textually, such as emojis, captions, byline material, and other non-relevant content. During the process of cleaning, a catalog file of metadata should be created that contains file names, article titles, source name, date written, author, or other information that is associated with each file.

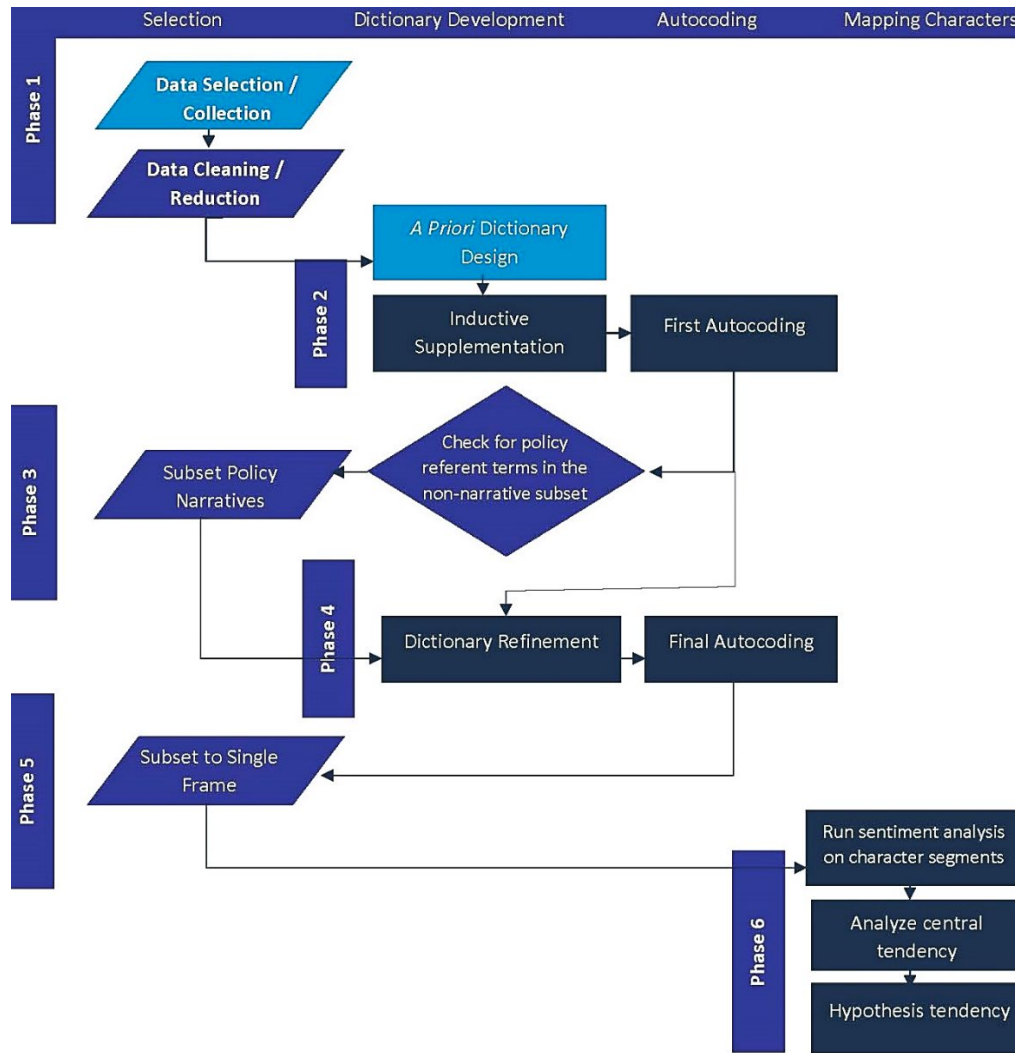
---

<sup>1</sup> MAXQDA Plus or MAXQDA Analytics Pro are the versions of the software that contains MAXDictio plugin for a number of functions listed here.

## Phase 2: Design and Refinement of Dictionaries

In this method, dictionaries are used for autocoding content and take the place of codebooks that researchers use in hand-coding activities. Dictionaries are used to identify and categorize text that is indicative of a theme or concept. Before beginning the discussion of the dictionary development process used in this study, we define several relevant terms and a few technical aspects of designing dictionaries for computer-assisted analysis. Based on Krippendorff (2004) and similar seminal works on content analysis methods, a quantitative content analysis schema contains multiple components: 1) a *text corpus*—the texts being analyzed, 2) a *dictionary*—a list of terms associated with theoretical concepts of interest, 3) an unlimited number of *categories*—groupings of the terms that have shared theoretical characteristics, and 4) rules of identification for each term or phrase such as those for capitalization, plurality, and/or verb forms.

Figure 1. General process outline



A common way to create dictionaries is in table format, such as the spreadsheet-based dictionary sample in Table 1 that may be read into autocoding software, such as MAXQDA. The



sample, taken from the character dictionary for this hydraulic fracturing study, contains a column for the search item, its associated category, and rules for its identification. The rules of identification—whole word, case sensitivity, and starting letters—are entered as on/off (1/0) switches that dictate which forms of a term will be tagged. For example, the search item in row 1, “anti-fracking group,” has been placed in the category *activist group*. Because “anti-fracking group” has the whole word and case sensitivity toggles turned off, the phrases “**anti-fracking groups**” as well as “Colorado **Anti-Fracking Group**” would be tagged during autocoding. The toggle for starting letters, which is irrelevant for the search items in this dictionary sample, becomes meaningful for words with prefixes or that are common portions of words. For example, the toggle would need to be activated when identifying terms like “sent” and “tree” to prevent words like “absent” and “street” from being tagged.

Considerations about word forms are also essential in dictionary design. Both contextual checks of search items as well as conceptual reviews are necessary to avoid coding errors and design problems. Contextual checks are part of the process of inductive dictionary design (described below) and are beneficial to confirming phrase use, particularly in the case of phrases that potentially have multiple meanings. Reviews of the dictionary for conceptual adherence, as well as consistency in the level of detail of the search items, should occur once at minimum between dictionary development phases.

*Table 1. A content analysis dictionary contains categories of search items and rules for term identification*

Category	Search item	Whole word	Case sensitivity	Starting letters
ACTIVIST_GROUP	anti-fracking group	0	0	1
ACTIVIST_GROUP	Aurora Citizens for Responsible Energy	1	1	1
ACTIVIST_GROUP	Citizen's Alliance	1	1	1
ALLIES	advocacy group	0	0	1
ALLIES	Advocate	0	0	1
ALLIES	Backer	0	0	1
EMPLOYEE	Assistant	0	0	1
EMPLOYEE	consult firm	0	0	1
EMPLOYEE	Employee	0	0	1
EMPLOYEE	Worker	0	0	1
ENVIRON	Conservationist	0	0	1
ENVIRON	Ecologist	0	0	1
ENVIRON	Environmentalist	0	0	1

### *A Priori Dictionary Development*

Typical coding schema development for qualitative content analysis includes several iterations, including initial theoretical development and then cyclical revision of codes. An initial, theory-based, deductive development of coding schema should reflect terms and categories that are directly relevant to the framework used. This stage of development will result in an a priori starting list of search terms and categories reflecting concepts or variables in research questions, the literature review, and hypotheses (Miles et al., 2020). Although every search term may not

occur in the analyzed content, the value of this exercise extends beyond the results of the analysis by helping the researcher define the bounds of their study. As is the case with the development of qualitative coding schema for research, the initial development of dictionaries will yield categories that need to be merged, eliminated, or further refined. An assessment of the categories includes consideration of the level of detail, reflection on theoretical concepts, and structural unity (Miles et al., 2020). This assessment can occur after the development of the initial dictionary as well as during contextual checks after the first round of tagging.

NPF Character and Policy Actor Dictionaries. Because NPF characters are identified per chosen unit of analysis (e.g., article, Tweet, paragraph, sentence), an automated method must allow for fluidity in character assignment. For instance, the Colorado Oil and Gas Association can be cast as a villain in one article, play no NPF character role in another, and be a hero in a third. The approach to investigate the portrayal of policy actors is to identify sentences that contain NPF character search terms, tag the terms, and then tag the associated actors. For example,

villain      industry

*That is, a good-looking representative of a villainous gas company would dupe the townspeople into selling him their mineral rights, only to repent after deciding that his employer was bad and fracking, as it is known, potentially worse. (Wines 2013)*

A priori dictionary development for NPF character identification will include categories named for the characters defined in the NPF literature. The categories can include the standard character list of *villain*, *victim*, and *hero* or an extended group (i.e., *opponents*, *allies*, and *experts*). The a priori dictionary may be populated with search terms for each of the NPF characters derived from synonyms related to their theoretical definitions. For example, “criminal,” “perpetrator,” and “crook” could be (and are in this study) search terms that may be used for the *villain* category (see Table 2 and Appendix D for more examples).

To be able to connect policy actors to NPF characters, a policy actor dictionary must also be developed. Although much of this dictionary will be populated with inductive supplementation as discussed in the next section, it is valuable to consider the types of terms that will fill relevant categories of actors in the policy domain. Consider categories of actors that are relevant to the policy domain and the degree to which these categories need to be divided. For instance, depending on the goal of the study, *government* actors might be divided into the subcategories of *international*, *federal*, *state*, and *local*.<sup>2</sup> A priori terms representing the categories could potentially include both general terms such as “resident” or “spokeswoman” as well as proper names such as “Flatirons Responsible Energy” or “National Wildlife Federation.”

Policy Referent Dictionary. To be effectively studied with the NPF, media must qualify as a policy narrative—which, by definition, should include a policy referent and at least one character. To ensure that the dataset only contains policy narratives, a policy referent dictionary may be

---

<sup>2</sup> If the researcher chooses to code using general and subdividing categories containing overlapping search terms such as “government” and “local government,” the more general term tag may be dropped when a double-tagging occurs. In *R*, this may be accomplished using decision-making structures (i.e., if-else or case statements). However, this can often be avoided by increasing specificity of the search term such as using “the government” and “local government” rather than “government,” as was done in this study.

constructed. This is an optional step as the sorting may otherwise be done manually when the media is in *Phase 1: Selection, Collection, and Cleaning*. Even if this is done manually, the construction of this dictionary will assist the researcher in the sorting process because they will have increased clarity about what media qualifies as a policy narrative. A policy referent dictionary will contain words that indicate the presence of a policy discussion, often reflecting processes or actions that citizens or government could take. Potentially, some overlap with the solution dictionary can occur if one is being designed.

*Table 2. Design approach and a priori dictionary examples*

Dictionary	Dictionary Design Approach Example	Examples of A Priori Dictionary Terms
NPF Characters and Policy Actors	NPF VILLAIN: Synonyms about those that cause the problems in stories (Jones and McBeth 2010).  INDUSTRY: General character phrases and proper names referring to characters that are part of the oil and gas industry.	criminal, perpetrator, crook  oil and gas industry, gas developer, industry official, Anadarko
Policy Referent	Terms that indicate the presence of a policy discussion, primarily reflecting policy processes or actions that citizens or government could take. This dictionary may heavily overlap with a solution dictionary.	tax, vote, political, ban, legislature, governor, deregulation
Solutions	Terms that suggest policy solutions, both general and specific. Some terms may be the same as those in the policy referent dictionary. Often indicated by the presence of certain types of words such as modal auxiliaries (“need to,” “should”), superlatives (“is better than”), and present tense verbs (“to improve”).	taxation, subsidies, regulation, law, mobilize voters, Paris Climate Accord
Issue Frames	CLIMATE CHANGE: Words or phrases that are indicative of the discussion of an issue frame within a policy domain, in this case, climate change.	greenhouse gas, global warming, carbon dioxide

**Solution Dictionary.** Researchers may seek to identify policy solutions, which the NPF literature also calls the “moral of the story” or “call to action” (Shanahan et al., 2018, p. 336). Solutions may include proposed specific or general policy fixes. The study design will define how these solutions are categorized. For instance, a basic codebook on solutions may solely seek to determine the absence or presence of solutions, in which it would not be necessary to categorize the solutions. However, for purposes of classifying solutions, the level of government, dominant

policy instrument, or other solution characteristics may serve as broad categories in which to capture a varying range of solutions. In conceptual a priori dictionary design, it may be important to consider attaching actions to general or specific policy solutions. To increase the precision of solution identification, these terms may need to include present tense verbs (“go vote”) but may also include modal auxiliaries (“need to,” “should”) and superlatives (“is better than”).

Frame Dictionary. Initial development of a frame dictionary is based on the conceptual definition of frames, the choice of which may vary depending on the theoretical basis for the investigation as discussed above. Terms in a frame dictionary should be indicative of the ways in which the frames are discussed, such as terms that invoke certain imagery or narrow a reader’s attention to a specific interpretive lens rather than just a topical focus. For example, the a priori construction of a frame dictionary sometimes employs aspects of a broader topic called news frames (e.g., climate change is a topic, but it could be narrated by using various frames such as *economic, health and safety, environment, political*). We chose to focus on issue frames in this analysis and developed the dictionary to contain narrower subcategories within the oil and gas policy domain (e.g., *property rights, jobs, worker safety, climate change, air quality*). There are other approaches to defining frames that could be substituted, but it is important to build the frame dictionary at a conceptual level. The terms included in the dictionary must be mutually exclusive in that they are indicative of only one frame.

