

Text Analysis and Social Simulation

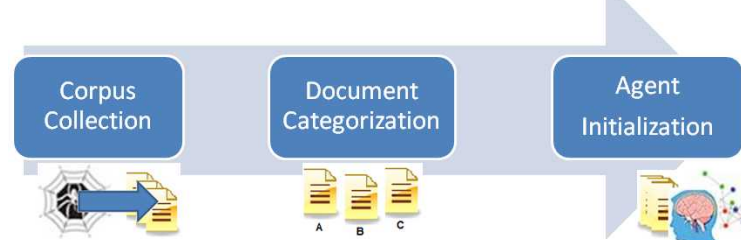
Zach Heath zheath@sandia.gov

JT McClain jtmccl@sandia.gov

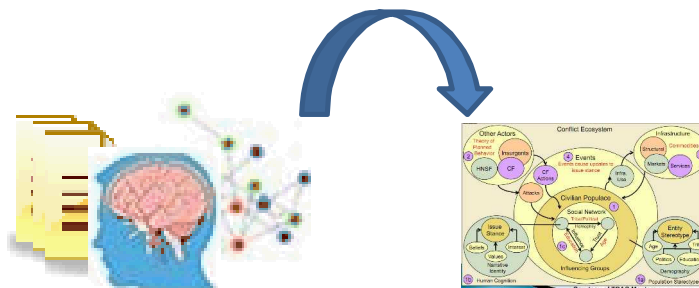
Sandia National Laboratories

TRAC-MTRY/Sandia National Laboratories Collaboration

- Two Main Tasks
 - Create a **modular pipeline** that uses **text analysis** for **social simulation model initialization**



- Augment TRAC's Cultural Geography (CG) Model with **cognitive agents** derived from **text analysis**



Motivation

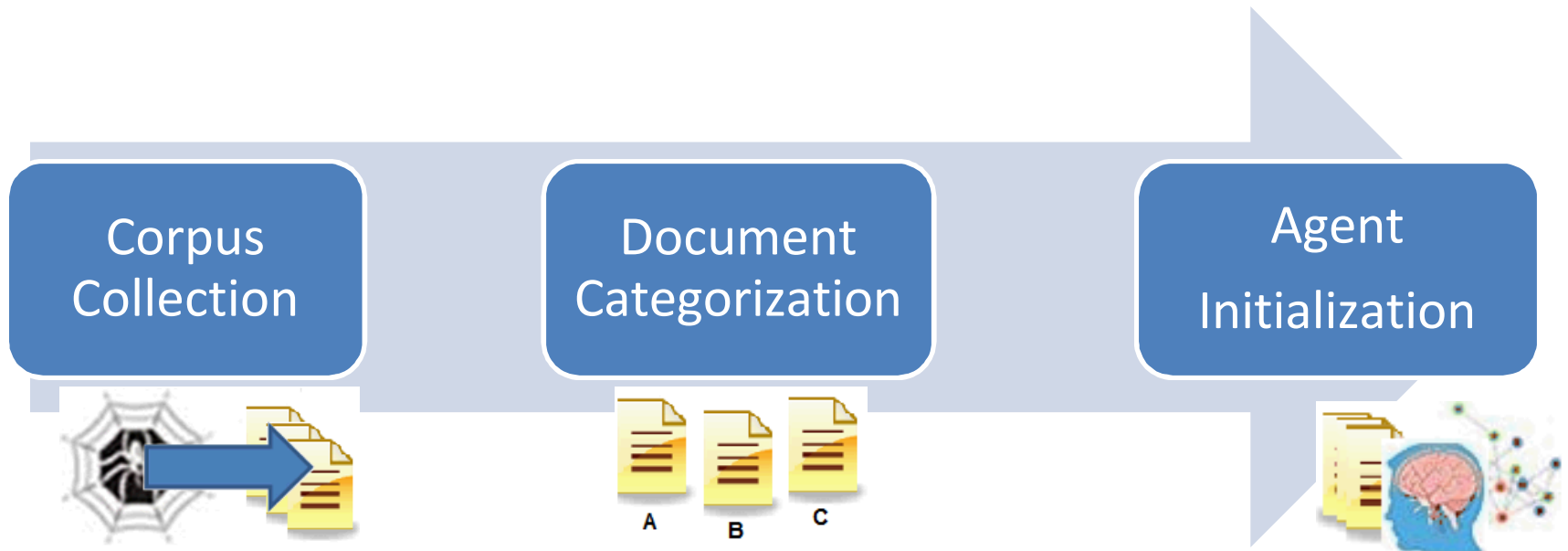
- Agent based model initialization is difficult
 - Manual source material collection and subject matter expertise is expensive
 - Subject matter expertise injects subjectivity
- Text analysis can **automate**, **augment**, and provide **objectivity** to much of this process

Caveats

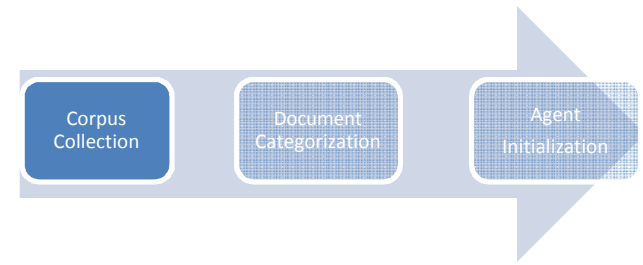
- This is in the proof of concept stage.
- This has not been validated
- Nowhere near being able to remove subject matter expertise from the process

Agent Initialization Through Text Analysis

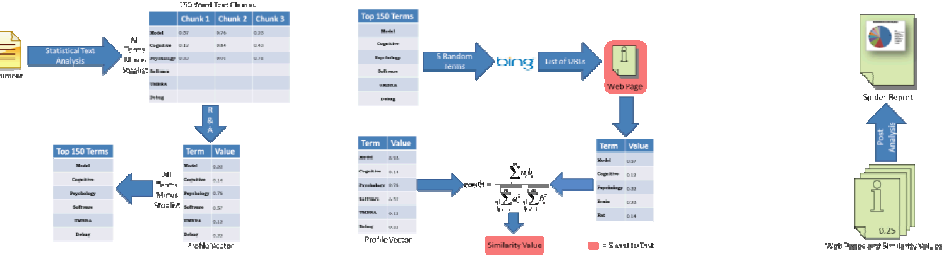
- Task 1 Goals
 - Create modular pipeline for agent initialization process
 - Provide initial test implementations for each module
 - Run and measure pipeline implementation using a test case



