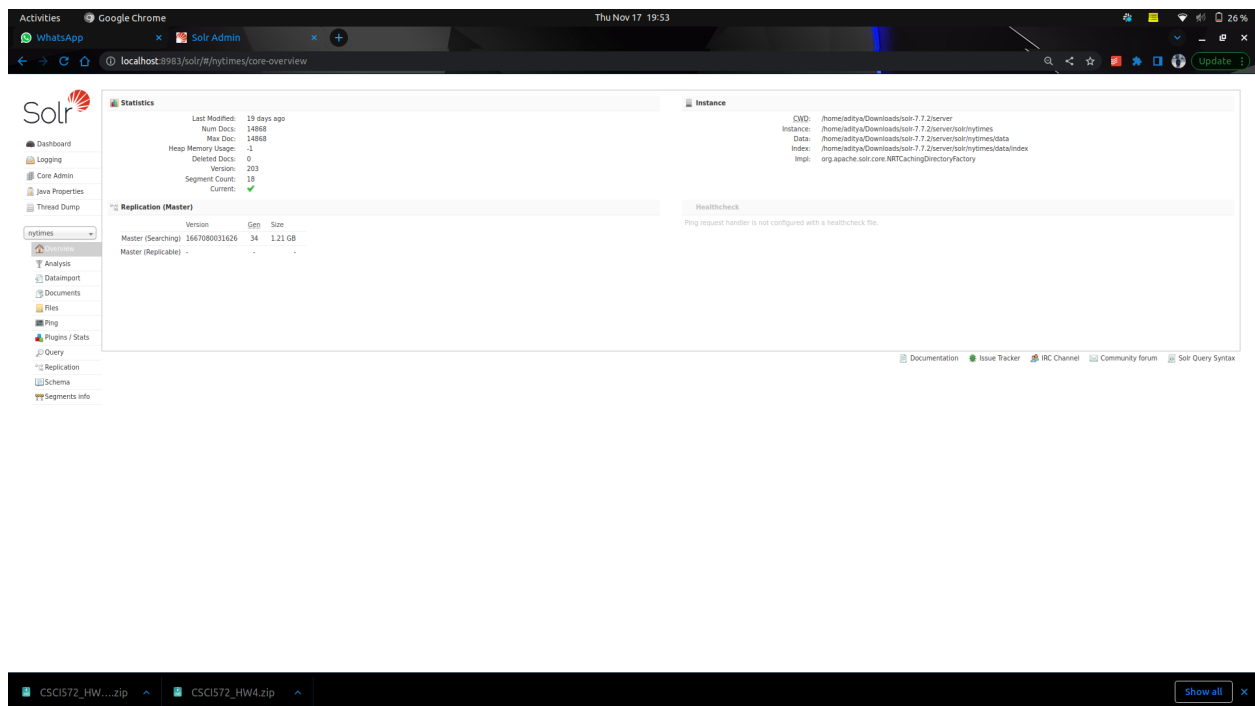# CSCI 572 - HW4 - Report

## Aditya Jain

### Steps to complete the assignment

1. I have used Ubuntu 22.04 to complete this assignment. I have used Solr version 7.7.2 for building the search engine. Started solr using command **bin/solr start.** By default Solr uses the Lucene algorithm to create a search engine. The first thing to do is to create a core in Solr (**bin/solr -c nytimes**) and then create an index using all crawled HTML files for NYtimes in that core (**bin/post -c nytimes -filetypes html /home/aditya//Downloads/solr_7.7.2/crawl_data**). Solr inherently uses Apache Tika to extract content from documents to be indexed. Opened browser and typed http://localhost:8993/solr/#/ to test queries.
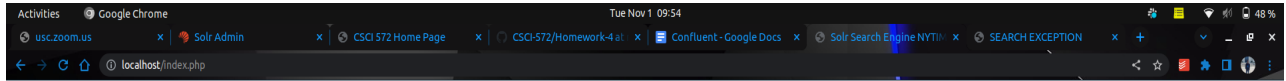


2. Once the index was created by Solr, I developed a web application to perform a search and display the top 10 results for the given query. I have used solr-php-client to interact with Solr and fetch results. To develop a web application I used HTML and **PHP 8.1.2** and made it run on an Apache2 Web server. The web application shows the query box and search button. After clicking the search button, the application sends the HTTP GET request to PHP which uses solr-php-client to fetch ranked results for the given query. In this application, I am fetching all results but displaying only the top 10 results.

3. Solr also supports ranking based on external files. Here I have calculated the PageRank scores of all web pages. For calculating the PageRank score, we need a count of all incoming and outgoing links of all web pages in the NYTimes HTML pages corpus. To extract all outgoing links of all HTML pages, I have used the external **JSoup 1.13.1** library in JAVA. After compiling and running this Java program, I created an edge list file, (**EdgeList.txt**) which contains a directed graph that is actually a mapping file. Each row in this file contains two HTML pages where there is a link from the first page to the second.

