

# Networks of Texts and People

Andrew Salway  
Uni Research, Bergen

# Networks of Texts and People (NTAP)

- New methods to analyse the distribution, flow and development of statements in online social networks
  - e.g. climate change discourse in the blogosphere
  - to contribute to social science research into framing, information diffusion and polarization
- Text analysis, network analysis and visualization
  - text analysis: corpus linguistics, text mining and IE
  - identify salient linguistic constructions / information structures
    - ➔ analyse their occurrence over the blogosphere and over time

# Networks of Texts and People (NTAP)

- Dept. of Information Science and Media Studies, UiB and Computational Language Unit, Uni Research
  - Nick Diakopoulos
  - Dag Elgesem
  - Knut Hofland
  - Andrew Salway
  - Lubos Steskal
  - Samia Touileb
- Funded by Research Council of Norway, 2012-2016.
- [www.ntap.no](http://www.ntap.no)

# Climate Change Corpus

January 1, 2012 to March 1, 2012

8,415  
authors

112,510  
posts

8,123  
blogs

(enter search terms here)

Go



## Statements

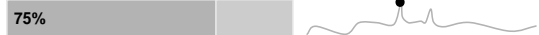
Sort by frequency

### Climate change

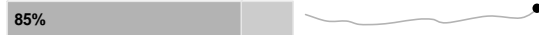
is\_caused\_by, "natural causes" 35



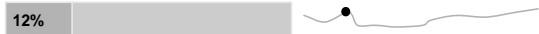
is\_caused\_by, "the excessive amount of carbon emissions poured into the atmosphere" 32



is, "changing the arctic" 25



is, "based on fraudulent science" 16



will be, "arrested by peak oil" 5



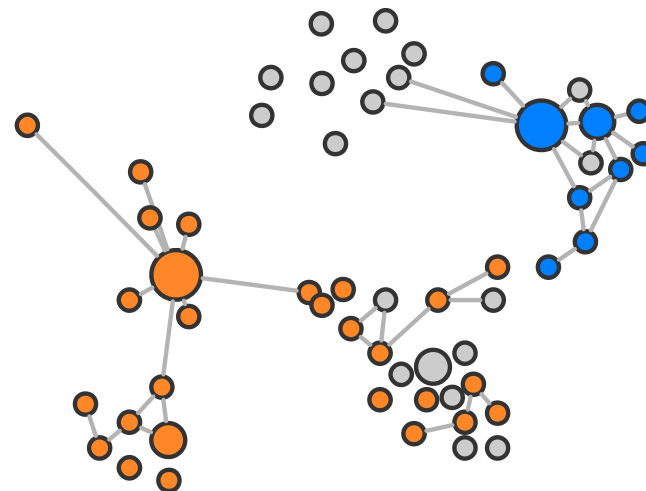
### Global warming

is\_caused\_by, "burning fossil fuels like coal" 45



## Network

[View Options](#)



# NTAP blog corpus

- English-language blogs that mention broad climate change issues across science, politics, environment, etc.
- Method:
  - Handpicked 5 seed blogs - well connected in the blogosphere, representing different viewpoints in the climate debate
  - Harvested all posts from the seed blogs and extracted key terms to be used to determine topical relevance in the crawl:
    - frequent words typical of the domain: e.g. *climate, global, carbon, emissions, temperature, sea, solar, greenhouse ...*
    - n-grams ( $2 \leq n \leq 5$ ) containing words, e.g. *climate change, climate science, carbon dioxide, emissions trading, sea level ...*
  - Breadth first crawl from the seed blogs:
    - harvest blog if English-language and has a key term on homepage
    - follow links from the homepage of each blog
    - limited to WordPress and Blogspot blogs