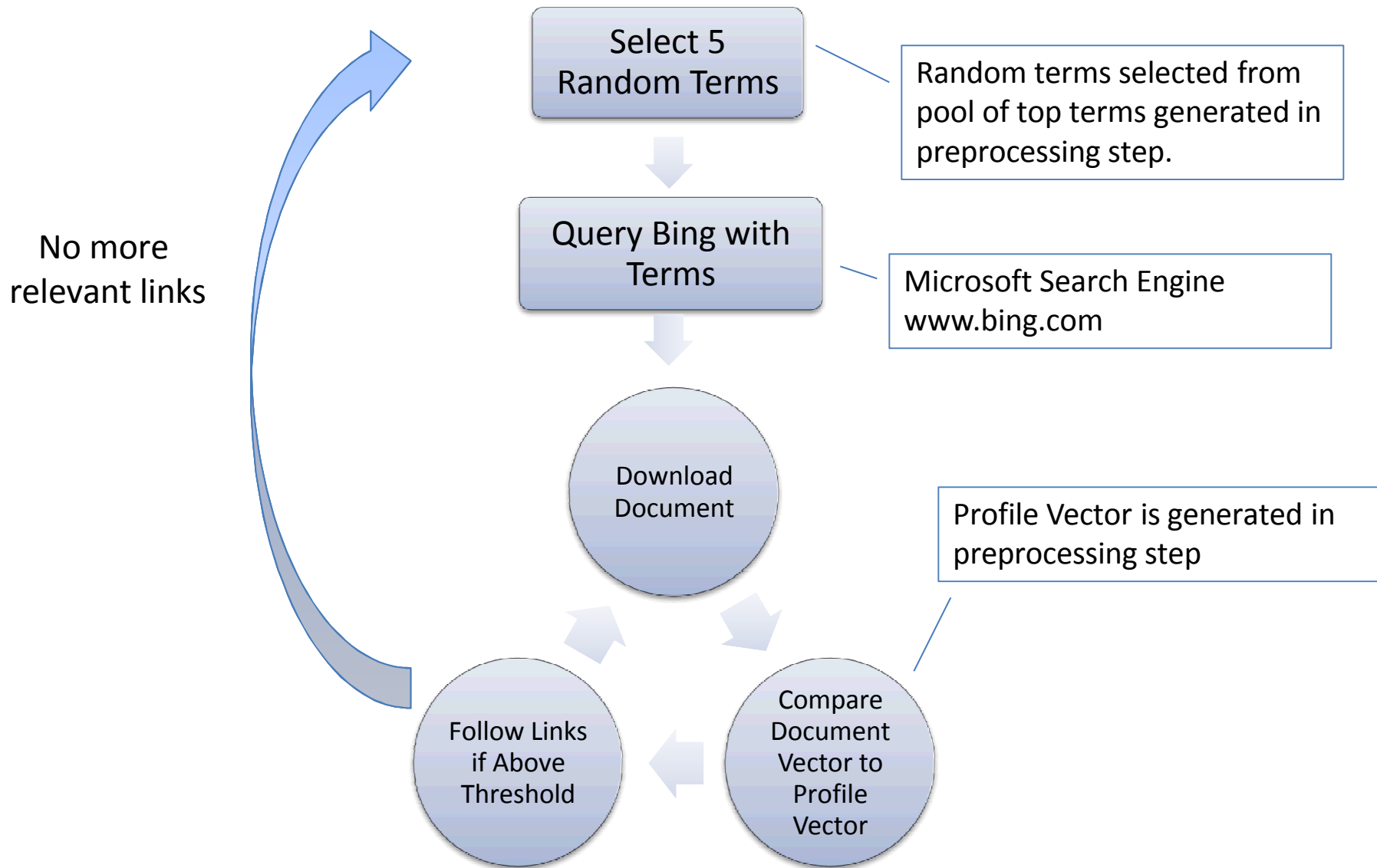
Corpus Collection



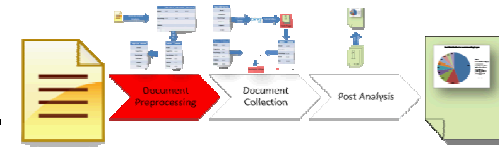
- Purpose:
 - Automatically or manually collect a set of documents that will be analyzed for agent based model initialization
- Inputs:
 - Example source document of interest and/or search keywords
 - Location to search
 - Directory or Database of Files
 - List of web domains
 - The internet as a whole
- Outputs:
 - Directory of documents
 - List of documents with model matching score



Basic Crawling Process



Document Preprocessing



150 Word Text Chunks



Statistical Text
Analysis

All
Terms
Minus
Stoplist

	Chunk 1	Chunk 2	Chunk 3
Model	0.57	0.76	0.23
Cognitive	0.12	0.84	0.43
Psychology	0.32	0.91	0.75
Software			
UMBRA			
Debug			

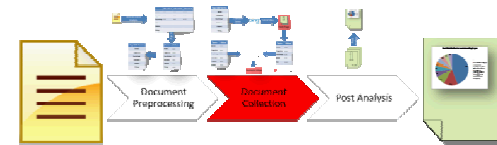
Profile Vector

Top 150 Terms
Model
Cognitive
Psychology
Software
UMBRA
Debug

All
Terms
Minus
Stoplist

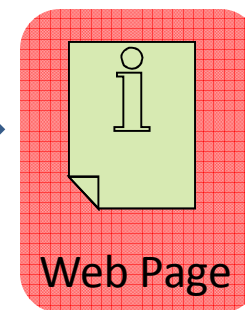
Term	Value
Model	0.23
Cognitive	0.14
Psychology	0.76
Software	0.57
UMBRA	0.12
Debug	0.32

