

# Once Upon A Time In Visualization: Understanding the Use of Textual Narratives for Causality

Arjun Choudhry, Mandar Sharma, Pramod Chundury, Thomas Kapler, Derek W. S. Gray,  
Naren Ramakrishnan, and Niklas Elmqvist, *Senior Member, IEEE*

**Abstract**—Causality visualization can help people understand temporal chains of events, such as messages sent in a distributed system, cause and effect in a historical conflict, or the interplay between political actors over time. However, as the scale and complexity of these event sequences grows, even these visualizations can become overwhelming to use. In this paper, we propose the use of textual narratives as a data-driven storytelling method to augment causality visualization. We first propose a design space for how textual narratives can be used to describe causal data. We then present results from a crowdsourced user study where participants were asked to recover causality information from two causality visualizations—causal graphs and Hasse diagrams—with and without an associated textual narrative. Finally, we describe CAUSEWORKS, a causality visualization system for understanding how specific interventions influence a causal model. The system incorporates an automatic textual narrative mechanism based on our design space. We validate CAUSEWORKS through interviews with experts who used the system for understanding complex events.

**Index Terms**—Causality visualization, natural language generation, data-driven storytelling, temporal data, quantitative studies.

## 1 INTRODUCTION

STORIES are a central part of what it means to be human [40, 53]. They teach, guide, and caution; they store, recall, and archive; they praise, spread joy, and inspire. In particular, stories are especially useful for encapsulating *causality*—the cause and effect of events in a plot—in a regular, understandable, and memorable format. This format is also surprisingly scalable. Examples abound of textual narratives representing complex chains of cause and effect ranging from the winding plots of G. R. R. Martin’s *A Song of Ice and Fire* and Neal Stephenson’s *The Baroque Cycle*, through shelf yards of history books laying out the intricacies of the Napoleonic Wars or the American Revolution in all their gritty detail, and all the way to quarterly reports telling the story of a company’s accomplishments over the last three months. However, despite all of this utility, little work exists on the use of textual narratives to represent causality in modern visualization tools. On the contrary, visualization and visual analytics researchers tend to view textual narratives with suspicion, often instead opting to apply text analytics and visualization methods to minimize their use.

In this paper, we attempt to remedy this gap in the literature by investigating how textual narratives can be used to represent causality. Our classification of narratives is primarily based on their utility as a complement to causality visualization techniques, such as dynamic graphs [5] and Hasse diagrams [17, 18]. Textual representations are generally much less compact than geometric ones (i.e., visualizations), and must thus be designed with specific questions in mind. We first propose and discuss a design space of causality representations, focusing in particular on textual narratives. We then report on a crowdsourced user study where we operationalized parts of our design space and asked participants to recover causality information from dynamic graphs versus Hasse diagrams, with and without an associated textual narrative. Our findings indicate that narratives can fill an important complementary role for key questions on causality, and thus serve as a “story-like” format to summarize a specific causal event chain.

To capitalize on these findings and demonstrate the use of our design space, we also present a textual narratives implementation in CAUSEWORKS, a causality visualization system on the Causal Exploration of Complex Operational Environments program [1] for understanding the impact of specific interventions in a causal model. These narratives

are based on best practices from our design space as well as the user study, and serve as a quick-reference textual summary of the selected interventions and objectives shown in a dynamic graph. We studied the utility of these narratives by interviewing several users with experience of causality, who used the system to understand climate change data. Our findings confirm many of the results from the crowdsourced study.

The contributions of our paper are the following: (1) a design space for complementary textual narratives in representing causality; (2) results from a crowdsourced study evaluating different causality visualizations with and without companion narratives; (3) an implementation of textual narratives in an existing causality analytics and visualization system (CAUSEWORKS); and (4) qualitative results from 5 experts using these narratives to understand climate change data.

## 2 BACKGROUND

Here we discuss the existing literature on causality, causality visualization, and data-driven storytelling.

### 2.1 Causality Visualization

Causal networks or directed acyclic graphs are commonly used to map relationships between variables [47]. Much of the work in causal visualization aims to encode aspects of causality such as temporal developments of cause and effect [20] using interactivity [64, 65] and animations [34] to improve the accuracy of causal inference.

Researchers have uncovered characteristics and shortcomings of specific visual representations of causality. For example, Bae et al. [4] find that multiple to/from connections from a particular node may influence how an analyst perceives indirect effects. Similarly, Hasse diagrams are used widely, but require the user to backtrace every effect, and can also introduce an overwhelming number of crossings in a large-scale causal system [19]. Wang et al. attempted to improve causal inference by overlaying salient statistical parameters such as p-values and regression co-efficients on 2D-graphs so that analysts can draw more reliable conclusions about causal relationships [64]. However, interpreting these parameters requires understanding statistical inference.

Research on perceptions of causality also show that inference is context-dependent, and a non-expert with regards to statistics or domain could see an illusion of causality in data [68]. In our work, we propose to mitigate misinterpretation by both experts and non-experts alike through the use of textual narratives to augment causal visualizations.

### 2.2 Narratives in Visualization

Historically spanning thousands of years [53, 62], *storytelling* conveys a series of events, usually involving characters and locations—*stories*—

- Arjun Choudhry, Mandar Sharma, and Naren Ramakrishnan are with Virginia Tech in Arlington, VA, USA. Email: {aj07lfc, mandarsharma}@vt.edu, naren@cs.vt.edu
- Pramod Chundury and Niklas Elmqvist are with the University of Maryland, College Park, MD, USA. Email: {pchundur, elm}@umd.edu
- Thomas Kapler and Derek W. S. Gray are with Uncharted Software in Toronto, ON, Canada. Email: {tkapler, dgray}@uncharted.software

using speech, sound, and visuals [26]. Generally, stories are communicated using visual media, such as illustrations, pictures, animations, video, and—now—visualization [16, 57]. Visualization, inherently, is inclined for communication by virtue of its graphical form, resulting in the notion of *communication-minded visualization* [63]. Combining the idea of *communication-minded visualization* with *storytelling* yields the notion of *data-driven storytelling*: narrative techniques for data [54].

We believe that data-driven storytelling naturally follows the idea of visualization for explanation (the latter). The production, presentation, and dissemination of analysis results is an important challenge in visualization and visual analytics [60]. Gershon and Page first proposed using storytelling for visualization [24], and their work has since been followed up by workshops [12, 13], surveys [31, 54], and even commercial tools [37]. Viégas and Wattenberg note the inclination of visualization for communication by virtue of its graphical form, and encourage focusing on so-called *communication-minded visualization* [63] for social analysis. In recent years, the use of textual data to aid visualization and vice versa have been explored [6, 39]. Furthermore, verbalization [10, 30, 55] has also been used for understanding machine learning models.

### 2.3 Causality and Causal Networks

The statistical and ML sciences have developed many formalisms to reason with both the structure and dynamics of causal networks [8]. To encapsulate causal structure, while there are many network formulations, one of the more popular ones is the Bayesian network formalism popularized by Pearl [47]. A Bayesian network is a directed acyclic graph (DAG) and can be thought of as a way to represent a factorization of the underlying joint distribution of random variables. However, interpreting such DAGs is difficult for humans and interpretation rules such as d-separation [23] and the ‘Bayes Ball’ algorithm [56] have been proposed. These rules essentially are ways to read or infer conditional independence relationships from the networks.

To overcome such interpretation difficulties, other representational formalisms have been proposed, e.g., dependency networks [27], which allow cycles, and Markov networks [36] (also called Markov random fields, or MRFs), which are undirected. In terms of dynamics, a causal representation must allow us to probe the effect of interventions and to posit and explore counterfactuals. Interventions are modeled using a calculus (e.g., Pearl’s do-calculus) that mutates the given network to propagate and understand the downstream consequences of the intervention. Counterfactuals allow us to ask more expressive questions and explore the progression of different variables in alternative worlds or situations. We assume in this paper that the underlying causal representation is fixed and a suitable interpretation of dynamics is available to probe the effect of interventions, and focus on the role of visualization in communicating cause-effect relationships.

## 3 DESIGN SPACE: TEXTUAL NARRATIVES FOR CAUSALITY

Visualizations are themselves considered as ways to tell stories with data [54] and, in this paper, we view textual narratives as an augmented form of storytelling that aims to increase insight [9, 45], comprehension, and decision making. We focus on textual narratives as a way to express causal information in event sequences, specifically as a complement to causality visualizations [33, 43], such as causal graphs or Hasse diagrams. For this reason, we tend to think of these textual narratives as a form of *data-driven storytelling* [50]—the use of traditional narrative methods to convey data—that relies on a textual, rather than a visual, medium. We believe that a textual narrative can also replace the visualization, at least to provide a high-level summary [22, 44].

### 3.1 Definitions

The causal relation  $\rightarrow$  is a relation that connects two elements (events)  $x$  and  $y$  as  $x \rightarrow y$  iff  $x$  is the cause of  $y$ . Sets of events are called processes  $P_1, \dots, P_N$ . Internal events are sequential and causally related. External events interconnect processes through messages. We denote the events for a process  $P_i$  as  $E_i = \{e_1^i, e_2^i, e_3^i, \dots\}$ . The causal relation is typically irreflexive, asymmetric, and transitive.

While some causal tasks are concerned with the entire causal model—i.e., the set of all processes and their associated events—many real-

world tasks use a more directed formulation. Of most interest is the ability to impose specific *interventions* or perturbations on a particular process and understanding the resulting impact on *objective* process(es). A causal model may involve a set of interventions  $I$  and objectives  $O$ , each corresponding to specific causes and effects.