#### *Inductive Dictionary Supplementation*

A second step of dictionary development employs an inductive search process in which meaningful words and phrases that were not identified in the initial design are added to the dictionary. With the help of text mining techniques, the analyst can search the corpus, selecting terms that fit the conceptual reasoning of the study. With a large text corpus, such as the one in this study that included the text of 5,708 articles, several techniques may assist the researcher including the generation of parts-of-speech lists, word and phrase frequencies, and topical associations. Inductive dictionary development can be accomplished using qualitative analysis software (such as MAXQDA Pro, which contains the MAXDictio program) or a programming language (such as R—as is done in this research). The generated lists are visually searched for words and phrases that fit the dictionary or categorical criteria, which are then added to the search items and assigned a category. If uncertain, viewing a word in context may be helpful with term selection. Table 3 includes examples of software tools that will be able to assist in these searches.

Table 3. Inductive supplementation assistance tools

Inductive Approach	MaxQDA/MaxDictio <sup>a</sup> Function	R Software <sup>b</sup>	Description of output
Parts-of-speech analysis	indirectly with word frequencies function (using case sensitivity or word ending rules <sup>c</sup> )	<i>openNLP</i>	identification of proper nouns, verbs, adjectives, superlatives, modal auxiliaries
Word and phrase frequency	word cloud (list mode) or word frequencies function, word combinations	<i>tm</i> , <i>tidytext</i>	frequent terms and term counts in text corpus
Topical associations	word tree (visual)	<i>tm</i> , <i>topicmodels</i>	relationships between terms
Term in context	keyword-in-context	<i>grep</i> (base), <i>which</i> (base), <i>str_which</i> (stringr)	phrase in segment
Regular expressions (regex)	N/A	<i>grep</i> (base), <i>str_extract_all</i> (stringr)	strings that fit a specified pattern

<sup>a</sup> A good resource for working with MAXQDA is Kuckartz, Udo, and Stefan Radiker. 2019. *Analyzing Qualitative Data with MAXQDA: Text, Audio, and Video*. Switzerland: Springer Nature.

<sup>b</sup> Numerous blogs and resources exist for using text analysis programs in R. A good place to start is this free e-book: Silge, Julia, and David Robinson. 2020. *Text Mining with R: A Tidy Approach*. O'Reilly Media. <https://www.tidytextmining.com/index.html>

<sup>c</sup> Lemmatization settings control whether words can be reduced to their root words (e.g., “running” to “run”). Word endings are often important to identifying phrases in various tenses.

One helpful method for generating lists of names and phrases that fit a general pattern is the use of regular expressions or regex, available to those programming in R. For instance, a list of all four-letter acronyms may be found using an expression such as “[A-Z]{4}” in a command like *str\_extract\_all*:

```
str_extract_all("[A-Z]{4}", dataframe$text, simplify= FALSE)
```

Because regex also allows coders to skip words in phrases, grab words of unknown length, or include or disregard punctuation, it can be quite helpful for inductive supplementation as well as tweaking dictionary phrases to capture only the desired forms.

As a supplemental method for frame dictionary refinement, R coders may use LDA to inductively investigate topical associations using the R package, *topic models*. LDA provides scores that show how closely words are positioned in the text. The primary reason that we

recommend utilizing LDA for frame dictionary supplementation is because it automatically detects the patterns and co-occurrences of words. The resulting topic groupings can be interpreted as ways of framing an issue, as LDA identifies and groups specific keywords used as framing devices (Jacobi et al., 2015). The granularity, or level of detail determined by the number of topics, may result in broader “themes” or finer “issues,” so this number should be conceptually driven. There are several general steps to the LDA process:

- 1) Convert text into a document-term matrix. A document-term matrix stores the frequency of terms in the collection of documents using *RWeka* and *tm*.
- 2) Optimize the number of topics. In order to correctly parameterize the model, an iterative “elbow” method may be used to identify optimal clustering of the data. With this approach, the researcher can iteratively increase the number of topics for an estimated range of topics, plot the sum of squared errors, and then observe the topic number at which an “elbow” or kink appears in the plotted data.<sup>3</sup>
- 3) Run the clustering analysis at the optimal number of topics, outputting the results to a spreadsheet for visual search.

Although frame dictionaries should be primarily conceptually designed to ensure construct validity and theoretical alignment, this exercise is useful for correctly associating terms as well as highlighting missing frame categories and terms, which enhances content validity—or the full representation of concept indicators. Construct validity—which addresses whether we are measuring what we intend to measure—may be more difficult to attain in automated content analysis. However, the researcher can make significant improvements with the following refinement activities in *Phase 4* such as manually checking the coded segments.

#### *First Autocoding*

The next step is a test-run of the dictionaries on the text corpus. MAXDictio is specifically designed for the application of dictionaries to a dataset of text, though *R* programmers may adjust the author’s software in Appendix F for their own use. The tagging of the text will result in the output of a spreadsheet (MAXDictio) or a dataframe with columns containing text segments, located terms, and/or their dictionary categories (Figure 2).

---

<sup>3</sup> Note that this method typically requires increased amounts of RAM to be able to analyze large text corpora. This study used 64GB of RAM.

Figure 2. Sample output after tagging for a) MAXDictio and b) R code that uses dataframes

a)

	A	B	C	E	H
1	Color	Document name	Code	Segment	Document group
	●	2007-02-04_Outdoor_groups_rally_behind	Char_CO\CHAR_GOVTS_STATE	A state lawmaker said he hopes to create a national model for balancing wildlife protection and energy development when he introduces a bill laying out guidelines for softening the impact of oil and gas drilling.	TimesCall
2					
3	●	2007-02-04_Outdoor_groups_rally_behind	Char_CO\CHAR_ALLIES	Supporters say 55 environmental, hunting and fishing groups are behind the proposal.	TimesCall
	●	2007-02-04_Outdoor_groups_rally_behind	Char_CO\CHAR_GOVTS	The guidelines include reducing the amount of land disturbed by development; speeding restoration; and encouraging consultation between energy companies, landowners and wildlife officials.	TimesCall
4					
	●	2007-02-04_Outdoor_groups_rally_behind	Char_CO\CHAR_INDUSTRY	The guidelines include reducing the amount of land disturbed by development; speeding restoration; and encouraging consultation between energy companies, landowners and wildlife officials.	TimesCall
5					

b)

date	paper	state	SorL	text	Category	word
2008-01-15	PittsburghPostGazette	Pennsylvania	L	"By putting the extraction plant in Pittsburgh, we are giving...	RESIDENT	farmer
2008-02-22	PittsburghPostGazette	Pennsylvania	L	The decision was hailed yesterday by environmental groups...	EMPLOYEE	manager
2008-02-22	PittsburghPostGazette	Pennsylvania	L	"At long last the Washington office has told the local Alleg...	EMPLOYEE	manager
2008-02-22	PittsburghPostGazette	Pennsylvania	L	The decision was hailed yesterday by environmental groups...	ENVIRON	environmental group
2008-02-22	PittsburghPostGazette	Pennsylvania	L	The decision settles 80 appeals filed about drilling in the for...	ENVIRON	environmental group
2008-02-22	PittsburghPostGazette	Pennsylvania	L	Drilling in the Allegheny National Forest is possible becaus...	GOVT_FED	federal government

### Phase 3: Subsetting to Only Policy Narratives

NPF analysis is limited to application on policy narratives. As such, a dataset may need to be reduced to only textual data that includes a policy referent. This may be done manually during the data collection. To this end, the researcher may develop a policy referent dictionary to assist in the subsetting process. As modeled in Figure 1, this process includes a priori dictionary development and inductive supplementation followed by the first autocoding run (discussed in *Phase 2* above). After the first coding run is complete, the researcher may subset the data to exclude those texts not containing policy referents. At this point, an additional step may be included that involves checking the excluded textual data for policy referents by running a word/term-frequency analysis on only the excluded texts. The resulting list may be searched for additional policy referent terms, which may then be added to the dictionary. If this optional step is taken, the data should be autocoded and subset again, potentially yielding a slightly larger dataset. However, if there are no policy referents in the word frequency list, the step of autocoding again may be skipped.

### Phase 4: Dictionary Refinement and Final Autocoding

Often manual qualitative coding relies on intercoder reliability to provide a measure of internal validity after a codebook has been refined and self-tested by the researchers. However, in an autocoding process using dictionary-based methods, the validity of the results is reliant on the validity of the dictionary. Dictionaries should be checked during construction for context and conceptual coherence, particularly if they are used to assign meaning. We briefly discuss two

different dictionary types—counting and interpretive dictionaries—and approaches to increasing internal validity before the final autocoding of the data.

### *Counting Dictionaries*

Our solutions and policy actor dictionaries are used for tagging nouns, verbs, and proper nouns to identify who, what, where, and counting those occurrences. They are designed to investigate the co-occurrence of elements rather than the symbolism, rhetoric, or messaging embedded in communications. To a much higher degree in the counting than in interpretive dictionary design, the terms are found inductively, and the categories assigned as a result of known information (e.g., policy actor job, solution policy tools). To ensure internal validity for this type of counting dictionary, the terms must fit appropriately into the categories, the terms within one category need to be of similar scale (e.g., all corporations and not individuals and corporations), and cutoff term frequencies must be consistent if every term will not be entered into the dictionary (i.e., terms must occur more than X number of times to appear in the dictionary). This last item, term-frequency cutoff, becomes important as the size of the dataset increases.

In automated content analysis, although every term may not be captured, the point of using a large size dataset is to capture more generalization (possibly at the expense of nuance). For example, the cutoff for a proper name to be added to the policy actor dictionary must be above 25 times in this study. This helps the researcher to balance time input into dictionary design as well as reduce the unintentional weighting of one category, such as might occur if adding every named individual. Additionally, in the NPF, the names of policy actors in the narratives must mean something to readers to be valuable to a characterization. However, due to the journalistic norm of adding an actor's role in the same sentence as their proper name (i.e., “Bob Jones, Lafayette resident...”), those actors are not missed in the counting. Thus, the less familiar actors should still be captured by terms for their general policy actor category (e.g., “resident,” “scientist,” “advocate”).

### *Interpretive Dictionaries*

With dictionaries that provide assignments of meaning to text, such as the NPF and frame dictionaries in this study, internal validity is important, and high intercoder reliability (>95%) is necessary. Assessing the internal validity of an interpretive dictionary includes meeting criteria of both appropriate meaning and mutual exclusivity: 1) Is the term meaningful to the category (i.e., does “criminal” always imply the presence of a *villain* for all coders)? 2) Does the term only belong to one category (i.e., does “enemy” fall into both *opponent* and *villain* categories)? With these criteria in mind, the most effective way of analyzing the level of agreement is to provide dictionary terms in random order removed from their assigned categories. Provide the categories available with their definitions as well as the category *none* to allow coders to reject the term or indicate its overlap into several categories. Researchers may be tempted towards designing larger dictionaries. However, small but meaningful and mutually exclusive interpretive dictionaries are valuable with the large datasets that are typically chosen for automated narrative analysis research—enough instances of a category should be found to be able to create reasonable generalizations.<sup>4</sup>

---

<sup>4</sup> Typically, the minimum is considered to be counts above 30.



Another check for internal validity is to run a correlation and clustering analysis after autocoding with the dictionaries. This check is designed to catch potential overlap of categories, which is especially useful in the case of defining frames. Multiple frames may occur in a single story, but clustering analysis is used to analyze frames at the article level. For NPF characters, articles typically include several character categories and thus clustering analysis may be at smaller unit (e.g., journalistic paragraph or Tweet). Clustering analysis requires the number of counted values for each category per unit of analysis provided—typically as a crosstab, matrix, or dataframe. An example of the results of such an analysis is provided from our study in Figure 4 in the *Hydraulic Fracturing Study Methodology* section below. Correlation values may be interpreted with the generally accepted ranges of the correlation coefficient in Appendix C.

### *Final Autocoding*

After the data have been subset and autocoded with the dictionaries, the coded segments are inspected for conceptual coherence as well as validity. Although conceptual coherence should be relatively clear from inspection of the dictionary categories in the dictionary design phase, terms may be catching segments that reveal unexpected results. This may result from a few design issues including 1) incorrect dictionary settings for word or phrase specificity (e.g., only code the whole word, use case sensitivity, or use only starting letters); 2) specifying the incorrect form of a word (e.g., plural, verb endings); or 3) incorrect regex pattern matching. Time spent on refining dictionary terms and settings will result in higher internal validity in the second, and potentially final, round of autocoding.

### **Phase 5: Subset to Single Frame**

In order to map identified policy actors (or groups of actors) onto NPF characters, we next restrict the dataset to a single frame. The reasoning for this lies in the potential for actors to change NPF character categories in various frames. The process of subsetting may significantly lower the size of the dataset, so the researcher may briefly investigate the descriptive properties of each category to be investigated. Primarily, is the number of data points per character category still adequate to be able to apply rules of normalcy and create reasonable generalizations?<sup>5</sup> Not every NPF character category may be fully represented in every frame as some frames are less contentious than others. However, for this single frame, the researcher may choose to use just the characters that are adequately represented.