4. Once we have links between all edges, we can create a graph. For this, I have used the **Python NetworkX library**. NetworkX helps in creating, modifying, and study of the complicated graph. It also supports functionality to compute PageRank. Using NetworkX we will read the edge list created in the above step and initialize a graph. We will use the PageRank functionality of the graph to compute the PageRank score of all HTML pages (edges) in the graph. I have used the default configuration to compute PageRank i.e. **100 max iterations**, **damping parameter (alpha) of 0.85**, **tolerance (tol) of 1e-6, weight='weight'** etc. After computing the PageRank score, I have written the PageRank score of all HTML pages in a file (**external_PageRankFile**). This file is used by Solr to support PageRank as one of the ranking methods. Also made few changes in **managed-schema** and **solrconfig.xml** in core's **conf** folder.

5. I have modified the web application to support both Lucene as well as Pagerank to rank HTML pages. Later I compared the results of various queries using both Lucene and Pagerank scores.

**Why some pages have higher PageRank values than others.** Pages that are being pointed by many pages or important pages with good PageRank scores will have higher page ranks than the pages which are not pointed by many pages or are being pointed by pages with lower PageRank scores. Also if page A is pointed by page B of a higher PageRank but that page B has too many out-links, then page B will contribute less to the PageRank of page A. So, the Pagerank of a page depends on various factors such as the number of pages that point to that page, PageRank of that page, and out-links of that page. Based on these factors some pages get higher PageRank than others.
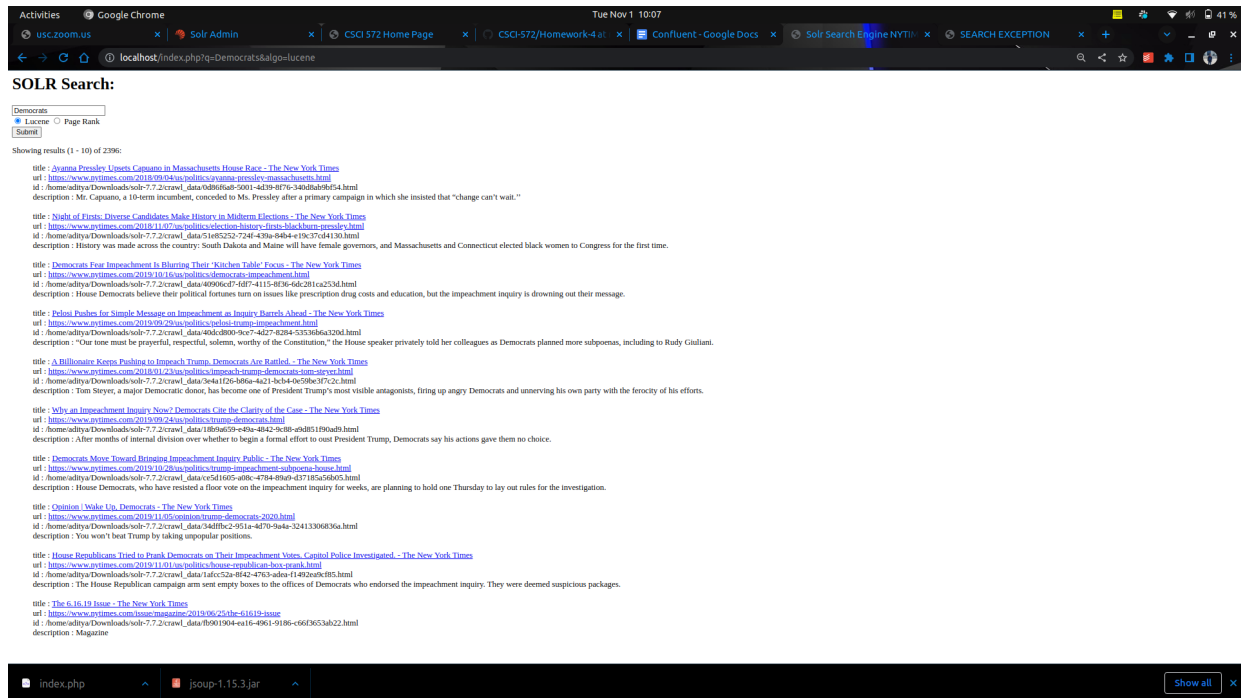
For query: **Democrats**