### 3.2 Narrative Rendering Pipeline

We view the representation of causal data using textual narratives as an interactive *rendering pipeline*, akin to a classic graphics rendering or visualization pipeline. In this model, we think of the *sentence clause* as the building block. Natural language generation (NLG) systems [48, 49] tend to consist of several stages:

1. **Content selection:** Determining the causality data to display;
2. **Document structuring and aggregation:** Prioritizing the order of data and merging sentences on similar causal data (same source or destination processes) to improve readability;
3. **Realization:** Generating the actual text for each piece of causality information to render in the summary; and
4. **Interaction:** Providing a feedback loop to allow interaction with the textual narratives, such as to drill down, link to related narratives, or brush to highlight items in associated views.

Using natural language to represent data as text is quite different from using visualization, which uses geometric shapes. Unlike visualization, natural language is typically precise. This leads to early fixation, as well as serial representations, which limits parallel processing. In practice, this means that natural language is better suited to presenting specific pieces of information rather than the holistic and parallel overviews that characterize data visualization.

### 3.3 Step 1: Extracting Causality Information

Our proposed text generation pipeline starts with identifying the specific causality information that users desire. This generally depends on the application, which we model using a degree-of-interest (DOI) function in the next language generation step. Thus, our treatment here includes all potentially useful causality information, given our data model. We organize this information into the following categories:

- **Cause and effect:** A central question when reasoning about causality tends to be the factors that caused a specific effect, which we capture as interventions and objectives. *Example:* a white cue ball striking the eight ball, sending it bouncing off the nearest wall of a pool table.
- **Correlation:** While correlation is not causation, many forms of causation have their roots in correlation. Depending on the causality model, the exact cause and effect may not be known; in such cases, correlations between nodes—i.e., a change in one node followed by a change in another node—can be used as a weaker form. *Example:* a medication administered to a patient followed by their blood pressure dropping.
- **Life cycle:** Processes may come and go, often as a result of them receiving an intervention or an internal event. Such life cycle information is commonly of interest in causal reasoning. *Example:* traffic in a computer network being directed around a faulty router that is no longer responding.
- **Connectivity:** Causality modeled as above is essentially a graph, which means that understanding a causal model requires understanding the topology and dynamic connectivity of the events passed in the system. *Example:* tracing infections of an airborne virus in a population based on their social contacts.

For all of the above causality information categories, we can also identify specific common metadata for them all: **path:** the processes on the path between source and destination; **weights:** the values or weights associated with each process; and **time:** the time stamps associated with each of the events. The types of causal information described above are included in the generated narratives of the crowdsourced study and in CAUSEWORKS (figure 1 shows an example).

### 3.4 Step 2: Calculating Order

Here we determine the structure and order of information that we will use for the textual narrative. The primary challenge is that even a moderately complex causal system will have a significant number of candidate causal information to convey.

To address this challenge, we use a degree-of-interest (DOI) function  $f_{DOI}(e) \subseteq \mathbb{R}$  based on user interest and task to prioritize each event  $e$  involved in the sequence of causality data extracted from the prior step. As a first level of prioritization, we propose limiting reports to the sets of interventions  $I$  and objectives  $O$ , as described in Section 3.1. We group items on a per-process basis, and then further prioritize events based on their occurrences, magnitude of change, and influence. To represent this information, we use a directed acyclic graph (DAG) as a scene graph to store the abstract data to render, where each causal process becomes a top-level container for associated causality data.



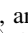


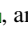
Generally speaking, generating a complete sentence for each clause—recall that a clause corresponds to an individual item of data—is the most clear and unambiguous approach. However, this often leads to significant repetition, which is often seen as clumsy and unnatural to a reader, as well as unnecessary verbosity, which is wasteful given that summaries are often limited in length. For this reason, we use *aggregation* to merge similar clauses that share the same source or destination process into a single sentence. We have aggregated events in our generated narratives (see Figures 2 and 7).

### 3.5 Step 3: Rendering Textual Narratives

We think of realizing the ordered causality data to be expressed as *rendering* the narrative, akin to how a computer graphics system may render a sorted list of triangles to generate a 3D scene. Since our focus is on generating summaries, the notion of a *character budget* is central to our approach: this is the maximum number of characters that we want to use to realize the textual narrative. This budget is not prescriptive, only restrictive; in other words, if space is not an issue, the budget can be set to infinity, resulting in exhaustive textual summaries.

The actual rendering process proceeds by iterating through the sorted data graph, where items are grouped based on top-level processes, as described above. By knowing the number of characters for each branch of the data graph, the renderer can determine how deeply to traverse while maintaining the character budget. Furthermore, we can also identify the available visual channels for conveying data using text:

- **Textual content:** The primary visual channel is obviously the written content that the text spells out.
- **Font size:** Most summaries will use a uniform font size, as changing the size of individual words or sentences throughout a text can be disruptive to reading as well as when calculating its space needs. However, it can be an effective way to show emphasis, particularly for titles and section headings.
- **Typographic emphasis:** A more common and typographically accepted practice is to use emphasis such as **boldface**, *italics*, or underlining to communicate data in the narrative, such as to mark processes, effects, or magnitudes. Additional such emphasis markers include SMALL CAPS, ALL CAPS, or the use of punctuation (!) or “quotation marks” in the text.
- **Color:** As for visualization, color can be an effective visual channel. We differentiate between the use of **font color** and **background color**, which can be used to convey different data (although, as always, care must be taken to avoid interference).
- **Hierarchical lists:** While not part of classic running prose, which tends to just use sentences and paragraphs as its typographic structures, we also consider lists—both enumerated and itemized ones—a useful visual channel. In particular, nesting lists can allow for showing hierarchical part-of relationships.
- **Word-scale graphics:** Despite our focus on written text, we cannot resist drawing on visualization, in the form of *word-scale graphics* [25]: small data-driven graphics that can be embedded into running text. Examples include mathematical symbols such

as  $\uparrow$ ,  $\Delta$ , and  $\propto$ , icons such as , , and , as well as micro visualizations such as , , and  (sparklines).

### 3.6 Step 4: Interacting with Narratives

Finally, since our intended output format almost always is on a computer screen—and not paper—we should also consider how to interact with these textual narratives. We propose the following possibilities:

- **Brushing:** Hovering over a process or an event in a narrative highlights all of its occurrences in related views (or narratives).
- **Hyperlinking:** Entities in the narrative are hyperlinks where clicking on one will navigate to it in related views (or narratives).
- **Drill-down/roll-up:** Dynamically changing the user’s degree of interest will allow for drilling down, e.g., by unfolding the items in a list or expanding suppressed elements in an enumeration.
- **Search:** Directly querying specific elements in a narrative by typing partial or complete search terms allows for quick access [21].

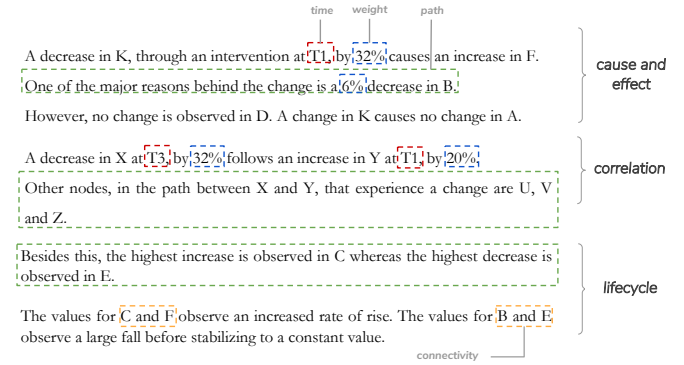


Fig. 1: Annotated example of a narrative conveying causal information about interventions and objectives.

### 3.7 Discussion: What Makes an Effective Narrative?

Textual narratives are slowly making their way into visualization systems, either as a way to generate data insights to accompany a visualization [58] or to structure a visualization for better communication [59]. Research into what makes an effective narrative is still in its infancy and is necessarily tied to the underlying analytical task and domain. For the causal networks domain considered here we identify four facets:

- **Language diversity:** More language diversity avoids monotony but might detract from conveying key messages and conclusions. Less language diversity supports comparison of generated narratives but might lead to ‘glossing over’ by analysts.
- **Level of detail:** Should the narrative capture an executive summary or provide in-depth access to the underlying data? We briefly discuss the preferred level of detail in our expert review.
- **Verbalizing numbers:** Verbalizing quantitative/probabilistic data (e.g., using Kent’s words of estimative probability [35] or the NIC/Mercyhurst standardization) is considered important in specific domains (e.g., intelligence analysis [29]) but other applications argue for direct access to the original numeric information.
- **Human performance aspects:** Understanding the characteristics of narratives that lead to improved human performance is an ongoing research problem [42]. Narratives provide increased comprehension, interest, and engagement and are known to contribute “distinct cognitive pathways of comprehension” with increased recall, ease of comprehension, and shorter reading times [11]. Conversely, the challenge of the written word implies slowness and error-prone behavior due to short-term memory limits.

In general, successful narrative research requires a standardization of both the generation and evaluation space, and an understanding of how a narrative fits into the larger comprehension process of the analyst.

## 4 CROWDSOURCED STUDY: NARRATION FOR CAUSALITY

We conducted a crowdsourced study to understand how narratives augment causal data exploration through visual analysis.

### 4.1 Participants

We recruited our participants through crowdsourcing from *Amazon Mechanical Turk (MTurk)* to complete visual analysis tasks that did not require prior training or data visualization expertise. Owing to the nature of MTurk, we had limited control over participant demographics, technology, and skill level. However, prior work indicates that simple tasks such as ours are flexible to a crowdsourced study design [28]. We planned to recruit 150 participants; all were drawn from within the United States due to tax and compensation restrictions by our Institutional Review Board (IRB). To ensure that our participants understood our task instructions, we screened our participants for working English knowledge. Participants were allowed to participate only once. We estimated our study completion time to be 20–30 minutes, and compensated our participants ethically at a rate of at least \$8/hour (similar to the U.S. federal minimum wage in 2019 of \$7.25).

