

# Using Python to organise my physical paper

Alex Chan · [alexwlchan.net](http://alexwlchan.net) · [they/she](https://they/she)

Slides: [alexwlchan.net/files/2021/campug.pdf](http://alexwlchan.net/files/2021/campug.pdf)

container





### How to vote

1. When you receive your postal vote, read the instructions carefully.
2. Your postal vote includes your ballot paper and a postal voting statement.
3. Complete both of these and return them immediately.
4. We need to receive your postal vote by 10pm on Thursday 12 December 2019.

If you need information in another format, please call our helpline below.

If you need help to vote, you can ask someone you know or get independent help by calling our helpline:

Helpline:

**01438 242213**

E-mail:

**electoral@stevenage.gov.uk**

Web:

**www.stevenage.gov.uk**

If you lose your postal vote or make a mistake

- Please phone the helpline immediately.
- We can only issue a replacement postal vote before 5pm on Thursday 12 December 2019.

If you would rather vote in person, or ask someone else to vote on your behalf, you must cancel your postal vote before 5pm on Tuesday 26 November 2019. For more information, please call the helpline.

It is an offence to:

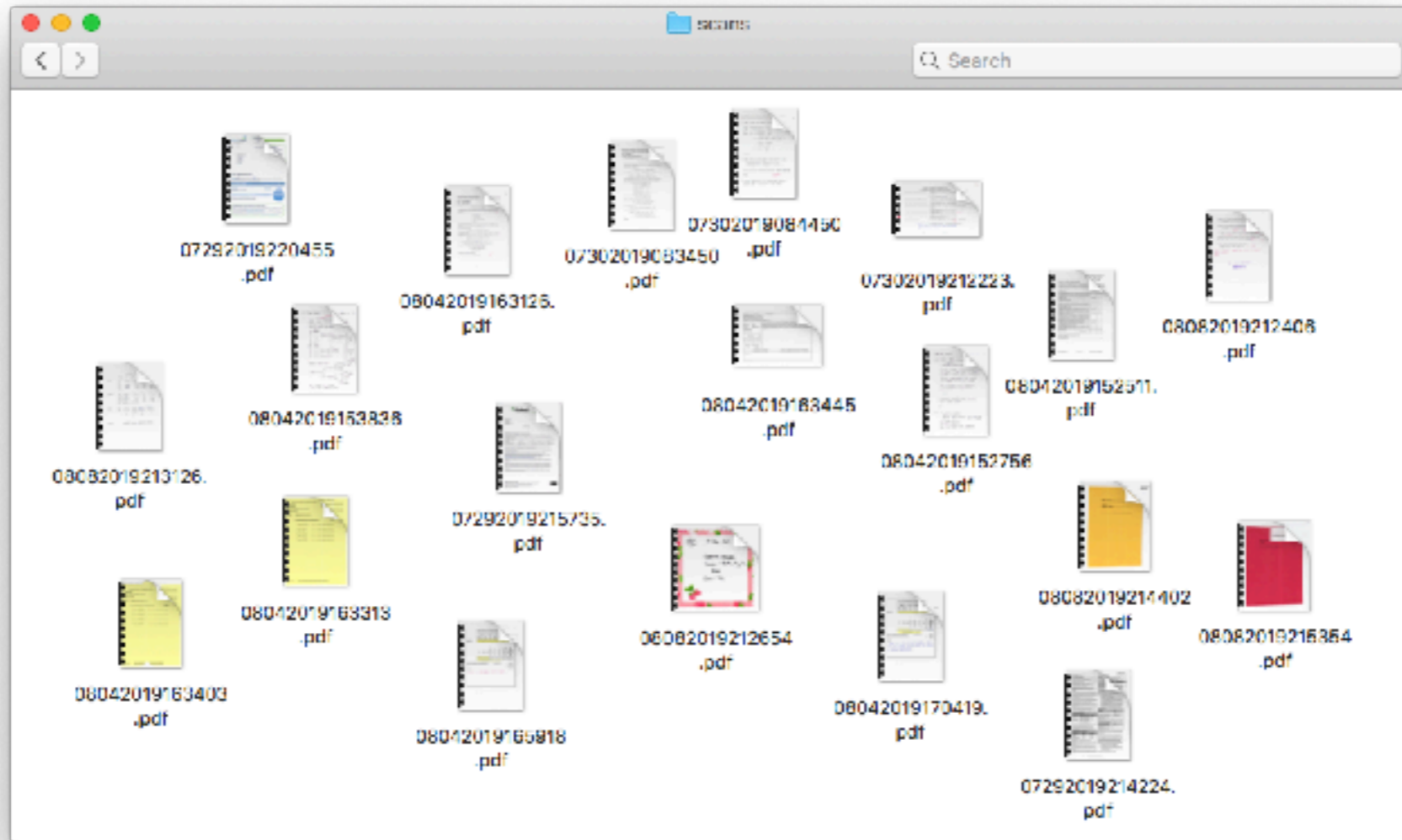
- vote using a ballot paper that was not sent for you, use or interfere with another voter's ballot paper
- vote more than once at this election, unless you are voting on your own behalf and as a proxy for another person
- vote as a proxy at this election for more than two people, unless you are their spouse, civil partner, parent, grandparent, brother, sister, and/or grandchild
- vote as a proxy for someone if you know that by law they are not allowed to vote

The Acting Returning Officer issued this card.