# NTAP blog corpus

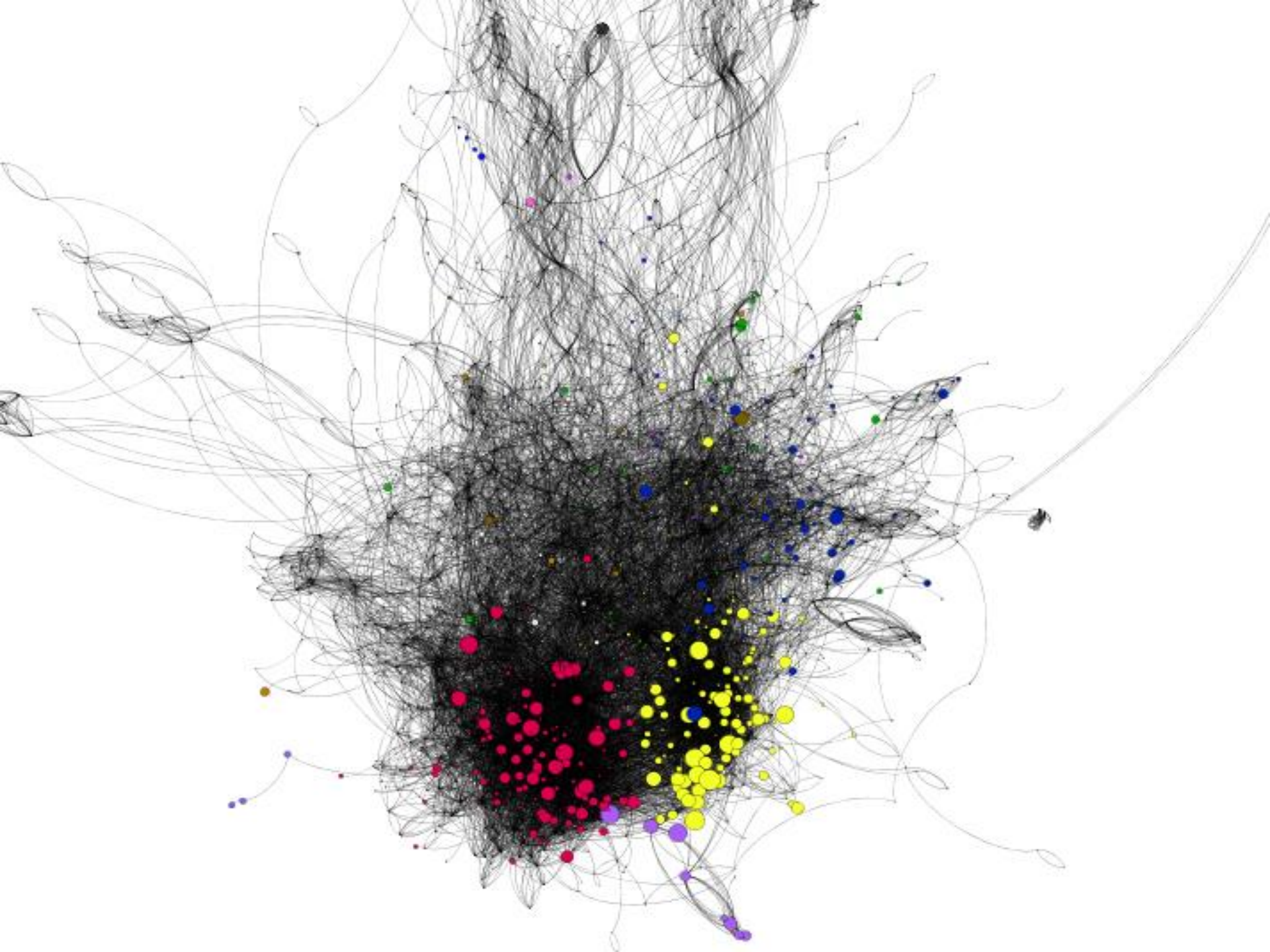
- Crawl carried out June-September 2012
- Harvested the complete content of about 3,000 English-language blogs, about 1.5m blog posts
- Text content of each post extracted using the Alchemy API and stored in a MongoDB database, with date
- Hyperlinks extracted links and stored these in a Neo4j graph database.