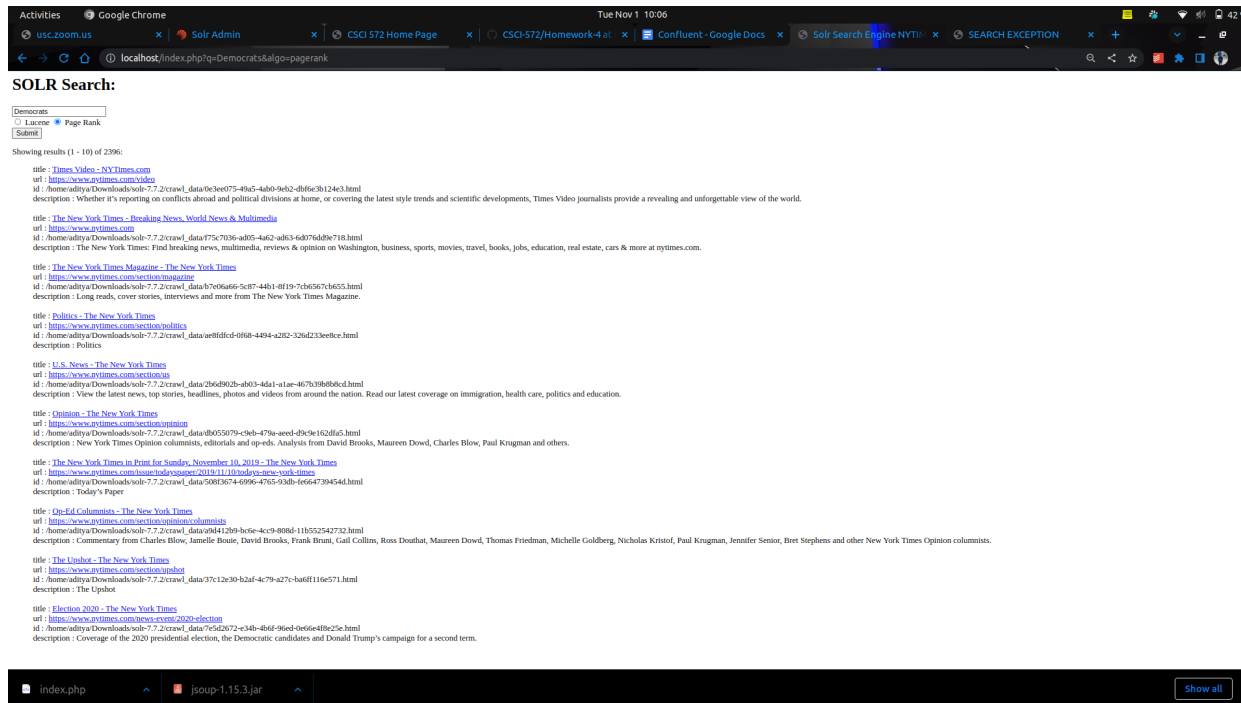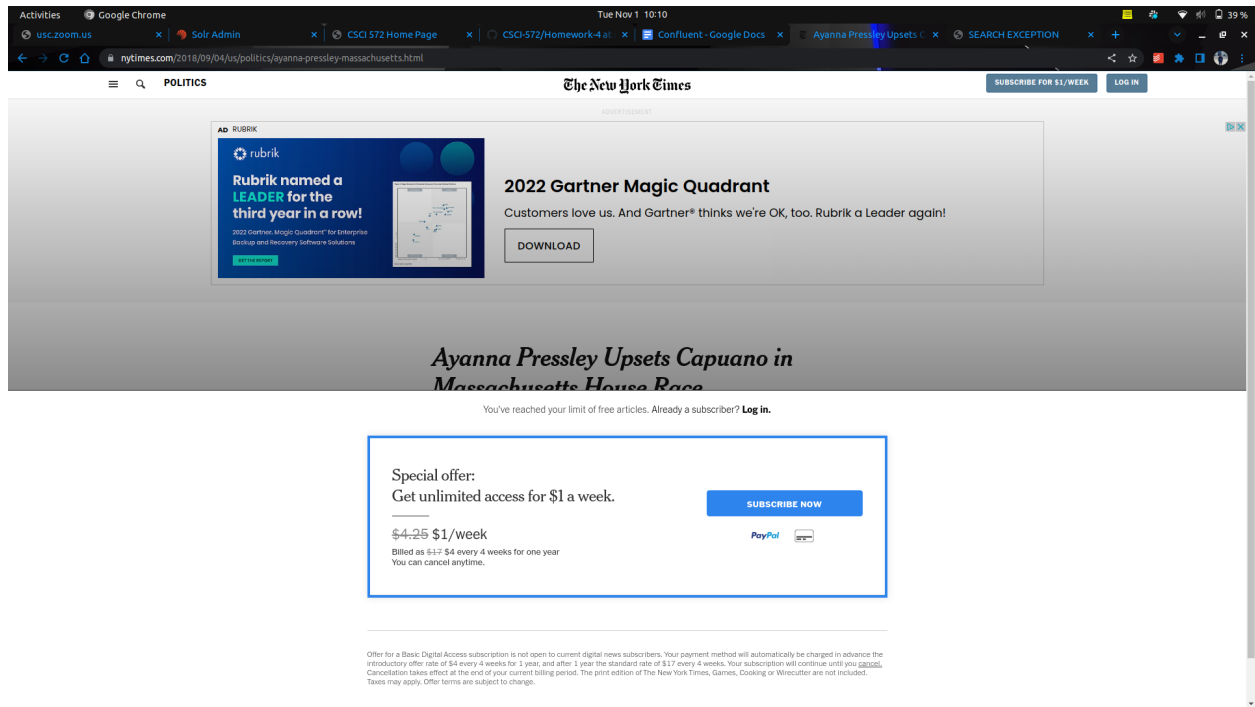## Screenshot of Initial Page where we enter the query