If undelivered return to:  
The Acting Returning Officer, Stevenage Borough Council

Canon





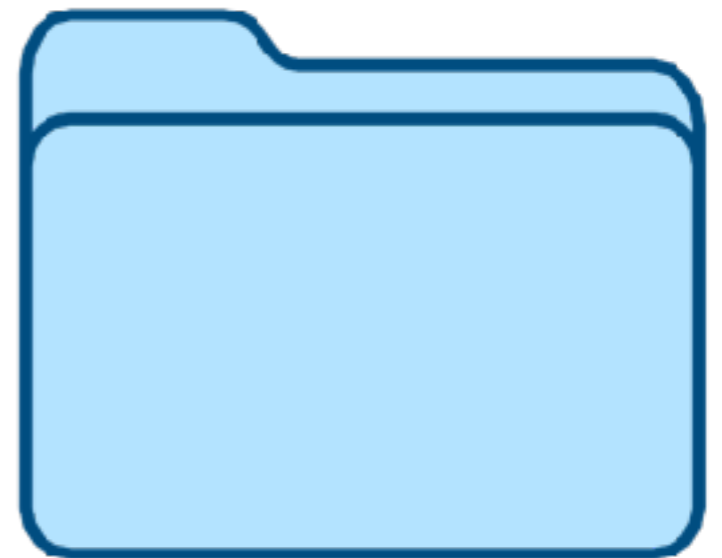
# How do I want to find documents?