Document Collection



Top 150 Terms

Model
Cognitive
Psychology
Software
UMBRA
Debug



Term	Value
Model	0.57
Cognitive	0.12
Psychology	0.32
Brain	0.23
Rat	0.14

Term	Value
Model	0.23
Cognitive	0.14
Psychology	0.76
Software	0.57
UMBRA	0.12
Debug	0.32

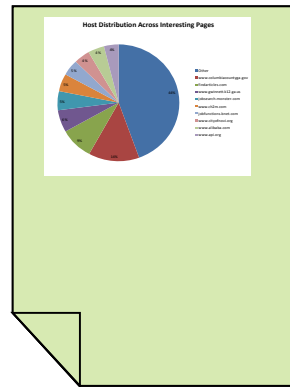
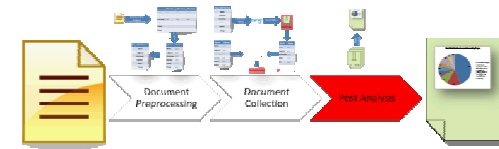
Profile Vector

$$\cos \Theta = \frac{\sum_{i=1}^m a_i b_i}{\sqrt{\sum_{i=1}^m a_i^2} \sqrt{\sum_{i=1}^m b_i^2}}$$

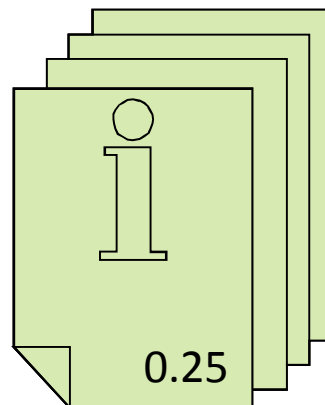
Similarity Value

 = Saved to Disk

Post Analysis



Crawler Report



Web Pages and Similarity Values

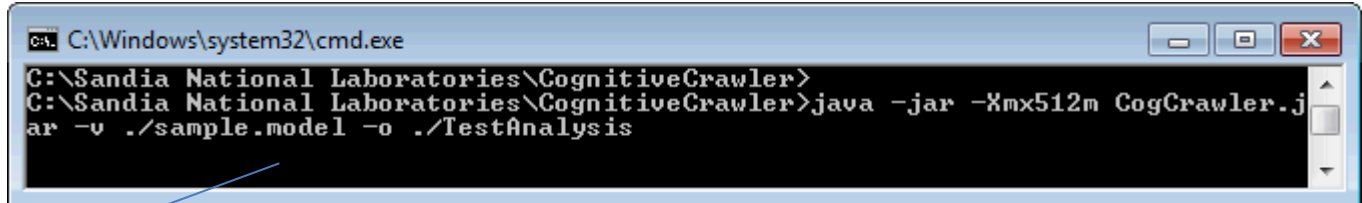
Produces aggregated list of collected documents with scores indicating level of match to source documents.

Collected documents can be used to produce a new profile for the corpus

Retrieved documents can be accessed through Google like search interface

Crawler report provides pie charts that indicate 'interesting' web domains

Cognitive Crawler



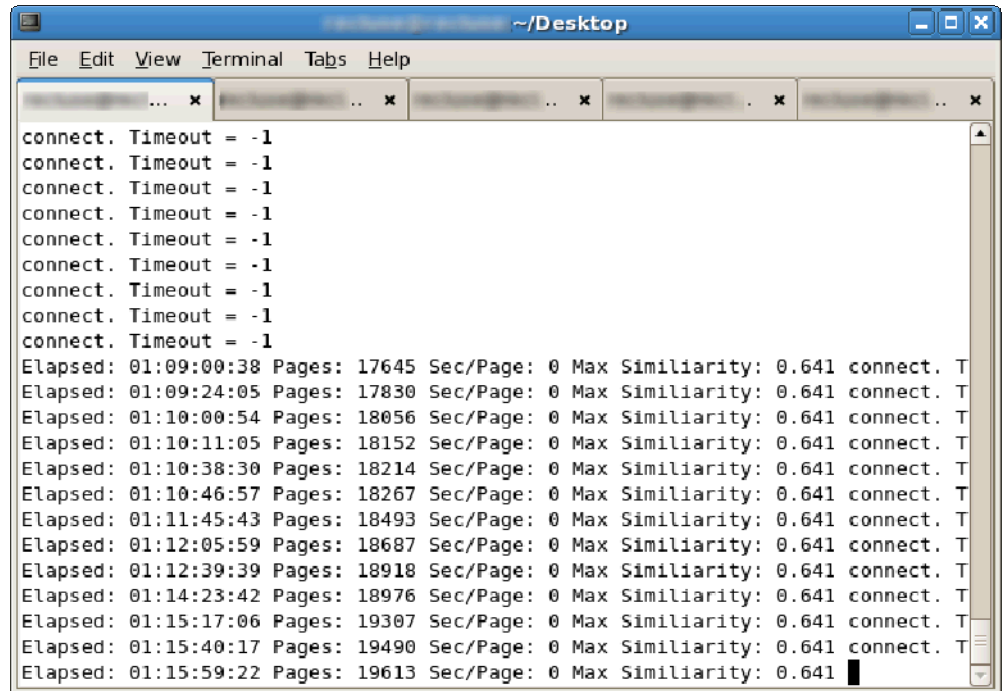
```
C:\Windows\system32\cmd.exe
C:\Sandia National Laboratories\CognitiveCrawler>
C:\Sandia National Laboratories\CognitiveCrawler>java -jar -Xmx512m CogCrawler.jar
ar -v ./sample.model -o ./TestAnalysis
```

Seed document model is built from example documents

Cognitive Crawler is launched through command line. Console output tracks crawler progress