### 4.2 Apparatus

We required our participants to use a desktop computer (no mobile devices), and the study was distributed through a web browser. We ensured that the visual representations, textual narratives, and their labels were legible for all common device formats. The testing platform was implemented as a *Qualtrics* survey with static trials saved as non-interactive mockups that were manually created using *Microsoft PowerPoint*, and were based on various factors such as polarity of links, link overlaps, and the number of intervening/objective nodes. The narratives were created manually based Section 3.2 (Figure 2).

### 4.3 Experimental Factors

Our goal was to first experimentally understand how the presence of a narrative augments causal analysis using visual representation. We chose a more familiar and less temporal causal representation (Causal Graph) and less familiar and more temporal causal representation (Hasse Diagram). We modeled four factors in our experiment:

- **Causality Visualization (VR):** The visual representation used for conveying causality. We chose two levels:
  - **Causal Graph (NL):** A Causal Graph is a node-link representation of the causal network.
  - **Hasse Diagram (HD):** We use a similar representation of Hasse diagrams as seen in previous work [19].
- **Textual Narrative (TN):** This is a key factor in our study: the 1) presence (ON) or 2) absence (OFF) of a textual narrative.
- **Difficulty (DL):** The difficulty of the trial is expressed in the size of the causal system involved in the event sequence. We chose three levels for this factor:
  - **Simple (S):** 3 or 5 nodes, and up to 4 time-hops (T1–T4)
  - **Medium (M):** 5 to 8 nodes, and up to 5 time-hops (T1–T5)
  - **Hard (H):** 9 to 12 nodes, and up to 5 time-hops (T1–T5)
 We settled on these values through pilot testing to ensure that our tasks can typically be completed in 30 minutes.
- **Narrative Scope (NS):** Inherently, the Hasse diagram affords the explicit showing of changes in a node across time intervals (e.g. T1–T2; T2–T3, etc.). On the other hand, the Causal Graph (NL) requires the user to follow the causal path between nodes to extrapolate temporal information. This means that the accompanying textual narrative could describe effect propagation between successive time-hops—**Instantaneous (IS)**—or can provide a **Cumulative (CU)** summary across all observed time.

This leads to 24 conditions. Since NS is only relevant for situations when TN is ON, this yields a total of 18 conditions. We presented a

total of 12 causal graph systems (CGS) to each of our participants. Although we do not include all aspects of our design space as experimental conditions, we use our narrative rendering pipeline in our mockups. We were also limited by the non-interactivity of our stimulus. Figure 2 shows a representation of the above mentioned factors, and also is annotated with applicable aspects of our narrative rendering pipeline. The nodes and edges in both visual representations of our abstract data were drawn manually with the aim of reducing edge crossing and length minimization. For larger and realistic datasets, we recommend using graph layout algorithms that minimize edge length and crossings. We can also note that the generated narratives are similar to those in Figure 1, which are a manifestation of the proposed design space.

### 4.4 Experimental Design

We used a mixed design in our study: between-subjects for VR, TN and NS; and within-subjects for DL (Table 1).

Table 1: Six groups with 25 participants per group; N=150 (25 × 6).

	H1	H2	H3	N1	N2	N3
VR	HD	HD	HD	NL	NL	NL
TN	OFF	ON	ON	OFF	ON	ON
NS	–	CU	IS	–	CU	IS
DL	S	S	S	S	S	S
	M	M	M	M	M	M
	H	H	H	H	H	H

Each participant saw all conditions of difficulty, but only one causality visualization. The relatively small total number of conditions enabled us to keep the session duration shorter than 30 minutes in duration to minimize fatigue and maximize attention for crowdworkers. In total, we planned to recruit 25 participants to each of six groups.

### 4.5 Analysis Tasks

Causal systems are complex structures that involve many processes (events) and messages propagating through a network of connections. In our user study, we use the words ‘node’ to mean a process and ‘link’ to mean a connection between ‘nodes’. The comprehensibility of a cause-effect relationship between, say, two nodes might also require an understanding of other effects that have propagated or will propagate through the system. Broadly, an understanding of *causality* might require a user to ask questions such as a) what factors caused a specific effect?, b) how does the effect on a node affect other connected nodes and to what extent?, c) what are the sequential and temporal impacts of this effect on the entire system?, d) how does this effect change or not change the earlier trend of the node. We adapted types of analysis tasks from previous work [19], and created 24 tasks for our participants.

Each task was of one of the following three types (QT): 1) Influence analysis (I), 2) Cause-effect analysis (C), and 3) Life-cycle Analysis (L). These task types are adaptations of causality information described in Section 3.3. In the tutorial, we explained to our participants that we choose certain *Intervention* and *Objective* nodes to analyze causal relationships. Table 2 shows the 9 types of tasks included in our study.

Within each group, each participant saw 12 graph systems (4 × [S, M, H]). A fixed order of increasing graph difficulty and tasks were used to improve familiarity by limiting chances of early task failure. Each graph system had 2 analysis sub-tasks. We distributed the first 8 task sub-types sequentially to each graph system within a DL, and alternated Trend (L4) with Spike (L5), in the event of a particular graph system showing a spike in a particular node (more on spikes in Section 5.2). Thus, each DL covered all the Task types (QT). Each sub-task required participants to read the question text and choose 1 out of 4 possible responses. Thus, for 150 participants, we planned to collect a total of 3,600 trials – 150 × 2 (questions) × 4 (graph systems) × 3 (DL).

These graph systems were modeled after abstract causal relationships with each node being labeled by alphabets (Figure 2). We avoided modelling real-world phenomena to avoid knowledge bias affecting

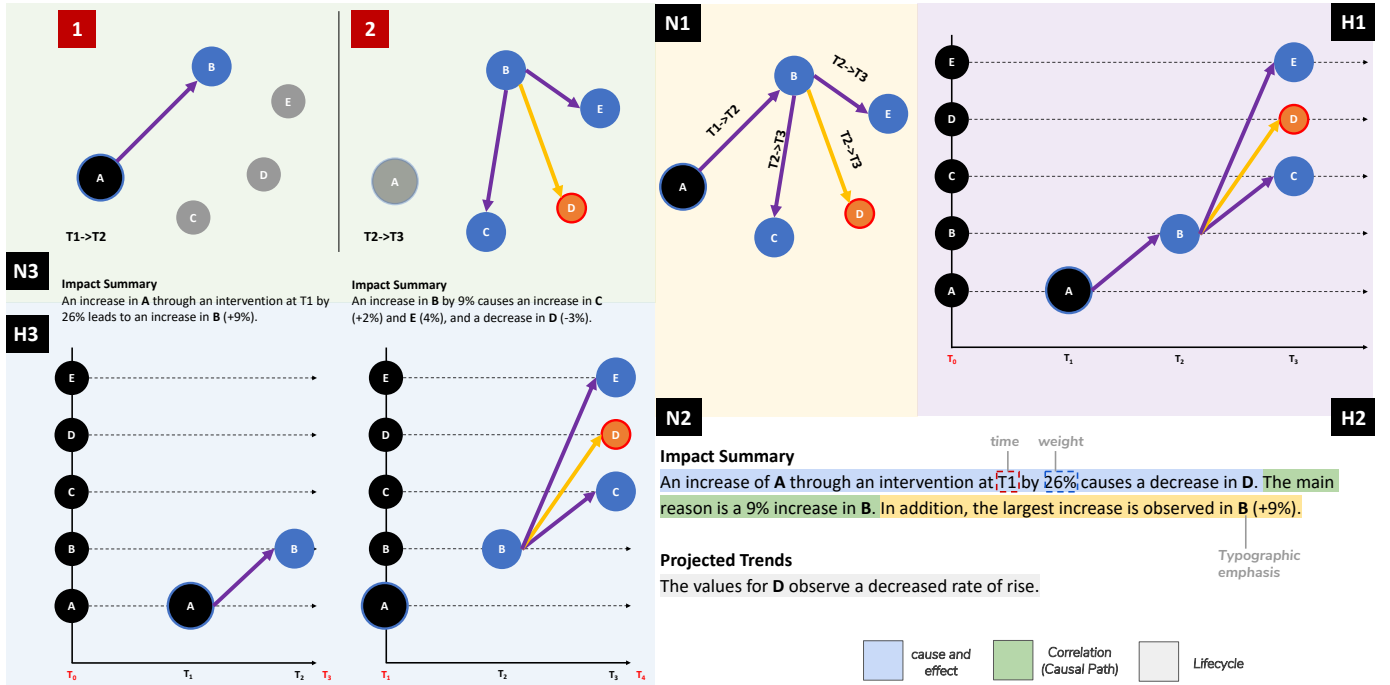


Fig. 2: Sample stimuli (DL: Simple) used in our 6 groups, with Hasse Diagram (HD) and Causal Graphs (NL). Sample narrative from groups N2 and H2 has been annotated with elements used from our design space.

performance. As described in Section 4.3, we assume that all 4 repetitions of a DL are equally simple or hard. Each survey page consisted of a chart (VR + TN + NS) corresponding to the experimental group.

#### 4.6 Collected Metrics

The tasks for all trials were controlled so that all participants saw the same graph systems, and were asked the same set of questions to allow comparison of participant performance between experimental groups.

