

# CS121: Assignment 2: Web Crawler

Jasper Doan, Maxwell Rehm, Sierra Martin

34455076, 71641940, 49838935

University of California, Irvine — October 31, 2024

## Question 1 (*Unique Pages*)

How many unique pages did you find? Uniqueness for the purposes of this assignment is ONLY established by the URL, but discarding the fragment part.

Our crawler found 12983 unique pages.

## Question 2 (*Longest Page*)

What is the longest page in terms of the number of words? (*HTML markup doesn't count as words*)

"<http://www.ics.uci.edu/~shantas/publications/20-secret-sharing-aggregation-TKDE-shantanu>": 147340 words (tokens)

## Question 3 (*Most Common Words*)

What are the 50 most common words in the entire set of pages crawled under these domains ? (Ignore English stop words, which can be found, for example, [here](#)Links to an external site.) Submit the list of common words ordered by frequency.

research: 26512  
2018: 16117  
student: 15673  
software: 15294  
2021: 15252  
information: 14962  
computer: 14927  
data: 14761  
2020: 13683  
2022: 13562  
2019: 13527  
2023: 13338  
will: 13004  
2016: 12999  
students: 12947  
2017: 12869  
uci: 12546

may: 12355  
2015: 11849  
ramesh: 11423  
values: 11120  
engineering: 11030  
10: 10806  
news: 10614  
informatics: 10550  
insert: 10381  
graduate: 10347  
undergraduate: 10052  
ics: 9982  
september: 9960  
us: 9701  
june: 9687  
october: 9651  
july: 9614

projects: 9392  
university: 9391  
ph: 9184  
design: 9098  
march: 8971  
people: 8939  
november: 8830  
april: 8824  
december: 8823  
read: 8811  
science: 8781  
courses: 8760  
department: 8713  
search: 8641  
events: 8624  
january: 8621

#### Question 4 (Subdomain)

How many subdomains did you find in the uci.edu domain? Submit the list of subdomains ordered alphabetically and the number of unique pages detected in each subdomain.