# Network analysis

- Used a community detection algorithm based on modularity maximization
  - blogs grouped to maximize inter-group hyperlinking (Louvain method, implemented in the Gephi tool)
- This suggested 11 major communities in the corpus, accounting for about 60% of all blogs
- A selection of blogs were inspected to manually code each community as “accepting”, “skeptical” or “neutral” regarding anthropogenic global warming

# Topic modelling

- Latent Dirichlet Allocation (LDA) used to identify topics within the corpus, using MALLET tool
- Two (out of 20) topics related strongly to climate change:
  - “climate change science”: *climate, warming, global, change, ice, data, temperature, years, science, scientists, carbon, sea, earth, year, ocean, time, temperatures, scientific, research*
  - “climate change politics”: *climate, change, countries, world, environmental, international, development, global, emissions, carbon, india, environment, people, government, nations, policy, china, issues, sustainable*
- Other topics: “energy”, “wildlife”, “legal”, “education”, “economic policy”, “transportation”, “American politics”, “storms and floods”, “farming”, “health”, “new age”, and some noise (incoherence topics and non-English words)



# Sub-corpora for two communities

- Our network analysis, topic modelling and manual coding suggested two communities concerned with climate change science, one broadly “accepting”, one “skeptical”

## ➔ sub-corpora

- “accepting” 204 blogs, 69k posts, 27m words
- “skeptical” 417 blogs, 290k posts, 127m words

# Sub-corpora for two communities

- Little difference between the 300 most frequent words (and the 300 top keywords) in the two sub-corpora
  - perhaps not surprising given role of topic modelling in selecting the sub-corpora, and a large volume of peripheral material?
- But, some lower frequency words, and word clusters, seem to be preferred by one community or the other...

# Sub-corpora for two communities

	“accepting”		“skeptical”	
	(204 blogs, 69k posts)		(417 blogs, 290k posts)	
	no. of blogs	no. of posts	no. of blogs	no. of posts
<b>acidification</b>	64	1782	75	412
<b>coral</b>	70	939	117	701
<b>ocean</b>	122	4627	206	6550
<b>species</b>	112	3455	196	5319
<b>tax</b>	100	2339	213	17,180
<b>Gore</b>	84	845	200	10,669

# Sub-corpora for two communities

	“accepting” (204 blogs, 69k posts)		“skeptical” (417 blogs, 290k posts)	
	no. of blogs	no. of posts	no. of blogs	no. of posts
<b>climate science</b>	99	2115	155	3551
<b>anthropogenic c c</b>	55	349	79	360
<b>human-caused c c</b>	37	124	57	201
<b>human-induced c c</b>	33	144	57	273
<b>man-made c c</b>	31	73	97	566
<b>climate change denial</b>	38	165	41	109
<b>climate alarmist</b>	6	6	47	216

# Sub-corpora for two communities

## Mentions of causes of climate change

	<b>“skeptical” (35k instances of climate change)</b>	<b>“accepting” (22k instances of climate change)</b>
cause [verb forms] climate change	147	49
cause(s) of climate change	117	34
contribute(s d) to climate change	68	34
affect [verb forms] climate change	18	7
lead to [verb forms] climate change	6	5
result in [verb forms] climate change	3	2
<b>TOTAL</b>	<b>359</b>	<b>131</b>