**Performance Measures.** *Correctness* (TRUE or FALSE) is our primary performance measure to interpret the effectiveness of narratives in augmenting visual exploration. We also recorded time spent on each trial (from when the two tasks were displayed until the participant submitted both the answers) to understand if and how *Completion Time* influences correctness. However, due to a limitation in Qualtrics and the need to maintain low session time, we recorded both tasks together.

**Subjective Responses.** We also asked our participants to rate the ease-of-understanding of both graphs, and narratives (when applicable) after each DL. This was measured on a 5-point Likert scale (1: extremely easy, 5: extremely difficult). In the conditions where TN was ON, we also asked participants to rate the usefulness of the narratives on a 5-point scale (1: extremely useful, 5: not at all useful). Participants also provided open-ended feedback about graphs and narratives.

#### 4.7 Procedure

All recruitment was conducted via MTurk. Participants that fit the eligibility criteria opened the Qualtrics survey in a separate browser window. At the end of their participation, they copied a unique completion code back into the MTurk interface, and were later paid as their work was checked. Each session started with a consent form with waived signed consent. Failing to give consent terminated the experiment. Participants were instructed that they could abandon their session at any time. Unfortunately, due to constraints in Qualtrics, we were not able to pay participants who only completed a partial session.

After consenting, we showed participants a tutorial to explain the visualization and narrative that they would see. We also explained causal relationships, and how to interpret them. The tutorial included 3 examples of simple causal relationships. There was also a separate page explaining the visual mappings that were used in our visualizations.

This legend information was also accessible across every survey page, along with the visualization and a sample narrative from the tutorial.

Then participants were shown a single illustrated page of instructions explaining the task. Additionally, we also introduced 3 “attention trials,” which involved 3 easy cause-effect (C) tasks. Each DL had one attention trial. The purpose of these attention trials was to eliminate responses from crowdworkers who did not pay attention to the task and only “clicked through” the experiment. Participants that spent less than 7 minutes (roughly 1/3rd the average study duration) were also discarded; these participants were also not paid. We informed participants in the consent form that they would be paid only after response validation.

Typical sessions lasted between 7 to 50 minutes in duration. A few participants took significantly longer to complete their sessions, but our logs indicate that these participants took significant breaks between trials (presumably due to interruptions halfway through). Some participants also contacted us with reasons for the delay, such as trials genuinely being hard, or network issues. We believe that the effective time spent on the experiment was no more than 23 minutes. Participants were also asked some demographic questions about their age, education level, and knowledge of statistical concepts and graph visualizations.

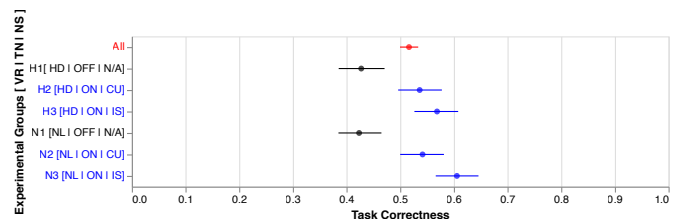


Fig. 3: Correctness comparison between 6 experimental groups (error bars represent 95% confidence intervals derived through bootstrapping).

#### 4.8 Results

We ran our crowdsourced graphical perception study on MTurk and collected a total of 4,824 responses from 201 unique respondents. This was higher than the 150 we had initially planned. During the recruitment process, we invalidated and rejected respondents (n=44) that just



Task type (QT)	Task sub-type	Question Structure
Influence (I)	Major Cause (I1)	Considering all the nodes, which node(s) caused the most influence on the system?
	Most Affected (I2)	Considering all the nodes, which node(s) were affected the most by changes in the system?
Causality (C)	Cause-Effect (C1)	Which statement best describes the cause-effect relationship between <I> and <O>?
	Major Factors (C2)	Choose all the nodes, including the objective that were affected by a change in <I>.
Lifecycle (L)	Max Increase (L1)	Excluding interventions/objectives, which node(s) goes through the greatest increase?
	Max Decrease (L2)	Excluding interventions/objectives, which node(s) goes through the greatest decrease?
	Time-Change (L3)	In the above system, at which time does node <X> increase/decrease the most from its initial level?
	Trend (L4)	Which statement best describes the trend that node <X> goes through?
	Spike (L5)	In the above system, which node goes through a sharp increase or decrease?

Table 2: List of task types and corresponding question structures for our user study. Each trial corresponded to a given task sub-type.

“clicked through” and completed the survey in less than 7 minutes. Thus, 157 participants were compensated for their time. During our analysis, we excluded data from participants (n=20) who spent less than 10 minutes on the survey. We expected a reasonable attempt to take 20 minutes based on our pilot, and believe that our complex perception tasks along with the tutorial required at least half of the estimated time. We also eliminated one participant that submitted a survey response after 3.8 hours. We present below results from the analysis of n=134 participants that completed 3,216 tasks (trials). The trials were distributed across experimental groups as follows: H1-480 (n=20) — H2-528 (n=22) — H3-528 (n=22) — N1-600 (n=25) — N2-552 (n=23) — N3-528 (n=22). We analyzed all our data using estimation methods to derive 95% confidence intervals (CIs). We employed non-parametric bootstrapping [15] with R = 1,000 iterations. This was done to follow current best practices for fair statistics in the field of HCI [14].

**Task Correctness.** Overall, we observed an accuracy of 51.6% (1,661/3,216) (Figure 3). First, we observed that the participants assigned to experimental groups that included a textual narrative — H2, H3, N2, N3 — outperformed those assigned to groups without a narrative — H1, N1. Secondly, we noted that participants that interacted with Causal Graphs (N1, N2, N3) performed better than those that used Hasse Diagrams (H1, H2, H3). We believe this to be a byproduct of participants being more familiar with node-link diagrams. We can infer that narratives providing Instantaneous NS (H3, N3) fare better across both VR, with the causal graphs (N3) outperforming all other groups.

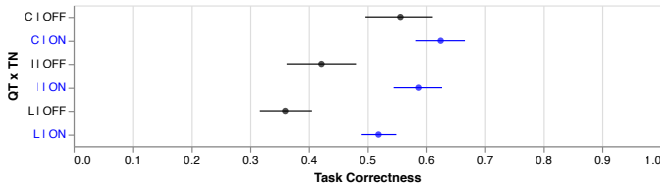


Fig. 4: Comparing effectiveness of the narrative in answering different question types (error bars represent 95% confidence intervals).

Figure 4 highlights the specific type of questions that the narratives were most effective in answering. Although the correctness increases for each condition that includes narratives, participants found the presence of narratives most helpful in answering the Influence (I) and Life-cycle Analysis (L) questions, for both type of visualizations. This became a useful insight while deciding on the modules in Section 5.

Figure 5 further highlights the improvement in the correctness of participant’s scores, across each difficulty level (S, M, H), when the visualizations are coupled with textual narratives. The comparatively lower task correctness improvements for the Hard (H) task type, in comparison to the Simple (S) and Medium (M) graph sets, can be attributed to the inherent added complexity, in terms of the added edges and number of nodes, within those datasets.

**Completion Time.** Completion time was measured per graph system. There were 4 repetitions of graph systems for every DL — S, M, H. In other words, completion time reflects time spent by a participant

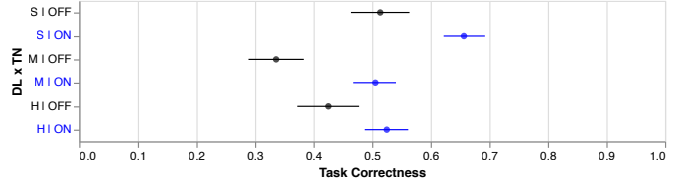


Fig. 5: Comparing the correctness between graph difficulty levels (error bars represent 95% confidence intervals).

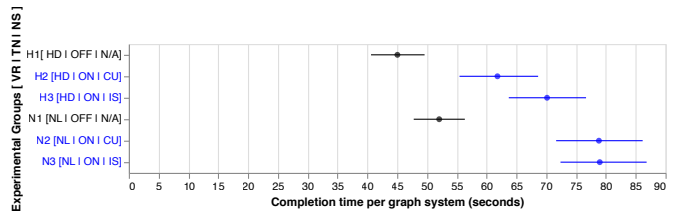


Fig. 6: Completion time for the test trials across different conditions (error bars represent 95% confidence intervals).

for two tasks. We eliminated outlier trials three standard deviations away from the mean for our analysis (Figure 6). We note that participants took much longer in groups where a narrative was present (H2, H3, N2, N3); with participants taking more time for Causal Graphs. Participants who were provided with narratives (H2, H3, N2, N3) took on an average 23.6 seconds more to answer 2 analysis tasks per graph as compared to participants without the narrative.

**Subjective Responses.** On an average, our participants ranked ease-of-understanding of the DL in the following order: Simple (mean=3.21), Medium (mean=3.68), Hard (mean=3.95). Graphs were rated as more easily understandable in the conditions where textual narratives were present versus when narratives were absent: Simple (mean=3.55 [OFF] vs. mean=3.0 [ON]); Medium (mean=3.82 [OFF] vs. mean=3.60 [ON]); Hard (mean=3.55 [OFF] vs. mean=3.0[ON]). The same trend was observed in the ease-of-understanding of narratives for each DL: Simple (mean=2.98), Medium (mean=3.38), Hard (mean=3.80). Additionally, the usefulness of the narratives decreased with increasing graph difficulty: Simple (mean=2.22), Medium (mean=3.38), Hard (mean=3.05).

Reviewing open-ended feedback showed us that difficulty understanding was attributed to unfamiliarity with a visualization: “It was somewhat challenging because I’m not familiar with this type of graph. Additionally, the abstract nature of our graph systems, and also novelty effects influenced difficulty. P111 says, “Their abstract nature was the most difficult to understand. Observing this with a real life example would make it easier to visualize and conceptualize.” Finally, our participants indicated that they used both visualization and narratives for causal inference: “I think having both the summary and the color coded chart makes it much easier to understand [...]”