### **Phase 6: Mapping of Policy Actors into NPF Character Sentiment Ranges**

With the expectation that the construction of policy actors and characters will differ depending on the frame, we propose a method of relating policy actors to their associated NPF character roles within each frame. The sentiment analysis of sentences containing characters leads to a large volume of scores for each actor and character, which can be related to each other using descriptive and inferential statistical tools. There are several steps to this method: 1) evaluate sentiment in segments tagged as including NPF characters and policy actors, 2) use descriptive statistics to

---

<sup>5</sup> This is generally considered to be  $n > 30$ .

investigate the distribution of the sentiment scores, and 3) relate measures of central tendency and hypothesis testing in order to “map” policy actors onto NPF characters.<sup>6</sup>

1) *Evaluate sentiment in sentences within a frame for all NPF characters and policy actors*

In the previous autocoding process, text segments have been tagged with character/actor categories and subset into a single frame. Sentiment analysis is most meaningful at the sentence level but can be evaluated at a slightly larger level of analysis such as the journalistic paragraph, Tweet, or caption. Beyond the paragraph level, sentiment analysis is not a useful exercise unless the analyst is just looking for understanding tone as positive or negative.

Sentiment analysis is used to evaluate the context around NPF characters and policy actors in this study. Sentiment analysis refers to the identification of emotions, positive and negative opinions, and evaluations within text. The reason for employing sentiment analysis in addition to our inductive dictionary design is to understand the emotion portrayed about each policy actor. One core assumption of the NPF is that “meaningful parts of policy reality are socially constructed” (Shanahan et al., 2018, p. 333). To analyze this social assignment of meaning to actors, objects, and processes in an automated fashion, an appropriate method needs to include evaluation of the words and phrases constructing the subject. Sentiment analysis addresses the socially constructed aspect of character interpretation by evaluating the context around a character or policy actor.

Although there are numerous sentiment analysis programs available, this analysis uses the psychology-oriented Harvard-IV dictionary as used in the General Inquirer software, which is a simple polarity analysis using a dictionary of words associated with negative (2,005 words) or positive (1,637 words) (Stone, 2018). These words are oriented towards psychological concepts, as well as evaluations such as *important*, *entrepreneurial*, and *failure*. The *Sentiment Analysis* package for R allows for the analysis of sentiment of a single word up to a whole article (this study assigned sentiment scores to sentences). The result of sentiment analysis on a segment of text with the General Inquirer dictionary yields a score that is between a negative pole (-1) and a positive pole (1). For example, a sentiment score of 0.133 (0.200 positivity minus 0.067 negativity) was assigned to this sentence: *Paul’s light protests have been part of his and the Boulder County Protectors’ effort to stop the oil and gas industry from drilling in Boulder County* (Arvenson, 2017).

This approach evaluates the linguistic context of the sentences that characters are found in. If characters are most often nested in sentences with negative sentiment, whether it is actually a direct construction about that character or not, it is arguable that the sentiment of that sentence is passed to the reader and associated with the character. Additionally, threats to the internal validity when using a dictionary approach (e.g., words may be inappropriately scored based on their context or include negation terms) may be minimized by the addition of sentiment analysis. Sentiment analysis also accounts for a larger portion of the identified words in a sentence, rather than depending on the correct interpretation of one or two dictionary words.

---

<sup>6</sup> For this phase, it is likely that the researcher will have to rely fully on R rather than on MAXQDA unless they are analyzing Tweets (MAXQDA does have a function for analyzing the sentiment of Tweets, but it will not assist with the mapping method described here).

2) *Use descriptive statistics to investigate the distribution of the sentiment scores*

After running the text through sentiment analysis and storing values with their respective sentences (in a new dataframe column in *R*), the next step is to investigate the descriptive statistics of each category. *FBasics* is an *R* package that may be used to summarize descriptive statistics of the sentiment for each policy actor and character category within a frame. Visual analysis using probability distribution functions allows the researcher to decide which measures of central tendency are most appropriate for comparison and whether statistical methods designed for normal distributions will be appropriate for application.

Distributions may very well peak off the central sentiment value of 0.0 or be heavily skewed, yet still be representative. However, if they are multimodal—which implies the presence of several populations—central distribution values may be misleading. The presence of multiple populations may indicate a shift in sentiment over time, poorly constructed dictionaries, or representation of multiple communities (e.g., mixing news sources that may construct policy actors differently). If multimodal distributions are present, consider investigating this as part of the analysis as it may yield increased nuance in your results.

3) *Relate measures of central tendency and hypothesis testing in order to “map” policy actors onto NPF characters*

To map the NPF characters to policy actors as would be the purpose in a manual NPF coding, Welch two-sample t-tests will establish whether coded *villains*, *victims*, and *heroes* (or, additionally, *allies*, *opponents*, and *experts* as we use in this study) hold significant relationships with particular policy actors. When NPF characters occupy a discrete space in the range of sentiment distributions, policy actor sentiment scores from the same group of articles may be mapped to them with some confidence. Based on this concept, this method identifies similarities of sentiment between policy actors and NPF characters by identifying those that hold the same sentiment (Welch t-tests that show an acceptance of the null hypothesis).<sup>7</sup> To demonstrate the mapping method on individual frames as well as analytical interpretation techniques, we walk through examples in the *Results* section of this study.

*Method Premise*

We use the above three steps to briefly demonstrate the statistical premise for the application of this method. Although the general method suggests subsetting to only one frame, we conduct this part of the statistical analysis across all frames to demonstrate that each NPF character’s distribution of sentiment scores occupies a distinct space that can be used to “map” onto the sentiment scores of policy actors. An investigation into the sentiment distributions of segments coded with NPF characters shows that most are relatively normal distributions, though *allies*, *experts*, and *opponents* show slightly bimodal distributions—multiple bumps rather than one hump (Figure 3). Across frames, multimodality is not totally unexpected for these three characters because they may have slightly varying constructions in different frames. In contrast, *heroes*, *villains*, and *victims* are strongly unimodal likely because they have more universally consistent constructions.

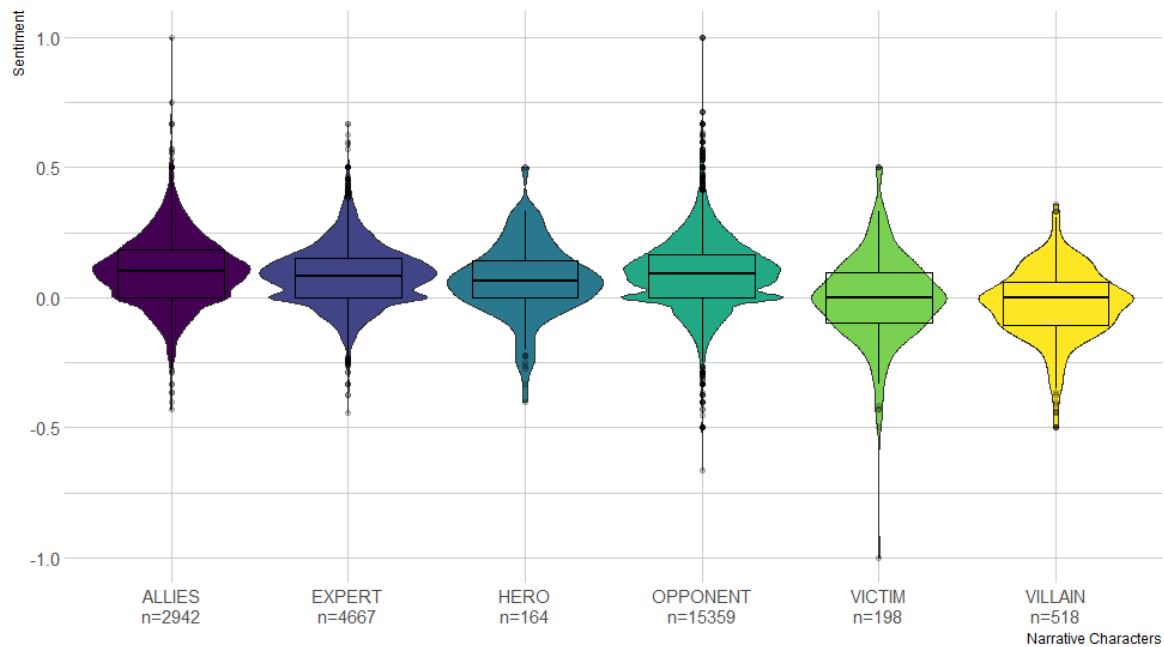
---

<sup>7</sup> For Welch t-tests, the null hypothesis is that the true difference in means is equal to 0.0.

If the sentiment distribution of scores of each NPF character is significantly different from each other, they may be considered distinct. Welch two-sided t-tests performed on each character category of sentiment scores relative to every other character demonstrate that NPF characters have statistically different distributions of the sentiment, except for *heroes* and *experts* above the 90% confidence level (Table 4). Considering the commitment of journalists to report with low levels of bias, it is conceivable that *experts* may be equivalent to *heroes* in news policy narratives, particularly when it comes to more technical, less value-driven frames, like those about wastewater injection or air quality in hydraulic fracturing for shale gas. Although the table shows a clear distinction between categories, within one frame the separation between characters has the potential to be even more clearly defined.

Additionally, the distinction between characters may be increased depending on the type of media. With our dataset of news articles, the frequencies of *allies*, *experts*, and *opponents* (n=2942, 4667, and 15,359 respectively) relative to the low frequencies of traditional NPF characters suggests three possibilities: 1) the dictionary-based terms for *heroes*, *villains*, and *victims* may not pick up all the subtleties of those characterizations; 2) journalists avoid using extreme characterizations; and 3) journalists frequently present two sides of an issue as well as an expert to provide “ground truthing” with evidence or expertise. However, if a researcher were using a sample containing op-eds, Tweets, or other intentionally biased text, the differences in sentiment score for characters would very likely be even more pronounced than found here.

Figure 3. Character sentiment distributions with mean and 1st and 3rd quartiles



*Table 4. Two-sided t-tests of NPF characters*

	Mean difference $\mu_1 - \mu_2$	Std. Error	t-statistic	Pr(> t )	Significance
ALLIES-EXPERT	0.023	0.003	7.194	0.000	***
ALLIES-HERO	0.037	0.012	3.038	0.003	**
ALLIES-OPPONENT	0.015	0.003	5.337	0.000	***
ALLIES-VICTIM	0.125	0.011	11.376	0.000	***
ALLIES-VILLAINS	0.148	0.007	19.856	0.000	***
EXPERT-HERO	0.015	0.012	1.204	0.230	No
EXPERT-OPPONENT	-0.008	0.002	-3.639	0.000	***
EXPERT-VICTIM	0.102	0.011	9.438	0.000	***
EXPERT-VILLAINS	0.126	0.007	17.339	0.000	***
HERO-OPPONENT	-0.022	0.012	-1.853	0.066	+
HERO-VICTIM	0.088	0.016	5.442	0.000	***
HERO-VILLAINS	0.111	0.014	7.950	0.000	***
OPPONENT-VICTIM	0.110	0.011	10.246	0.000	***
OPPONENT-VILLAIN	0.133	0.007	18.792	0.000	***
VICTIM-VILLAIN	0.023	0.013	1.812	0.071	+

Note: = $p < 0.1$ ; \* $p < 0.05$ ; \*\* $p < 0.01$ ; \*\*\* $p < 0.001$

## HYDRAULIC FRACTURING STUDY METHODOLOGY

To illustrate how the use of automated tools—and specifically the method described here—can advance policy narrative research, we autocoded 5,708 state and local news articles focusing on hydraulic fracturing of oil and gas. We used this automated method to investigate the relationship between narrative components—namely, characters and proposed solutions—and the frames that tie policy narrative elements to one another. In this section, we provide details of our use of the textual analysis process described above.

## Data

The dataset utilized in this paper consists of news articles collected from six state and six local print news sources from January 1, 2007 to December 31, 2017.<sup>8</sup> The initial year of 2007 was chosen because that is the year that horizontal hydraulic fracturing technology allowed for rapid expansion of oil and gas development in the United States. The news sources were chosen based on the circulation level but limited to only those newspapers that included online access or have archives available through Newsbank or Nexis (online news repositories) for the entire sampling period. The dataset is composed only of articles related to the development of unconventional oil and gas. The search terms “hydraulic fracturing” and “fracking” were used to filter the articles to be able to limit the sample to only those containing references to shale oil and gas development using hydraulic fracturing. Six states—comprising the top producing states in the United States—were selected based on their estimated total production of oil and gas in 2017 as reported by the US Energy Information Administration. These states are Texas, Pennsylvania, Oklahoma, Louisiana, North Dakota, and Colorado. With filtering by date and search terms, 20,417 articles were collected from the 12 newspapers.

Several data filtering and cleaning processes were implemented: opinion pieces and letters to the editor were removed and data were screened for duplicate articles. Byline material, captions, and titles were used to catalog as metadata information in a separate spreadsheet but were removed from the text for automated analysis. After these cleaning processes were applied, 6,008 articles were available for the analysis. However, after the policy narrative subsetting process (described below), the total number used in the analysis decreased to 5,708. Appendix A lists the newspapers chosen to represent state and local news sources (based on circulation and archive availability) as well as the final article count from each source.

The unit of analysis in this research is the article, though the unit of observation is the sentence within a policy narrative. The primary reason for aggregating up to the article level is that frames are typically presented at the article level. However, multiple frames can occur within a single policy narrative. Although characters are identified and their sentiment evaluated at the sentence level, results are aggregated to the article level. When several characters occur within one sentence, both characters are assigned the same sentiment.

