link to home page) s	on about this spreadsheet, please visit https://github.com/archivers-space/research/tree/master/web.archiving General information Users & developer features								**************************************																	Other notes					
	Open Source Primary dev. Target				Target				Paralle	Scheduled	Crawl storage	Capture raw	"Out of the		wling capabili nced Extra		Run	Extract links	:		ranced data harv Manual form			Archive n	anagement	eatures Full-text		Other notes	Date/versi Evaluation	sion checked Version	
	source?	repo	Operating system(s)	language	audlence	CLI (I API		crawlin	g crawling	format(s)	responses	Follow links UR				vaScript Handle I	eact from Flash	Run Flash		interaction		data	Browse	Playback	search	Notable users	Notes and comments	date	examined
Whole site archiving sy	ystems																														
rchive-It	16	n/a	web app	Java	enterprise	*	* ·	*	*		~	ARC, WARC		~	~		•	× ×	~	*		*	×		-	~	~	Run by IA. Used by 100's of institutions.	Paid service. Core software is Heritrix, with some extensions.	2017-05-21	
Prozzler A	Apache	GitHub	Lin, Win, macOS	Python	user	~	* ~	*	*	~	*	WARC	~	~	* ;	•	~	· ·	*	*	*	*	*	*	~	~	*	IA	Uses Chromium to fetch pages. Uses warcprox and pywb.	2017-07-14	1611
rawler	~	GitHub	Lin, Win, macOS	PHP	user	*	* v	*	* *	×	×	MySQL	*	~	*		*	* *	*	*	*	*	*	*	*	*	*	Made by the FCC.	Bare bones crawler. License not stipulate. Purpose unclear.	2017-06-01	2012-06-04
Crawler4j A	Apache	GitHub	Lin, Win, macOS	Java	user	×	× ×	~	■ Java classes	·	×	files on disk	*	~	* ;	•	*	* *	*	*	*	*	*	*	×	×	×			2017-05-23	
<u>Orawljax</u> A	Apache	GitHub	Lin, Win, macOS	Java	user	~	* v	~	≭ Plug-ins	~	limited	log file; plug-ins can do more	*	~	* .	•	•	, ,	*	*	~	~	*		~	*	×		Several papers written about the implementation.	2017-06-01	3.5
		GitHub	Lin, Win, macOS	Python	user	~	* ·	×	IPC Python	*	×	WARC		~	~		~	* *	*	*		*	*	×	*	×	*		Uses wpull internally.	2017-05-15	
ecco	MIT	GitHub	Lin, Win, macOS	Java	user	×	× ×	~	■ Java classes	-	*		*	~	*		~	, ,	*	*	~	~	*	*	*	*	×		Code comments are in Chinese.	2017-06-01	
		GitHub	Lin	Java	user	~	* ·	*	MX Java classes	-	*	ARC, WARC	~	~	-		~	* *	-	*	*	×	×		*	*	*	Used by IA.		2017-05-15	
fTTrack	GPL	GitHub	Lin, Win	С	user	-	*	~	C callbacks	-	×	files on disk	×	~	,	,	V	* *	*	*	*	*	×	×	*	×	*			2017-07-03	
Sucks	GPL	<u>SF</u>	Lin, Win, macOS	Java	user	-	/ ×	~	■ Java classes	-	×	files on disk	*	~	,	,	*	* *	*	*		*	*		*	×	×			2017-05-15	
letarchiveSuite I	LGPL	GitHub	Lin	Java	user, enterprise	~	* v	×	Java classes	·	~	ARC, WARC	~	~	•		~	* *	*	*	*	*	*:		~		•	Netarkivet at The Royal Library of Denmark	Uses Heritrix for crawling,	2017-07-14	
lutch A	Apache	Apache	Lin, Win, macOS	Java	user, enterprise	v .	× ×	~	Plug-ins	~	~	several db options	×	~		,	~	· ·	×	*	~	~	×	~	*	*	×			2017-06-01	
Octoparse	*	*	Win	NET	user, enterprise	*	× *	* 1	EST #	~	~	database, CSV, Excel, files on disk	~	~		,	~	, ,	*	*	~	~		~	,				Seems to run in the cloud, but there's a downloadable console or something.	2017-09-03	6.4.3
ageFreezer	×	×	web app	n/a	user, enterprise	×	* v	×	*	~	~	web pages	*	~					*	*		*	×	×	,	~		EDGI web monitors use it		2017-10-04	
implecrawler	BSD	GitHub	Lin, Win, macOS	Nodejs	user	v	× ×	~	Node module:		*	files on disk	~	~	, ,	,	v .	* *	*	*	*	v	×	*	*	*	×			2017-09-03	115
iquidwarc G	GPLv3	Github	Lin, macOS	Nodejs	user, enterprise	v	× ×	~	Node module:		*	WARC	~	~	* ,		·	, ,	*	*	*	*	×	*	*	×;	×		A high fidelity archival crawler that uses Chrome or Chrome Headless	2017-07-21	d4ca0b8
itormCrawler A	Apache	GitHub	Lin, Win, macOS	Java	user,		× ×	~	■ Java classes	~	limited	several db options	×	~	, ,			, ,		*	~	~	×	~	×	*	×	Several companies, apparently.	https://github.com/DigitalPebble/storm- crawler/wiki/Presentations	2017-06-01	
VAIL (Electron) G	GPLv3	GitHub	Lin, Win, macOS	Node js (Flectron)	user	*	× ×	~	* *	~	~	WARC	~	~	,		,	, ,	*	*	*	*	×	*	,	v	×		Uses Chrome Browser and Heritrix for crawling, Pywb, Twitter Monitoring And Automatic Archival	2018-12-13	120-Beta2
VAIL (py)	MIT	GitHub	Win, macOS	Python	user	×		×	× ×	~	×	WARC	~	~	× ;		×	× ×	×	×	×	*	×	×	,	v	×		OpenWayback and Heritrix, cf. WAIL (Electron)	2018-12-13	v0.2016.0709
VebMagic A	Apache	GitHub	Lin, Win, macOS	Java	user	*	× ×	v	≭ Java		×	files on disk		~	, ,	,	~	v *	*	*	~	~	*	~	×	*	×		Intended as a programming framework, not end-user app. Comments are in Chinese	2017-10-04	0.7.3