## 4.9 Discussion

The main takeaway from our crowdsourced study is that narratives complement visualizations by providing descriptions to explain changes in the causal system. Based on our analysis of difficulty level and subjective responses, we believe that narratives will be more useful as the complexity of the system increases. Our participants also indicated that interactivity would have eased task difficulty—a known limitation in our study. We also strongly believe that interactivity can be leveraged to facilitate details on-demand in the narratives, especially when system complexity is bound to increase verbosity. The fact that Causal Graphs had a higher accuracy score than Hasse Diagrams further corroborated the prioritization of DAGs during the DOI step, and drove us to use them as the visualization medium in the CAUSEWORKS system. The study also encourages us to allocate a separate paragraph to talk about the trends followed by important nodes, owing to the high accuracy gains observed in the Lifecycle (L) task.

We also observed that the experimental groups that had higher accuracy also demonstrated higher completion times. We believe that the additional time stems from having to read the narratives before making an inference. This supports the aggregation DOI prioritization feature, wherein to reduce the verbosity of the textual snippets, nodes experiencing similar trends should be combined. This corroborates with our prediction that narratives aid in causal inference by providing descriptive texts that explain the changes occurring in the network.

## 5 APPLICATION: THE CAUSEWORKS SYSTEM

CAUSEWORKS, a system for intelligence analysis [51, 66, 67], integrates a range of network analysis, natural language generation (NLG), and data analytics techniques to develop coherent, concise, and explainable causal visualizations augmented by narratives for use by analysts. Drawing on our design space, our visualizations and narratives provide two main mechanisms (Figure 7): (1) a summary of changes and their impact on the objectives; and (2) additional projected trends.

### 5.1 System Overview

Figure 8a shows a screenshot of the CAUSEWORKS system. Various visual aspects of the design space (Section 3.5 and 3.6), such as ‘Node Coloring’, ‘Hyperlinking’ and ‘Brushing’, form an integral part of the narrative rendering process and its subsequent interactivity. Moreover, the performance metrics helped determine the usefulness of the various types of information snippets (Section 3.3), whereas the subjective responses played an important role in deciding the order in which the various information snippets are bundled together (Section 3.4). The left pane displays the *whiteboard*, a drawing space for causal graphs that allows the user to create and edit the network by adding, deleting, and modifying nodes as well as the edges amongst nodes, thus defining the semantics of the network. The *whiteboard* itself is unbounded, which allows the pane to incorporate a large number of nodes, and can be navigated using the *scrolling wheel* and the *magnifying scope* tools on the top left. Furthermore, the system also displays the chosen *objective* nodes, *intervention* nodes as well as the generated *narrative*, the placeholders for which can be seen in the right pane. The interactive GUI allows the user to select multiple intervention nodes and multiple objective nodes, and displays an explanatory narrative in real-time.

### 5.2 Extracting Causality

The impact summary elucidates how interventions introduced over one or more nodes propagate through the network and change target nodes (*Effect*) and the major nodes that help propagate that change (*Major Effect*). Note that the interventions could be made over one or multiple source nodes, and, further, they could be point interventions introduced at a specific timestamp or a sustained intervention introduced over a time period. The precise differences in how such interventions create observable changes in the target nodes is dependent on the causal model semantics (e.g., whether it is an ODE-based model or a discrete time-stamped Bayesian model), which is beyond the scope of this paper. Irrespective of the causal semantics, the impact summary encapsulates the cumulative effect of the interventions and identifies nodes in the causal path that depict the highest and least changes.

**Impact Summary.** Generating a summary of causal impacts is non-trivial due to the multitude of paths between source and target nodes. An effective narrative must reduce the number of words utilized to describe the associated effects. Below we detail how changes made on a set of interventions propagate through the network and affect the target nodes (*Effect*), the major nodes that help propagate that change (*Major Effect*), the intervening nodes that have no observable variation on the target nodes (*No Effect*) and, finally, the nodes that experience the most impact (*Max Effect*).

- **Effect Module:** The ‘Effect Module’ usually contributes the first sentence of the narrative and provides information on the propagation effect of each intervention on the specified target nodes. The set of source nodes and target nodes are grouped together based on common nodes in their causal path via a dictionary of  $\langle \text{key}, \text{value} \rangle$  pairs. Then, paths are grouped by merging source nodes that have at least one common node for each target node and then subsequently merging together target nodes. This often requires multiple passes and merge operations over the dictionary constructed. Figure 7 depicts a sample snippet detailing the effect of decreasing *Fossil Fuel Consumption* on *Quality of Marine Ecosystem*.
- **Major Effect Module:** Each ‘Effect’ sentence in the above module may or may not be followed by a sentence from the ‘Major Effect’ module. This module tries to capture the important nodes along the causal path between the set of source and target nodes, thus shining light on those causal path nodes that experienced the highest variation in either direction. This enlightens the user regarding the nodes that were the highest contributors to the effect propagation. Sample text snippets are shown in Figure 7 to list out the major factors that propagate the effect between the chosen intervention and target nodes.
- **No Effect Module:** This module articulates the specific source nodes that aren’t responsible for the change observed in the target node as well as the target nodes that remain unaffected due to the combined effect of all the interventions imposed on the network. This step can potentially lead to highly verbose sentences, thus affecting readability. To address this problem we introduce another grouping over the (input, output) node pairs, create n-grams of the source nodes, and articulate the most frequently occurring tuples amongst all the target nodes. These groups of tuples are then plugged into sentences, which are then included in the narrative. Figure 7 shows sample text snippets showing the non-impact of the interventions on ‘Government Policies against Climate Change’ as well as describes the  $\langle \text{Intervention}, \text{Objective} \rangle$  pairs that are not connected.
- **Maximum Effect Module:** The narrative generated until now focuses only on a subset of the whole network. This subset covers the edges and the nodes that lie in the causal paths between the set of *Intervention* and *Objective* nodes. However, there may still be nodes that might have been affected by the interventions but may have not been considered before. These nodes may provide interesting insights to the user and thus are worth adding to the final narrative. Hence, this module traverses through all the nodes in the system, instead of only the causal path nodes, and finds the nodes experiencing the maximum variation along both the positive and negative axis. Finally, it wraps both the nodes in a well structured sentence and attaches it to the end of the *Impact Summary* narrative. Figure 7 points out the most positively impacted node, (*Risk of Diseases*), as well as the most negatively impacted node, (*Atmospheric CO<sub>2</sub>*).

**Projected Trends.** While the *Impact Summary* articulates the overall influence of the source nodes on the target nodes, it leaves out information such as the temporal patterns observed by the nodes, or spikes in values that may have occurred in the course of the intervention, or other contextual information from external data sources (e.g., Wikipedia). We outline these parts of the narrative below.

### Impact Summary:

A decrease in 'Fossil Fuel Consumption', through an intervention 2 months from now, by 32% and an increase in 'Ozone Layer Depletion', through an intervention 3 months from now, by 20% causes an increase in 'Quality of Marine Ecosystem'. One of the major reasons behind the change is a 6% decrease in 'Sea Level above Baseline'. An increase in 'Land Degradation' by 21% causes a decrease in 'Food Availability'. One of the leading reasons behind the change is a 15% decrease in 'Agricultural Land'. However, no change is observed in 'Government Policies against Climate Change'. A change in 'Fossil Fuel Consumption' and 'Ozone Layer Depletion' causes no change in 'Food Availability'. Furthermore, a change in 'Land Degradation' leads to no change in 'Quality of Marine Ecosystem'. Besides this, the highest increase is observed in 'Risk of Diseases(+23%)' whereas the highest decrease is observed in 'Atmospheric CO<sub>2</sub>(-17%)'.

### Projected Trends:

A large number of factors were affected by the interventions. The top detected trends include the following: The values for 'Land Degradation' and 'Deforestation' observe an increased rate of rise. 'Land Degradation' is a process in which the value of the biophysical environment is affected by a combination of human-induced processes acting upon the land. The values for 'Methane Emissions' and 'Greenhouse Effect' observe an increased rate of fall. 'Greenhouse Effect' is the process by which radiation from a planet's atmosphere warms the planet's surface to a temperature above what it would be without this atmosphere. Radiatively active gases(i.e. greenhouse gases) in a planet's atmosphere radiate energy in all directions. Part of this radiation is directed towards the surface, warming it. The values for 'Quality of Marine Ecosystem' observe an increased rate of rise followed by a small decline before stabilizing to a constant value. For 'Quality of Marine Ecosystem', a sharp negative spike is projected 6 months from now.



Fig. 7: Narrative explaining a detailed causal network. The 'processes' or 'events' are depicted within single-quotes in this figure.

- **Time Series Module:** The Time Series module parses over the temporal information for entities in the causal path between the source and target nodes and captures key change trajectories observed over those nodes. Following a k-means clustering [41] over the temporal progressions across a 12 month period (with number of clusters selected using the Silhouette coefficient [52]), the clusters are sorted based on the number of nodes in each cluster and the high volume clusters are verbalized in the narrative. To limit the description length, a Pagerank score [46] is used as a filtering criterion to determine the most important nodes. Figure 7 details the time series patterns observed by 'Land Degradation', 'Deforestation', 'Methane Emissions' and 'Greenhouse Effect'.
- **Spike Detection Module:** Important nodes found in the previous step are further analyzed to check for presence of spikes or troughs during the timespan in consideration. This provides information to the user of any key abnormalities or milestones that might have occurred in these nodes. Each sentence of the *Time Series module* may or may not be followed by an output from the 'Spike Detection' module. For detecting spikes in the timeseries, the concept of a moving window is used to distinguish between gradual and sudden rises or falls in the time series value. Example text snippet showing the spikes in the time series value for 'Quality of Marine Ecosystem' is shown in Figure 7.
- **Wikification Module:** This constitutes the final module of the generated narrative. However, similar to the scenario with previous module, the *Wikification* module is also interleaved with the *Time Series* module to provide a seamless and continuous reading experience to the user. This module involves parsing through the summary paragraphs in the corresponding Wikipedia pages for important nodes mentioned in the *Timeseries* module and attaching key descriptive information to provide context to the narrative. If a Wikipedia page does not exist for the node mentioned, this module is skipped in its entirety.