home

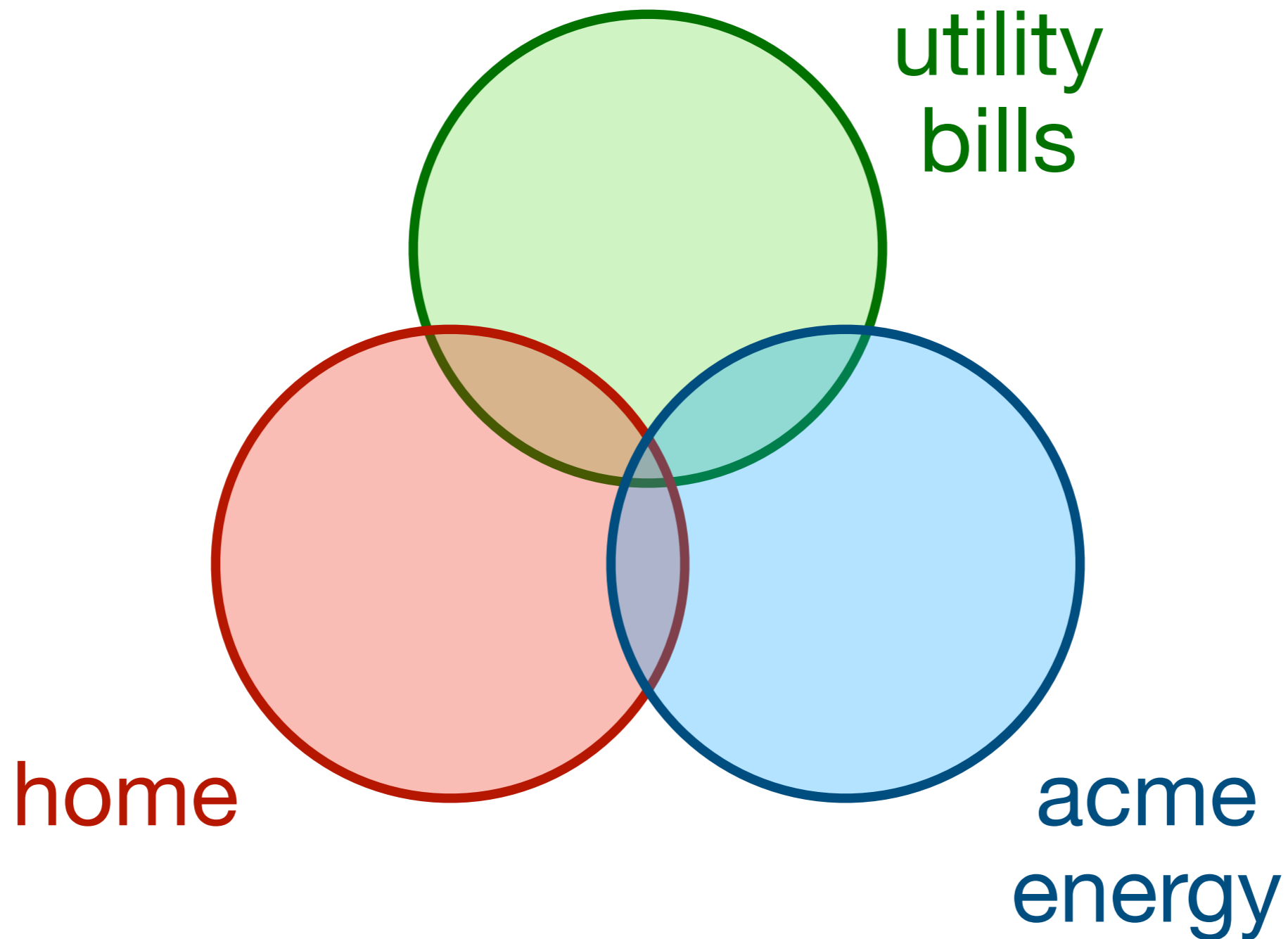


utility  
bills



acme  
energy

# How do I want to find documents?



```
docstore add 'BillJul21.pdf' \  
  --source_url='https://acme.co/gas' \  
  --title='2021-07: Gas bill for 2021' \  
  --tags='home,utility-bills,acme-energy'
```



showing documents 1-54 of 54.  
▶ tag list  
▶ tag cloud



### For All Mankind Lunar Surface Access Module (LSAM) LEGO instructions, by [SpaceXplorer](#)

date saved: 6 Jun 2021  
source: [rebrickable.com](#)  
tagged with: for-all-mankind lego sci-fi



### Journey into our Underworld

date saved: 17 Apr 2021  
source: [www.docspstore.com](#)  
tagged with: [archives](#)



### Volcano pattern, by [CrossStitchMari](#)

date saved: 26 Feb 2021  
source: [www.fastmail.com](#)  
tagged with: cross-stitch-patterns



### JPEG 2000 spec

date saved: 4 Oct 2019  
tagged with: [file-format-specs](#) [welcome](#)



### Papers from "Making IT Work" conference

date saved: 16 Aug 2016

showing documents 1-54 of 54.

▶ tag list

▼ tag cloud

archives aws by:CrossStitchMari by:Happinest by:James Clear by:Mathysphere by:SpaceXplorer by:The Duckbill Group by:ZAnnaCrossStitch chinese  
 community-safety conferences **cross-stitch-patterns** elasticsearch esu family file-format-specs for-all-mankind  
 from:Gendered Intelligence gsk languages lego **lgbtq-network** music org:mhfa-england org:open-university  
 org:prospect-union personal pride-planets recipes rpg sci-fi self-improvement sketchnotes **stonewall** talks **trans-awareness-training**  
**trans-inclusion-policies** wellcome write-the-docs



### For All Mankind Lunar Surface Access Module (LSAM) LEGO Instructions, by SpaceXplorer

date saved: 6 Jun 2021  
source: [rebrickable.com](http://rebrickable.com)  
tagged with: [for-all-mankind](#) [lego](#) [sci-fi](#)



### Journey Into our Underworld

date saved: 17 Apr 2021  
source: [www.deepstore.com](http://www.deepstore.com)  
tagged with: [archives](#)



### Volcano pattern, by CrossStitchMari

date saved: 26 Feb 2021  
source: [www.fastmail.com](http://www.fastmail.com)  
tagged with: [cross-stitch-patterns](#)



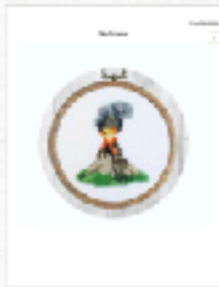
### JPEG 2000 spec

date saved: 4 Oct 2019  
tagged with: [file-format-specs](#) [wellcome](#)

showing documents 1-10 of 10.

by:CrossStitchMari by:Happinesst by:Mathysphere by:ZannaCrossStitch **cross-stitch-patterns** music pride-planets

filtering to tag **cross-stitch-patterns** [x]



**Volcano pattern, by [CrossStitchMari](#)**

date saved: 26 Feb 2021  
source: [www.fastmail.com](http://www.fastmail.com)  
tagged with: cross-stitch-patterns



**Math Sampler pattern, by [Mathysphere](#)**

date saved: 2 Jan 2020  
source: [www.fastmail.com](http://www.fastmail.com)  
tagged with: cross-stitch-patterns



**Solar Eclipse cross-stitch pattern, by [Mathysphere](#)**

date saved: 2 Jan 2020  
source: [www.fastmail.com](http://www.fastmail.com)  
tagged with: cross-stitch-patterns