## Dictionaries

For this study, dictionaries were designed for the identification of policy narratives, NPF characters, frames, and solutions. The dictionary designs followed the processes detailed in the previous sections using a priori theoretical design, inductive supplementation, and refinement. Due to the large volume of words with the text corpus, all terms that occurred with a frequency of 0.4% and over (occurring in more than ~25 articles) were considered for use in the dictionaries, though some were included that occurred at lower frequencies because they were generated in more specific lists (i.e., lists of proper nouns or acronyms). After an initial autocoding of the text, words

---

<sup>8</sup> The news articles were originally collected for a larger study on policy debates related to shale oil and gas development as described in (Berardo et al., 2020). Funding support for the larger study came from the National Science Foundation Award Foundation, Decision, Risk, and Management Sciences Program (grant SES-1734310 and SES-1734294).

and phrases were examined in the context of their sentences for an additional determination of their validity for the dictionary and further dictionary refinement.

### *Policy Referent*

To ensure that our dataset of articles only included policy narratives—which must include a policy referent and at least one character—a policy referent dictionary was created. The dictionary was designed to recognize policy discourse, attempting to detect both general and specific policy solutions as well as political behaviors, actors, or venues. The initial dictionary was a product of brainstorming dimensions of politics (e.g., “vote,” “election,” “tax,” “senator”). The terms were then supplemented inductively by using a word/phrase frequency program, yielding more general terms, but also proper names of specific policies (e.g., “Safe Drinking Water Act”). This dictionary contained a number of overlapping terms with the solution dictionary. However, unlike the solution dictionary, because its purpose was solely to determine *whether* a policy discussion was occurring to either include or exclude it from the study, no categories were necessary.

### *Issue Frame Dictionary*

*As discussed earlier in this chapter, for this study we selected issue frames as the frame-type of our focus. The initial issue frame dictionary was developed with general guidance from the theoretical definition provided in Chong and Druckman (2007). Initial terms represented the key expected areas of discussion related to energy production and development. These terms, derived from brainstorming, were used to create categories and key terms. Next, the dictionary was inductively supplemented using frequency lists, and an LDA analysis was used to check our topical associations. In order to optimize the number of topics to parameterize the model, the elbow method was used to find an optimal number of 13 topics. While the results of the LDA model did not determine our final categories, this exercise was useful in associating terms as well as showing us issue frame categories and terms that were missing from our inductively created dictionary. After dictionary category and term refinement with LDA, we found a total of 18 categories. Table 6 provides example terms for each issue frame (category or subcategory) and identifies how the issue frames align with broader “news frames” (as defined earlier in this chapter) for the group.*

### *NPF Characters and Policy Actors*

In this study, NPF characters and policy actors have different emphases—the character dictionary assists with interpretation while the actor dictionary is primarily for counting. However, they both contribute to the same end goal: relating the two to understand how policy actors are characterized in policy narratives.

The general approach to the dictionary design along with examples of each category are detailed in Table 5. The NPF dictionary categories and terms were based on synonyms of the NPF character definitions (Table 5), which were then inductively supplemented from word frequency lists. The policy actor a priori dictionary was created by brainstorming and categorizing the actors involved in the issue area of hydraulic fracturing of oil and gas. The policy actor types used in this study are entity actors, general groups, and named people. Entity actor category codes include *industry*, *industry ally*, *local government*, *state government*, *federal government*, and *activist group*. The terms in these categories are proper names of companies, government agencies, or groups. Individuals associated with these groups were included as named people with a code

associated with their role (e.g., Governor Hickenlooper would be categorized as a *named state government* actor). General group category codes include *environment*, *resident*, *government*, and *employee*. As the type implies, these categories are populated with general terms rather than proper names such as *environmentalist*, *rancher*, and *manager*.

All search terms in this character and policy actor dictionary are mutually exclusive; however, it is possible that one character may be tagged with a title (captured here in the entity categories) and named people category code within a single sentence. For instance, a sentence containing a named person with their title “US Senator for Colorado, Cory Gardner, said...” would be tagged as both *state government* and *named state government*. Thus, care was used with interpretation of the named people category code analysis.



Table 5. NPF Character and policy actor dictionary development

Group	Dictionary design approach	Character Code	Examples of dictionary terms
<b>ENTITY ACTORS</b>	General character phrases and proper names referring to characters that are part of the oil and gas industry.	INDUSTRY	Anadarko, oil and gas industry, gas developer, industry official
	General character phrases and proper names referring to characters that are allies of the oil and gas industry.	INDUSTRY ALLY	trade group, Flatirons Responsible Energy, Oil and Gas Association
	General character phrases and proper names referring to characters that can occur at the local level of government	GOVT LOCAL	city council, county commissioner, mayor
	General character phrases and proper names referring to characters that can occur at the state level of government.	GOVT STATE	health department, state lawmaker, state legislator, Water Resource Board
	General character phrases and proper names referring to characters that can occur at the federal level of government.	GOVT FED	Obama, Bush administration, federal land manager
	Only proper names of activist groups identified to be opposed to oil and gas development.	ACTIVIST GROUP	National Wildlife Federation, Wilderness Society, Stop Fracking Wayne County
<b>GENERAL GROUPS</b>	Descriptions of general characters specifically referencing environmental alignment.	ENVIRON	environmentalist, preservationist, tree hugger
	General character phrase referencing resident or community members.	RESIDENT	tenant, locals, rancher, Cherokees, community member, homeowner
	General terms referencing government. These terms could apply to various levels of government.	GOVT	Senator, governor, bureaucrat, legislative committee
	General character phrases referencing employee or workers.	EMPLOYEE	manager, spokesman, spokeswoman, political consultant
<b>NPF CHARACTERS</b>	General character phrases referencing experts.	EXPERT	scientist, geologist, professor
	Synonyms about those that cause the problems in stories (Jones and McBeth, 2010).	VILLIAN	criminal, perpetrator, crook
	Synonyms about those that fix or attempt to fix the problems in stories (Jones and McBeth, 2010).	HERO	hero, winner
	Synonyms about those that are harmed by the problems in stories (Jones and McBeth, 2010).	VICTIM	victim, casualty, fatality, injured party
	An individual, organization, or government entity that is not explicitly blamed, but is identified as holding a policy position with which the author disagrees (Merry, 2016).	OPPONENT	competitor, adversary, antagonist
	An individual, organization, or governmental entity that is not explicitly praised, but is identified as holding a policy position with which the author agrees (Merry, 2016).	ALLIES	proponent, supporter, enthusiast
<b>NAMED PEOPLE</b>	Named people from the following groups: GOVT FED, ACTIVIST GROUP, RESIDENT, EXPERT, GOVT STATE, GOVT LOCAL, GOVT FED, INDUSTRY, INDUSTRY ALLY	Add "SP" at end of group	Chad Warmington, Ken Salazaar

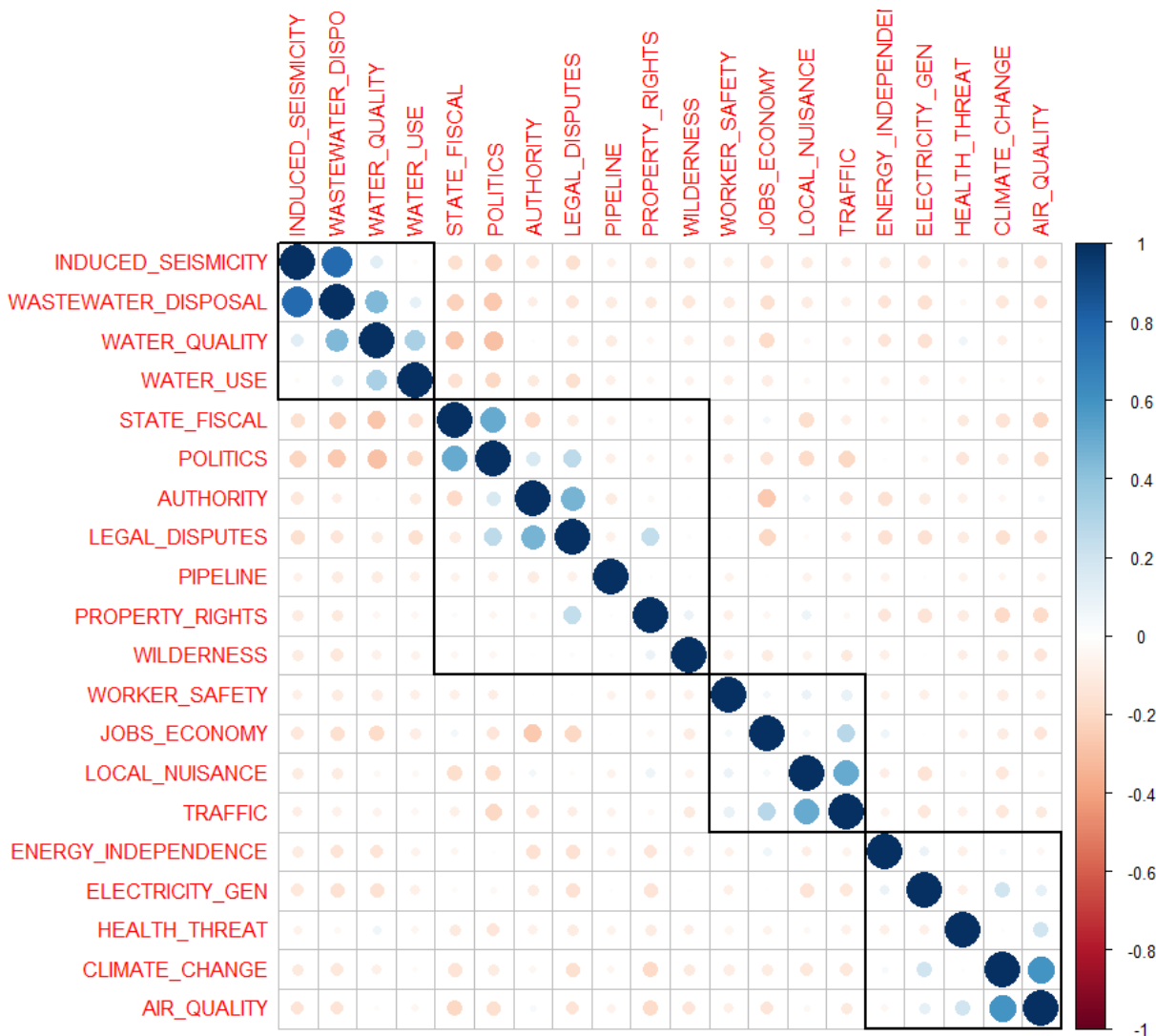
Table 6. Final frame dictionary design

News Frames	Dictionary Categories/Issue Frames	Examples of Dictionary Terms
Water/Disposal	INDUCED SEISMICITY	quake, infrastructure damage, faults
	WASTEWATER DISPOSAL	reinjection, injection well, wastewater
	WATER QUALITY	brackish, toluene, radioactive
	WATER USE	water consumption, amount of water used, drought
Political/Legal	STATE FISCAL	appropriations, budget, fiscal
	POLITICS	protest, public disapproval, policy
	AUTHORITY	moratorium, extension, hearing
	LEGAL DISPUTES	legal, lawsuit, court
	PIPELINE	natural gas pipes, right-of-way, energy hub
	PROPERTY RIGHTS	royalties, mineral owner, landowner
Local Effects	WILDERNESS	wilderness, woodland, prairie chicken
	WORKER SAFETY	OSHA, deaths per, workplace safety
	JOBS ECONOMY	jobs, employment, hiring
	LOCAL NUISANCE	odor, noise, dust
Debate about the Use of Energy Types	TRAFFIC	trucks, highway, driver, collision
	ENERGY INDEPENDENCE	energy security, Middle East, Saudi Arabia
	ELECTRICITY GENERATION	renewable credit, wind production, electricity
	HEALTH THREAT	immune system, respiratory, cancer
	CLIMATE CHANGE	greenhouse gas, methane leak, warming potential
	AIR QUALITY	particulate, ozone, air pollution

Some of these categories co-occurred so often they could potentially be considered as aspects of the same issue frame. After autocoding with a refined version of the frame dictionary, we ran a correlation and clustering analysis to investigate the distinctness of our dictionary

categories. This method uses the tagged category counts, representing issue frames, at the article level. We found a high degree of positive correlation between a number of issue frames (Figure 4) suggesting that we could combine several of these issue categories. For instance, *wastewater disposal* and *induced seismicity* are correlated, as are *politics* and *state fiscal*, *climate change* and *air quality*, and *local nuisance* and *traffic*. However, for the results presented here, we decided to keep these issue frames separate to add a more nuanced investigation into the relatedness of aspects of issues. Figure 4 also reveals the distinctness of inversely-correlated frames (implying that when

Figure 4. Correlation and clustering analysis of frames in autocoded articles



one frame increases the other decreases), such as *jobs and economy* and *authority, politics* and *water quality*, and *wastewater disposal* and *state fiscal*.

### Solution Dictionary

An initial solution dictionary was designed using the NPF concept of “moral of the story”—a policy solution that may end in a call to action (Shanahan et al., 2018). As with the policy referent dictionary, which contains some similarly derived terms, the solution dictionary

may include general steps to achieving a policy goal, a stated policy preference, or named policies. To reduce the complexity posed by the large number of proposed solutions in the news articles, we employed a categorical coding design by policy tool (e.g., *regulatory*, *subsidies*, *taxation*). Though this categorization focused on proposed government solutions, we added a category for those that were oriented towards citizen action. To capture general statements of solutions to policy problems, we created the category, *general solutions*. Identified general solutions included phrases that suggest the need for a solution. For example, statements including a modal auxiliary (such as *need to*, *should*), indicated that they are proposing solutions. Additionally, phrases about the future (i.e., *for our future*), superlatives (i.e., *is better than*), and verbs in the present tense (i.e., *to improve*) usually pointed to solution statements. The dictionary was supplemented with named policy solutions by generating lists of proper names and acronyms.