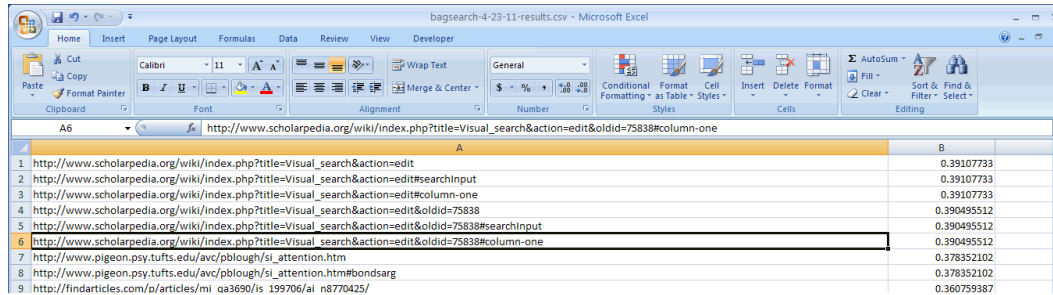
Crawler can run for weeks but generally finds most relevant material in first 1-2 days

Produces 15-20 GB of data



```
File Edit View Terminal Tabs Help
connect. Timeout = -1
connect. Timeout = -1
connect. Timeout = -1
connect. Timeout = -1
connect. Timeout = -1
connect. Timeout = -1
connect. Timeout = -1
connect. Timeout = -1
connect. Timeout = -1
Elapsed: 01:09:00:38 Pages: 17645 Sec/Page: 0 Max Similiarity: 0.641 connect. T
Elapsed: 01:09:24:05 Pages: 17830 Sec/Page: 0 Max Similiarity: 0.641 connect. T
Elapsed: 01:10:00:54 Pages: 18056 Sec/Page: 0 Max Similiarity: 0.641 connect. T
Elapsed: 01:10:11:05 Pages: 18152 Sec/Page: 0 Max Similiarity: 0.641 connect. T
Elapsed: 01:10:38:30 Pages: 18214 Sec/Page: 0 Max Similiarity: 0.641 connect. T
Elapsed: 01:10:46:57 Pages: 18267 Sec/Page: 0 Max Similiarity: 0.641 connect. T
Elapsed: 01:11:45:43 Pages: 18493 Sec/Page: 0 Max Similiarity: 0.641 connect. T
Elapsed: 01:12:05:59 Pages: 18687 Sec/Page: 0 Max Similiarity: 0.641 connect. T
Elapsed: 01:12:39:39 Pages: 18918 Sec/Page: 0 Max Similiarity: 0.641 connect. T
Elapsed: 01:14:23:42 Pages: 18976 Sec/Page: 0 Max Similiarity: 0.641 connect. T
Elapsed: 01:15:17:06 Pages: 19307 Sec/Page: 0 Max Similiarity: 0.641 connect. T
Elapsed: 01:15:40:17 Pages: 19490 Sec/Page: 0 Max Similiarity: 0.641 connect. T
Elapsed: 01:15:59:22 Pages: 19613 Sec/Page: 0 Max Similiarity: 0.641
```

Cognitive Crawler



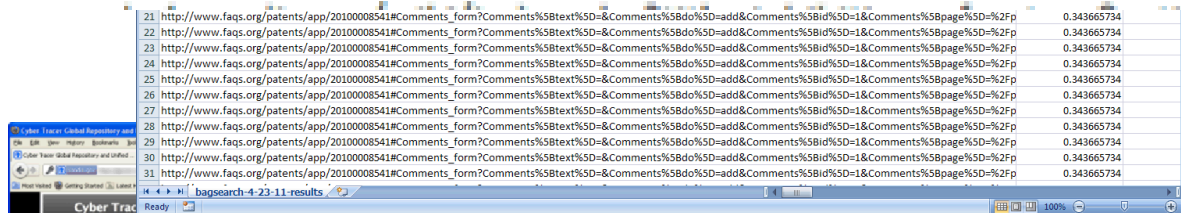
bagsearch-4-23-11-results.csv - Microsoft Excel

A	B
1 http://www.scholarpedia.org/wiki/index.php?title=Visual_search&action=edit	0.39107733
2 http://www.scholarpedia.org/wiki/index.php?title=Visual_search&action=edit#searchinput	0.39107733
3 http://www.scholarpedia.org/wiki/index.php?title=Visual_search&action=edit#column-one	0.39107733
4 http://www.scholarpedia.org/wiki/index.php?title=Visual_search&action=edit&oldid=75838	0.39045512
5 http://www.scholarpedia.org/wiki/index.php?title=Visual_search&action=edit&oldid=75838#searchinput	0.39045512
6 http://www.scholarpedia.org/wiki/index.php?title=Visual_search&action=edit&oldid=75838#column-one	0.39045512
7 http://www.pigeon.psy.tufts.edu/avc/pblough/si_attention.htm	0.378352102
8 http://www.pigeon.psy.tufts.edu/avc/pblough/si_attention.htm#bondsarg	0.378352102
9 http://findarticles.com/p/articles/mi_qa3690/is_199706/ai_n8770425/	0.360759387

Score represents how well collected documents matches text profile generated in preprocessing step

http://www.pigeon.psy.tufts.edu/avc/pblough/si_attention.htm#bondsarg
http://findarticles.com/p/articles/mi_qa3690/is_199706/ai_n8770425/
http://findarticles.com/p/articles/mi_qa3690/is_199706/ai_n8770425/#talkback

0.378352102
0.360759387
0.356080931
0.343753779
0.343753779
0.343753779
0.343753779
0.343665734