**Bi Planet pattern, by [Mathysphere](#)**

showing documents 1-6 of 6.

by:Mathysphere cross-stitch-patterns [pride-planets](#)

filtering to tags **cross-stitch-patterns** [x] **by:Mathysphere** [x]



### Math Sampler pattern, by Mathysphere

date saved: 2 Jan 2020  
source: [www.fastmail.com](http://www.fastmail.com)  
tagged with: cross-stitch-patterns



### Solar Eclipse cross-stitch pattern, by Mathysphere

date saved: 2 Jan 2020  
source: [www.fastmail.com](http://www.fastmail.com)  
tagged with: cross-stitch-patterns



### Trans Pride Planet pattern, by Mathysphere

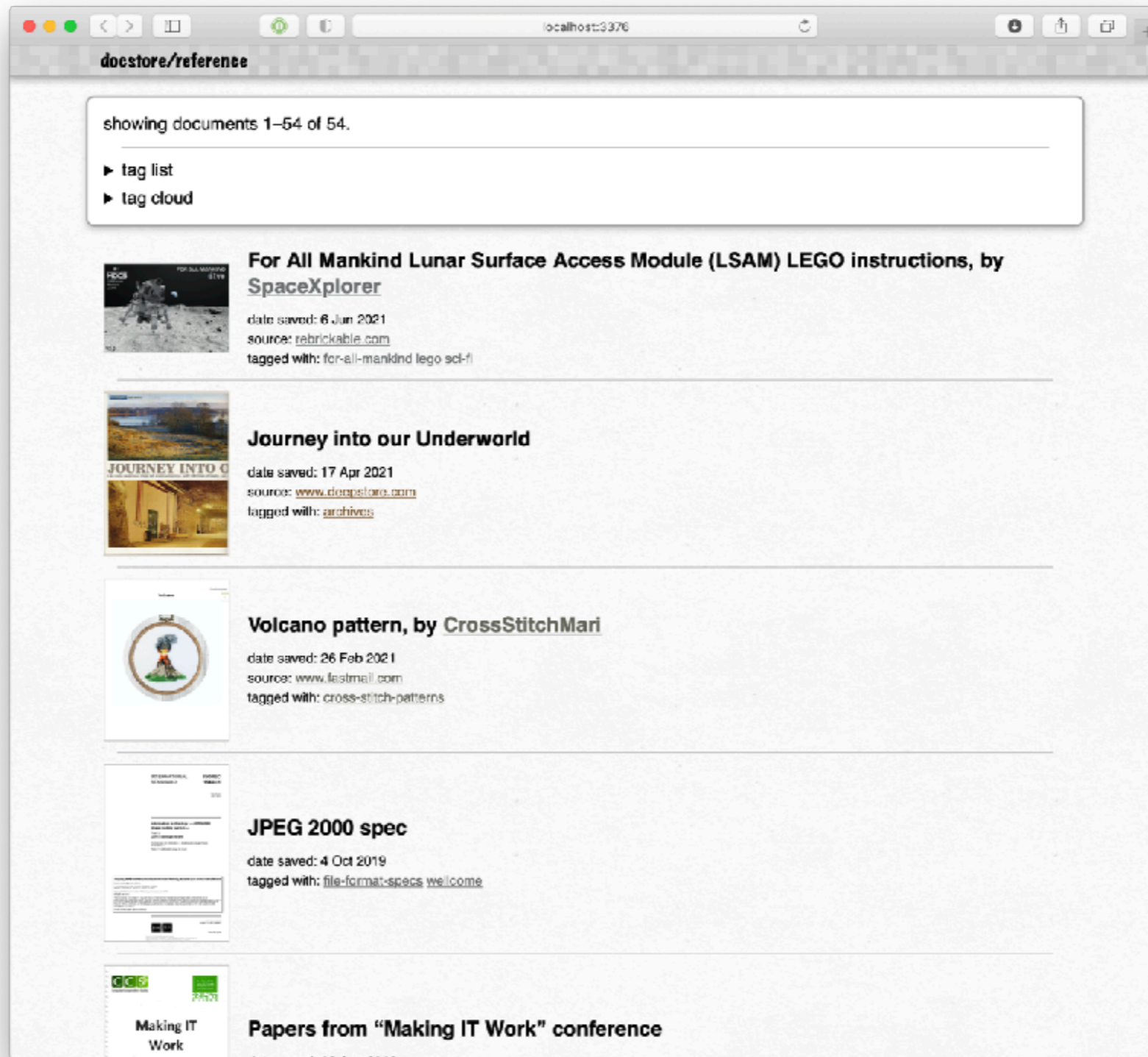
date saved: 19 Apr 2019  
tagged with: cross-stitch-patterns [pride-planets](#)



### Bi Planet pattern, by Mathysphere







alexwlchan/docstore

Storing the files

Storing the metadata

Previewing the files

Wrangling the tags



1/

Storing the files

What should I save  
these files as?