# Sub-corpora for two communities

## Mentions of effects of climate change

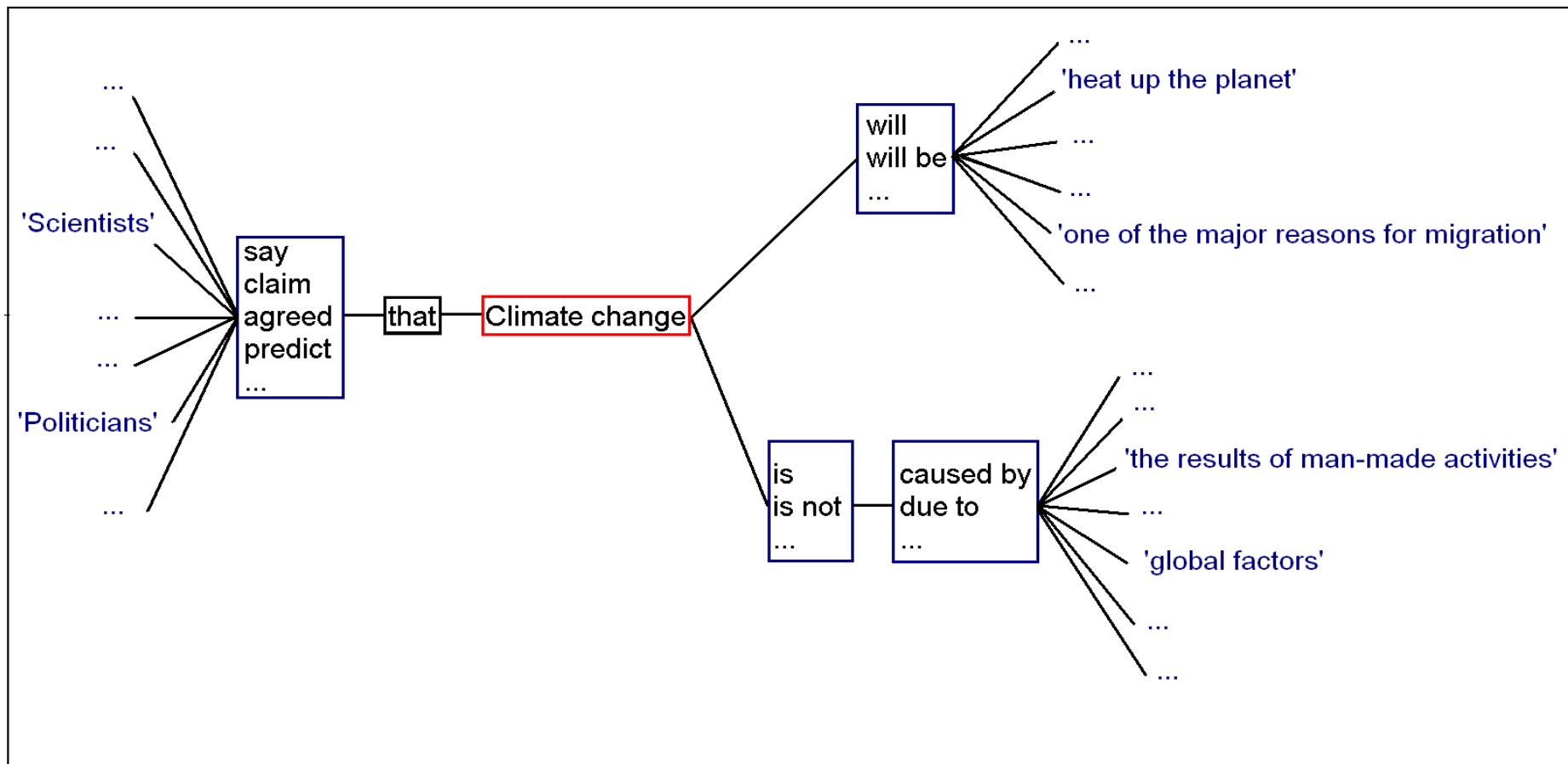
	<b>“skeptical” (35k instances of climate change)</b>	<b>“accepting” (22k instances of climate change)</b>
result   effect(s)   impact(s)   consequence(s) of climate change	1,412	1,034
due to climate change	179	148
climate change cause   affect   lead to   result in   contribute to [verb forms]	68	34
<b>TOTAL</b>	<b>1,659</b>	<b>1,216</b>