### Subsetting, Mapping, and Analysis

After a final autocoding, we subset coded segments into separate issue frame categories. After running the segments through sentiment analysis, the resulting data contains segments that are each associated with categorical tags and sentiment scores.

We investigated the relationships between the frames and coded narrative elements using the mapping technique presented above as well as regression. Ordinary least squares (OLS) regression analysis was performed with the package *lm* in *R* in order to address the research question that asks how frames, characters, and solutions are related. There is no assumption of causality with the use of regression in this study, rather it is used to investigate the relationship between frame, character code, and solution code. The general theoretical relationship between variables is  $\text{frame} \approx f(\text{characters} + \text{solutions})$ . The regressions of each frame were calculated using frequency count matrices, which are large tables containing a row for each article with a column for each frame. Because more than one frame identifier can occur within one article, these counts are whole numbers equal to or greater than zero. Similarly, these large matrices were constructed for characters and solutions. The probability distributions of characters, solutions, and frames are each approximately normal. For these reasons, OLS is suitable for this analysis.

Measures of central tendency, such as statistical mean, kurtosis, skew, and quartiles, were calculated for character and policy actor sentiment scores with the R package, *BasicStats*. We also analyzed the Welch two-sample t-tests for characters and policy actors with the base R function, *t.test*. These scores allow for mapping policy actors onto NPF characters for each frame.

## RESULTS

The results of this analysis suggest that frames are often correlated with certain characters and solutions. Using OLS regression analysis, sentiment analysis, descriptive statistics, and hypothesis testing, we find relationships between frames and characters and solutions. We also relate certain policy actors to NPF characters. The general form of the regression is  $\text{FRAME} \sim \beta_0 + \beta_1\text{CHAR1} + \beta_2\text{CHAR2} \dots + \beta_3\text{SOLN1} + \beta_4\text{SOLN2} + \dots$

To investigate the relationship of frames with solutions and characters, we first calculated OLS regression statistics of a single frame against all solutions and characters using frequency

count matrices. As a second step, we evaluate the sentiment associated with characters and policy actors to show how policy actors relate to NPF characters in this semi-automated method of textual analysis. The results of this investigation show that all frames have significant relationships to characters and solutions (Table 7). We briefly discuss several examples of significant findings from all those regressions (see Appendix A for regressions).

An interesting result we mentioned in the *General Method Description* section, is that the NPF characters *hero* and *experts* hold a similar overall sentiment signature, whereas all other NPF characters have a significantly different sentiment distribution (Table 4). This finding suggests that the heroes of news articles may be experts, potentially influenced by the journalistic norm of presenting unbiased information. In the examples presented below for individual frames, we analyzed these characters separately, as they are discussed in NPF literature.

### **Frame: State Fiscal**

An example of a finding that suggests high predictability of elements in a certain frame is the *state fiscal* frame. The predictors explain much of the variance of the issue frame ( $R^2_{\text{adj}}=0.777$ ,  $F(7,5696)=2,834.88$ ,  $p<0.001$ ). The government policy actor codes (significant above the 99.9% confidence level) were relatively predictive of the frame, suggesting that policy narratives containing frames about state budgets have a predictable formula, even across six state and local papers. Based on an understanding of how news articles are constructed, it is highly unlikely that journalists start with a policy actor and *then* create a frame. The interpretation is of co-occurrence, but not of causation (i.e., the independent variables of characters and solutions likely did not determine the chosen frame (the dependent variable)). The significant policy actors in the *state fiscal* frame were found to be *activist group*, *general govt*, *local government*, and *named people* (categories ending in SP in Table 7) at the *state* and *local* levels. However, the coefficients in this regression suggest that as the frequency of some actors increase in a narrative, the frame terms decrease. For example, the coefficients (regression results in Appendix B) for *activist group* and *local government* ( $\beta_{\text{ACTIVISTGROUP}}=-0.296$ ,  $\beta_{\text{LOCALGOVT}}=-0.065$ ) have an inverse relationship with this frame. Therefore, we focus on those with positive coefficients.

Table 7. Regression Variables and Significance

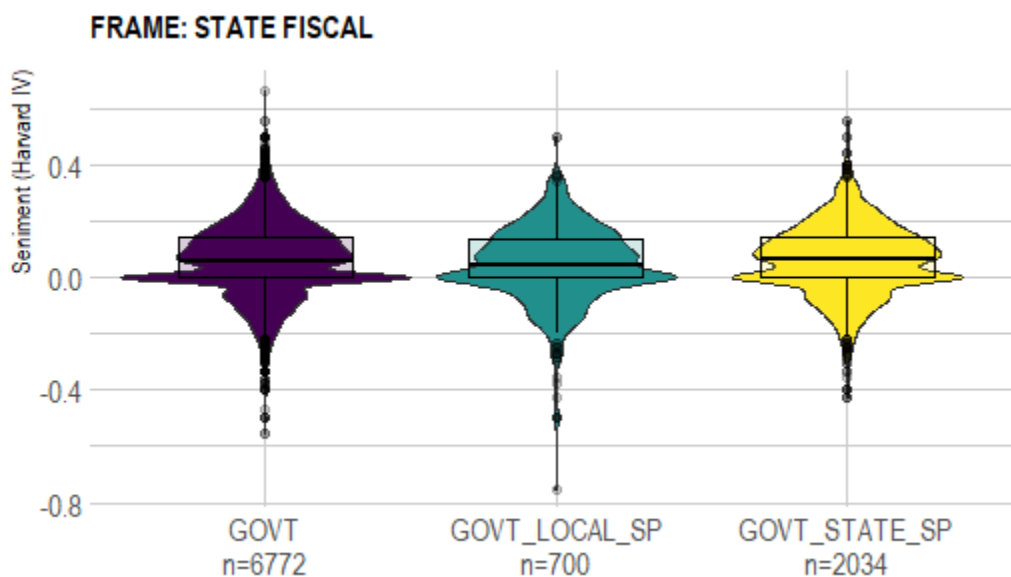
	SOLUTIONS																				
	PROPERTY RIGHTS	STATE FISCAL	JOBS ECONOMY	WORKER SAFETY	CLIMATE CHANGE	AIR QUALITY	HEALTH THREAT	LOCAL NUISANCE	INDUCED SEISMICITY	WATER QUALITY	WATER USE	WASTE WATER DISPOSAL	ENERGY INDEPENDENCE	PIPELINE	POLITICS	AUTHORITY	LEGAL DISPUTES	ELECTRICITY GEN	TRAFFIC	WILDERNESS	
BAN	+					***	+	***					***	***	***	***	***				
CITIZEN ACTION		*	***		+	***	***		*			*		**	***						
GENERAL	***	***	***	***	**	***			***	*				+	***	**				***	
GOALS					***	***		**								**	***				
GOALS SP														+	***						
INFORMATION		*		+	**	***	***		***	***	*	***			***					***	
INFORMATION SP				***			***		**												
PERMIT	**	*				*	**			***	***		**		+					***	
PERMIT SP						+			***		+			***		***				***	
REGULATION		***	**	*	*	***	***	***		*	+	*		***	***	***				***	
REGULATION SP			*						***		***					***				***	
SUBSIDIES	*													+	+	+	***				
SUBSIDIES SP			**		**							*									
TAX	**	***	**		*	*			**			+	*	***	***		*			**	
TRAINING		*	***	*					**											+	
TRAINING SP			+																	+	
INDUSTRY	***	**	***	**			*	*	***	**	**	***		***	+	+	**			***	
INDUSTRY ALLY	*			***					**	+	+		***			***	**				
GOVT LOCAL		***				*			***	*	**	**		**	***	***					
GOVT STATE			**	***	*		***	***	***	***	*		***	***	***	***				**	
GOVT FED	*	**	*		***	***	**	***	***	***	***	***			***					***	
ACTIVIST GROUP	*	***	*		***		**			+	+				*	***		+		***	
ENVIRON	*	**		+	***	***		**	**	***	*				*	*	**			***	
GOVT	**	***	***		*	***		+		**				***	**		**	**		**	
RESIDENT	***	**	***	***	+	+	***	***	**	+				***	***		***	**	**	***	
EMPLOYEE		+	***	***	*			***			**			***		+	+	***	**	**	
INDUSTRY SP							***								+						
GOVT LOCAL SP		***													+	+					
ACTIVIST GROUP SP		*			*									+	*						
GOVT STATE SP		***					+	*	*	+	*		***			**					
INDUSTRY ALLY SP	**			*		***		+	+	+	*		*		**						
EXPERT SP			+		*			***	*		*		***		+						
EXPERT		***	***	***	***	***	***	***	***	***	***		*	**		+	*			*	
ALLIES	**	*	+				***		***	*				***	**	***	***				
OPPONENT			***		***			*	***	**	***	***	***	***		***	***				
VICTIM			+	***			***			**						***	***		***		
VILLAIN		**		+						***			**		***	***					
HERO						***				+		***				***	***				
R <sup>2</sup>		0.15	0.78	0.24	0.19	0.06	0.09	0.02	0.15	0.22	0.08	0.02	0.10	0.05	0.03	0.55	0.51	0.22	0.14	0.33	0.09

Note: = $p < 0.1$ ; \* $p < 0.05$ ; \*\* $p < 0.01$ ; \*\*\* $p < 0.001$

As a first step in relating the policy actors to NPF characters, sentiment analysis was conducted on character-containing sentences within articles that have a *state fiscal* frame. From this process, each sentence containing a character was associated with a sentiment score. The

resulting distributions with mean and 1<sup>st</sup> and 3<sup>rd</sup> quartiles are shown in Figure 5. The boxplots show that in the articles that contain this frame, the *named local government official* has the lowest sentiment ( $M_{GOVTLOCALSP}=0.0450$ ) and *named state government official* ( $M_{GOVTSTATESP}=0.0628$ ) has the highest. Comparatively, the values for the theoretical characters tagged in these articles have a low range suggesting that there is very little polarization in these articles. An interesting finding in this work and Crow and Wolton (2020) is that *allies* have a higher sentiment than *heroes* ( $M_{HERO}=0.0700 < M_{ALLIES}=0.0927$ ). This is a consistent finding in sentiment across all the articles and many other frames as well. Policy actors are not as villainized in the *state fiscal* frame as they are in the whole set of articles as evidenced by the mean sentiment that is closer to zero ( $M_{VILLAIN\_ALL}=-0.0401 < M_{VILLAIN\_FISCAL}=-0.0172$ ). These findings suggest there is less polarization and overt depiction of *heroes* and *villains* in this local frame.

Figure 5. Sentiment distributions for significant characters



The second step in relating policy actors to NPF characters was hypothesis testing of the significant characters against the NPF characters (Table 8). Welch two-sample, two-sided t-tests were performed between the sentiment distributions of each significant actor and the NPF characters from the *state fiscal* frame. This maps the significant actors in a frame, identified at the sentence level with dictionary tagging, to the closest NPF character. Some actors will not be significantly associated with an NPF character in every frame. The hypothesis testing relies on the rejection of the null hypothesis in the case of testing the true difference in means, meaning that the populations cannot be interpreted as significantly different (on the chart this is marked by a *No*—we cannot reject the null—for significance. Table 8 shows that *named state government*, *government*, and *named local government* actors are most similar to the *hero* sentiment population, with the *named state government* actor having the most significant relationship. In this frame, these government actors are associated with moderately constructed *heroes*, while there are no clear actor associations with other characters.