VolcanoPattern.pdf

Alex Chan\_5312.pdf

Statement.pdf

alex.chan@metaswitch.com>Alexander  
Chan>Payslip November 2014-2015.PDF

V5C:3 scrappage note.pdf

UUIDs?

VolcanoPattern.pdf

↳ volcanopattern.pdf

I don't store files under  
their original filenames

Alex Chan\_5312.pdf

↳ alex\_chan\_5312.pdf

Statement.pdf

↳ statement.pdf, statement\_aea1.pdf, ...

alex.chan@metaswitch.com>Alexander Chan>Payslip  
November 2014-2015.PDF

↳ alex\_chan\_metaswitch\_com\_alexander\_chan\_payslip\_  
november\_2014-2015.pdf

V5C:3 scrappage note.pdf

↳ v53\_c\_scrappage\_note.pdf

# unicodecode

```
>>> import unicodecode
```

```
>>> unicodecode.unicodecode("Statement.pdf")  
"Statement.pdf"
```

```
>>> unicodecode.unicodecode("ACME™ Corp Shipping note")  
"ACME(tm) Corp Shipping note"
```

```
>>> unicodecode.unicodecode("Alex Chan>Payslip Nov.PDF")  
"Alex Chan>Payslip Nov.PDF"
```

```
>>> unicodecode("北京")  
"Bei Jing "
```

# open()

```
>>> with open("greeting.txt", mode="w") as outfile:  
...     outfile.write("Hello world!")
```

12

```
>>> with open("greeting.txt", mode="r") as infile:  
...     print(infile.read())
```

"Hello world!"

```
>>> if not os.path.exists("important.txt"):  
...     with open("important.txt", mode="w"):
```

```
...
```

# open(..., "x")

```
>>> with open("name.txt", mode="x") as outfile:  
...     outfile.write("My name is Alex!")
```

16

```
>>> with open("name.txt", mode="x") as outfile:  
...     outfile.write("My name is Tibs!")
```

Traceback (most recent call last):

File "<stdin>", line 1, in <module>

FileExistsError: [Errno 17] File exists: 'name.txt'

<a>

<a href="/files/beijing.pdf">

↳ downloads "beijing.pdf"

Content-Disposition: attachment; filename="北京.pdf"

<a href="/files/beijing.pdf">

↳ downloads "北京.pdf"

<a href="/files/beijing.pdf" download="北京.pdf">

↳ downloads "北京.pdf"

# Storing the files:

- You don't have to store a file under its original name
- Exclusive creation mode "x" can avoid filename collision
- `unicodecode` can create ASCII-safe strings
- The Content-Disposition header or download attribute can suggest a filename



# 2/

## Storing the metadata

- How do I model the metadata?
- How do I save the metadata to disk?

# attrs

```
>>> import attr
```

```
>>> @attr.s
```

```
... class Document:
```

```
...     path = attr.ib()
```

```
...     tags = attr.ib()
```

```
>>> doc = Document(
```

```
...     path="scanned_doc.pdf",
```

```
...     tags=["home", "bills", "acme energy"]
```

```
... )
```

# attrs

```
>>> import attr
```

```
>>> @attr.s
```

```
... class Document:  
...     path = attr.ib()  
...     tags = attr.ib()
```

```
>>> doc = Document(  
...     path="scanned_doc.pdf",  
...     tags=["home", "bills", "acme energy"]  
... )
```

```
>>> doc.path  
"scanned_doc.pdf"
```

```
>>> doc.tags  
["home", "bills", "acme energy"]
```

# attrs

```
>>> import attr
```

```
>>> @attr.s
```

```
... class Document:  
...     path = attr.ib()  
...     tags = attr.ib()
```

```
>>> doc = Document(  
...     path="scanned_doc.pdf",  
...     tags=["home", "bills", "acme energy"]  
... )
```

```
>>> repr(doc)
```

```
Document(path="scanned_doc.pdf", tags=["home",  
"bills", "acme energy"])
```

# attrs

```
>>> import attr
```

```
>>> @attr.s
```

```
... class Document:  
...     path = attr.ib()  
...     tags = attr.ib()
```

```
>>> doc = Document(  
...     path="scanned_doc.pdf",  
...     tags=["home", "bills", "acme energy"]  
... )
```

```
>>> eval(repr(doc)) == doc
```

```
True
```

```
>>> doc == Document(path="cat.jpg", tags=["pets"])
```

```
False
```

# attrs

```
>>> import attr
```

```
>>> @attr.s
```

```
... class Document:
...     path = attr.ib()
...     tags = attr.ib()
```

```
>>> doc = Document(
...     path="scanned_doc.pdf",
...     tags=["home", "bills", "acme energy"]
... )
```

```
>>> attr.asdict(doc)
{"path": "scanned_doc.pdf", "tags": ["home", "bills",
"acme energy"]}
```