# Local grammar induction

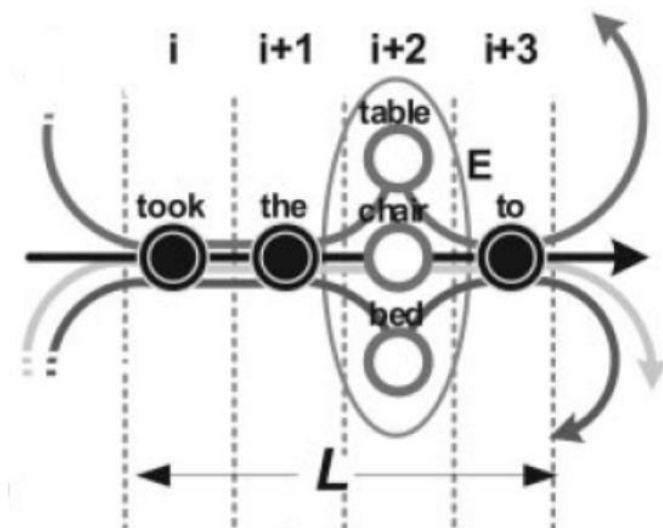
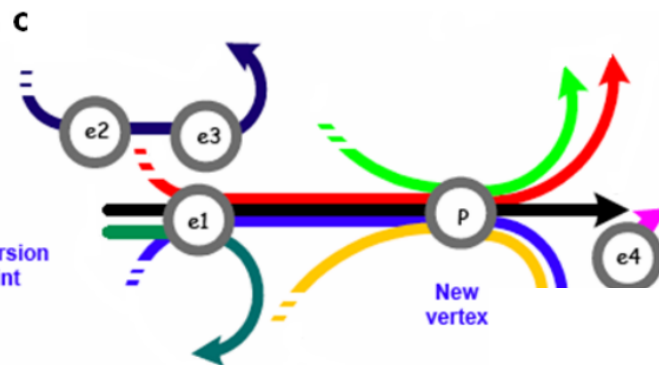
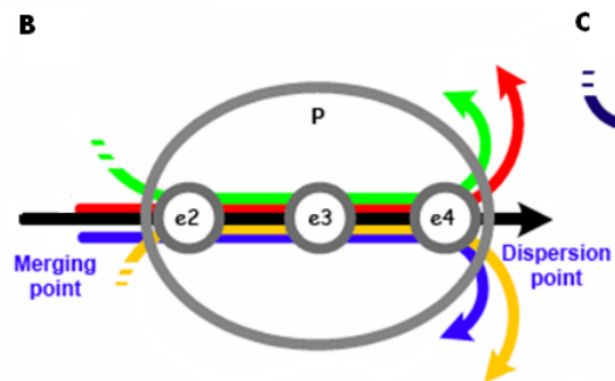
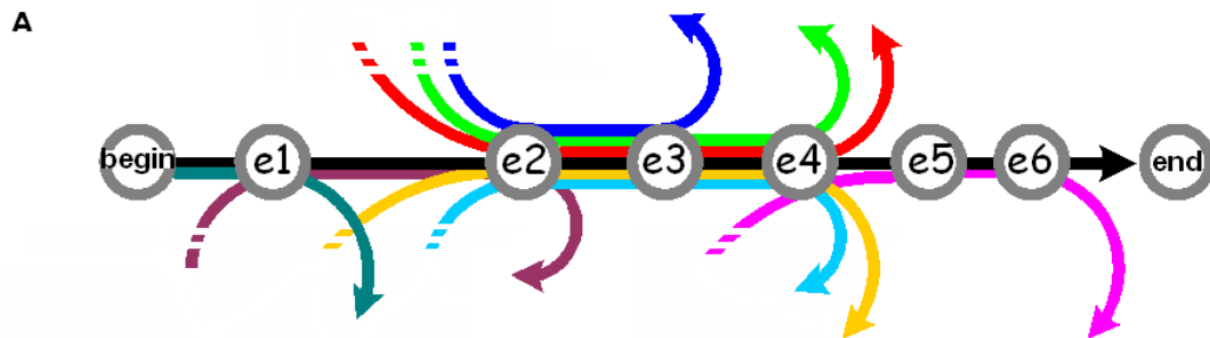
- For exploratory purposes, we need an overview of what is typically written about key terms in the corpus
  - 100,000's lines for “climate change”, 10,000's lines for “sea levels”
  - manual analysis of concordance lines is not feasible (?)
- Can we exploit relatively restricted and repetitive language use in the corpus (i.e. same text genre, topic, function) to induce and visualize interesting patterning?
  - cf. collocations, clusters, phrases, constructions, lexical bundles, n-grams, collocation frameworks, formulaic units, multiword expressions, phrase frames
  - cf. The Sketch Engine, The Word Tree