Table 8. T-test results on significant actors against NPF characters

$\mu_1 - \mu_2$	Mean difference	Std. Error	t-statistic	Pr(> t )	Similar sentiment
GOVT STATE SP – HERO	-0.007	0.016	-0.452	0.652	Yes
GOVT – HERO	-0.008	0.016	-0.486	0.629	Yes
GOVT LOCAL SP – HERO	-0.025	0.017	-1.506	0.136	Yes
GOVT LOCAL SP – VICTIM	0.063	0.022	2.847	0.006**	No
GOVT – VICTIM	0.081	0.022	3.716	0.000***	No
GOVT STATE SP – VICTIM	0.081	0.022	3.717	0.000***	No
GOVT LOCAL SP – VILLAIN	0.062	0.014	4.490	0.000***	No
GOVT STATE SP – VILLAIN	0.080	0.013	6.076	0.000***	No
GOVT – VILLAIN	0.080	0.013	6.136	0.000***	No
GOVT STATE SP – EXPERT	-0.025	0.004	-5.880	0.000***	No
GOVT LOCAL SP – EXPERT	-0.043	0.006	-7.091	0.000***	No
GOVT – EXPERT	-0.026	0.004	-7.084	0.000***	No
GOVT STATE SP – ALLIES	-0.036	0.005	-7.684	0.000***	No
GOVT LOCAL SP – ALLIES	-0.054	0.006	-8.459	0.000***	No
GOVT – ALLIES	-0.036	0.004	-8.958	0.000***	No
GOVT LOCAL SP – OPPONENT	-0.051	0.005	-9.346	0.000***	No
GOVT STATE SP – OPPONENT	-0.033	0.003	-9.960	0.000***	No
GOVT – OPPONENT	-0.033	0.002	-14.175	0.000***	No

Note: = $p < 0.1$ ; \* $p < 0.05$ ; \*\* $p < 0.01$ ; \*\*\* $p < 0.001$

The solutions highly related to the *state fiscal* frame include *tax*, *regulation*, and *general solutions* based on regression coefficients ( $\beta_{\text{TAX}}=1.875$ ,  $\beta_{\text{REGULATION}} = -0.128$ ,  $\beta_{\text{GENERAL}} = -0.059$ ). However, while the positive coefficient for *tax* predicts the co-occurrence of the *state fiscal* frame, both *regulation* and *general* solutions have a relatively strong negative relationship with this frame. This relationship implies that as the occurrence of the *regulation* and *general* solutions decrease and the *tax* solution increases, the likelihood of the frame being *state fiscal* increases. Although the regression coefficient implies positivity about the *tax* solution, the coding method used here is not evaluative of whether the solution is accepted. The support for solutions in the coded sentences can vary:

*Mr. Rendell's prediction that the tax would bring in \$1.7 million in its first year, Mr. Rhoads said, is "dead wrong."* (Eagle, 2009)



*Oil and gas companies, however, argue that the tax will have the opposite effect, harming Pennsylvanians by destroying the opportunity for thousands of well-paying jobs. (Eagle, 2009)*

*While state legislators continue to debate the need for a tax on natural gas drilling, proponents of the tax yesterday concerned themselves with where the revenue should go. (Eagle, 2009)*

As previously shown in Figure 4, the *state fiscal* frame is moderately correlated (0.3) with the *political* frame and mildly anti-correlated (-0.09) with the *authority* frame, a result supported by the coefficient findings for *tax* and *regulations*. The *state fiscal* frame was found in 1,937 articles—though not necessarily the dominant frame—indicating that is an important theme in the discussion of shale development in the newspapers evaluated. The high occurrence of this frame category and the high coefficient of determination associated with the use of these actors and solutions suggest a similarity of the frame and the use of actors and solutions across states and over the 10-year period.

### **Frame: Wastewater Disposal**

The OLS regression results also indicate that some issues are associated with a larger number of narrative elements, including NPF characters. The *wastewater disposal* issue, which is highly correlated (0.51) with *induced seismicity* and moderately correlated with *water quality* (0.28) and *water use* (0.19), contains significant relationships with the NPF characters of *villains*, *victims*, and *experts*. Three solutions variables (*information*, *permit*, and *named regulation*) and actors (*state government*, *experts*, *government*, *government local*, *government federal*, and *villain*) mildly explain the variance of the issue frame at the 99.9% confidence level ( $R^2_{adj} = 0.082$ ,  $F(9,5698) = 56.514$ ,  $p < 0.001$ ). The strong positive correlation of the wastewater disposal frame to other frames (see Figure 4) suggests that there are various ways of talking about this issue. Comparing this frame to the relatively compact *state fiscal* frame suggests that some frames are associated with greater variability in character use.

Additionally, the regression results show the use of more NPF characters, which indicate more polarization in these narratives. In fact, the NPF *villain* character has the highest coefficient of characters that explain variability ( $\beta_{VILLAIN} = 0.227$ ,  $\beta_{EXPERT} = 0.163$ ,  $\beta_{STATEGOVT} = 0.106$ ,  $\beta_{FEDGOVT} = 0.064$ ,  $\beta_{GOVT} = -0.032$ ,  $\beta_{LOCALGOVT} = -0.051$ ) at the 99.9% confidence level. The sentiment scores for the characters support the idea of increased narrativity of *wastewater disposal* policy narratives. Sentiment analysis was conducted on character-containing sentences within articles that have a *wastewater disposal* frame. The results are shown in Figure 6. In the articles that contain this frame, the figure shows there is often a relatively negative *villain* compared to the *villains* in all of the articles ( $M_{VILLAIN\_ALL} = -0.0261 > M_{VILLAIN\_WWD} = -0.0485$ ).

Yet, at the 99.9% significance level, *villains* were not associated with the *industry* or *opponent* actor, which is significant to the model at the 95% confidence level. Interestingly, *villains* in these narratives are not highly associated with any particular policy actor. A partial explanation may be that policy actors are reluctant to name *industry*, which is constructed very positively in the dataset ( $M_{INDUSTRY\_ALL} = 0.1020$ ) and has a similar construction to the character *allies* ( $M_{INDUSTRY\_ALL} = 0.1020 \sim M_{ALLIES\_ALL} = 0.1083$ ). Examples of segments containing the

*villain* character show the reluctance for villainization of the industry that is producing the wastewater but often portray the haulers of the waste as deviant individuals:

*Local officials in the Oil Patch say they know some truck drivers are dumping liquid waste in remote areas to save time and money, but it's rare to be able to pursue criminal charges. (Dalrymple, 2014)*

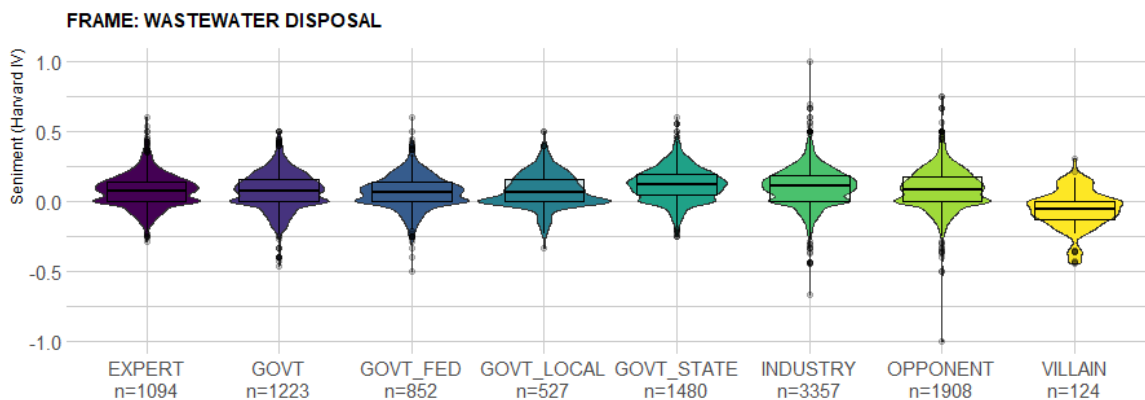
*Mckenzie County, which has the highest concentration of drilling activity, has not yet had any criminal cases for illegal dumping, yet Samuelson estimates he receives two to three reports a month. (Dalrymple, 2014)*

*A Western Pennsylvania waste-hauler, Robert Allan Shipman, was convicted of illegally dumping waste in 2012, and sentenced to serve seven years of probation and 1,750 hours of community service, and to pay \$382,000 in restitution and fines. The attorney general has appealed the sentence, arguing that Shipman deserved jail time. (Maykuth, 2013)*

*Nathan Garber of Kalispell, Mont., is accused of ordering the illegal dumping of saltwater waste into a former oil well, putting drinking water in Stark County at risk for contamination. He faces up to five years in prison and a \$5,000 fine if convicted of the felony. (Associated Press, 2012)*

Figure 6 shows that of the most positively constructed character group in this frame is *state government* actors, which are constructed most closely with *allies* ( $M_{ALLIES}=0.1027 < M_{STATEGOVT}=0.1197$ ). Again, *allies* have a higher sentiment than *heroes* ( $M_{HERO}=0.0976 < M_{ALLIES}=0.1027$ ).

Figure 6. Significant policy actors and NPF characters in the wastewater disposal frame



Among the articles using this frame, there were few references to *victim* (n=10) or *hero* (n=22), which prevented hypothesis testing, though *victims* were significantly related to this frame in the regression at the 99% confidence level (i.e., when it is used, this is one of the frames it is

highly associated with). The results in Table 6 showed that *local government* and *government* are most associated with *experts* and that *state government* and *industry* are most associated with *allies*.

The solutions presented in this frame of highest significance (99.9% confidence level) are *named regulations*, *permit*, and *information*. Based on an evaluation of the regression coefficients ( $\beta_{\text{NAMEDREGULATION}} = 0.243$ ,  $\beta_{\text{PERMIT}} = 0.104$ ,  $\beta_{\text{INFORMATION}} = -0.129$ ), *named regulation* is the solution co-occurring most often. When investigated, these were found to be environmentally-related laws (e.g., *Clean Streams Law*, *Ground Water Protection Act*, *Clean Water Act*). While the positive coefficients for *named regulation* and *permit* show the prediction of the co-occurrence of the *wastewater disposal* frame, *information* solutions have a relatively strong negative relationship with this frame. Because it is negative, the regression coefficient implies that as the occurrence of *information* solutions increase, the likelihood of the frame being *wastewater disposal* decreases. Negative coefficients can occur in the regression results because the regression was done across all frames, characters, and solutions.

Table 9. T-test results on significant actors against NPF characters

$\mu_1 - \mu_2$	Mean difference	Std. Error	t-stat	Pr(> t )	Similar sentiment
GOVT LOCAL – EXPERT	-0.001	0.006	-0.198	0.843	Yes
GOVT – EXPERT	-0.002	0.005	-0.479	0.632	Yes
GOVT STATE – ALLIES	0.006	0.008	0.740	0.460	Yes
INDUSTRY – ALLIES	-0.006	0.008	-0.745	0.457	Yes
GOVT LOCAL – OPPONENT	-0.013	0.006	-2.121	0.034*	No
GOVT FED – EXPERT	-0.012	0.005	-2.178	0.030*	No
GOVT – OPPONENT	-0.014	0.005	-2.868	0.004**	No
GOVT LOCAL – ALLIES	-0.036	0.009	-3.812	0.000***	No
GOVT – ALLIES	-0.037	0.009	-4.266	0.000***	No
INDUSTRY – OPPONENT	0.017	0.004	4.302	0.000***	No
GOVT FED – OPPONENT	-0.024	0.005	-4.524	0.000***	No
GOVT FED – ALLIES	-0.047	0.009	-5.250	0.000***	No
GOVT STATE – OPPONENT	0.029	0.004	6.583	0.000***	No
INDUSTRY – EXPERT	0.029	0.004	6.845	0.000***	No
GOVT FED – VILLAIN	0.126	0.013	9.390	0.000***	No
GOVT STATE – EXPERT	0.041	0.005	8.775	0.000***	No
GOVT LOCAL – EXPERT	0.136	0.014	9.888	0.000***	No
GOVT – VILLAIN	0.135	0.013	10.156	0.000***	No
INDUSTRY – VILLAIN	0.166	0.013	12.850	0.000***	No
GOVT STATE – VILLAIN	0.179	0.013	13.631	0.000***	No

Note: +p<0.1; \* p<0.05; \*\* p<0.01; \*\*\* p<0.001

*Villains* in the wastewater disposal issue frame have lower than average negative sentiment and *state government* are *allies* with greater than average positive sentiment, suggesting this is a contentious issue. Interestingly, *named regulations* are frequently associated with this topic, while *information* is a solution rarely posed. Likely, when policy actors are characterized as *villains*, we don't imagine that the soft policy tool of *information* would make an impact in the situation. Here,

state government is seen as an ally in the fight against dumping using existing laws and regulations to solve the problem.

## DISCUSSION

One of the main objectives of this chapter has been to investigate the relationship between frames, characters, and solutions. Our results show significant relationships between frames and those narrative elements, confirming some definitional assumptions about frames that had not, up to this point, been investigated. Here we have found aspects of Entman's (1993) definition of frames to hold empirical validity: frames narrow our lens by promoting particular presentations of reality as well as treatment recommendations. Our findings suggest that not only are policy actors often closely tied to frames, but they are also accompanied by specific policy solutions. These findings have implications for advancing our understanding of the structure that frames can impose on a narrative, including elements such as characters and solutions.

While we only investigated a small portion of the characters, solutions, and frames in the hydraulic fracturing analysis, the application of this method provided us with findings about the narratives of the issue in the six state and six local newspapers as well as policy narratives in general. The *state fiscal* frame has a very high predictability of characters and solutions, but has a very low range of sentiment, marking it as a consistently used, yet uncontentious issue frame of hydraulic fracturing. These findings suggest that frames can carry consistent elements regardless of their level of NPF characterization. More contentious issue frames, like *wastewater disposal*, involve an increased number of NPF characterizations that are constructed with a wider range of sentiment than others. Importantly to the advancement of the NPF, this method provides a quantitative measure of the relative level of good or bad associated with a character. The relative level of characterization clues the researcher into the aspects of issues that are the thorniest and can potentially highlight the solutions that may be acceptable to the public.

We have demonstrated that NPF research will benefit from including frames in a standard analysis because they allow for advances in understanding narrative structure and when characters are expected to be associated with particular solutions. Research on large datasets of policy narratives could increase our understanding of other factors that likely impact characters and solutions, such as urgency or proximity of the problem. The use of sentiment analysis in conjunction with inductive methods may also lead us to insights on devil-shift or angel-shift patterns. NPF research could use the relationship between frames and policy design we have demonstrated here to address scope of conflict strategies in an automated fashion as well, particularly since we found that policy actors are associated with certain designs that hold relatively fixed cost/benefit distributions.