WebRecorder.io A	Apache	GitHub	Lin, Win, macOS	Python	user		<i>,</i> ,	~	* *	×	×	WARC	~		* ,		~	v v	*	*	*	~	*	×	~	v			Interactive, high-fidelity web archiving tool	2017-07-03	eccea96
rget	GPL	Savannah	Lin, Win, macOS	С	user	~	× ×	×	* *	×	×	WARC, files on disk	~	~	✓ lim	ited	×	* *	×	*	*	×	×	×	*	*	×		Use option –save-headers to save HTTP headers	2017-05-17	119
rpull	GPL	GitHub	Lin, Win, macOS	Python	user	~	* *	V	✔ plug-ins, scrip	s ×	×	WARC	~	~	~		v .	* *	*	*	×	*	×	*	*	*	×			2017-05-20	201
Single page snapshot/	/archivin	na system	10																												
	*	n/a	web app	n/a	user	*	× v	×	EST #	v	×	web page	×	×	* ;		×	V	×	v	×	limited	×	*	×	v	×		Very good quality page captures.	2017-05-15	
		GitHub	Lin, Win, macOS	С	user	v	× ×	,	× ×	*	*	files on disk	~	*	· ,		*	* *	*	*	*	limited	×	*	*	*	*	Standard with many	Yesy good quarry page captures.	2017-05-21	
	×	n/a	web app	n/a	user	×	* v	×	* *	*	*	web page	×	*	* ;		×	v *	×	*	*	*	×	*	,	~	×	operating systems	Seems to be free for use, but not open source.	2017-06-25	
aparazzi!	×	n/a	macOS	n/a	user		v *	×	* *	×	×	PDF, PNG, JPG, TIFF	*	*	* ·		*	v *	*	~	*	limited	*	×	×	*	×	mhucka uses this all the time	Very good quality full-page captures.	2017-05-24	
Perma.cc Mil	IT + GPL	GitHub	web app	Python	user,	×	* v	*	EST Django	*	×	WARC, PDF, PNG	×	*	x ,		×	~	×	*	*	×	×	×	,	limited	×			2017-05-24	
VARCreate	MIT	GitHub	Chrome extension	JavaScrint	user	*	/ ×		* *	*	*	WARC			. ,		×	v *		*	*		*	*	×	×	×			2017-07-14	
		GitHub	Lin, Win, macOS	Python	user	v	* *	×	* *	*	*	PNG	*	*	* ,		×	v *	*	v	×	*	×	×	*	*	*			2017-05-27	
Data scraping systems	ıc																														
		GitHub	Lin, Win, macOS, EMR	Java	user,	v	× ×	v	■ Java classes	~	*	files	*	v	v		×	* *	*	*	*		×	*	×	×	×		Author seems to be creator of Krugle. Uses Apache	2017-06-25	
	*	n/a	EMR web app	n/a	enterprise user,	*	* v	*	EST #		V	JSON, CSV,	×	,	,		v li	mited	*	*	_	v	×		*	*	×		Nutch, Hadoop, Tika, others. Offers edu & charity discounts. Has a pretty active	2017-05-24	
RobotSoft.com	×	*	Win	n/a	enterprise user,	_	· *		× ×	,	,	Gdocs, Tableau files on disk, user	*	~	, ,		·	~	*		_	,	*	~	*	*			user forum.	2017-10-04	282
	Affero	GitHub	Lin, Win, macOS,	Ruby	enterprise user,	<i>-</i>	* ·		Plug-ins (suppo	ts	limited	db SQLlite		,	,		,	,		*	,	· ·		·		~ ×	×		Offers cloud-based scraping. Active user forum.	2017-05-29	202
			Docker, web app Lin, Win, macOS,	Python	enterprise		* '		many langs;		limited		*	,	<i>y</i>			, ,	*	* _	,	<i>y</i>	*		* ~	*	* *			2017-05-29	
ortia		GitHub	Docker	.,		*		,	Python			files, MySQL, git		-													-	A version of this is used by	Has visual scraping definition editor. This is a fork; the original (by Martins Balodis) has		
	r:DLv2	GitHub	Chrome extension	JavaScript	user	×	× v	×	× ×	*	×	CSV, CouchDB	×	~	,		×		×	×	-		×	~	*	×	×	WebScraperio (commercial)	not changed since 2014.	2017-10-04	0.3.1
WebScraper.io (fork)																															
WebScraper.io (fork)		n/a	Win	n/a	user		× ×	*	* *	*	*	files, user db	*	~		,	~		*	*	~	~	*	~	*	*	×		Can extract data & store in database. May no longer be supported.	2017-05-23	

For more information a	bout this spi	out this spreadsheet, please visit https://github.com/archivers-space/research/tree/master/web_archiving. General information Users & developer features																													
Name			General information	on		L	Jsers & develop	er features		**************************************											lata harvesting capabilities Archive management features Other notes al form Auto form Extract file Full-facet Full-facet							ion checked			
Name (link to home need)	Open	Source		Primary dev.	Target	CII CIII	work Extensi	ibility	Parallel Scheduled	Crawl storage	Capture raw	Follow links URL filtering	Advanced	Extract links	Run	Handa Dar	Extract links	Pun Flore	Targeted	Manual form	Auto form	Extract file data	Brown	Full-to Playback sear	xt Notable users	Notes and comments	Evaluation date	Version examined			
80legs.com	source/	repo (Operating system(s)	language	audience	CLI GUI W	IUI API A	UPI trames	work c	crawling crawling	format(s)	responses	Follow links URL filtering	tiltering	from JavaScript	JavaScript	Handle Heact	from Flash	Mun Flash	scraping	Interaction	Interaction	data	Browse	Playback sear	h Notable users	Notes and comments	date	examined		