## 5.3 Generating Textual Narratives

Figure 7 shows an example narrative<sup>1</sup>. Most importantly, the narrative highlights important aspects of the cumulative causal impacts caused by the interventions on the specified objectives. It also explains the effect that the interventions made on 'Fossil Fuel Consumption(-31%)', 'Land Degradation(+21%)' and 'Ozone Layer Depletion(+20%)' had on the target nodes, 'Quality of Marine Ecosystem', 'Food Availability' and 'Government Policies against Climate Change'.

Beyond the basic cause and effect relationships, the narrative also accounts for the changes to the entire system by mentioning the nodes

experiencing the highest rise and decline across the whole network. Furthermore, it clusters together nodes experiencing the same value patterns, and details the time-series patterns as well as spikes for the important nodes in those clusters. It also provides additional insights, by attaching the associated Wikipedia summary, for the available nodes.

## 5.4 Expert Review: CAUSEWORKS Narratives

We conducted an expert review [61] to validate the narrative engine.

**Method.** We recruited 5 experts who worked with causal systems in varying capacities (P1 and P5 were developers working on building system frameworks and analytics for causal systems; P2 and P4 were usability experts working on user research and visualization design for causal systems; and P3 was a visualization expert working with causal systems). We engaged with our experts in an hour-long semi-structured feedback session in a remote video call where they indirectly interacted with the CAUSEWORKS system. We encouraged the participants to think aloud, and interrupt at any point to ask questions, make, and share observations. These experts all work with planners who model and infer causality with the aid of visualization. P1, P3, and P5 were familiar with causal visualization and modelling, but were relatively unfamiliar with narrative generation. We gave a brief tutorial that explained both the visualization and the narrative structure. We created a scenario to demonstrate the system features. Two researchers collaborated with experts in the review: one researcher and the expert performed pair-analytics [2] while exploring the scenario, while the second researcher asked questions. We explained to our experts that the focus of the session was to critique the generated narratives in the system, and not features such as the visualization or other user-interface elements.

**Results.** Overall, our experts were impressed with how the narrative augments the causal graph in the system, especially to tackle large-scale causal systems "with multiple factors" (P2), and the narrative "puts everything in context" (P2) when beginning causal exploration. Referring to our design space, we summarize these results below:

### • Content Generation:

P3 suggested making the narrative more robust by including a 'model summary' of the underlying causal model: "Narrative should try (very) hard to have model scope (temporal and geographic scope)." Experts also noted that the narrative compensates for information loss from visual mapping by showing absolute values. Said P2: "[my] first instinct was to look at the graph because that would tell the specific percentage." P1 noted that the reliability of the narrative can be improved by "including the baseline trends."

- **Document Structuring and Aggregation:** Experts were satisfied with the presentation order of causal information. P3 commented that the order of presentation can remain the same for

<sup>1</sup>Please refer to our supplementary material for the reference graph.



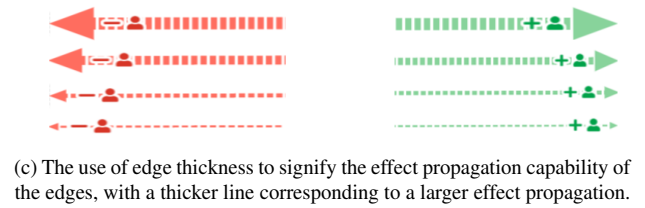
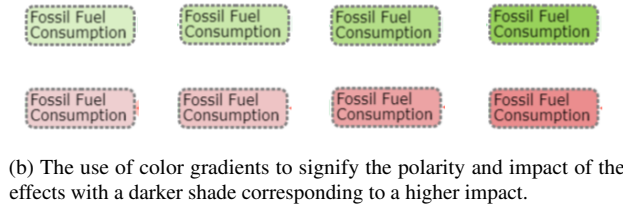
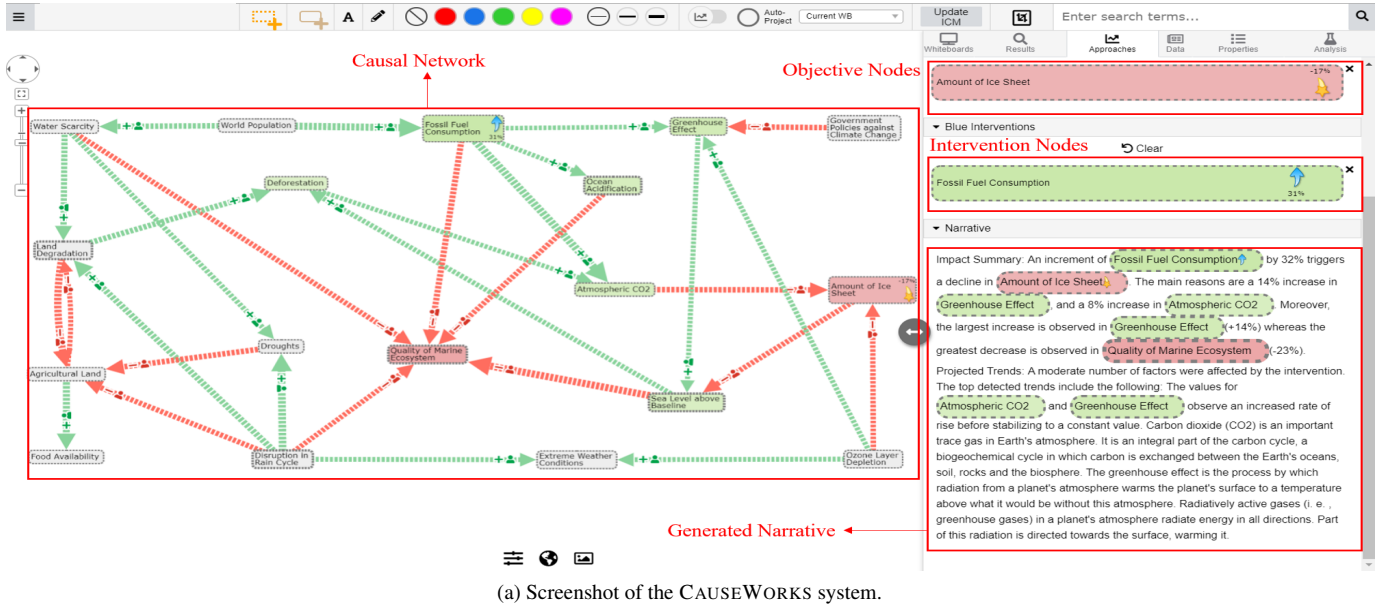


Fig. 8: Overview of interface and features of the CAUSEWORKS system.

a particular type of narrative and may change if more types are introduced. For example, the impact summary, projected trends, and model summary may have different causal information—aggregated and structured in different order.

- **Realization and Interaction:** All 5 experts acknowledged the idea of a character budget, and that the current rendering of the textual narratives can be improved with the use of rendering effects such as hierarchical lists (all 5) and interactivity such as Brushing (P2, P4), Search (P5), and Sliders (P5).

## 6 DISCUSSION

Our work indicates that textual narratives help users infer information from causal networks. More specifically, the user study shows our method faring favorably for time, correctness, and confidence of information absorption. These narratives also reinforce the ‘narrative intelligence’ viewpoint proposed by Blair and Meyer [7]. Narratives can be used to generate quick, precise, and informative reports (or subsections of reports) [38] owing to their structured representations.

The CAUSEWORKS system adds another layer of abstraction to narratives. Furthermore, the use cases presented here support past findings on visualization rhetoric [32] in combining interactivity with organized information presentations to enhance the decision-making process for the end-user. The findings are also in line with graph comics [3], which explored the effectiveness of using textual snippets with graphical images for communicating changes in dynamic networks.

As we mentioned earlier, we did not experimentally validate all of our design space in the crowdsourced study. However, we used our 4-step narrative rendering pipeline broadly in both the mockups of the crowdsourced study (e.g., causality information extraction, typographic emphasis, calculating order, and font size) and CAUSEWORKS (causality information extraction, calculating order, textual rendering and color, word-scale graphics, and interactivity through brushing). Our expert review results also validate our design choices.

We believe that future evaluations on the effectiveness of our narrative design space can help expand the space for causal systems, and eventually other systems. To test this hypothesis, we also plan on conducting another user study with the CAUSEWORKS system to evaluate the performance benefits offered by the system in a more interactive setting. Limitations in our work pertain to the scope of questions about causality that it can answer, e.g., we have not focused on communicating dynamics of the system as a whole. The current methodology is also designed in-line with the temporal datasets. Future work will be targeted towards a more generic approach that ingests non-temporal datasets as well as incorporating a visual DOI function to focus the user’s attention on important nodes and links.