Gaining a deeper understanding of how frames bound the presence of characters and solutions in narratives will contribute to a stronger tie between the framing and NPF literature. With our initial findings pointing to a relationship between frame and narrative elements such as characters and solutions, we recognize the potential for some shared properties of NPF plots and frames. For instance, frames may act in similar ways to an NPF plot, "situat[ing] the characters and their relationship in time and space" (Shanahan et al., 2017, p. 175). More deeply

understanding the overlap between plots and frames may provide routes to increased consistency in the coding of plots, which has historically been challenging.

As researchers consider using semi-automated methods, they must also consider the tradeoffs. Automated textual analysis, even with the use of iteratively refined dictionaries, may not detect some nuances of language. In some cases, researchers may consider this an acceptable compromise with the increased generalizability of the findings from having a larger dataset and the gains in reliability from automation. Additionally, while these methods appear to be faster because they are automated, they take a similarly deep consideration of research design as do hand-coding studies. Significant time must be invested in conceptualizing, populating, checking the validity of, and refining dictionaries for this process. And if one is a novice programmer, time must be spent in understanding terminology, data structures, and available software.

The narrative research in this study was only a first step towards understanding automated tools that may increase the ability to address larger datasets systematically. Our research has limitations that could be addressed in future research. For instance, while context around characters is important, is a whole sentence a good predictor of emotional response about one paragraph? Perhaps it is only in certain parts of speech in that sentence, as suggested by Shanahan et al. (2018), or can it be contained in one evaluative word? We also make assumptions in our theoretical definitions that could be refined by future research. For instance, although we did use some entities as policy actors, such as *industry* and *levels of government*, we did not include non-human characters and abstract principles as entities. For example, although we did not treat it as such, *the environment* could be constructed as a victim or *greenhouse gas emissions* could be a villain in a frame about climate change.

The results from this research have implications for the way that policy narrative researchers do business. While this chapter only discusses several of the frames, policy actors/characters, and solutions, it should be clear that these methods add autocoding and computer-assisted tools to the policy narrative researcher's toolkit. Additionally, advances in computation linguistics research suggest there are more automated tools to come. For instance, Jacobs (2019) shows the potential for vector space models aided by sentiment analysis to predict the emotional potential of fictional characters—a method that could aid NPF researchers, particularly in the realm of character identification with large volumes of text. While there are limitations to automated methods, we have shown that an iterative context-considered approach can open up the potential for addressing larger, more complex datasets with relatively high levels of confidence. The iterative, inductive approach we have presented also means that researchers hold greater levels of understanding of their own methods and results.

## REFERENCES

- Adgate, John L., Bernard D. Goldstein, and Lisa M. McKenzie. 2014. "Potential Public Health Hazards, Exposures and Health Effects from Unconventional Natural Gas Development." *Environmental Science & Technology* 48 (15): 8307-8320.  
<https://doi.org/10.1021/es404621d>. <https://doi.org/10.1021/es404621d>.
- Arvenson, Amelia. 2017. "David Paul Pleads not Guilty in Case of Beaming Light onto Boulder Courthouse." *Longmont Times Call*, December 11, 2017.  
<https://www.timescall.com/2017/12/11/david-paul-pleads-not-guilty-in-case-of-beaming-light-onto-boulder-courthouse/>.
- Associated Press. 2012. "Mont. Oil Executive Pleads not Guilty to Crime in N.D." *Bismarck Tribune*, November 28, 2012.
- Bamberger, Michelle and Robert E. Oswald. 2015. "Long-term Impacts of Unconventional Drilling Operations on Human and Animal Health." *Journal of Environmental Science and Health, Part A* 50 (5): 447-459.  
<https://doi.org/10.1080/10934529.2015.992655>.
- Berardo, Ramiro, Federico Holm, Tanya Heikkila, Christopher M. Weible, Hongtao Yi, Jennifer Kagan, Catherine Chen, and Jill Yordy. 2020. "Hydraulic Fracturing and Political Conflict: News Media Coverage of Topics and Themes across Nine States." *Energy Research & Social Science* 70: 101660.  
<https://doi.org/10.1016/j.erss.2020.101660>.
- Blair, Benjamin, Tanya Heikkila, and Christopher M. Weible. 2016. "National Media Coverage of Hydraulic Fracturing in the United States: Evaluation Using Human and Automated Coding Techniques: Hydraulic Fracturing Media Coverage." *Risk, Hazards & Crisis in Public Policy* 7 (3): 114-128.  
<https://doi.org/10.1002/rhc3.12097>.
- Chong, Dennis and James N. Druckman. 2007. "Framing Theory." *Annual Review of Political Science* 10 (1): 103-126. <https://doi.org/10.1146/annurev.polisci.10.072805.103054>.
- Crow, Deserai A., Lydia A. Lawhon, John Berggren, Juhi Huda, Elizabeth Koebele, and Adrienne Kroepsch. 2017. "A Narrative Policy Framework Analysis of Wildfire Policy Discussions in Two Colorado Communities." *Politics & Policy* 45 (4): 626-656.  
<https://doi.org/10.1111/polp.12207>.
- Crow, Deserai A., and Andrea Lawlor. 2016. "Media in the Policy Process: Using Framing and Narratives to Understand Policy Influences." *Review of Policy Research* 33 (5): 472-491.  
<https://doi.org/10.1111/ropr.12187>.
- Crow, Deserai A. and Laura Wolton. 2020. "Talking Policy in Congressional Campaigns: Construction of Policy Narratives in Electoral Politics." *Politics & Policy* 48 (4): 658-699.  
<https://doi.org/10.1111/polp.12369>.
- Dalrymple, Amy. 2014. "Truckers Dumping Illegal Waste Rarely Caught." *Grand Forks Herald*, February 6, 2014.

- Eagle, Jess. 2009. "Supporters of Gas Drilling Tax Suggest How it Should be Spent." *Pittsburgh Post Gazette*, June 13, 2009.
- Entman, Robert M. 1993. "Framing: Toward Clarification of a Fractured Paradigm." *Journal of Communication* 43 (4): 51-58.  
<https://doi.org/10.1111/j.1460-2466.1993.tb01304.x>.
- Gallegos, Tanya J., Brian A. Varela, Seth S. Haines, and Mark A. Engle. 2015. "Hydraulic Fracturing Water Use Variability in the United States and Potential Environmental Implications." *Water Resources Research* 51 (7): 5839-5845.  
<https://doi.org/10.1002/2015WR017278>.
- Hand, Eric. 2015. "Oil and Gas Operations Could Trigger Large Earthquakes." *Science*, April 23, 2015.  
<https://www.sciencemag.org/news/2015/04/oil-and-gas-operations-could-trigger-large-earthquakes>.
- Heikkila, Tanya, and Christopher M. Weible. 2017. "Unpacking the Intensity of Policy Conflict: a Study of Colorado's Oil and Gas Subsystem." *Policy Sciences* 50 (2): 179.  
<https://doi.org/10.1007/s11077-017-9285-1>.
- Humphreys, Ashlee, and Rebecca Jen-Hui Wang. 2017. "Automated Text Analysis for Consumer Research." *Journal of Consumer Research* 44 (6): 1274-1306.  
<https://doi.org/10.1093/jcr/ucx104>.
- Jacobi, Carina, Wouter van Atteveldt, and Kasper Welbers. 2015. "Quantitative Analysis of Large Amounts of Journalistic Texts Using Topic Modelling." *Digital Journalism* 4 (1): 89-106.  
<https://doi.org/10.1080/21670811.2015.1093271>.
- Jacobs, Arthur M. 2019. "Sentiment Analysis for Words and Fiction Characters From the Perspective of Computational (Neuro-)Poetics." *Frontiers in Robotics and AI* 6 (53).  
<https://doi.org/10.3389/frobt.2019.00053>.
- Jones, Michael D. 2013. "Cultural Characters and Climate Change: How Heroes Shape Our Perception of Climate Science." *Social Science Quarterly* 95 (1): 1-39.
- Jones, Michael D. and Mark K. McBeth. 2010. "A Narrative Policy Framework: Clear Enough to Be Wrong?" *Policy Studies Journal* 38 (2): 329.
- Jones, Michael D., Elizabeth A. Shanahan, and Mark K. McBeth. 2014. *The Science of Stories: Applications of the Narrative Policy Framework in Public Policy Analysis*. London: Palgrave MacMillan.
- Jones, Michael D. and Geoboo Song. 2014. "Making Sense of Climate Change: How Story Frames Shape Cognition." *Political Psychology* 35 (4): 447-476.  
<https://doi.org/10.1111/pops.12057>.
- Joyce, Stephanie and Jordan Wirfs-Brock. "The Rising Cost of Cleaning Up After Oil and Gas." *Inside Energy*, October 1, 2015.



- Konkel, Lindsey. 2016. "Salting the Earth: The Environmental Impact of Oil and Gas Wastewater Spills." *Environmental Health Perspectives* 124 (12): A230-A235. <https://doi.org/doi:10.1289/ehp.124-A230>.
- Krippendorff, Klaus. 2004. *Content Analysis: An Introduction to its Methodology*. Thousand Oaks, Calif: Sage.
- Lawlor, Andrea. 2015. "Framing Immigration in the Canadian and British News Media." *Canadian Journal of Political Science* 48 (2): 329. <https://doi.org/10.1017/S0008423915000499>.
- Lawlor, Andrea and Deserai A. Crow. 2018. "Risk-Based Policy Narratives." *Policy Studies Journal* 46 (4): 843-867. <https://doi.org/843-867>.
- Lawton, Ricky N. and Murray A. Rudd. 2014. "A Narrative Policy Approach to Environmental Conservation." *Ambio* 43 (7): 849-857. <https://doi.org/10.1007/s13280-014-0497-8>.
- Mason, Krystal L., Kyla D. Retzer, Ryan Hill, and Jennifer M. Lincoln. 2015. "Occupational Fatalities During the Oil and Gas Boom -- United States, 2003 - 2013." *Morbidity and Mortality Weekly Report (MMWR)* 64 (20): 551-554.
- Maykuth, Andrew. 2013. "Shale Criminal Charges Stun Drilling Industry." *Philadelphia Inquirer*, September 12, 2013.
- McGinnis, Michael D. and Elinor Ostrom. 2014. "Social-ecological System Framework: Initial Changes and Continuing Challenges." *Ecology and Society* 19 (2): 30. <https://doi.org/10.5751/es-06387-190230>.
- Merry, Melissa. 2016. "Making Friends and Enemies on Social Media: The Case of Gun Policy Organizations." *Online Information Review* 40 (5): 624-642. <https://doi.org/10.1108/OIR-10-2015-0333>.
- Merry, Melissa K. 2018. "Narrative Strategies in the Gun Policy Debate: Exploring Proximity and Social Construction." *Policy Studies Journal* 46 (4): 747-770. <https://doi.org/10.1111/psj.12255>.
- Miles, Matthew B., A. Michael Huberman, and Johnny Saldaña. 2020. *Qualitative Data Analysis: A Methods Sourcebook*. Thousand Oaks: Sage Publications, Inc.
- Moskowitz, Peter. 2015. "New Report Estimates Enough Natural Gas is Leaking to Negate Climate Benefits." *The Guardian*, June 24, 2015. <http://www.theguardian.com/environment/2015/jun/24/natural-gas-leaks-methane-environment>.
- Olofsson, Kristin L., Christopher M. Weible, Tanya Heikkila, and J. C. Martel. 2018. "Using Nonprofit Narratives and News Media Framing to Depict Air Pollution in Delhi, India." *Environmental Communication* 12 (7): 956-972. <https://doi.org/10.1080/17524032.2017.1309442>.
- Poirier, William, Catherine Ouellet, Marc-Antoine Rancourt, Justine Béchar, and Yannick Dufresne. 2020. "(Un)Covering the COVID-19 Pandemic: Framing Analysis of the Crisis

- in Canada." *Canadian Journal of Political Science* 53 (2): 365-371.  
<https://doi.org/10.1017/S0008423920000372>.
- Ratner, Bruce. 2009. "The Correlation Coefficient: Its Values Range Between +1/-1, or Do They?" *Journal of Targeting, Measurement and Analysis for Marketing* 17 (2): 139-142.  
<https://doi.org/10.1057/jt.2009.5>.
- Scheufele, Dietram A. and Shanto Iyengar. 2014. "The State of Framing Research: A Call for New Directions." In *The Oxford Handbook of Political Communication*, edited by Kate Kenski and Kathleen Hall Jamieson. New York: Oxford University Press.
- Schäfer, Mike S. and Saffron O'Neill. 2017. "Frame Analysis in Climate Change Communication: Approaches for Assessing Journalists' Minds, Online Communication and Media Portrayals." In *Oxford Encyclopedia of Climate Change Communication*, edited by Matthew Nisbet, Shirley Ho, Ezra Markowitz, Saffron O'Neill, Mike S. Schäfer and Jagadish Thaker. New York: Oxford University Press.
- Scott, Tyler A., Nicola Ulibarri, and Ryan P. Scott. 2020. "Stakeholder Involvement in Collaborative Regulatory Processes: Using Automated Coding to Track Attendance and Actions." *Regulation & Governance* 14 (2): 219-237.  
<https://doi.org/10.1111/rego.12199>.
- Shanahan, Elizabeth A., Michael D. Jones, and Mark K. McBeth. 2018. "How to Conduct a Narrative Policy Framework Study." *The Social Science Journal* 55 (3): 332-345.  
<https://doi.org/10.1016/j.soscij.2017.12.002>.
- Shanahan, Elizabeth A., Michael D. Jones, Mark K. McBeth, and Claudio M. Radaelli. 2017. "The Narrative Policy Framework." In *Theories of the Policy Process*, edited by Christopher M. Weible and Paul A. Sabatier. Boulder, CO: Westview.
- Stone, Philip J., Robert F. Bales, J. Zvi Namenwirth, and Daniel M. Ogilvie. 1962. "The General Inquirer: A Computer System for Content Analysis and Retrieval Based on the Sentence as a Unit of Information." *Behavioral Science* 7 (4): 484.
- Stone, Phillip. 2018. "Descriptions of Inquirer Categories and Use of Inquirer Dictionaries." Accessed December 8, 2018.  
<http://www.wjh.harvard.edu/~inquirer/homecat.htm>.
- U.S. Energy Information Administration. 2018. "How Much Shale Gas is Produced in the United States?" Accessed April 10, 2018.  
<https://www.eia.gov/tools/faqs/faq.php?id=907&t=8>.
- Wines, Michael. 2013. "Gas Drilling Not Imminent, but Debate Roils N.Y. Region." *Pittsburgh Post-Gazette*, January 13, 2013.
- Witze, Alexandra. 2015. "Artificial Quakes Shake Oklahoma: Earthquakes Linked to Oil and Gas Operations Prompt Further Research into Human-induced Seismic Hazards." *Nature* 520 (7548): 418.
- Yu, Jingyuan, Yanqin Lu, and Juan Muñoz-Justicia. 2020. "Analyzing Spanish News Frames on Twitter during COVID-19—A Network Study of El País and El Mundo." *International*

*Journal of Environmental Research and Public Health* 17 (15): 5414.  
<https://doi.org/10.3390/ijerph17155414>.