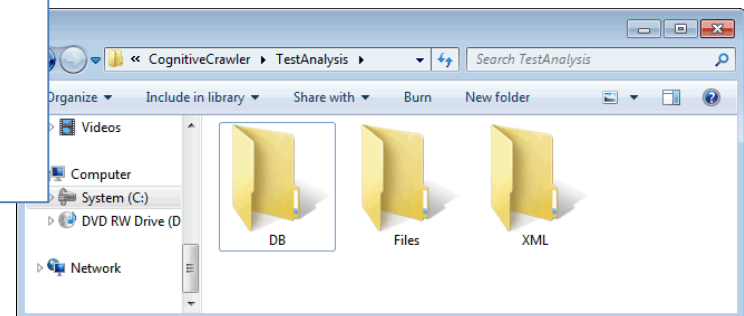


Cyber Tracer Global Repository and Unified Natural Language Search

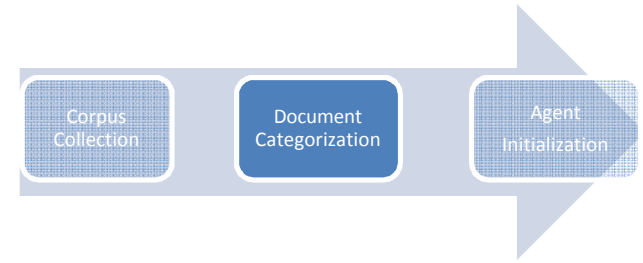
bagsearch-4-23-11-results

21 http://www.faa.gov/patents/app/20100008541#Comments_form?Comments%5Btext%5D=&Comments%5Bdo%5D=add&Comments%5Bsid%5D=1&Comments%5Bpage%5D=2Fp	0.343665734
22 http://www.faa.gov/patents/app/20100008541#Comments_form?Comments%5Btext%5D=&Comments%5Bdo%5D=add&Comments%5Bsid%5D=1&Comments%5Bpage%5D=2Fp	0.343665734
23 http://www.faa.gov/patents/app/20100008541#Comments_form?Comments%5Btext%5D=&Comments%5Bdo%5D=add&Comments%5Bsid%5D=1&Comments%5Bpage%5D=2Fp	0.343665734
24 http://www.faa.gov/patents/app/20100008541#Comments_form?Comments%5Btext%5D=&Comments%5Bdo%5D=add&Comments%5Bsid%5D=1&Comments%5Bpage%5D=2Fp	0.343665734
25 http://www.faa.gov/patents/app/20100008541#Comments_form?Comments%5Btext%5D=&Comments%5Bdo%5D=add&Comments%5Bsid%5D=1&Comments%5Bpage%5D=2Fp	0.343665734
26 http://www.faa.gov/patents/app/20100008541#Comments_form?Comments%5Btext%5D=&Comments%5Bdo%5D=add&Comments%5Bsid%5D=1&Comments%5Bpage%5D=2Fp	0.343665734
27 http://www.faa.gov/patents/app/20100008541#Comments_form?Comments%5Btext%5D=&Comments%5Bdo%5D=add&Comments%5Bsid%5D=1&Comments%5Bpage%5D=2Fp	0.343665734
28 http://www.faa.gov/patents/app/20100008541#Comments_form?Comments%5Btext%5D=&Comments%5Bdo%5D=add&Comments%5Bsid%5D=1&Comments%5Bpage%5D=2Fp	0.343665734
29 http://www.faa.gov/patents/app/20100008541#Comments_form?Comments%5Btext%5D=&Comments%5Bdo%5D=add&Comments%5Bsid%5D=1&Comments%5Bpage%5D=2Fp	0.343665734
30 http://www.faa.gov/patents/app/20100008541#Comments_form?Comments%5Btext%5D=&Comments%5Bdo%5D=add&Comments%5Bsid%5D=1&Comments%5Bpage%5D=2Fp	0.343665734
31 http://www.faa.gov/patents/app/20100008541#Comments_form?Comments%5Btext%5D=&Comments%5Bdo%5D=add&Comments%5Bsid%5D=1&Comments%5Bpage%5D=2Fp	0.343665734

Collected documents stored as flat files and/or in a database. Files can be accessed through Google like search interface



Document Categorization

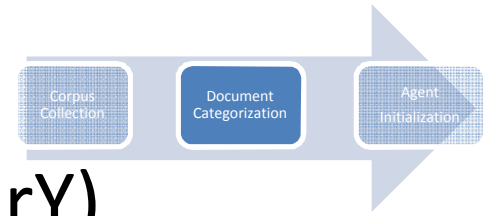


- Purpose:
 - Given a collection of documents, group or categorize those documents based on some set criteria
- Inputs:
 - Document Corpus
 - Categorization criteria and pre-categorized documents
 - Categorization thresholds
- Outputs:
 - Document Categorizations
 - Categorization Statistics

Proof of concept: STANLEY

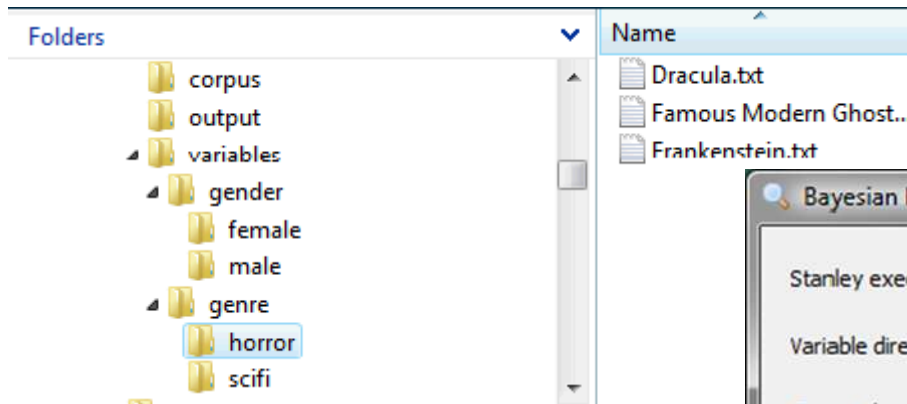
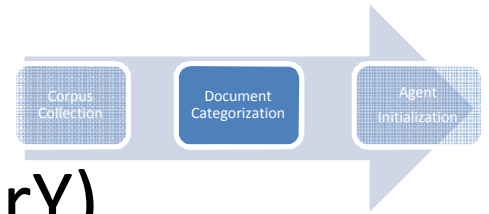
(Sandia Text ANaLysis Extensible library)