## 7 CONCLUSION

We have presented a review of the design space for a specialized form of data-driven storytelling: the use of natural language narratives for causal network data. Based on the review, we isolated several interesting questions about the role of textual narratives for this purpose. To answer these questions, we conducted a large-scale crowdsourced user study where participants saw causal systems of increasing complexity. The data was displayed using one of two visualization techniques, causal graphs and Hasse diagrams, with and without the presence of textual narratives. The main finding is that the coupling of causality visualization techniques with textual narratives significantly increases accuracy and acts as a pivotal complement to information visualization.

## ACKNOWLEDGMENTS

This work was partially supported by the Defense Advanced Research Projects Agency (DARPA) under Contract Number FA8650-17-C-7720. The views, opinions and/or findings expressed are those of the authors and should not be interpreted as representing the official views or policies of the Department of Defense or the U.S. Government. The authors wish to thank all DARPA Causal Exploration collaborators for their support and encouragement.

## REFERENCES

- [1] DARPA, Causal Exploration, <https://www.darpa.mil/program/causal-exploration>, 2018.
- [2] R. Arias-Hernández, L. T. Kaastra, T. M. Green, and B. D. Fisher. Pair analytics: Capturing reasoning processes in collaborative visual analytics. In *Proceedings of the Hawaii International International Conference on Systems Science*, pp. 1–10. IEEE Computer Society, 2011. doi: 10.1109/HICSS.2011.339
- [3] B. Bach, N. Kerracher, K. W. Hall, S. Carpendale, J. Kennedy, and N. Henry Riche. Telling stories about dynamic networks with graph comics. In *Proceedings of the ACM Conference on Human Factors in Computing Systems*, pp. 3670–3682, 2016. doi: 10.1145/2858036.2858387
- [4] J. Bae, T. Helldin, and M. Riveiro. Understanding indirect causal relationships in node-link graphs. *Computer Graphics Forum*, 36(3):411–421, 2017. doi: 10.1111/cgf.13198
- [5] F. Beck, M. Burch, S. Diehl, and D. Weiskopf. The state of the art in visualizing dynamic graphs. In R. Borgo, R. Maciejewski, and I. Viola, eds., *State of the Art Reports for the Eurographics Conference on Visualization*. Eurographics Association, 2014. doi: 10.2312/eurovisstar.20141174
- [6] F. Beck and D. Weiskopf. Word-sized graphics for scientific texts. *IEEE Transactions on Visualization and Computer Graphics*, 23(6):1576–1587, 2017. doi: 10.1109/TVCG.2017.2674958
- [7] D. Blair and T. Meyer. Tools for an interactive virtual cinema. In *Creating Personalities for Synthetic Actors*, pp. 83–91. Springer, 1997. doi: 10.1007/BFb0030572
- [8] M. Bunge. *Causality and Modern Science*. Dover Publications, 1979.
- [9] R. Chang, C. Ziemkiewicz, T. M. Green, and W. Ribarsky. Defining insight for visual analytics. *IEEE Computer Graphics and Applications*, 29(2):14–17, 2009. doi: 10.1109/MCG.2009.22
- [10] A. Chatzimparmpas, R. M. Martins, I. Jusufi, K. Kucher, F. Rossi, and A. Kerren. The State of the Art in Enhancing Trust in Machine Learning Models with the Use of Visualizations. *Computer Graphics Forum*, 2020. doi: 10.1111/cgf.14034
- [11] M. F. Dahlstrom. Using narratives and storytelling to communicate science with nonexpert audiences. *Proceedings of the National Academy of Sciences*, 111(Supplement 4):13614–13620, 2014. doi: 10.1073/pnas.1320645111
- [12] N. Diakopoulos, J. DiMicco, J. Hullman, K. Karahalios, and A. Perer. Telling stories with data: The next chapter—a visweek 2011 workshop, 2011.
- [13] J. DiMicco, M. McKeon, and K. Karahalios. Telling stories with data—a VisWeek 2010 workshop, 2010.
- [14] P. Dragicevic. *Fair Statistical Communication in HCI*, pp. 291–330. Springer International Publishing, Cham, 2016. doi: 10.1007/978-3-319-26633-6\_13
- [15] B. Efron. Bootstrap methods: another look at the jackknife. In *Breakthroughs in statistics*, pp. 569–593. Springer, 1992. doi: 10.1007/978-1-4612-4380-9\_41
- [16] W. Eisner. *Graphic Storytelling and Visual Narrative*. W. W. Norton & Company, New York, NY, USA, 2008.
- [17] N. Elmqvist and P. Tsigas. Causality visualization using animated growing polygons. In *Proceedings of the IEEE Symposium on Information Visualization*, pp. 189–196, 2003. doi: 10.1109/INFVIS.2003.1249025
- [18] N. Elmqvist and P. Tsigas. Growing squares: Animated visualization of causal relations. In *Proceedings of the ACM Symposium on Software Visualization*, pp. 17–26, 2003. doi: 10.1145/774833.774836
- [19] N. Elmqvist and P. Tsigas. Animated visualization of causal relations through growing 2d geometry. *Information Visualization*, 3:154–172, 07 2004. doi: 10.1057/palgrave.ivs.9500074
- [20] F. Elwert. *Graphical Causal Models*, pp. 245–273. 03 2013. doi: 10.1007/978-94-007-6094-3\_13
- [21] M. Feng, C. Deng, E. M. Peck, and L. Harrison. The effects of adding search functionality to interactive visualizations on the web. In R. L. Mandryk, M. Hancock, M. Perry, and A. L. Cox, eds., *Proceedings of the ACM Conference on Human Factors in Computing Systems*, pp. 137:1–137:13. ACM, New York, NY, USA, 2018. doi: 10.1145/3173574.3173711
- [22] M. Gambhir and V. Gupta. Recent automatic text summarization techniques: a survey. *Artificial Intelligence Review*, 47:1–66, 2016. doi: 10.1007/s10462-016-9475-9
- [23] D. Geiger, T. Verma, and J. Pearl. d-separation: From theorems to algorithms. In *Machine Intelligence and Pattern Recognition*, vol. 10, pp. 139–148. Elsevier, 1990. doi: 10.1016/B978-0-444-88738-2.50018-X
- [24] N. D. Gershon and W. Page. What storytelling can do for information visualization. *Communications of the ACM*, 44(8):31–37, 2001. doi: 10.1145/381641.381653
- [25] P. Goffin, J. Boy, W. Willett, and P. Isenberg. An exploratory study of word-scale graphics in data-rich text documents. *IEEE Transactions on Visualization and Computer Graphics*, 23(10):2275–2287, 2017. doi: 10.1109/TVCG.2016.2618797
- [26] J. Gottschall. *The Storytelling Animal: How Stories Make Us Human*. Mariner Books, New York, NY, USA, 2012. doi: 10.1075/ssol.2.2.07bor
- [27] D. Heckerman, D. M. Chickering, C. Meek, R. Rounthwaite, and C. M. Kadie. Dependency networks for inference, collaborative filtering, and data visualization. *Journal of Machine Learning Research*, 1:49–75, 2000. doi: 10.1162/153244301753344614
- [28] J. Heer and M. Bostock. Crowdsourcing graphical perception: Using mechanical turk to assess visualization design. In *Proceedings of the ACM Conference on Human Factors in Computing Systems*, pp. 203–212. ACM, New York, NY, USA, 2010. doi: 10.1145/1753326.1753357
- [29] R. J. Heuer Jr. Analysis of competing hypotheses. *Psychology of Intelligence Analysis*, pp. 95–110, 1999.
- [30] F. Hohman, A. Srinivasan, and S. M. Drucker. TeleGam: Combining visualization and verbalization for interpretable machine learning. In *Proceedings of the IEEE Conference on Visualization*, pp. 151–155, 2019. doi: 10.1109/VISUAL.2019.8933695
- [31] J. Hullman and N. Diakopoulos. Visualization rhetoric: Framing effects in narrative visualization. *IEEE Transactions on Visualization and Computer Graphics*, 17(12):2231–2240, 2011. doi: 10.1109/TVCG.2011.255
- [32] J. Hullman and N. Diakopoulos. Visualization rhetoric: Framing effects in narrative visualization. *IEEE Transactions on Visualization and Computer Graphics*, 17(12):2231–2240, 2011. doi: 10.1109/TVCG.2011.255
- [33] J. Hullman, S. Drucker, N. H. Riche, B. Lee, D. Fisher, and E. Adar. A deeper understanding of sequence in narrative visualization. *IEEE Transactions on Visualization and Computer Graphics*, 19(12):2406–2415, 2013. doi: 10.1109/TVCG.2013.119
- [34] N. Kadaba, P. Irani, and J. Leboe. Visualizing causal semantics using animations. *IEEE Transactions on Visualization and Computer Graphics*, 13:1254–61, 11 2007. doi: 10.1109/TVCG.2007.70528
- [35] S. Kent. *Words of Estimative Probability*. 1964.
- [36] D. Koller and N. Friedman. *Probabilistic Graphical Models: Principles and Techniques - Adaptive Computation and Machine Learning*. The MIT Press, 2009. doi: 10.5555/1795555
- [37] R. Kosara. Story points in Tableau Software. Keynote at Tableau Customer Conference, Sept. 2013.
- [38] S. Latif and F. Beck. Interactive map reports summarizing bivariate geographic data. *Visual Informatics*, 3(1):27 – 37, 2019. Proceedings of PacificVAST. doi: 10.1016/j.visinf.2019.03.004
- [39] S. Latif and F. Beck. Vis author profiles: Interactive descriptions of publication records combining text and visualization. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):152–161, 2019. doi: 10.1109/TVCG.2018.2865022
- [40] T. M. Leitch. *What Stories Are: Narrative Theory and Interpretation*. Pennsylvania State University Press, University Park, PA, 1986.
- [41] J. MacQueen et al. Some methods for classification and analysis of multivariate observations. In *Proceedings of the Berkeley Symposium on Mathematical Statistics and Probability*, vol. 1, pp. 281–297. Oakland, CA, USA, 1967.
- [42] R. Metoyer, Q. Zhi, B. Janczuk, and W. Scheirer. Coupling story to visualization: Using textual analysis as a bridge between data and interpretation. In *Proceedings of the ACM Conference on Intelligent User Interfaces*, pp. 503–507, 2018. doi: 10.1145/3172944.3173007
- [43] M. Moezzi, K. B. Janda, and S. Rotmann. Using stories, narratives, and storytelling in energy and climate change research. *Energy Research & Social Science*, 31:1–10, 2017. doi: 10.1016/j.erss.2017.06.034
- [44] A. Nenkova and K. McKeown. A survey of text summarization techniques. In *Mining Text Data*, 2012. doi: 10.1007/978-1-4614-3223-4\_3
- [45] C. North. Toward measuring visualization insight. *IEEE Computer Graphics & Applications*, 26(3):6–9, 2006. doi: 10.1109/MCG.2006.70
- [46] L. Page, S. Brin, R. Motwani, and T. Winograd. The PageRank citation ranking: Bringing order to the web. In *Proceedings of the ACM Conference on the World Wide Web*, 1999. doi: 10.1.1.38.5427
- [47] J. Pearl. *Causality: Models, Reasoning and Inference*. Cambridge University Press, USA, 2nd ed., 2009. doi: 10.5555/1642718
- [48] E. Reiter and R. Dale. Building applied natural language generation systems. *Natural Language Engineering*, 3:57–87, 1997. doi: 10.1017/