accessibility.ics.uci.edu: 7	grape.ics.uci.edu: 467	ngs.ics.uci.edu: 1779
acoi.ics.uci.edu: 109	graphics.ics.uci.edu: 12	oai.ics.uci.edu: 6
aiclub.ics.uci.edu: 3	graphmod.ics.uci.edu: 16	observium.ics.uci.edu: 1
archive.ics.uci.edu: 121	hack.ics.uci.edu: 2	omni.ics.uci.edu: 1
asterix.ics.uci.edu: 14	hai.ics.uci.edu: 7	onboarding.ics.uci.edu: 1
asterixdb.ics.uci.edu: 2	hana.ics.uci.edu: 4	password.ics.uci.edu: 1
auge.ics.uci.edu: 2	helpdesk.ics.uci.edu: 2	pastebin.ics.uci.edu: 2
awareness.ics.uci.edu: 1	hobbes.ics.uci.edu: 12	pgadmin.ics.uci.edu: 1
cbcl.ics.uci.edu: 58	hombao.ics.uci.edu: 2	phpmyadmin.ics.uci.edu: 2
cert.ics.uci.edu: 2	honors.ics.uci.edu: 1	prometheus-infra-blue.ics.uci.edu: 1
cgvw.ics.uci.edu: 1	hpi.ics.uci.edu: 6	psearch.ics.uci.edu: 6
checkmate.ics.uci.edu: 1	http:www.ics.uci.edu: 1	radicle.ics.uci.edu: 7
chenli.ics.uci.edu: 11	hub.ics.uci.edu: 4	redmiles.ics.uci.edu: 8
cherry.ics.uci.edu: 1	i-sensorium.ics.uci.edu: 6	riscit.ics.uci.edu: 2
chime.ics.uci.edu: 1	iasl.ics.uci.edu: 2	rstudio-hub.ics.uci.edu: 1
cloudberry.ics.uci.edu: 49	icde2023.ics.uci.edu: 49	satware.ics.uci.edu: 1
cml.ics.uci.edu: 170	ics45c-hub.ics.uci.edu: 2	sdcl.ics.uci.edu: 206
code.ics.uci.edu: 15	ics45c-staging-hub.ics.uci.edu: 2	se.ics.uci.edu: 1
codeexchange.ics.uci.edu: 2	ics46-hub.ics.uci.edu: 2	seal.ics.uci.edu: 41
computableplant.ics.uci.edu: 89	ics46-staging-hub.ics.uci.edu: 2	seraja.ics.uci.edu: 2
containers.ics.uci.edu: 1	ics53-hub.ics.uci.edu: 1	sherlock.ics.uci.edu: 9
coronavirustwittermap.ics.uci.edu: 2	ics53-staging-hub.ics.uci.edu: 2	sli.ics.uci.edu: 388
courselisting.ics.uci.edu: 4	ieee.ics.uci.edu: 6	sourcerer.ics.uci.edu: 1
cradl.ics.uci.edu: 12	industryshowcase.ics.uci.edu: 23	speedtest.ics.uci.edu: 1
create.ics.uci.edu: 5	informatics.ics.uci.edu: 1	sprout.ics.uci.edu: 2
cs.ics.uci.edu: 15	insite.ics.uci.edu: 9	staging-hub.ics.uci.edu: 1
cs260p-hub.ics.uci.edu: 2	intranet.ics.uci.edu: 3	stairs.ics.uci.edu: 5
cs260p-staging-hub.ics.uci.edu: 1	ipubmed.ics.uci.edu: 2	statconsulting.ics.uci.edu: 5
cwicsocal18.ics.uci.edu: 13	isg.ics.uci.edu: 219	statistics-stage.ics.uci.edu: 14
cyberclub.ics.uci.edu: 8	jgarcia.ics.uci.edu: 3	student-council.ics.uci.edu: 16
dataguard.ics.uci.edu: 1	jujube.ics.uci.edu: 2	summeracademy.ics.uci.edu: 2
dataprotector.ics.uci.edu: 1	julia-hub.ics.uci.edu: 2	swiki.ics.uci.edu: 82
dblp.ics.uci.edu: 2	kdd.ics.uci.edu: 2	tad.ics.uci.edu: 3
dejavu.ics.uci.edu: 2	kpassword.ics.uci.edu: 1	tastier.ics.uci.edu: 2
dgillen.ics.uci.edu: 33	luci.ics.uci.edu: 6	tippers.ics.uci.edu: 2
ds4all.ics.uci.edu: 4	mailman.ics.uci.edu: 9	tippersweb.ics.uci.edu: 3
duke.ics.uci.edu: 1	malek.ics.uci.edu: 2	transformativeplay.ics.uci.edu: 53
duttgroup.ics.uci.edu: 123	mapgrid.ics.uci.edu: 1	tutoring.ics.uci.edu: 6
dynamo.ics.uci.edu: 21	mcs.ics.uci.edu: 12	tutors.ics.uci.edu: 4
eli.ics.uci.edu: 5	mdogucu.ics.uci.edu: 8	ugradforms.ics.uci.edu: 1
emj-pc.ics.uci.edu: 1	mds.ics.uci.edu: 30	unite.ics.uci.edu: 11
emj.ics.uci.edu: 37	metaviz.ics.uci.edu: 1	vision.ics.uci.edu: 215
esl.ics.uci.edu: 6	mhcid.ics.uci.edu: 22	wearablegames.ics.uci.edu: 3
evoke.ics.uci.edu: 6	mondego.ics.uci.edu: 18	wics.ics.uci.edu: 996
flamingo.ics.uci.edu: 44	motifmap-rna.ics.uci.edu: 3	wiki.ics.uci.edu: 32
fr.ics.uci.edu: 14	motifmap.ics.uci.edu: 3	www-db.ics.uci.edu: 38
frost.ics.uci.edu: 1	mover.ics.uci.edu: 25	xtune.ics.uci.edu: 7
futurehealth.ics.uci.edu: 122	mse.ics.uci.edu: 1	yarra.ics.uci.edu: 1
gitlab.ics.uci.edu: 667	mswe.ics.uci.edu: 12	
grafana-infra-blue.ics.uci.edu: 1	nalini.ics.uci.edu: 8	