Document Categorizer

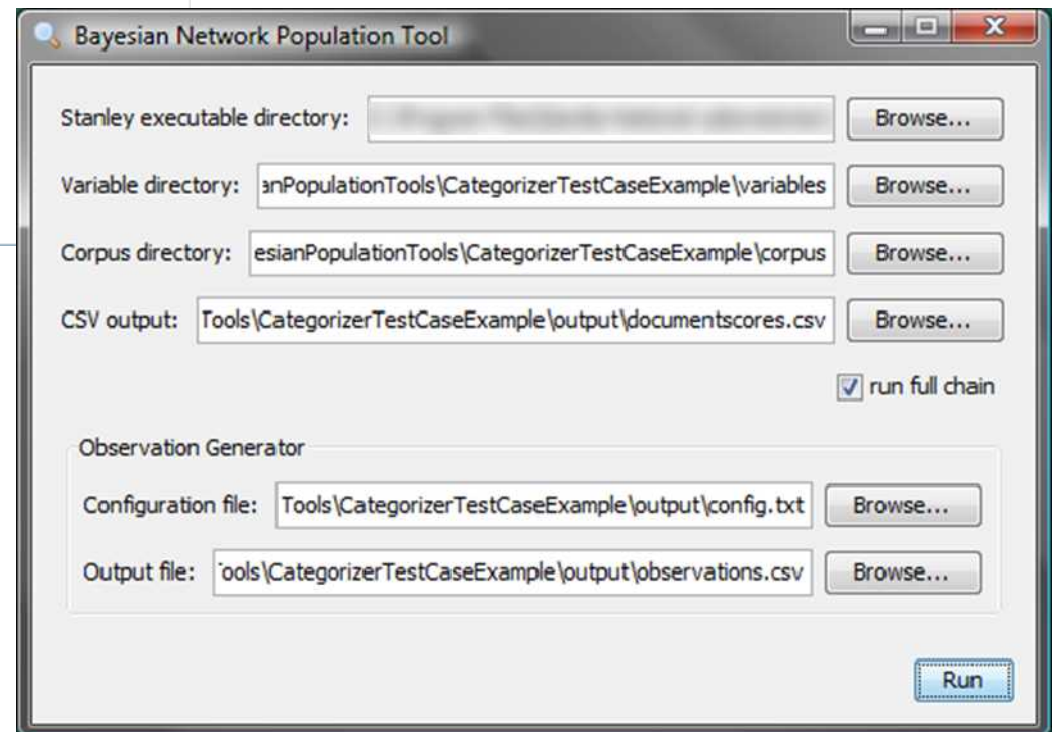


- Builds textual profile for each set of pre-categorized documents
- Matches and scores new documents from corpus against all profiles using same text analysis tools explained in Cognitive Crawler section
- Marks document categorizations based on thresholds set by user
- To be replaced with Cortext (java based text analysis library)

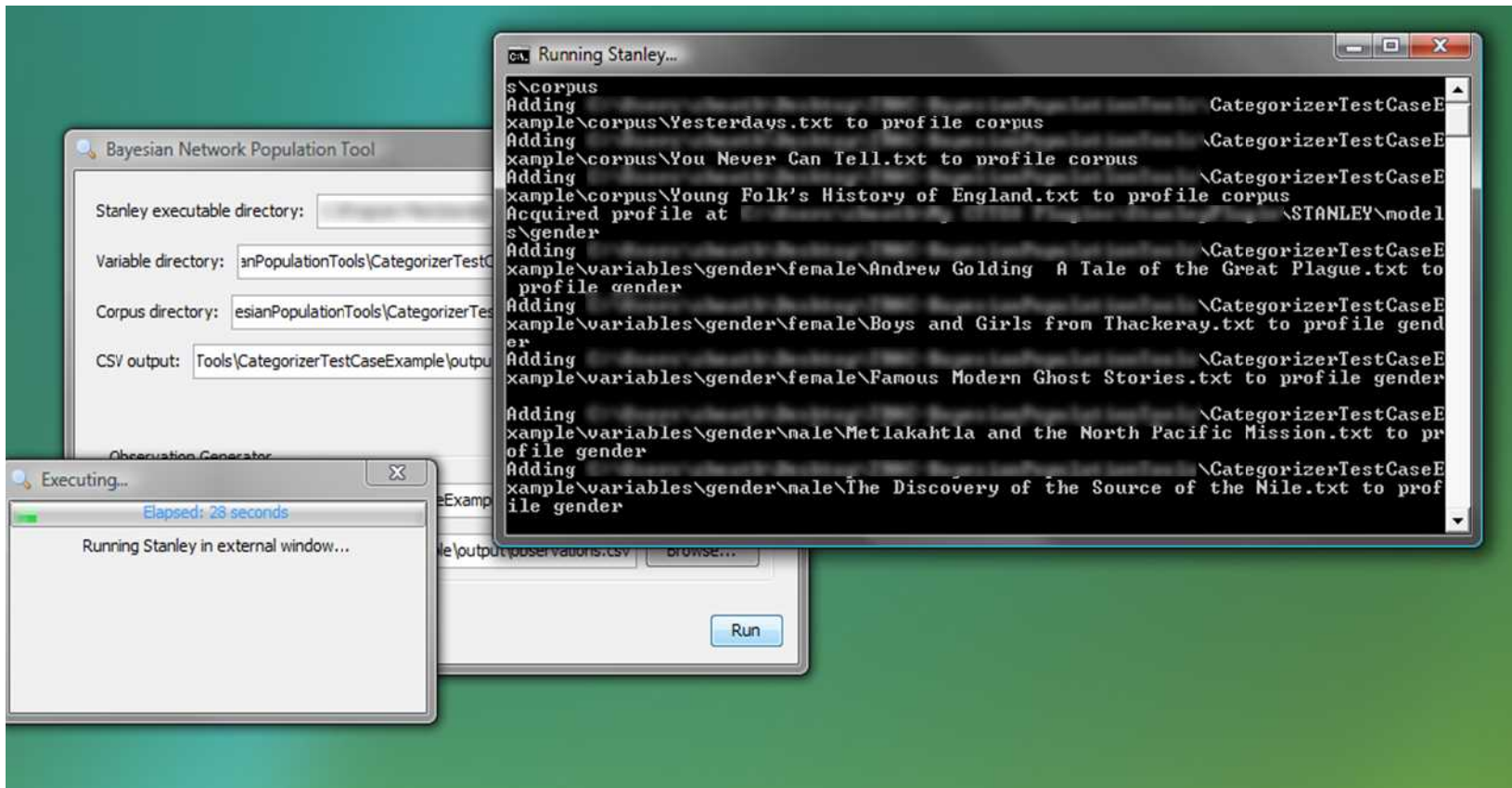
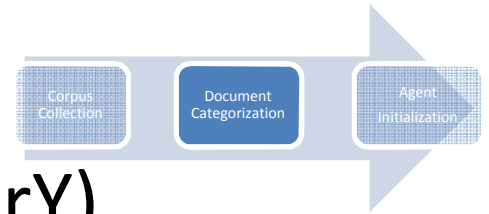
Proof of concept: STANLEY (Sandia Text ANaLysis Extensible librarY) Document Categorizer



