UC Berkeley School of Information (South Hall)

- bcarver@ischool.berkeley.edu
- @brianwc



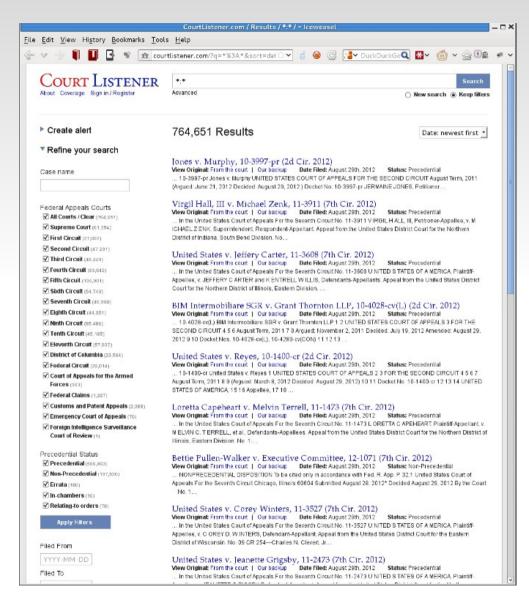




CourtListener.com

Michael Lissner





 "PLEASE IN THE NAME OF GOD, I WOULD PREFER YOU PROTECT MY PRIVACY AND REMOVE THIS CASE FROM THE WEBB SITE. I HAVE ENOUGH INJUSTICE AND HUMILIATION IN REGARDS TO THIS CASE, I THEREFORE WOULD PREFER IT OFF THE WEB SITE, WHERE THE WHOLE WORLD WOULD SEE IT. THANK YOU VERY MUCH. YOUR COOPERATION IS HIGHLY APPRECIATED IN THIS MATTER. SINCERELY,"

• "i was really disgusted with this article my brother is [Full Name]. Thanks, and i hope all your indecency gets aired to the world. People make mistakes my brother is paying for his, but he is still a good person. You don't know all the facts unless you were there. GO TO HELL!!!!!"

 "SIR, PLEASE REMOVE MY CASE FROM UNIVERSAL GOGGLE POSTING. THE FACTS ARE INCORRECT AND WHAT IS CONTAINED IN THIS VIOLATES MY PRIVACH AND HAS CAUSED MANY PROBLEMS IN MY LIFE, I HAD TO HIRE AN ATTORNEY TO ADDRESS THESE ISSUES. PLEASE RESPOND; PLEASE UNDERSTAND IF YOU DO NOT REMOVE THE "RAILROAD JOB" I RECIEVED AT THE HANDS OF THE US GOV'T., SPECIFIC LAWYERS, ARBITRATORS ETC. AND THE INACCCURATE ACCOUT OF THIS CASE...

 ... [case name] I WILL BE FORCED TO BRING APPROPRIATE LEGAL ACTION AGAINST YOU. SOMEONE HAS BEEN POSTING THIS ON ALL SITES FOR YEARS TO DISCREDIT MY REPUTATION. YOUR COMPLICATE PART IN THIS TRAVESTY OF JUSTICE WILL BE MEET WITH."

CourtListener.com

Motivations:

- Daily Awareness Tool (researchers, journalists, lawyers)
- Public Access to the Law
 - The *complete* case law corpus is simply not online without enormous expense (yet!)
- Survivability/Persistence (The "If I get hit by a bus" test)
 - All the code is under open source licenses on bitbucket.org
 - All the documents can be downloaded by anyone
- General purpose case law research tool and citator
- Give it all away for free

Features of our corpus

- 766,000+ documents (with more added daily)
- Coverage of the federal appellate courts, *i.e.*, the 13 circuit courts within the federal system and the Supreme Court of the United States (with more jurisdictions, e.g., the 50 states, to be added soon).
- The entire Supreme Court corpus from 1754 to the present. (From 1 U.S. 1 on; 61,256 documents)
- A (mostly complete) Circuit Court corpus from the late 1940s to present.
- Inter-linked citations.

Collecting Legal Decisions

- In August 2012 we added 2,195 documents, or about 100 documents added each weekday.
- This is from the federal appellate courts. We are in the process of adding all the state appellate courts and the federal district courts and that will increase the document count significantly.
- Unique web scrapers for each court, often multiple scrapers for the same court!
- Unified in our "Juriscraper" project at bitbucket.org

- Stated reasons for requesting removal/blocking:
- (Preliminary data; review is still ongoing)
 - 10% Damage to business / career
 - 40% Privacy
 - 4% Disputes the court's statement of the facts
 - 17% Reputation
 - 19% Safety
 - 10% No reason given

Career:

• "I'm requesting that you remove all information from your site containg my name, [name]. The info dated [date] regarding my discrimination litigation could be very detrimental to my career, namely future business opportunities. I can be reached at the above email address if you need to contact me.; I understand you can publish public material. However, please block any articles containing my name, [name] from the search engines."