# Screenshot of page which shows the results for Lucene(Default).

localhost/index.php?q=Democrats&algo=lucene

**SOLR Search:**

Democrats
○ Lucene  ○ Page Rank
Submit

Showing results (1 - 10) of 2396:

title : Ayanna Pressley Upsets Capuano in Massachusetts House Race - The New York Times
url : https://www.nytimes.com/2018/09/04/us/politics/ayanna-pressley-massachusetts.html
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/0d86f6a8-5001-4d39-8f76-340d8ab9bf54.html
description : Mr. Capuano, a 10-term incumbent, conceded to Ms. Pressley after a primary campaign in which she insisted that "change can't wait."

title : Night of Firsts: Diverse Candidates Make History in Midterm Elections - The New York Times
url : https://www.nytimes.com/2018/11/07/us/politics/election-history-firsts-blackburn-pressley.html
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/51e85252-724f-439a-84b4-e19c37cd4130.html
description : History was made across the country: South Dakota and Maine will have female governors, and Massachusetts and Connecticut elected black women to Congress for the first time.

title : Democrats Fear Impeachment Is Blurring Their 'Kitchen Table' Focus - The New York Times
url : https://www.nytimes.com/2019/10/16/us/politics/democrats-impeachment.html
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/40906cd7-fdf7-4115-8f36-6dc281ca253d.html
description : House Democrats believe their political fortunes turn on issues like prescription drug costs and education, but the impeachment inquiry is drowning out their message.

title : Pelosi Pushes for Simple Message on Impeachment as Inquiry Barrels Ahead - The New York Times
url : https://www.nytimes.com/2019/09/29/us/politics/pelosi-trump-impeachment.html
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/40dcd800-9ce7-4d27-8284-53536b6a320d.html
description : "Our tone must be prayerful, respectful, solemn, worthy of the Constitution," the House speaker privately told her colleagues as Democrats planned more subpoenas, including to Rudy Giuliani.

title : A Billionaire Keeps Pushing to Impeach Trump. Democrats Are Rattled. - The New York Times
url : https://www.nytimes.com/2018/01/23/us/politics/impeach-trump-democrats-tom-steyer.html
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/3e4a1f26-b86a-4a21-bcb4-0e59be3f7c2c.html
description : Tom Steyer, a major Democratic donor, has become one of President Trump's most visible antagonists, firing up angry Democrats and unnerving his own party with the ferocity of his efforts.

title : Why an Impeachment Inquiry Now? Democrats Cite the Clarity of the Case - The New York Times
url : https://www.nytimes.com/2019/09/24/us/politics/trump-democrats.html
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/18b9a659-e49a-4842-9c88-a9d851f90ad9.html
description : After months of internal division over whether to begin a formal effort to oust President Trump, Democrats say his actions gave them no choice.

title : Democrats Move Toward Bringing Impeachment Inquiry Public - The New York Times
url : https://www.nytimes.com/2019/10/28/us/politics/trump-impeachment-subpoena-house.html
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/ce5d1605-a08c-4784-89a9-d37185a56b05.html
description : House Democrats, who have resisted a floor vote on the impeachment inquiry for weeks, are planning to hold one Thursday to lay out rules for the investigation.

title : Opinion | Wake Up, Democrats - The New York Times
url : https://www.nytimes.com/2019/11/05/opinion/trump-democrats-2020.html
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/34dffbc2-951a-4d70-9a4a-32413306836a.html
description : You won't beat Trump by taking unpopular positions.

title : House Republicans Tried to Prank Democrats on Their Impeachment Votes. Capitol Police Investigated. - The New York Times
url : https://www.nytimes.com/2019/11/01/us/politics/house-republican-box-prank.html
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/1afcc52a-8f42-4763-adea-f14492ea9cf85.html
description : The House Republican campaign arm sent empty boxes to the offices of Democrats who endorsed the impeachment inquiry. They were deemed suspicious packages.

title : The 6.16.19 Issue - The New York Times
url : https://www.nytimes.com/issue/magazine/2019/06/25/the-61619-issue
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/fb901904-ea16-4961-9186-c66f3653ab22.html
description : Magazine

# Screenshot of page which shows the results for PageRank.

localhost/index.php?q=Democrats&algo=pagerank

**SOLR Search:**

Democrats
○ Lucene  ● Page Rank
Submit

Showing results (1 - 10) of 2396:

title : Times Video - NYTimes.com
url : https://www.nytimes.com/video
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/0e3ee075-49a5-4ab0-9eb2-dbf6e3b124e3.html
description : Whether it's reporting on conflicts abroad and political divisions at home, or covering the latest style trends and scientific developments, Times Video journalists provide a revealing and unforgettable view of the world.

