Starrydata: a plot mining web system for literature data

Yukari Katsura^{1,2}, Masaya Kumagai^{3,4}, Mitsunori Kaneshige⁵, Yuki Ando^{2,3}, Sakiko Gunji^{1,2}, Yoji Imai^{2,3}, Takushi Kodani^{1,2}, Riku Sato^{1,2}, Kaoru Kimura¹, Koji Tsuda^{1,2,3}

Graduate School of Frontier Sciences, the University of Tokyo¹, NIMS-MI²I², RIKEN-AIP³, SAKURA Internet Inc.⁴, X-Ability Co., Ltd.⁵

Plot mining from published papers on materials science



Our goals

1. Make every published data accessible

2. Make research efficient by data science



Our approach

- 1. Develop an efficient web system for data collection
- 2. Design an economic system (community) for data collection

東京大学 CRIKEN AP SAKURA internet

3. Share the data freely for worldwide researchers



Starrydata web system

We developed this reference-manager-like web system to share numerical data extracted from plot images written in JavaScript and Python, by using Django, MongoDB, VEGA.js and WebPlotDigitizer (https://automeris.io/WebPlotDigitizer).

Free access

X-Ability



300 320 340 360 380 400 420 440 46

9. Visualization of all data in the mylist







11. Reading data in python

										rawc	la	ita (a	allı	num	neri	c)						
											1	figureid	pa	perid	prop	pertyid	_x p	propertyid_y	sampleid		x	
										73272		623		145			0	3	904	820.6	9060	1.25754
iman		1.0								73273	•	623		145			0	3	904	837.79	9285	1.25056
import pandas as pd							73274		623		145			0	3	904	869.04	18740	1.25003			
										73275		623		145			0	3	904	894.03	34353	1.23622
=op	en	'cl	ath	rate.json	1','r')					73276		623		145			0	3	904	915.93	36873	1.26264
lict	ct_all=json.load(f)							73277		623	145				0	3	904	933.13	8024	1.27575		
df_rawdata=pd.DataFrame(dict_all["rawdata"])							73278		sam	uple compositio		sitior	n paj	paperid	sampleid	sample	ename	3259	1.29565			
df_f	igu	ire=	pd.	DataFrame	(dict_all	dict_all["fi	gure"]	, כני		73279		884	Ba8Cu6Si1		6Ge24	4	140	884		4GPa	1913	1.31548
df_s	amp	ole=	pd.	.DataFrame(dict_all["sample"])						73281		885	Ba80	Cu6Si16Ge24		4	140	885		5GPa	747	1.35516
ui _p	1.01		.y-v	u.bucurru	meturce_u		prope	ar cy 17		/3201		886		Ba8Cu	16Si40	0	141	886	Ba8C	u6Si40	,,41	1.55510
pa	pe	erau	thor	author full		doi	issue	iournal	iour	mal full		887	Ba80	Cu6Ge2	205i20	D	141	887	Ba8Cu6Ge	20Si20	Ves	ar
_		Sara	mat	A Saramat			10040	Journal	Jou			888	Ba8	BCu6Ge	85132	ron	142	888		x=0		-
145	^.	G	. S	G. S	10.1063/1.216	3979	2	of Appli	50	Appli	C	889	Ba8	Cu6Ge	, P	pro	perty	y id	prope	rtyname	,	ı
146	J.	Marti	n, S. Er	J. Martin, S. Er	10.1063/1.217	1775	4	Journal of	Jo	ournal of Appli	c	890	Ba80	Cu6Ge1	"	0		0	Tem	perature)	
								Appli								1		2 E	Electrical con	ductivity	/ ohn	n^(-1)*m^
147	Ch	r f	igu	ire				caption	on 1	figureid	fig	gurenan	ne (paperid	1	2		1	Seebeck co	oefficien	t	V*K^
			619	Total them	mal conductivit	y of p	ourified	Ba8Ga16	i	619			5	144	• ;	3		5	Pow	er facto	r W⁺r	n^(-1) * K^
148			620	Figure of me	erit ZT of purifie	d Ba8	8Ga160	Ge30 sam		620			6	144	4	4		3	Thermal con	ductivity	/ W⁺r	n^(-1) * K^
149	A		621	a Se	ebeck coefficie	ent S a	as a fui	nction of t		621		2	(a)	145	i 1	5		8 Dimensi	noless figure	e of meri	t	
		١.	622	2. a S	Seebeck coeffic	cient \$	S as a f	function o)	622		2((b)	145		6		4	Electrical r	resistivity	/	ohn
			623	a Se	ebeck coefficie	ent S a	as a fui	nction of t		623		2	(C)	145		7		7	Inverse tem	perature	•	K**

Example data: Thermoelectric properties

Thermoelectric materials are the materials that interconvert temperature difference and electricity. The efficiency increases with increasing $ZT=S^2\sigma T/\kappa$ (or increasing power factor $PF=S^2\sigma$), where S is Seebeck coefficient [V/K], σ is electrical conductivity [$\Omega^{-1}m^{-1}$], and κ is thermal conductivity [Wm⁻¹K⁻¹] (*T* is temperature [K]). Even though all these