Abot																															
AbotX									- 1											1											
Andjing																															
Anemone																															
<u>Aperture</u>																															
Apifier																															
Arachnid (Java)																															
Arachnid (PHP)																															
<u>arachniweb</u>																															
Arale																															
ArchiveBot																															
ARCOMEM																															
ASPseek																															
Bingol																															
blekko	×																														
CCBot																															
cl-web-crawler		-	b					,																							
CrawlBot		-	web app					•												1						-					
crawler.js																				1											
crawwwler																															
<u>DataparkSearch</u>																				1						1					
DeepArc																															
	×		macOS																												
Django Dynamic Sc	raper																														
dryscrape																															
EIS Archiver																															
Ex-Crawler																															
F(b)arc																															
Gungho																															
Hounder																															
html-snapshots																															
html2warc																															
HyperSpider (JS)																															
icrawler																															
iwebcrawler							,																								
jedi-crawler																															
Jspider																															
JWAT																															
Knowlesys																															
LARM									_																						
Lassie																				1						-					
Lentil																				1						1					
LinkGrabber																															
<u>METIS</u>																															
Miru																															
mnoGoSearch																															
mummif.it																															
Newspaper																															
Norconex HTTP Colle	ector																														
NutchWAX																				1											
OpenWayback									-+															~	~						
OpenWebSpider																				1											
									-+											1											
OutbackCDX																				1											
PageVault																				1						-					
Panscient																															
ParseHub																															
eeep.us																															
hantom-crawler																															
inboard																															
Preservica																															
PromptCloud																															
uppeteer																											1				
appeteet																				1							1				

For more information at	more information about this spreadsheet, please visit https://github.com/archivers-space/research/tree/master/web_archiving_ General information Users & developer features																											211				
			General information								1				"Out of the box									rvesting capa		Archive r	managemen			Other notes		rsion checked
Name (link to home page)	Open source?	Source	Operating system(s)	Primary dev.	Target audience	cu e				Extensibility framework			Crawl storage format(s)		Follow links URL filts		from JavaSci			Extract links from Flash	Run Flash			Auto form interaction		Browse	Playback	Full-text search	Notable users	Notes and comments	Evaluation date	Version examined
scrala			, , , , , ,															,				, ,										
screen-scraper																																
ScrapBook																																
ScrapeKit																																
Scrapy																																
SeimiCrawler																																
Sentry	Affero	GitHub	Lin, Win, macOS	Go		~					~	~																		l		
Shine																																
Smart Cache Loader																																
Sparkler	Apache	GitHub																														
Sphider																																
Spiderman																																
Spindle																																
twarc																																
UbiCrawler																																
warc-discovery																																
warcbase																													1			
