

# CS121: Assignment 2: Web Crawler

Jasper Doan, Maxwell Rehm, Sierra Martin

34455076, 71641940, 49838935

University of California, Irvine — October 27, 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 25843 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: 79802  
2021: 79622  
2020: 74071  
2018: 68213  
2019: 66521  
student: 65335  
2022: 65179  
informatics: 61792  
2023: 60318  
2017: 60215  
2016: 58203  
software: 58028  
2015: 57093  
may: 54239  
june: 54066  
ics: 53617  
uci: 53496

september: 50816  
july: 50133  
engineering: 49205  
march: 46643  
january: 45892  
april: 45699  
february: 45641  
november: 45054  
december: 44680  
graduate: 44214  
august: 43906  
undergraduate: 42781  
october: 41426  
ph: 41210  
computer: 38023  
department: 36995  
students: 35704

support: 35330  
will: 35081  
design: 33334  
information: 32451  
projects: 29757  
2014: 29261  
data: 28827  
profiles: 28815  
news: 26431  
read: 25697  
continue: 25072  
reading: 24842  
university: 24613  
years: 23961  
courses: 23575  
alumni: 23371

#### 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: 497	ngs.ics.uci.edu: 127
acoi.ics.uci.edu: 110	graphics.ics.uci.edu: 18	oai.ics.uci.edu: 6
aiclub.ics.uci.edu: 3	graphmod.ics.uci.edu: 15	observium.ics.uci.edu: 1
archive.ics.uci.edu: 91	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: 3	password.ics.uci.edu: 1
auge.ics.uci.edu: 2	helpdesk.ics.uci.edu: 2	pastebin.ics.uci.edu: 2
cbcl.ics.uci.edu: 45	hobbes.ics.uci.edu: 11	pgadmin.ics.uci.edu: 1
cert.ics.uci.edu: 26	hombao.ics.uci.edu: 2	phpmyadmin.ics.uci.edu: 2
cgvw.ics.uci.edu: 1	honors.ics.uci.edu: 1	prometheus-infra-blue.ics.uci.edu: 1
checkin.ics.uci.edu: 5	hpi.ics.uci.edu: 6	psearch.ics.uci.edu: 6
checkmate.ics.uci.edu: 1	hub.ics.uci.edu: 4	pypmyadmin.ics.uci.edu: 1
chenli.ics.uci.edu: 11	i-sensorium.ics.uci.edu: 6	radicle.ics.uci.edu: 7
cherry.ics.uci.edu: 1	iasl.ics.uci.edu: 2	redmiles.ics.uci.edu: 8
chime.ics.uci.edu: 1	icde2023.ics.uci.edu: 49	riscit.ics.uci.edu: 2
cloudberry.ics.uci.edu: 49	ics45c-hub.ics.uci.edu: 2	rstudio-hub.ics.uci.edu: 1
cml.ics.uci.edu: 119	ics45c-staging-hub.ics.uci.edu: 2	satware.ics.uci.edu: 1
code.ics.uci.edu: 15	ics46-hub.ics.uci.edu: 2	sdcl.ics.uci.edu: 120
codeexchange.ics.uci.edu: 2	ics46-staging-hub.ics.uci.edu: 2	se.ics.uci.edu: 1
computableplant.ics.uci.edu: 87	ics53-hub.ics.uci.edu: 2	seal.ics.uci.edu: 46
containers.ics.uci.edu: 1	ics53-staging-hub.ics.uci.edu: 2	sherlock.ics.uci.edu: 9
coronavirustwittermap.ics.uci.edu: 2	ieee.ics.uci.edu: 3	sli.ics.uci.edu: 408
courselisting.ics.uci.edu: 4	industryshowcase.ics.uci.edu: 23	sourcerer.ics.uci.edu: 2
cradl.ics.uci.edu: 1	informatics.ics.uci.edu: 2	speedtest.ics.uci.edu: 1
create.ics.uci.edu: 5	insite.ics.uci.edu: 9	sprout.ics.uci.edu: 2
cs.ics.uci.edu: 17	instdav.ics.uci.edu: 1	staging-hub.ics.uci.edu: 2
cs260p-hub.ics.uci.edu: 2	intranet.ics.uci.edu: 3	stairs.ics.uci.edu: 5
cs260p-staging-hub.ics.uci.edu: 2	ipubmed.ics.uci.edu: 2	statconsulting.ics.uci.edu: 5
cwicsocal18.ics.uci.edu: 13	isg.ics.uci.edu: 180	statistics-stage.ics.uci.edu: 1
cyberclub.ics.uci.edu: 8	jgarcia.ics.uci.edu: 3	student-council.ics.uci.edu: 15
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	support.ics.uci.edu: 1
dblp.ics.uci.edu: 2	kdd.ics.uci.edu: 1	svn.ics.uci.edu: 1
dejavu.ics.uci.edu: 2	kpassword.ics.uci.edu: 1	swiki.ics.uci.edu: 148
dgillen.ics.uci.edu: 33	luci.ics.uci.edu: 6	tad.ics.uci.edu: 3
ds4all.ics.uci.edu: 4	mailman.ics.uci.edu: 7	tastier.ics.uci.edu: 2
duke.ics.uci.edu: 1	malek.ics.uci.edu: 2	tippers.ics.uci.edu: 2
duttgroup.ics.uci.edu: 105	mapgrid.ics.uci.edu: 1	tippersweb.ics.uci.edu: 3
dynamo.ics.uci.edu: 21	mcs.ics.uci.edu: 12	transformativeplay.ics.uci.edu: 53
eli.ics.uci.edu: 5	mdogucu.ics.uci.edu: 8	tutoring.ics.uci.edu: 6
elms.ics.uci.edu: 1	mds.ics.uci.edu: 30	tutors.ics.uci.edu: 4
emj-pc.ics.uci.edu: 1	metaviz.ics.uci.edu: 1	ugradforms.ics.uci.edu: 1
emj.ics.uci.edu: 37	mhcid.ics.uci.edu: 22	unite.ics.uci.edu: 11
esl.ics.uci.edu: 6	mondego.ics.uci.edu: 16	vision.ics.uci.edu: 215
evoke.ics.uci.edu: 6	motifmap-rna.ics.uci.edu: 1	wearablegames.ics.uci.edu: 3
flamingo.ics.uci.edu: 43	motifmap.ics.uci.edu: 3	wics.ics.uci.edu: 352
fr.ics.uci.edu: 13	mse.ics.uci.edu: 1	wiki.ics.uci.edu: 87
futurehealth.ics.uci.edu: 122	mswe.ics.uci.edu: 12	www-db.ics.uci.edu: 38
gitlab.ics.uci.edu: 54	mysqladm.ics.uci.edu: 1	xtune.ics.uci.edu: 7
grafana-infra-blue.ics.uci.edu: 1	nalini.ics.uci.edu: 8	yarra.ics.uci.edu: 1