title : The New York Times - Breaking News, World News & Multimedia
url : https://www.nytimes.com
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/f75c7036-ad05-4a62-ad63-6d076dd9e718.html
description : The New York Times: Find breaking news, multimedia, reviews & opinion on Washington, business, sports, movies, travel, books, jobs, education, real estate, cars & more at nytimes.com.

title : The New York Times Magazine - The New York Times
url : https://www.nytimes.com/section/magazine
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/b7e06a66-5c87-44b1-8f19-7cb6567cb655.html
description : Long reads, cover stories, interviews and more from The New York Times Magazine.

title : Politics - The New York Times
url : https://www.nytimes.com/section/politics
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/ae8fdfcd-0f68-4494-a282-326d233ee8ce.html
description : Politics

title : U.S. News - The New York Times
url : https://www.nytimes.com/section/us
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/2b6d902b-ab03-4da1-a1ae-467b39b8b8cd.html
description : View the latest news, top stories, headlines, photos and videos from around the nation. Read our latest coverage on immigration, health care, politics and education.

title : Opinion - The New York Times
url : https://www.nytimes.com/section/opinion
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/db055079-c9eb-479a-aeed-d9c9e162dfa5.html
description : New York Times Opinion columnists, editorials and op-eds. Analysis from David Brooks, Maureen Dowd, Charles Blow, Paul Krugman and others.

title : The New York Times in Print for Sunday, November 10, 2019 - The New York Times
url : https://www.nytimes.com/issue/todayspaper/2019/11/10/todays-new-york-times
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/508f3674-6996-4765-93db-fe664739454d.html
description : Today's Paper

title : Op-Ed Columnists - The New York Times
url : https://www.nytimes.com/section/opinion/columnists
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/a9d412b9-bc6e-4cc9-808d-11b552542732.html
description : Commentary from Charles Blow, Jamelle Bouie, David Brooks, Frank Bruni, Gail Collins, Ross Douthat, Maureen Dowd, Thomas Friedman, Michelle Goldberg, Nicholas Kristof, Paul Krugman, Jennifer Senior, Bret Stephens and other New York Times Opinion columnists.

title : The Upshot - The New York Times
url : https://www.nytimes.com/section/upshot
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/37c12e30-b2af-4c79-a27c-ba6ff116e571.html
description : The Upshot

title : Election 2020 - The New York Times
url : https://www.nytimes.com/news-event/2020-election
id : /home/aditya/Downloads/solr-7.7.2/crawl_data/7e5d2672-e34b-4b6f-96ed-0e66e4f8e25e.html
description : Coverage of the 2020 presidential election, the Democratic candidates and Donald Trump's campaign for a second term.

**Screenshot of the actual webpage that opens when clicking on one of the result's link.**



**A table containing list of URLs of the top ten results produced by both ranking methods for each query**

| index | query | lucene | pagerank |
|---|---|---|---|
| 0 | | https://www.nytimes.com/2019/05/25/movies/cannes-film-festival-winners-parasite.html | https://www.nytimes.com/interactive/2019/10/15/t-magazine/nick-cave-artist.html |
| 1 | | https://www.nytimes.com/2019/05/25/movies/cannes-film-festival-winners-parasite.html | https://www.nytimes.com/interactive/2019/10/15/t-magazine/rachel-weisz-acting-movies.html |
| 2 | Cannes | https://www.nytimes.com/2018/04/12/movies/cannes-lineup-netflix-boycott.html | https://www.nytimes.com/interactive/2019/10/15/t-magazine/nicolas-ghesquiere-louis-vuitton.html |
| 3 | | https://www.nytimes.com/2019/05/24/movies/cannes-almodovar-kechiche.html | https://www.nytimes.com/interactive/2019/10/15/t-magazine/the-greats.html |
| 4 | | https://www.nytimes.com/2019/05/16/movies/cannes-film-festival-jim-jarmusch.html | https://www.nytimes.com/interactive/2019/10/15/t-magazine/shigeru-ban.html |
| 5 | | https://www.nytimes.com/2019/05/21/movies/cannes-film-festival-antonio-banderas-pedro-almodovar.html | https://www.nytimes.com/2017/10/05/us/harvey-weinstein-harassment-allegations.html |
| 6 | | https://www.nytimes.com/2019/05/27/movies/ | https://www.nytimes.com/2019/09/06/movies/ |