# Part of an idealised local grammar for “climate change”



# Grammar Induction

- Grammar induction algorithms generate a grammar from an unannotated corpus, i.e. based on formal distributional properties of words(cf. Harris 1954)
- For example, ADIOS (Automatic Distillation of Structure, Solan et al. 2005)
  - Organises all input sequences (sentences) on a graph with a node for each unique vocabulary item
  - In each iteration, simultaneously forms patterns (syntagmatic units) and equivalence classes (paradigmatic units)



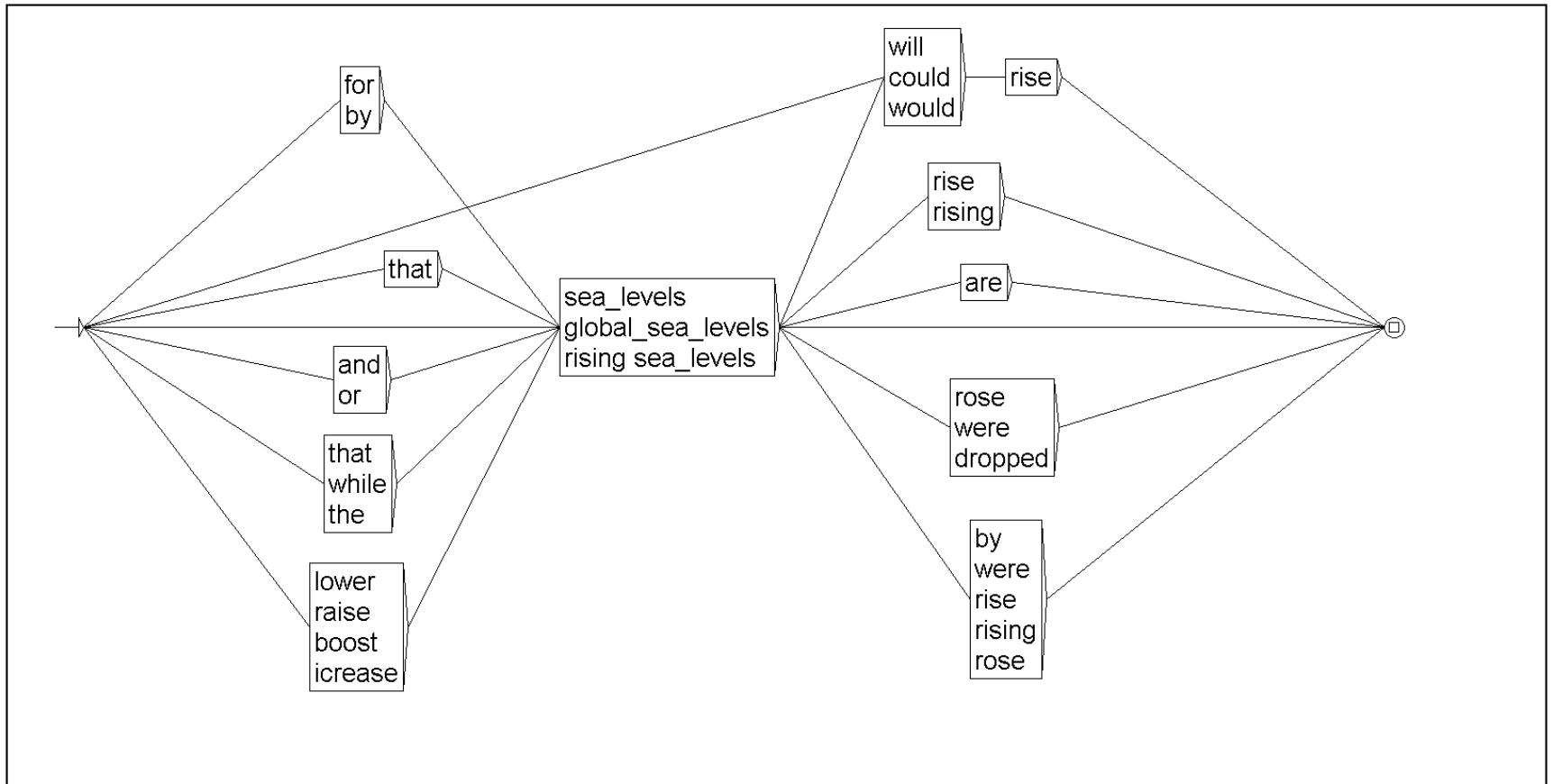
# Adapting ADIOS

(PhD work of Samia Touileb)

- Input is a set of sentences that contain the same key term
- The sentences are turned into sets of snippets, of varying sizes, around the key term
- Sets of snippets are presented to algorithm in size order, e.g. N iterations for the smallest snippets, then N for the next size, and so on
  - to focus the algorithm on patterns local to the key term
- After each iteration the 'best' patterns and equivalence classes are selected and inserted into the input file
  - to make further patterning more explicit

P\_o (of climate\_change)  
P\_1 (of the)  
P\_2 (to climate\_change)  
P\_3 (on climate\_change)  
P\_4 (climate\_change is)  
P\_5 ((to|with) the)  
P\_6 (**the (affects|effect|effects|impact|impacts) P\_o**)  
P\_7 ((while|and|induced) climate\_change)  
P\_8 (on the)  
*P\_9 (the (dangers|science|problem|risk)s) P\_o)*  
*P\_10 (to (meet|tackle|take|redirect))*  
P\_11 (to (escape|address) climate\_change)  
P\_12 (**the (psychological|inevitable|worst|visible|negative) effects P\_o**)  
P\_13 (**to (prevent|combat) climate\_change**)  
P\_14 (**((threat|threats|risk|risk)s|challenges) posed by**)  