# attrs

```
>>> import attr
```

```
>>> @attr.s
```

```
... class Document:
```

```
...     path = attr.ib()
```

```
...     tags = attr.ib()
```

```
>>> doc = Document(
```

```
...     path="scanned_doc.pdf",
```

```
...     tags=["home", "bills", "acme energy"]
```

```
... )
```

```
@attr.s
```

```
class Document:
```

```
    id = attr.ib()
```

```
    title = attr.ib()
```

```
    date_saved: datetime.datetime = attr.ib()
```

```
    tags: List[str] = attr.ib()
```

```
    files: List[File] = attr.ib()
```

```
@attr.s
```

```
class File:
```

```
    id = attr.ib()
```

```
    filename = attr.ib()
```

```
    path = attr.ib()
```

```
    size: int = attr.ib()
```

```
    checksum = attr.ib()
```

```
    thumbnail: Thumbnail = attr.ib()
```

```
    source_url = attr.ib()
```

```
    date_saved: datetime.datetime = attr.ib()
```



```
{
  "date_saved": "2020-10-21T23:22:56.878358",
  "files": [
    {
      "checksum": "sha256:3a6d3bc8641bc43530bb3e6cd183d26a5a05c69b3b74c38e99d5ef9fb2c01253",
      "date_saved": "2020-10-21T23:22:56.878358",
      "filename": "THE_HAUNTING_OF_TRAM_CAR_015_MOBI.mobi",
      "id": "0eb55943-26c7-4a94-99d0-c1fe39318fdf",
      "path": "files/t/the-haunting-of-tram-car-015-mobi.mobi",
      "size": 2267910,
      "source_url": "https://ebookclub.tor.com/",
      "thumbnail": {
        "dimensions": {
          "height": 400,
          "width": 250
        },
        "path": "thumbnails/t/the-haunting-of-tram-car-015-mobi.mobi.png",
        "tint_color": "#9b642c"
      }
    }
  ],
  "id": "9f9bb56e-35ab-451c-81a6-f5ee95fce912",
  "tags": [
    "fantasy",
    "historical-fantasy",
    "fiction",
    "by:P. Dj\u00e8l\u00eded Clark",
    "format:mobi"
  ],
  "title": "The Haunting of Tram Car 015"
}
```



attrs

cattrs



.json

# cattr

```
>>> import cattr
```

```
>>> doc = Document(  
...     path="scanned_doc.pdf",  
...     tags=["home", "bills", "acme energy"]  
... )
```

```
>>> cattr.unstructure(doc)  
{"path": "scanned_doc.pdf", "tags": ["home", "bills",  
"acme energy"]}
```

```
>>> cattr.structure(  
...     {"path": "cat.jpg", "tags": ["pets"]},  
...     Document)  
Document(path="cat.jpg", tags=["pets"])
```

# Storing the metadata:

- How do I model the metadata?  
Use the `attrs` library.
- How do I save the documents to disk?  
JSON is fine for small N. Often better than a traditional database.
- How do I serialise `attrs` models to JSON?  
Use `cattrs`.

3/

Previewing the files

showing documents 1-54 of 54.  
▶ tag list  
▶ tag cloud



### For All Mankind Lunar Surface Access Module (LSAM) LEGO Instructions, by [SpaceXplorer](#)

date saved: 6 Jun 2021  
source: [rebrickable.com](#)  
tagged with: [for-all-mankind](#) [lego](#) [sci-fi](#)



### Journey into our Underworld

date saved: 17 Apr 2021  
source: [www.deepstore.com](#)  
tagged with: [archives](#)



### Volcano pattern, by [CrossStitchMari](#)

date saved: 26 Feb 2021  
source: [www.fastmail.com](#)  
tagged with: [cross-stitch-patterns](#)



### JPEG 2000 spec

date saved: 4 Oct 2019  
tagged with: [file-format-specs](#) [wellcome](#)



### Papers from "Making IT Work" conference

date saved: 18 Aug 2019

showing documents 1-49 of 49.  
▶ tag list  
▶ tag cloud



**The Night Circus, by Erin Morgenstern**

date saved: 6 Jul 2020  
tagged with: [fantasy](#) [fiction](#) [format:epub](#)



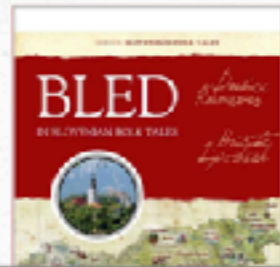
**Porting to Python 3 (phone)**

date saved: 27 Jul 2019  
source: [gumroad.com](#)  
tagged with: [format:odf](#) [programming](#) [python](#)