| | | | |
|---|---|---|---|
| | | cannes-oscars.html | fall-movies.html |
| 7 | | https://www.nytimes.com/2019/05/20/arts/cannes-film-festival-black-director.html | https://www.nytimes.com/2019/09/06/movies/fall-movies.html |
| 8 | | https://www.nytimes.com/2019/05/22/movies/quentin-tarantino-margot-robbie.html | https://www.nytimes.com/news-event/awards-season |
| 9 | | https://www.nytimes.com/2019/05/22/movies/dicaprio-pitt.html | https://www.nytimes.com/2017/10/05/us/harvey-weinstein-harassment-allegations.html |
| 10 | Congress | https://www.nytimes.com/interactive/2019/us/politics/trump-impeachment-congress-list.html | https://www.nytimes.com |
| 11 | | https://www.nytimes.com/2019/11/05/magazine/congress-president-impeachment.html | https://www.nytimes.com/section/magazine |
| 12 | | https://www.nytimes.com/2019/11/05/magazine/congress-president-impeachment.html | https://www.nytimes.com/section/t-magazine |
| 13 | | https://www.nytimes.com/2019/07/25/us/politics/mueller-impeachment.html | https://www.nytimes.com/issue/todayspaper/2019/11/10/todays-new-york-times |
| 14 | | https://www.nytimes.com/2019/07/26/us/politics/donald-trump-impeachment.html | https://www.nytimes.com/section/jobs |
| 15 | | https://www.nytimes.com/2019/09/07/us/politics/congress-assault-weapons-ban.html | https://www.nytimes.com/video/us/politics/100000006067682/president-trump-impeach-questions-answers-michael-cohen.html |
| 16 | | https://www.nytimes.com/2018/11/21/obituaries/james-billington-dead.html | https://www.nytimes.com/video/us/politics/100000005524412/paul-manafort-fraud-laundering-scandal.html |
| 17 | | https://www.nytimes.com/2018/09/17/us/politics/bloomberg-president-2020-democrat.html | https://www.nytimes.com/interactive/2018/obituaries/notable-deaths.html |
| 18 | | https://www.nytimes.com/2019/08/07/us/politics/foreign-aid-freeze-congress.html | https://www.nytimes.com/video/us/politics/100000006639814/presidential-candidates-campaigns.html |
| 19 | | https://www.nytimes.com/2018/07/30/obituaries/ron-dellums-forceful-liberal-in-congress-for-27-years-dies-at-82.html | https://www.nytimes.com/video/us/politics/100000006812398/warren-bill-gates.html |
| 20 | Democrats | https://www.nytimes.com/2018/09/04/us/politics/ayanna-pressley-massachusetts.html | https://www.nytimes.com/video |
| 21 | | https://www.nytimes.com/2018/11/07/us/politics/election-history-firsts-blackburn-pressley.html | https://www.nytimes.com |
| 22 | | https://www.nytimes.com/2019/10/16/us/politics/democrats-impeachment.html | https://www.nytimes.com/section/magazine |
| 23 | | https://www.nytimes.com/2019/09/29/us/politics/pelosi-trump-impeachment.html | https://www.nytimes.com/section/politics |

| | | | |
|---|---|---|---|
| 24 | | https://www.nytimes.com/2018/01/23/us/politics/impeach-trump-democrats-tom-steyer.html | https://www.nytimes.com/section/us |
| 25 | | https://www.nytimes.com/2019/09/24/us/politics/trump-democrats.html | https://www.nytimes.com/section/opinion |
| 26 | | https://www.nytimes.com/2019/10/28/us/politics/trump-impeachment-subpoena-house.html | https://www.nytimes.com/issue/todayspaper/2019/11/10/todays-new-york-times |
| 27 | | https://www.nytimes.com/2019/11/05/opinion/trump-democrats-2020.html | https://www.nytimes.com/section/opinion/columnists |
| 28 | | https://www.nytimes.com/2019/11/01/us/politics/house-republican-box-prank.html | https://www.nytimes.com/section/upshot |
| 29 | | https://www.nytimes.com/issue/magazine/2019/06/25/the-61619-issue | https://www.nytimes.com/news-event/2020-election |
| 30 | | https://www.nytimes.com/2019/10/01/travel/what-you-need-to-know-the-tokyo-2020-olympics.html | https://www.nytimes.com/video |
| 31 | | https://www.nytimes.com/2018/02/15/learning/should-technology-in-sports-be-limited.html | https://www.nytimes.com |
| 32 | | https://www.nytimes.com/2018/02/15/learning/should-technology-in-sports-be-limited.html | https://www.nytimes.com/section/magazine |
| 33 | | https://www.nytimes.com/2017/12/05/sports/olympics/ioc-russia-winter-olympics.html | https://www.nytimes.com/section/obituaries |
| 34 | | https://www.nytimes.com/2018/04/23/sports/breakdancing-olympics.html | https://www.nytimes.com/section/politics |
| 35 | Olympics 2020 | https://www.nytimes.com/2016/08/02/sports/olympics/rare-show-of-discord-between-ioc-and-world-anti-doping-agency-over-russian-scandal.html | https://www.nytimes.com/section/realestate |
| 36 | | https://www.nytimes.com/2019/09/23/sports/olympics/russia-doping-wada.html | https://www.nytimes.com/section/us |
| 37 | | https://www.nytimes.com/video/sports/olympics/100000005750317/johnny-weir-tara-lipinski-figure-skating.html | https://www.nytimes.com/section/movies |
| 38 | | https://www.nytimes.com/video/sports/olympics/100000005750317/johnny-weir-tara-lipinski-figure-skating.html | https://www.nytimes.com/video/opinion |
| 39 | | https://www.nytimes.com/video/sports/olympics/100000005752444/south-korea-speedskaters.html | https://www.nytimes.com/section/arts/television |
| 40 | Patriot Movement | https://www.nytimes.com/2019/08/14/opinion/joe-walsh-trump-primary.html | https://www.nytimes.com/video |