WARClight																													1			
warcprox																													1			
wayback				Java																									1	An older, unbranded version of IIPC OpenWayback		
Web Harvest																																
Web::Scraper																													1			
Web2Warc																													1			
WebArchivePlayer																																
WebCite																																
webcitation.org																																
WebEater																																
WebLech																																
	Annaha	CHILL	1 to 180 000	Python,																	~											
Webrecorder Player	мраспе	Giunub	LIT, WIT, MIACUS	Electron		,											~	~														
webshot																																
website-scraper																																
WebSPHINX																																
wiki-crawler																																
WWW-Mechanize																																
WWW:Crawler:Lite																																
Xapian																																
XArch																																
xidel																																
																											_					
- Dead or deprecate	ea softwa	are (nome p	pages no longer ex	ust, etc.) –																												
Combine											-											-				-		-	-		-	
Hyperix											-											-				-		-	1		-	
Nalanda											-															-					-	
- Insufficient inform	nation (e.	g., no expla	nations or descrip	tions, or de:	scriptions	written in	an unfa	miliar lang	guage) –	-																						
spiderframework																															2017-07-03	
																										-						
											1															1						
- Too specialized for	or a speci	ific purpose	·-																													
ArchiveFacebook	~		Firefox Add-On	IncoConta	user	,	, .,	*		×	×	*	Original Web	*	* *	*	*	×	*	*	×	v	×	*	v	~	v	×		Firefox add-on for archiving a Facebook acct. Uses	2017-07-03	
ni ci iiver acebook	•	GIGHUD	rireiux Add-UN	JavaScript	user	- '	•	•	•	•			Resources	•	- *			•	•	•	•	,	•	•	,			- 1		old XUL add-on format, no longer maintained	2017-07-03	
HamarCarld - (1)	Public	SF.net	Lin, Win, macOS	Java		,	× ×	*	×	*	*		CSV, txt, DOT, RDF, Prolog,	*	v *	*	*	*	*	*	*	×	×	*	*	×	*	×	1	Only collecte the link etayst of	2017-09-03	2002.02.0
HyperSpider (Java)	domain	SEnet	Lin, Win, macus	Java	user			*	*	•	*	•	HDF, Prolog,	*	, x	*	*			*	*		*					*	1	Only collects the link structure of a website	2017-09-03	2003-02-0
python-sitemap	GPLv3	GitHub	Lin, Win, macOS	Python	user	v :	× ×	*	×	*																				Only collects the link structure of a website	2017-09-03	1868340
											-															-						
pysitemap	Apache	GitHub	Lin, Win, macOS	Python	user	v :	× ×	×	×	*																×	×	×	-	Only collects the link structure of a website	2017-10-04	
python-goose	Apache	GitHub	Lin, Win, macOS	Python	user	v :	××	×	×	*	×	*	files on disk	*						*	*							1	1	Specialized for extracting text of news articles	2017-10-04	10.25
SyncUThink	GPL	GitHub	Lin, Win, macOS	Java	user																								1	Crawls CiteULike accounts to download PDFs	2017-07-03	
				5040	usci																					-		-		Does not archive or scrape. Might also be known as		
Visual Sitemap Gener	×	×				*	× /	×	×	*																		J	1	"DYNO Mapper".	2017-10-04	
WERA	~			PHP								×	WARC		* *							*	×	×	×	~	,			Tool for searching and navigating archived web		

For more information about this spreadsheet, please visit https://github.com/archivers-space/research/tree/master/web_archiving_																															
	General information Users & developer features											"Out of the box" website crawling capabilities										nced data hari	esting capab	ilities	Archive management features			Other notes			sion checked
Name	0	реп	Source		Primary dev.	Target		LI	ibrary Network	Extensibility	Parallel Scheduled	Crawl storage	Capture raw		Advanced	Extract links	Run	E	Extract links		Targeted	Manual form	Auto form	Extract file			Full-text			Evaluation	Version
(link to home pay	ge) so	ource?	repo	Operating system(s)	language	audience	CLI GUI	WUI	API API	framework	crawling crawling	format(s)	responses	Follow links URL filtering	filtering	from JavaScript	JavaScript	Handle React	from Flash	Run Flash	scraping	Interaction	Interaction	data	Browse	Playback	search	Notable users	Notes and comments	date	examined