- [49] E. Reiter and R. Dale. *Building Natural Language Generation Systems*. Studies in Natural Language Processing. Cambridge University Press, 2000. doi: 10.1017/CBO9780511519857
- [50] N. H. Riche, C. Hurter, N. Diakopoulos, and S. Carpendale. *Data-Driven Storytelling*. A. K. Peters, Ltd., USA, 1st ed., 2018.
- [51] C. Rooney, S. Attfield, B. L. W. Wong, and S. Choudhury. Invisque as a tool for intelligence analysis: The construction of explanatory narratives. *International Journal of HumanComputer Interaction*, 30(9):703–717, 2014. doi: 10.1080/10447318.2014.905422
- [52] P. J. Rousseeuw. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, 20:53 – 65, 1987. doi: 10.1016/0377-0427(87)90125-7
- [53] R. C. Schank and R. P. Abelson. Knowledge and memory: The real story. In J. Robert S. Wyer, ed., *Advances in Social Cognition*, vol. 8, pp. 1–85. Lawrence Erlbaum Associates, Hillsdale, NJ, USA, 1995.
- [54] E. Segel and J. Heer. Narrative visualization: Telling stories with data. *IEEE Transactions on Visualization and Computer Graphics*, 16(6):1139–1148, 2010. doi: 10.1109/TVCG.2010.179
- [55] R. Sevastjanova, F. Beck, B. Ell, C. Turkay, R. Henkin, M. Butt, D. A. Keim, and M. El-Assady. Going beyond visualization: Verbalization as complementary medium to explain machine learning models. In *Workshop on Visualization for AI Explainability at IEEE VIS*, 2018.
- [56] R. D. Shachter. Bayes-ball: The rational pastime (for determining irrelevance and requisite information in belief networks and influence diagrams). *arXiv preprint arXiv:1301.7412*, 2013. doi: 10.5555/2074094.2074151
- [57] D. Sless. *Learning and Visual Communication*. Wiley, New York, NY, USA, 1981.
- [58] A. Srinivasan, S. M. Drucker, A. Endert, and J. Stasko. Augmenting visualizations with interactive data facts to facilitate interpretation and communication. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):672–681, 2018. doi: 10.1109/TVCG.2018.2865145
- [59] A. Srinivasan, H. Park, A. Endert, and R. C. Basole. Graphiti: Interactive specification of attribute-based edges for network modeling and visualization. *IEEE Transactions on Visualization and Computer Graphics*, 24(1):226–235, 2017. doi: 10.1109/TVCG.2017.2744843
- [60] J. J. Thomas and K. A. Cook. *Illuminating the Path: The Research and Development Agenda for Visual Analytics*. IEEE Computer Society, addrIEEECS, 2005.
- [61] M. Tory and T. Möller. Evaluating visualizations: Do expert reviews work? *IEEE Computer Graphics and Applications*, 25(5):8–11, 2005. doi: 10.1109/MCG.2005.102
- [62] J. Vansina. *Oral Tradition as History*. University of Wisconsin Press, Madison, WI, USA, 1985. doi: 10.2307/3601125
- [63] F. Viégas and M. Wattenberg. Communication-minded visualization: A call to action. *IBM Systems Journal*, 45(4):801–812, 2006. doi: 10.1147/sj.454.0801
- [64] J. Wang and K. Mueller. The visual causality analyst: An interactive interface for causal reasoning. *IEEE Transactions on Visualization and Computer Graphics*, 22(1):230–239, 2016. doi: 10.1109/TVCG.2015.2467931
- [65] C. Ware. Perceiving complex causation through interaction. In *Proceedings of the Symposium on Computational Aesthetics*, pp. 29–35. Association for Computing Machinery, New York, NY, USA, 2013. doi: 10.1145/2487276.2487279
- [66] W. Wright, D. Schroh, P. Proulx, A. Skaburskis, and B. Cort. The sand-box for analysis: Concepts and methods. In *Proceedings of the ACM Conference on Human Factors in Computing Systems*, pp. 801–810. Association for Computing Machinery, New York, NY, USA, 2006. doi: 10.1145/1124772.1124890
- [67] W. Wright, D. Sheffield, and S. Santosa. Argument mapper: Countering cognitive biases in analysis with critical (visual) thinking. In *Proceedings of the International Conference on Information Visualisation*, pp. 250–255, 2017.
- [68] C. Xiong, J. Shapiro, J. Hullman, and S. Franconeri. Illusion of causality in visualized data. *IEEE Transactions on Visualization and Computer Graphics*, 26(1):853–862, 2019. doi: 10.1109/TVCG.2019.2934399

# Once Upon A Time In Visualization: Understanding the Use of Textual Narratives for Causality

## Supplementary Materials

### A ADDITIONAL FIGURES

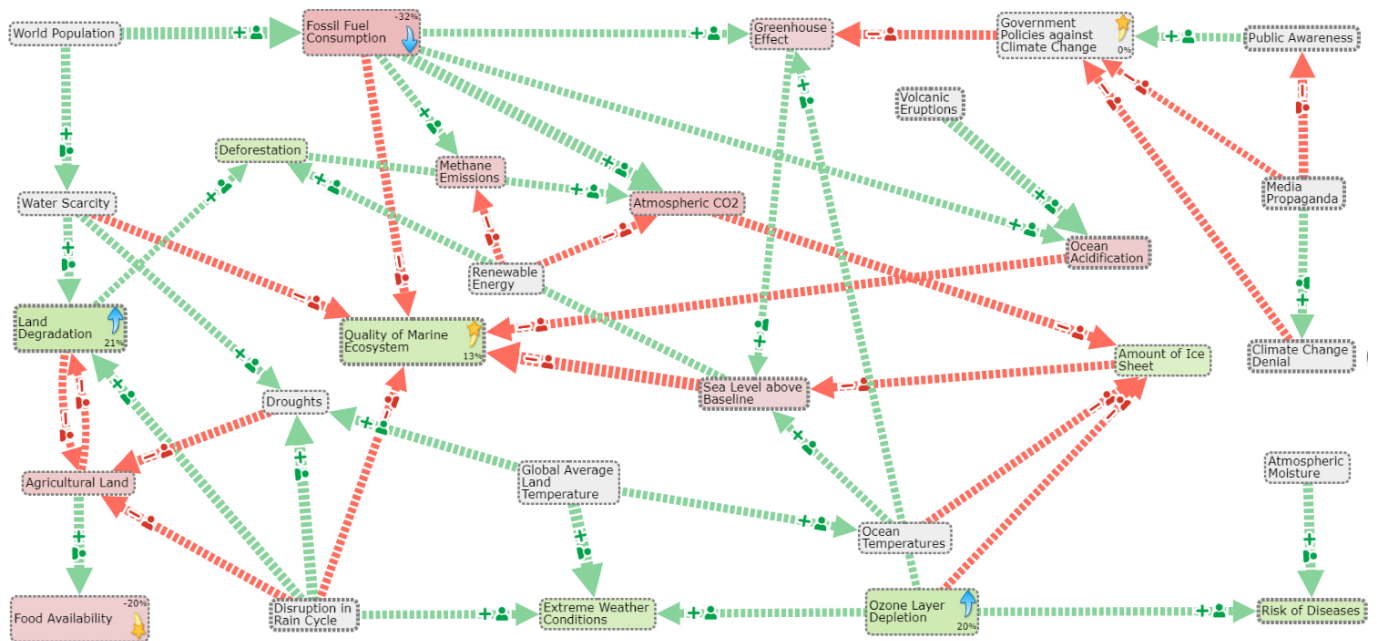


Fig. A: Sample causal network visualized in CAUSEWORKS.

### B VIDEO

We have attached a companion video showcasing the CAUSEWORKS system.

### C GRAPH SAMPLES

We have also attached images of the graph systems that we used in our crowdsourced study.

The naming convention is: **[Experimental condition][Difficulty Level][Repetition]**

Example: [H1][N][1] — The file name will be 'H1N1.png'

### D ANONYMIZED PERFORMANCE DATA

We have attached anonymized correctness, time, and subjective responses.