P\_15 (in (attempts|order) to)  
P\_16 (the (details|issue) P\_o)  
P\_17 (**climate\_change (conference|summit) in**)  
*P\_18 (climate\_change will (have|reduce|increase) the)*  
P\_19 (**((scientific|academic) literature)**)

# A step towards a local grammar fragment for “sea levels”



# Closing Remarks

- Network analysis, e.g. community detection, may offer a fruitful way to add a dimension to the analysis of social media corpora; perhaps need to fuse network analysis and text analysis to detect communities
- Early results from local grammar induction show promise, i.e. as a way to elucidate interesting patterning in large sets of concordance lines

# Ongoing / future work

- Norwegian and French blog corpora
- Comparison of sub-corpora, and temporal comparisons, with regards to the framing of key terms
  - How to determine significant differences between sub-corpora? (NB. skewed distribution across blogs, variable amounts of peripheral material due to permissive crawl)
- Analyses of causality, modality, proposed climate solutions
- Development of local grammar induction
  - Optimising parameters
  - How to create local grammar fragments from P's and E's
  - Evaluation

# A few patterns and equivalence classes for “sea levels” snippets

P_0	P2339	(E_0,E_1)	360
E_0	E2340	{and,or}	
E_1	E2334	{P_1,sea_levels}	
P_1	P2316	(rising,sea_levels)	1535
P_2	P2323	(E_2,rise)	282
E_2	E2324	{will,could,would}	
P_3	P2324	(E_3,are)	253
E_3	E2325	{global_sea_levels,sea_levels,P_1}	
P_4	P2329	(sea_levels,P_2)	234
P_5	P2344	(E_4,sea_levels)	219
E_4	E2343	{that,while,the}	
P_6	P2358	(E_5,E_6)	198
E_5	E2359	{lower,raise,boost,increase}	
E_6	E2360	{global_sea_levels,sea_levels}	
.			
.			
.			