Privacy:

• "I am writing to you because on your website there is my wife asylum case published. in this document my wife and friends personal information like compleate name are displayed without any consultaition. i would like to request a name removal or change in this document as is a private and personal matter. my wife name is [name] and i belive that the documet is a thesis of a master student. again i would like to request that the names on the published document were either changed or erased in order to protect my wife and our privacy."

Rebuttal:

• "Please let me know how I can remove or submit my side of the story."

Reputation:

• "When my name is Googled this is the first link that comes up. I am not very happy about this. I would like this link to be removed. How can I do that or who do I need to talk to to make this happen."

Safety:

• "some important information about me is in your internet page, i would like this information be removal or erase, because of my personal security. I am living in [Country Name], and due this information my personal security is at risk. Thank"

- I am compiling not only their stated reasons and attempting to categorize those, but also am gathering the following data:
 - Subjective reason for request
 - Follow-up email?
 - Plaintiff/Defendant, Appellant/Appellee status
 - Jurisdiction; filing date; request date
 - Pro se or represented by counsel?
 - Precedential or non-precedential case? # of citations?
 - Nature of suit; resolution of suit

- Preliminary results:
 - ~90% are appellants/petitioners, i.e., they lost their case below. (Most go on to lose their appeal).
 - ~20% pro se litigants & ~80% represented by counsel
 - ~20% precedential cases & ~80% non-precedential
 - ~34% asylum cases
 - Many of these documents describe facts embarrassing to the individual requesting it be removed/blocked, such as drug convictions or failed employment discrimination claims.
 - Some contain information about a minor.

- Asylum cases involving countries such as:
 - Albania, Columbia, Indonesia, Jamaica, Mauritania,
 Poland, South Africa, Yugoslavia, and more
- Some search results:
 - Asylum and "petition granted" = 380 results
 - Asylum and "petition denied" = 2,815 results
 - < ~12% of these appeals succeed (rough estimate)

Reasons for Public Access

- Promotes judicial integrity
- Enables public confidence in judiciary
- Due process; equal protection; free speech
- Democracy

Reasons for Online Access

- Convenience
- Reach
- Efficiency
- Lower costs (for providers and users)

Proactive Approach

• When a new document is added to the site, we automatically search it for strings that appear to be social security numbers, tax ID numbers, or alien ID numbers and replace the numbers with 'x' in our display version of the document and try to block search engines from indexing the unredacted version.

Balancing the factors

- Favors Privacy
 - Non-precedential cases
 - Some safety concerns seem especially compelling

- Favors Access
 - You lost your case. So?
 - You had counsel

Policy Decision

- CourtListener (inspired by resource.org's policy)
 - will not remove a document without a court order directing us to remove that document, but
 - will try to block search engines from indexing the document upon written request.
- Visitors to our site can still search for and find the document, but those submitting queries to a search engine will (typically) not find the document.
- We believe this typically strikes an appropriate balance between public access and privacy interests.

Relevant Technologies

- Technological fixes:
 - robots.txt file
 - HTML robots noindex metatag
 - The X-Robots-Tag HTTP header
 - sitemap.xml file
 - Search engine-provided webmaster tools
- (Extra special thanks to Michael Lissner for researching all of this!)

robots.txt file

- Pros:
 - Easy
- Cons:
 - Creates a convenient list of pages people want blocked.
 - Addresses only what pages are "crawled" not "indexed"
 - Even search engines that respect robots.txt files may index pages that they don't crawl if, e.g., a third-party site links to the page. Google does this.

HTML robots noindex metatag

- Pros:
 - Easy
 - Addresses indexing
- Cons:
 - Not all bots respect these tags (DuckDuckGo)
 - Only works for HTML documents, not PDFs

The X-Robots-Tag HTTP header

- Pros:
 - Addresses any file type, including PDFs
- Cons:
 - Less easy (e.g., use Apache's mod_headers)
 - Not all bots respect these tags (DuckDuckGo, I Archive)

sitemap.xml file

An XML file that lists the URLs for a site.

Pros:

- Can be used to encourage a search engine to browse URLs you want blocked so that they will encounter your HTTP header or noindex metatag
- Unlike robots.txt, pages to be blocked are not singled out

Cons:

 Unless you block some bots entirely, you will also be encouraging such search engines to index the very pages you want to block

Search engine-provided webmaster tools

Pros:

Search engines that provide these tools (Google/Bing)
 will typically do what you ask them to do

Cons:

- Laborious manual entry process for each URL
- These requests can expire

CourtListener's Approach

- Our robots.txt file does not list individual documents, but does block non-cooperative search engines (either completely or by doc types not supported by noindex tag/headers)
- We use the HTML robots noindex metatag
- We use the X-Robots-Tag HTTP header (and use it on our sitemap.xml file itself)
- Our sitemap.xml file encourages crawling of URLs
- We use search engine-provided webmaster tools

Future Policy Changes?

• When I complete my review of all of these requests, we may conclude that asylum cases represent a special case of a compelling safety concern that justifies efforts to proactively block search engines from indexing these. Stay tuned.

Questions?

bcarver@ischool.berkeley.edu