**CHAPTER APPENDIXES***Appendix A. Included news sources and article counts*

Newspaper Name	Publishing City	State	Pre-exclusion Count	Source	Article count post-exclusion
Houston Chronicle	Houston	TX	2935	Newsbank	592
Dallas Morning News	Dallas	TX	1543	Newsbank	441
Pittsburgh Post-Gazette	Pittsburgh	PA	3992		1170
Philadelphia Inquirer	Philadelphia	PA	2397	Newsbank	691
Oklahoman	Oklahoma City	OK	2909		393
Tulsa World	Tulsa	OK	1868		187
The Advocate	Baton Rouge	LA	663	Newsbank	233
Times-Picayune	New Orleans	LA	682	Newsbank	170
Grand Forks Herald	Grand Forks	ND	1003	Newsbank	360
Bismarck Tribune	Bismarck	ND	2425	Newsbank	580
Denver Post	Denver	CO	1595	Newsbank	567
Daily Times-Call	Longmont	CO	1466	Newsbank	315
<b>Total Count</b>					<b>5708</b>

*Appendix B. Regression of significant characters and solutions at 99.9% confidence*

	<i>Dependent variable:</i>
	STATE_FISCAL
GENERAL	0.079*** (0.013)
REGULATION	-0.124*** (0.015)
TAX	1.868*** (0.016)
GOVT_LOCAL_SP	0.220*** (0.037)
ACTIVIST_GROUP	-0.296*** (0.063)
GOVT	0.156*** (0.011)
GOVT_LOCAL	-0.065*** (0.017)
GOVT_STATE_SP	0.098*** (0.020)
Constant	0.220*** (0.045)
Observations	5,704
R <sup>2</sup>	0.778
Adjusted R <sup>2</sup>	0.778
Residual Std. Error	2.319 (df = 5695)
F Statistic	2,492.394*** (df = 8; 5695)
<i>Note:</i>	+ p<0.1; * p<0.05; ** p<0.01; *** p<0.001

	<i>Dependent variable:</i>
	WASTEWATER_DISPOSAL
SOL_INFORMATION	-0.127*** (0.036)
SOL_PERMIT	0.103*** (0.017)
SOL_REGULATION_SP	0.248*** (0.051)
GOVT_STATE	0.106*** (0.010)
EXPERT	0.163*** (0.012)
GOVT	-0.032*** (0.007)
GOVT_LOCAL	-0.051*** (0.011)
GOVT_FED	0.064*** (0.011)
VILLAIN	0.227*** (0.046)
Constant	0.135*** (0.031)
Observations	5,708
R <sup>2</sup>	0.082
Adjusted R <sup>2</sup>	0.080
Residual Std. Error	1.619 (df = 5698)
F Statistic	56.514*** (df = 9; 5698)
<i>Note:</i>	+ p<0.1; * p<0.05; ** p<0.01; *** p<0.001

Appendix C. Accepted interpretation of correlation coefficient

Correlation Coefficient Range	Relationship Interpretation (e.g., Ratner 2009)
0	no linear relationship
0 to 0.3 or 0 to -0.3	weak positive or negative relationship
0.3 to 0.7 or -0.3 to -0.7	moderate positive or negative relationship
0.7 to 1.0 or -0.7 to -1.0	strong positive or negative linear relationship
+1 or -1	perfect positive or negative relationship

*Appendix D. Character dictionary*

Category	Search item	Whole word	Case sensitivity	Starting letters
CHAR_ALLIES	advocacy group	0	0	0
CHAR_ALLIES	advocate	0	0	0
CHAR_ALLIES	backer	0	0	0
CHAR_ALLIES	campaigner	0	0	0
CHAR_ALLIES	champion	0	0	0
CHAR_ALLIES	counsel	0	0	0
CHAR_ALLIES	defender	0	0	0
CHAR_ALLIES	enthusiast	0	0	0
CHAR_ALLIES	friend	0	0	0
CHAR_ALLIES	guardian	0	0	0
CHAR_ALLIES	leader	0	0	0
CHAR_ALLIES	promoter	0	0	0
CHAR_ALLIES	proponent	0	0	0
CHAR_ALLIES	protector	0	0	0
CHAR_ALLIES	spokesperson	0	0	0
CHAR_ALLIES	supporter	0	0	0
CHAR_EXPERT	analyst	0	0	0
CHAR_EXPERT	chemist	0	0	0
CHAR_EXPERT	economist	0	0	0
CHAR_EXPERT	educator	0	0	0
CHAR_EXPERT	examiner	0	0	0
CHAR_EXPERT	expert	0	0	0
CHAR_EXPERT	geologist	1	0	0
CHAR_EXPERT	industry analyst	1	1	0
CHAR_EXPERT	inspector	0	0	0
CHAR_EXPERT	panelist	0	0	0
CHAR_EXPERT	physicist	0	0	0

CHAR_EXPERT	professor	0	0	0
CHAR_EXPERT	researcher	0	0	0
CHAR_EXPERT	scholar	0	0	0
CHAR_EXPERT	scientist	0	0	0
CHAR_EXPERT	seismologist	0	0	0
CHAR_EXPERT	specialist	0	0	0
CHAR_EXPERT	technician	0	0	0
CHAR_HERO	celebrity	0	0	0
CHAR_HERO	hero	0	0	0
CHAR_HERO	protagonist	0	0	0
CHAR_HERO	superstar	0	0	0
CHAR_HERO	winner	0	0	0
CHAR_OPPONENT	activist	0	0	0
CHAR_OPPONENT	adversary	0	0	0
CHAR_OPPONENT	agitator	0	0	0
CHAR_OPPONENT	antagonist	0	0	0
CHAR_OPPONENT	candidate	0	0	0
CHAR_OPPONENT	challenger	0	0	0
CHAR_OPPONENT	competitor	0	0	0
CHAR_OPPONENT	demonstrator	0	0	0
CHAR_OPPONENT	enem	0	0	1
CHAR_OPPONENT	foe	0	0	0
CHAR_OPPONENT	heckler	0	0	0
CHAR_OPPONENT	litigant	0	0	0
CHAR_OPPONENT	opponent	0	0	0
CHAR_OPPONENT	opposition	0	0	0
CHAR_OPPONENT	protester	0	0	0
CHAR_OPPONENT	radical	0	0	0
CHAR_OPPONENT	rival	0	0	0



CHAR_VICTIM	casualt	0	0	1
CHAR_VICTIM	fatalit	0	0	1
CHAR_VICTIM	injured part	0	0	1
CHAR_VICTIM	victim	0	0	0
CHAR_VILLAIN	convict	0	0	0
CHAR_VILLAIN	criminal	0	0	1
CHAR_VILLAIN	crook	0	0	0
CHAR_VILLAIN	delinquent	0	0	0
CHAR_VILLAIN	deviants	0	0	0
CHAR_VILLAIN	felon	1	0	0
CHAR_VILLAIN	guerrilla	0	0	0
CHAR_VILLAIN	guilty party	0	0	0
CHAR_VILLAIN	lawbreaker	0	0	0
CHAR_VILLAIN	perpetrator	0	0	1
CHAR_VILLAIN	rebel	0	0	0
CHAR_VILLAIN	suspect	0	0	0
CHAR_VILLAIN	terrorist	0	0	0
CHAR_VILLAIN	thug	0	0	0
CHAR_VILLAIN	villain	0	0	1

---

### *Appendix E. Autocoding example in R*

The following lines of code are for use on a dataframe that contains source data and clean text.

The way this code is designed is that each row of the dataframe input text is segmented into the desired unit for the analysis. The following chunk of code breaks the text from a dataframe into sentences (the unit of analysis for this study). Here, we use the *unnest\_sentences* function to change each whole block of text into a larger dataframe with the number of rows equal to the number of sentences in all the text. For example, if 100 news articles were originally read into a dataframe (*text\_df*), the dimensions would be the column of text (*text\_df\$text*) plus columns for metadata (6 information columns shown below) by 100 (a row for each of the 100 articles). The dimension of the post-unnesting dataframe (*text\_sen*) will be the number of sentences in all the documents.

```

library(lexRankr)
#text_df=dataframe with each cell in the "text" column contains the text
#from a whole document
dim(text_df)
#[1] 100 7
colnames(text_df)
#[1] "title" "paper" "date" "section" "byline" "text" "highlight"

#the column title contains each document name
#sents in the function call is a dataframe column that will be created containing the sentences
text_sen<-unnest_sentences(text_df,sents,text,doc_id=title)

#the number of rows in the new dataframe will equal the number of sentences
#in all the text documents
dim(text_sen)
#11220 8
colnames(text_sen)
#[1] "title" "paper" "date" "section" "byline" "highlight" [7] "sent_id" "sents"

```

For the dictionary to find all occurrences, the text must be cleaned. Remove all punctuation and multiple spaces.

```

#remove a little punctuation so the dictionary can find all occurrences
#this will create a new column of clean text in your dataframe (clntxt)
library(mgsub)
text_sen$clntxt<-mgsub(text_sen$sents, c("[:punct:]", "\r\n"), c("", " "))
text_sen$clntxt<-mgsub(text_sen$clntxt, "\\s+", " ")

```

Read in the dictionary designed similarly to the one shown in Figure 1. Then change the values in the columns of switches to be *grep* keyword friendly.

```

library(readxl)
#read in frame dictionary
#this example is from a spreadsheet
#read in as a tibble or dataframe for compatibility with the following code
dict<-read_xlsx('C:/Users/vegan/Dropbox/Campaign Book Project/Dictionary/pol_ref_dict.xlsx',col_names=TRUE,col_types=NULL)

#rename cols
colnames(dict)[c(1:5)]<-c("Category","word","ww","cs","sl")

#make the columns friendly the grep search
#have to swap 0 to 1 and 1 to 0
#this is for the ignore.case keyword in grep
dict$cs<-replace(dict$cs,dict$cs==1,"FALSE")
dict$cs<-replace(dict$cs,dict$cs==0,"TRUE")

```

```
#this is for the starting letter
ind<-which(dict$sl==1)
#add a single space to the beginning of the word
dict$word[ind]<-paste("",dict$word[ind])
#remove any multiple spaces
dict$word<-mgsub(dict$word,"\\s+", " ")

#if the word is the whole word add a space
#on the front and back of the word
ind=which(dict$ww==1)
dict$word[ind]<-paste("",dict$word[ind],"
```

Create a new empty dataframe for the coded segments. The dimensions of the new dataframe will contain two additional columns for category and search term.

```
####create an empty dataframe to save the coded segments
#you'll want to have enough columns to contain
#the original data info (e.g., title of article, source, date)
#plus columns that are needed for the category, search word
#ncol should be the number of columns in text_sen dataframe+2
coded_segs<-data.frame(matrix(ncol=10,nrow=0))

#####
for (wind in 1:length(dict$word)){
  #look for the word in tweets, tag it
  ind<-grep(dict$word[wind],text_sen$clntxt,value=FALSE,ignore.case=dict$cs[wind])
  #add to the output dataframe
  coded_segs<-rbind(coded_segs,cbind(text_sen[ind,c("date","paper","title","section","byline","sents","sent_id")],cbind(cbind(rep(dict$Category[wind],length(ind)),rep(dict$word[wind],length(ind))),text_sen$clntxt[ind]))
}
colnames(coded_segs)<-c(colnames(text_sen)[c(1:6)],"Code","SearchWord","clntxt")
coded_segs$clntxt<-as.character(coded_segs$clntxt)
```

The resulting dataframe of coded segments will contain a row for every tag. This dataframe can then be used for various analysis purposes.