# Developing an improved focused crawler for the IDEAL project

Ward Bonnefond, Chris Menzel, Zack Morris, Suhas Patel, Tyler Ritchie, Mark Tedesco, Franklin Zheng

# Refresher - IDEAL project

▶ Integrating Digital Event Archiving and Library

▶ Finding webpages related to an event (i.e. natural disaster)

▶ Store found webpages locally for parsing and analysis

# Refresher - Enhanced focus crawler

▶ Extract key words and key concepts (i.e. date, location, type of disaster)

▶ Construct trees based on these words and concepts

▶ Develop algorithm to compare different trees and their relationships

▶ Make this process accessible via a web application

# Refresher - Project components

1. Tree construction and visual representation

2. Event representation (i.e. key words and key concepts) versus actual event (i.e. webpage)

3. Integrating updated modules into the existing focused crawler

# Refresher - Original Implementation

| Start with a list of seed URLs | → | Web-crawler crawls through list of URLs | → | Outputs a score for each URL based on keyword matchings | → | Searches the webpage for other URLs | → | Adds any good URLs found to the list |

# Final Results

- Front-End
  - User can enter multiple seed URLS and other important attributes (date, location, type of disaster, etc.)
  - Constructs a visual tree representation of the articles found via the web application
  - Loads back-end results using JavaScript (previously PHP)
- Back-end
  - Constructs a full event tree from query created by the user
  - Intelligently extracts location, date, type of disaster, etc. using natural language processing as opposed to keyword analysis

# Scoring Algorithm

▶ Date: 20%

    ▶ Year (365) + Month (30) + Day (1)

    ▶ i.e. Year has 365 times the weight of a Day

▶ Location, Type of Event, Event Name, etc.: 80%

    ▶ Replaces the previously used tree-edit distance

▶ New scoring algorithm provides more accurate articles that using only tree-edit distance did not find

# Final Implementation

Start with a list of seed URLs → Web-crawler crawls through list of URLs → Outputs a score for each URL based on custom scoring algorithm → Searches the webpage for other URLs → Adds any good URLs found to the list

# Current Back End Example Data

- **Base Focused Crawler**

- 0.479449956362|0.555185525105|http://www.tmz.com/2014/04/27/donald-sterling-photos-magic-johnson-matt-kemp-v-stiviano-racist-audio/2/0.479449956362|0.555185525105|http://www.tmz.com/2014/04/27/donald-sterling-photos-magic-johnson-matt-kemp-v-stiviano-racist-audio/1/0.479449956362|0.555185525105|http://www.tmz.com/2014/04/27/donald-sterling-photos-magic-johnson-matt-kemp-v-stiviano-racist-audio/20/0.479449956362|0.596680123725|http://www.tmz.com/2014/04/27/donald-sterling-photos-magic-johnson-matt-kemp-v-stiviano-racist-audio/14/0.479449956362|0.614986314544|http://www.tmz.com/2014/04/27/donald-sterling-photos-magic-johnson-matt-kemp-v-stiviano-racist-audio/11/0.479449956362|0.603562290076|http://www.tmz.com/2014/04/27/donald-sterling-photos-magic-johnson-matt-kemp-v-stiviano-racist-audio/10/0.479449956362|0.555185525105|http://www.tmz.com/2014/04/27/donald-sterling-photos-magic-johnson-matt-kemp-v-stiviano-racist-audio/12/0.479449956362|0.682182274472|http://www.tmz.com/2014/04/27/donald-sterling-photos-magic-johnson-matt-kemp-v-stiviano-racist-audio/13/0.479449956362|0.621091837905|http://www.tmz.com/2014/04/27/donald-sterling-photos-magic-johnson-matt-kemp-v-stiviano-racist-audio/15/0.479449956362|0.588410106633|http://www.tmz.com/2014/04/27/donald-sterling-photos-magic-johnson-matt-kemp-v-stiviano-racist-audio/16/

Visited: 100Accepted: 81

**Extended Focused Crawler**

0.469920800878|0.52|http://www.dailymail.co.uk/tvshowbiz/article-2511646/Kim-Kardashian-eBays-old-clothes-help-Philippines-typhoon-victims--donates-just-10-proceeds-charity.html0.469920800878|1.0|http://www.fao.org/emergencies/crisis/philippines-typhoon-haiyan/en/0.469920800878|0.901299136822|http://www.fao.org/emergencies/crisis/philippines-typhoon-haiyan/crop-damages-map/en/0.469920800878|0.919740571108|http://www.fao.org/emergencies/crisis/philippines-typhoon-haiyan/impact-assessment-map/en/0.469920800878|0.855014256832|http://www.fao.org/emergencies/crisis/philippines-typhoon-haiyan/input-distribution-map/en/0.469920800878|1.0|http://www.fao.org/emergencies/crisis/philippines-typhoon-haiyan/seed-distribution-map/en/

Visited: 100Accepted: 73

# Current Front-End Example

# Current Front End – Tree View
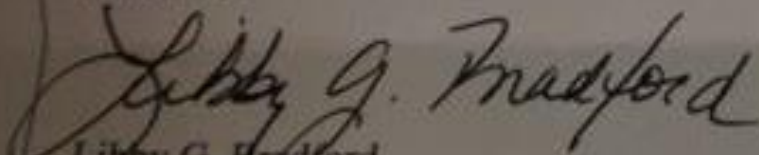
# Extras – Poster

# Extras – Second Place Letter

Dear Franklin,

The Department of Computer Science at Virginia Tech is pleased to award your group Second Place in the Capstone Category. For this award, you will share equally a $300 award with your teammates, which will be credited to your university account.

Since these funds are routed through both the Financial Aid and Bursar's Offices, they will review your accounts before crediting your award. If you have an outstanding balance, they will apply the award towards it before disbursing the remaining amount. If you have direct deposit, you will receive the funds in your account as directed. If you do not have direct deposit, you should receive a check in the mail.

You should be receiving the financial award soon. If you have not received it within 30 days, please let me know.

Sincerely,

Libby G. Bradford

Libby G. Bradford
Director, External Relations and Undergraduate Studies

# Questions?