**Intentional Friendships: A Gentle Guide (landscape)**

date saved: 27 Jul 2019  
source: [gumroad.com](#)  
tagged with: [format:odf](#) [zines](#)



**Bled in Slovenian folk tales, by Dusica Kunaver**

date saved: 1 Jul 2019  
source: [www.amazon.co.uk](#)

showing documents 1-36 of 36.  
filtering to tag council-tax [x]



**2016-03: Council tax bill 2016/2017**

date saved: 8 Mar 2019  
tagged with: council-tax



**2018-03: Council Tax bill 2018/2019**

date saved: 24 Feb 2019  
tagged with: council-tax



**2017-06: Council tax bill 2017/2018**

date saved: 24 Feb 2019  
tagged with: council-tax



**2015-03: Council tax bill 2015/2016**

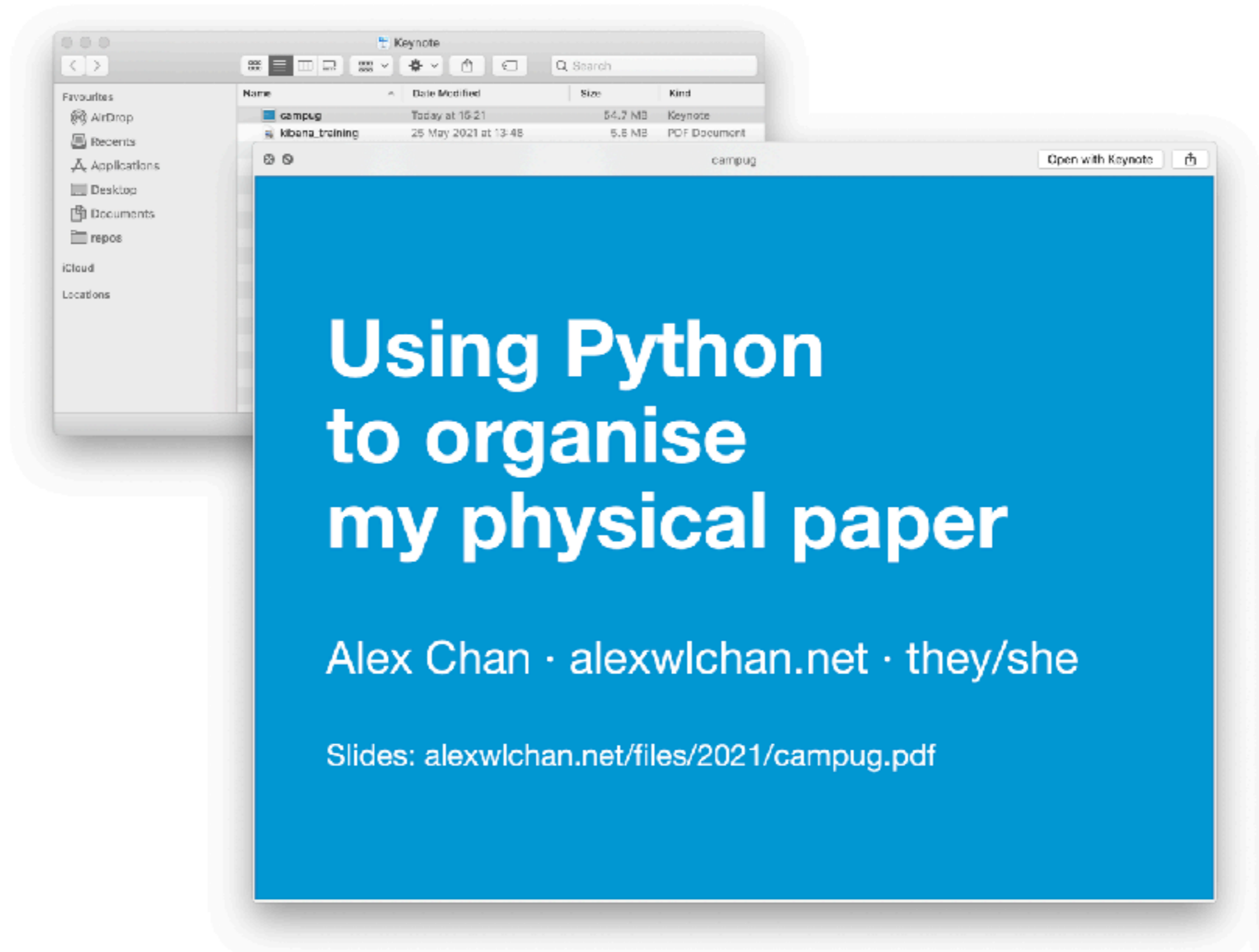
date saved: 23 Feb 2019  
tagged with: council-tax





```
pdftocairo document.pdf \  
-jpeg \  
-singlefile \  
-scale-to-x 400
```

```
qlmanage -t -s 400 document.jpeg
```



```
ffmpeg -i animated.gif \  
-movflags faststart \  
-pix_fmt yuv420p \  
-vf scale=400:400  
out.mp4
```



localhost:3373

docstore/books

showing documents 1-49 of 49.