Directory structure represents variables to categorize on. User must fill these directories with example documents



Proof of concept: STANLEY (Sandia Text ANaLysis Extensible librarY) Document Categorizer



Proof of concept: STANLEY (Sandia Text ANaLysis Extensible library) Document Categorizer



documentScores.csv - Microsoft Excel

Document	Variable	Score	Score	Score	Score
95 Theses.txt	corpus	0.51786308			
95 Theses.txt	authorGender	0.531130992	0.525042539	0.541182	
95 Theses.txt	genre	0.45101695	0.619124029	0.292756	0.525762
95 Theses.txt	timeWritten	0.583643893	0.562309539	0.548094	0.596231
95 Theses.txt	violence	0.57429249	0.571978532	0.548111	

How well document matches the example corpus as a whole

How well document matches violence variable 1 and 2 example documents

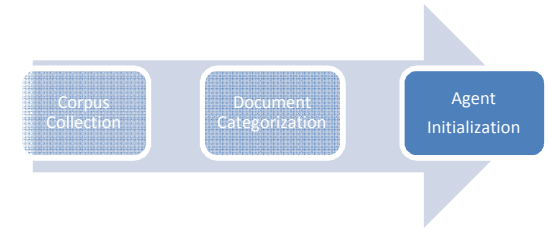
How well document matches all violence variable example documents

Agent Initialization



- Purpose:
 - Analyze a set of document categorizations or other text analysis operations to produce data for agent based model initialization
- Inputs:
 - Document Categorizations
 - Agent Initialization parameters
- Outputs:
 - Initialization files for Agent based model of interest
 - Additional Statistics

Proof of concept: Observation Generator



- Analyzes document categorization scores to produce list of document observations
- Document observations are turned into a Bayesian Network for use in the Cultural Geography Model

Proof of concept: Observation Generator



Categorization scores
from previous step

	A	B	C	D	E	F	G
1	% Document Score File						
2	% Corpus Documents:						
3	% Variable Documents:						
4	% Date: 10/29/2010 11:57 AM						
5							
6	@Variables						
7	corpus						
8	authorGender	female	male				
9	genre	biography	cooking	horror			
10	timeWritten	1700sAndEarlier	1800s	1900s			
11	violence	acceptable	unacceptable				
12							
13	@Data						
14	95 Theses.txt	corpus					
15	95 Theses.txt	authorGender					
16	95 Theses.txt	genre					
17	95 Theses.txt	timeWritten					
18	95 Theses.txt	violence					
19	A Confederate Girl's Diary.txt	corpus					
20	A Confederate Girl's Diary.txt	authorGender					
21	A Confederate Girl's Diary.txt	genre					
22	A Confederate Girl's Diary.txt	timeWritten					
23	A Confederate Girl's Diary.txt	violence					
24	A Connecticut Yankee in King Arthur's Court.txt	corpus					
25	A Connecticut Yankee in King Arthur's Court.txt	authorGender					

Observation Generator

Setup Input Config Output

Input file: tools\ClassifierTestExample\output\documentscores.csv

Configuration file: nTools\ClassifierTestExample\output\config.txt

Output file: iTools\ClassifierTestExample\output\observations.csv

```
cfBlais9.txt - Notepad
File Edit Format View Help
# Configuration threshold file for converting document scores to observations
# Format: VAR_NAME, VAR_OPTION_NAME, VAR_SCORE, VAR_OPTION_SCORE, \
# VAR_NAME: Name of the variable to be configured (ex: genre)
# VAR_OPTION_NAME: The name of the variable output value to be configured
# VAR_SCORE: The minimum score value that a document must have for the variable to be configured
# VAR_OPTION_SCORE: The minimum score value that a document must have for the variable option to be configured
# VAR_OPTION_DIFF: The minimum score difference that the document must have for the variable option to be configured
# "*" matches any variable name or any variable output value
# *,*,.6,.01 Sets the threshold for all variables and all variable options
# genre,*,.6,.01 Sets the threshold for all variable options of genre
# genre,fiction,.6,.01 Sets the threshold for the variable option of genre
# Rules are matched by the order they are specified within this file

violence, unacceptable,.3,.3,.001
violence, acceptable,.8,.85,.05

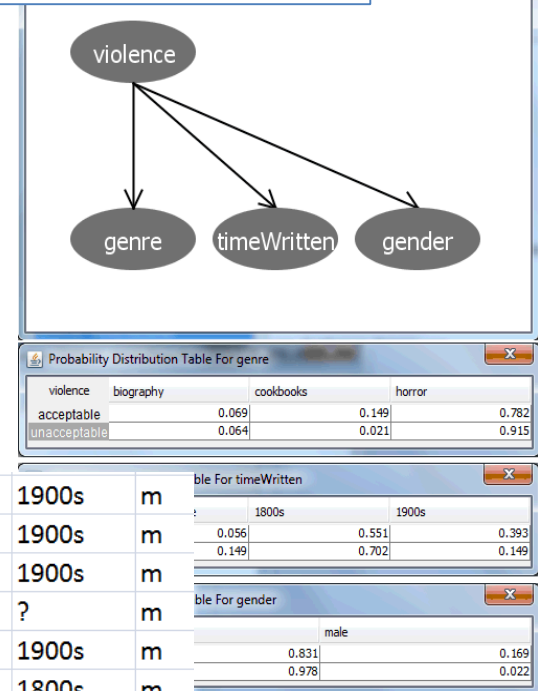
*,*,.3,.3,.005
```