| | | | |
|---|---|---|---|
| 41 | | https://www.nytimes.com/1976/08/27/archives/lewis-michaux-92-dies-ran-bookstore-in-harlem.html | https://www.nytimes.com |
| 42 | | https://www.nytimes.com/2019/05/15/us/politics/us-iraq-embassy-evacuation.html | https://www.nytimes.com/video/opinion |
| 43 | | https://www.nytimes.com/2019/01/07/obituaries/moshe-arens-dead.html | https://www.nytimes.com/section/t-magazine |
| 44 | | https://www.nytimes.com/2019/10/27/obituaries/john-conyers-jr-dead.html | https://www.nytimes.com/section/opinion |
| 45 | | https://www.nytimes.com/2018/12/29/obituaries/shehu-shagari-dead.html | https://www.nytimes.com/section/technology |
| 46 | | https://www.nytimes.com/2019/09/23/opinion/democrats-republicans.html | https://www.nytimes.com/issue/todayspaper/2019/11/10/todays-new-york-times |
| 47 | | https://www.nytimes.com/2018/09/04/us/politics/ayanna-pressley-massachusetts.html | https://www.nytimes.com/section/lens |
| 48 | | https://www.nytimes.com/interactive/2018/10/31/magazine/yemen-war-saudi-arabia.html | https://www.nytimes.com/interactive/2018/obituaries/notable-deaths.html |
| 49 | | https://www.nytimes.com/2019/08/15/t-magazine/tropical-brutalism.html | https://www.nytimes.com/interactive/2019/obituaries/notable-deaths.html |
| 50 | | https://www.nytimes.com/2019/10/24/us/politics/senate-republicans-trump.html | https://www.nytimes.com |
| 51 | | https://www.nytimes.com/2019/11/04/us/politics/trump-republicans-impeachment.html | https://www.nytimes.com/section/politics |
| 52 | | https://www.nytimes.com/2019/11/04/us/politics/trump-republicans-impeachment.html | https://www.nytimes.com/section/us |
| 53 | | https://www.nytimes.com/by/catie-edmondson | https://www.nytimes.com/issue/todayspaper/2019/11/10/todays-new-york-times |
| 54 | Republicans | https://www.nytimes.com/interactive/2019/08/14/magazine/republicans-racism-african-americans.html | https://www.nytimes.com/section/upshot |
| 55 | | https://www.nytimes.com/2019/10/04/us/politics/republicans-trump-defense.html | https://www.nytimes.com/news-event/2020-election |
| 56 | | https://www.nytimes.com/2019/10/31/us/politics/house-impeachment-partisan-vote.html | https://www.nytimes.com/video/us/politics/100000006067682/president-trump-impeach-questions-answers-michael-cohen.html |
| 57 | | https://www.nytimes.com/2019/11/09/us/politics/republican-strategy-impeachment-trump.html | https://www.nytimes.com/video/us/elections/100000006810114/kentucky-election-bevin-beshear.html |
| 58 | | https://www.nytimes.com/2019/11/06/us/trump-senate-republicans-courts.html | https://www.nytimes.com/video/us/politics/100000006812485/hillary-clinton-medicare.html |
| 59 | | https://www.nytimes.com/2019/11/06/us/trum | https://www.nytimes.com/video/election-2016 |