Thresholds for matching to be
set by user

Proof of concept: Observation Generator



Observations can be used to initialize a Bayesian network



Generated observations using matching thresholds

obsCBPaper.xlsx - Microsoft Excel									
Home Insert Page Layout Formulas Data Review View Team									
G19 f									
			D	E	F	G			
5	adventure buster	horror	unacceptable	1900s	male	1900s	?		
6	adventure club	horror	acceptable	1900s	male	1900s	?		
7	adventures of kathlyn	horror	acceptable	1900s	male	1700s	m		
8	Adventures of Reddy Fox	horror	unacceptable	1900s	female	1800s	f		
9	airplane	?	acceptable	1900s	female	1900s	m		
10	Alec Forbes of Howglen	horror	acceptable	1900s	female	1900s	m		
11	alexander	biography	acceptable	?	male	1800s	m		
12	amatuer	horror	acceptable	?	female	1800s	m		
13	american fairy	horror	unacceptable	?	female	?	m		
14	andrew	5	adventure buster	horror	unacceptable	1900s	female	1900s	m
15	animal	6	adventure club	horror	acceptable	1900s	?	1900s	m
16	Armour	7	adventures of kathlyn	horror	acceptable	1900s	female	1900s	m
17	arnold	8	Adventures of Reddy Fox	horror	unacceptable	1900s	female	?	m
18	Aussie	9	airplane	?	acceptable	1900s	male	1900s	m
19	Autumn	10	Alec Forbes of Howglen	horror	acceptable	1900s	female	1800s	m
20	bible	11	alexander	biography	acceptable	?	male	1900s	m
21	black a	12	amatuer	horror	acceptable	?	female	1800s	m
22	blind n								
23	books								
24	bracele								
25	bunny!								
26	captive								
27	castle								
28	cattle								
29	chocolate								

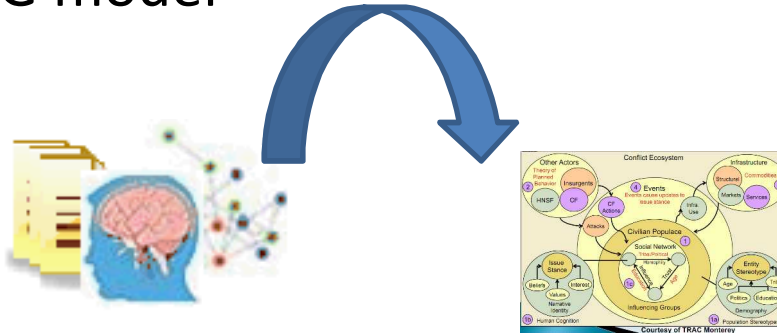
Agent Initialization Through Text Analysis

- Proof of concept implementations requiring subject matter expert guidance/analysis
- Pipeline will be modular to allow for testing of alternative/improved implementations
- Continual improvement of components will relieve burden on subject matter expert



Integration of Cognitive Model based Agents in CG

- Task 2 Goals
 - Integrate concept map based cognitive model agents into the CG model



- Run and measure implementation using a test case comparing with other agent implementations

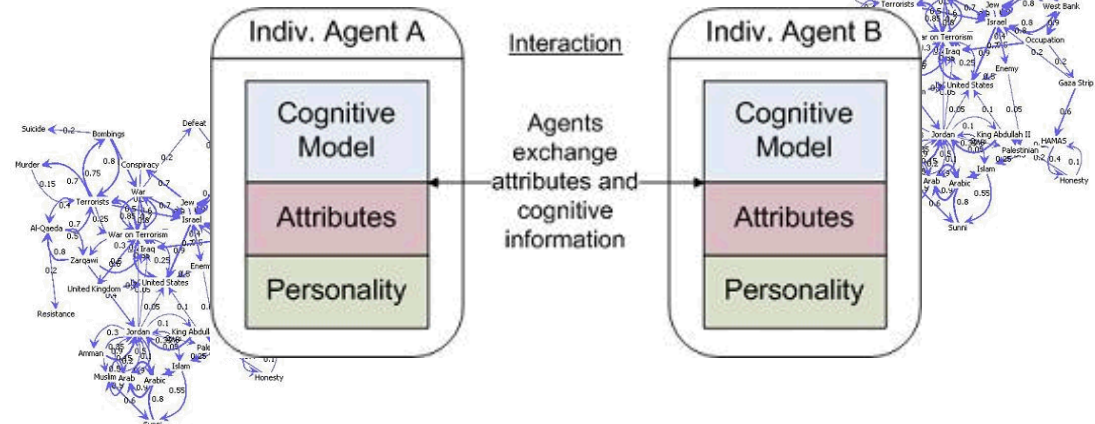
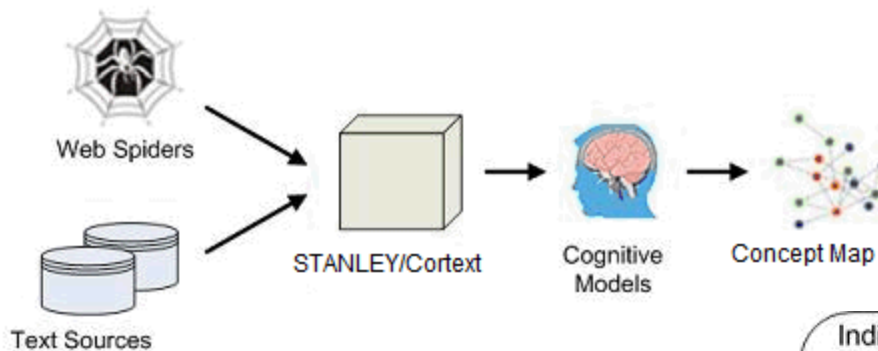
Concept map agents

Individuals have cognitive models created from text

Cognitive information is exchanged during interactions



Media Agents also inject cognitive information to their subscribers



Concept map agents