| | | | |
|---|---|---|---|
| | | p-senate-republicans-courts.html | |
| 60 | | https://www.nytimes.com/2019/10/01/politics/senate-trump-impeachment.html | https://www.nytimes.com/section/us |
| 61 | | https://www.nytimes.com/by/carl-hulse | https://www.nytimes.com/issue/todayspaper/2019/11/10/todays-new-york-times |
| 62 | | https://www.nytimes.com/2019/10/26/us/trump-senate-presidency.html | https://www.nytimes.com/section/climate |
| 63 | | https://www.nytimes.com/2019/10/24/us/politics/senate-republicans-trump.html | https://www.nytimes.com/news-event/2020-election |
| 64 | Senate | https://www.nytimes.com/2019/11/08/us/impeachment-senate-democrats.html | https://www.nytimes.com/video/opinion/100000006064751/coal-miner-to-trump-coal-mining-isnt-coming-back.html |
| 65 | | https://www.nytimes.com/2019/08/03/us/politics/senate-votes-mcconnell.html | https://timesjourneys.nytimes.com/tour-type/ |
| 66 | | https://www.nytimes.com/2019/08/31/us/politics/senate-race-2020.html | https://www.nytimes.com/times-journeys/tour-type/small-group/ |
| 67 | | https://www.nytimes.com/2019/08/22/us/politics/john-hickenlooper-senate-2020.html | https://www.nytimes.com/times-journeys/trip-type/politics-and-perspective/ |
| 68 | | https://www.nytimes.com/2019/10/18/us/politics/mcconnell-impeachment-senate.html | https://www.nytimes.com/video/us/politics/100000006639814/presidential-candidates-campaigns.html |
| 69 | | https://www.nytimes.com/2019/11/06/us/trump-senate-republicans-courts.html | https://www.nytimes.com/video/us/politics/100000006812485/hillary-clinton-medicare.html |
| 70 | | https://www.nytimes.com/interactive/2017/10/04/us/bump-stock-las-vegas-gun.html | https://www.nytimes.com/section/world |
| 71 | | https://www.nytimes.com/2019/05/02/business/tesla-stock-fundraising.html | https://www.nytimes.com/section/business |
| 72 | | https://www.nytimes.com/2018/12/16/technology/tech-workers-company-stock-shareholder-activism.html | https://www.nytimes.com/issue/todayspaper/2019/11/10/todays-new-york-times |
| 73 | Stock | https://www.nytimes.com/2019/05/10/technology/uber-stock-price-ipo.html | https://www.nytimes.com/section/automobiles |
| 74 | | https://www.nytimes.com/2019/11/04/business/under-armour-stock-investigation.html | https://www.nytimes.com/section/upshot |
| 75 | | https://www.nytimes.com/2019/10/24/style/rebag-clair-handbag-stock-market.html | https://www.nytimes.com/section/fashion |
| 76 | | https://www.nytimes.com/2019/05/30/technology/uber-stock-earnings.html | https://www.nytimes.com/section/world/middleeast |
| 77 | | https://www.nytimes.com/2019/09/20/technology/airbnb-employees-ipo-payouts.html | https://www.nytimes.com/video/opinion/100000006616246/the-king-of-fish-and-chips.html |
| 78 | | https://www.nytimes.com/2018/11/12/obituari | https://www.nytimes.com/by/vanessa-friedma |

| | | es/david-pearson-nascar-dead.html | n |
|---|---|---|---|
| 79 | | https://www.nytimes.com/2019/10/28/business/stock-market-record.html | https://www.nytimes.com/section/business/media |
| 80 | | https://www.nytimes.com/2019/11/09/science/seals-distemper.html | https://www.nytimes.com/section/us |
| 81 | | https://www.nytimes.com/2019/10/25/health/predict-usaid-viruses.html | https://www.nytimes.com/section/world |
| 82 | | https://www.nytimes.com/2019/10/28/us/politics/kay-hagan-dead.html | https://www.nytimes.com/section/science |
| 83 | | https://www.nytimes.com/2019/02/05/us/politics/trump-hiv-aids-plan.html | https://www.nytimes.com/section/health |
| 84 | | https://www.nytimes.com/by/donald-g-mcneil-jr | https://www.nytimes.com/es/ |
| 85 | Virus | https://www.nytimes.com/2019/09/27/well/live/when-is-the-best-time-to-get-your-flu-shot.html | https://www.nytimes.com/section/climate |
| 86 | | https://www.nytimes.com/video/opinion/100000006210828/russia-disinformation-fake-news.html | https://www.nytimes.com/2019/07/15/t-magazine/most-important-contemporary-art.html |
| 87 | | https://www.nytimes.com/video/opinion/100000006210828/russia-disinformation-fake-news.html | https://www.nytimes.com/2019/07/15/t-magazine/most-important-contemporary-art.html |
| 88 | | https://www.nytimes.com/2019/01/23/neediest-cases/daughter-cmv.html | https://www.nytimes.com/2019/09/03/t-magazine/walter-van-beirendonck.html |
| 89 | | https://www.nytimes.com/by/tariro-mzezewa | https://www.nytimes.com/section/well |

**An overlap graph showing the amount of overlap between the two ranking strategies for all queries. By overlap we mean the number of results that appear in the top ten under both ranking methods. Axes should be labelled correctly to indicate query name and overlap number. If no overlaps at all, do not plot an empty graph, but you must mention it in your report.**

| Query | Overlap |
|---|---|
| Cannes | 0 |
| Congress | 0 |
| Democrats | 0 |
| Olympics 2020 | 0 |

| | |
|---|---|
| Patriot Movement | 0 |
| Republicans | 0 |
| Senate | 0 |
| Stock | 0 |
| Virus | 0 |

There is no overlap at all for all queries.