- ▶ tag list
- ▶ tag cloud

---

**A Visit from Eldritchbot, by [Thomas Heasman-Hunt](#)**

date saved: 3 Oct 2020  
source: [www.patreon.com](http://www.patreon.com)  
tagged with: [format:pdf](#) [sci-fi](#) [smolrobots](#)

---

**Finn Family Moomintroll, by [Tove Jansson](#)**

date saved: 2 Oct 2020  
tagged with: [fiction](#) [moomins](#) [format:pub](#)

---

**Artificial Condition, by [Martha Wells](#)**

date saved: 17 May 2020  
tagged with: [fiction](#) [format:pub](#) [murderbot](#) [sci-fi](#)

---

**Bled in Slovenian folk tales, by [Dusica Kunaver](#)**

date saved: 1 Jul 2019  
source: [www.amazon.co.uk](http://www.amazon.co.uk)  
tagged with: [folk-stories](#) [format:epub](#)

---

**The Life-Changing Magic of Tidying Up, by [Marie Kondo](#)**



## Finn Family Moomintroll, by [Tove Jansson](#)

date saved: 2 Oct 2020

tagged with: [fiction](#) [moomins](#) [format:epub](#)

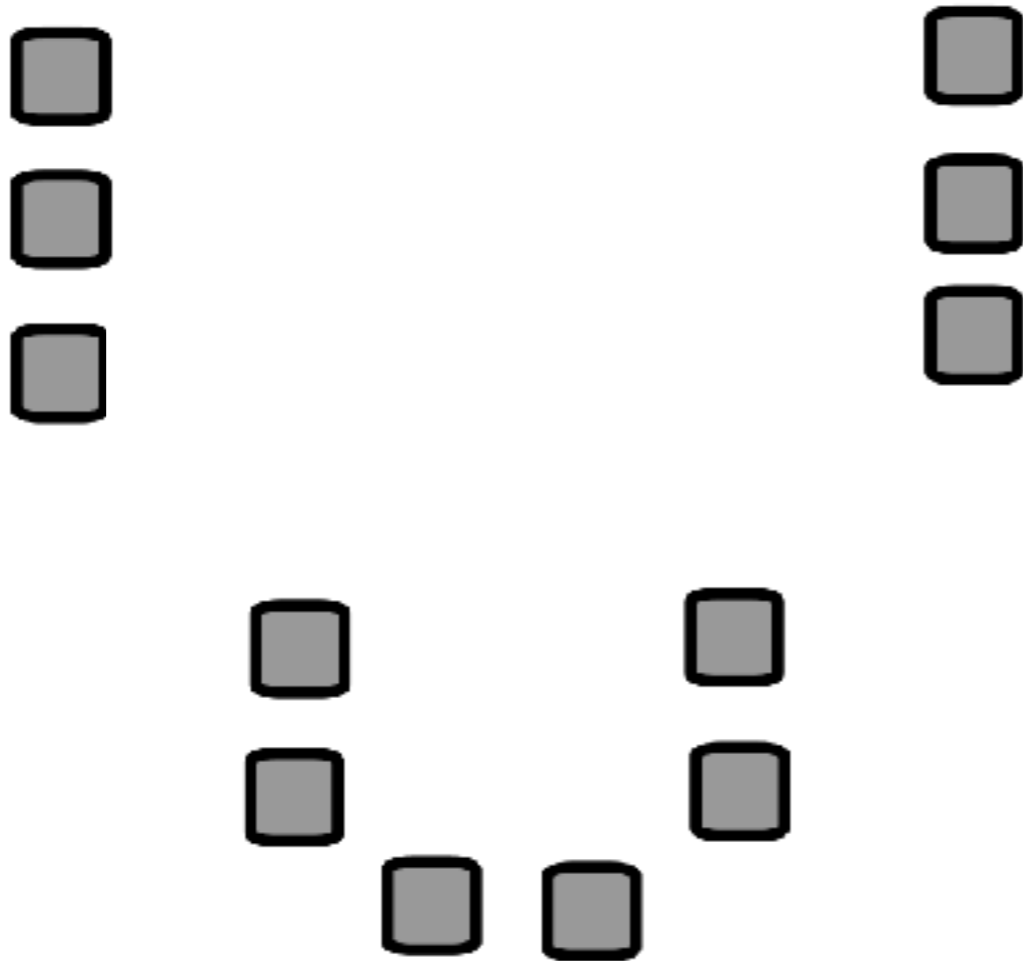
```
>>> from PIL import Image
>>> im = Image.open("cover.jpg")

>>> tally = collections.Counter(im.getdata())

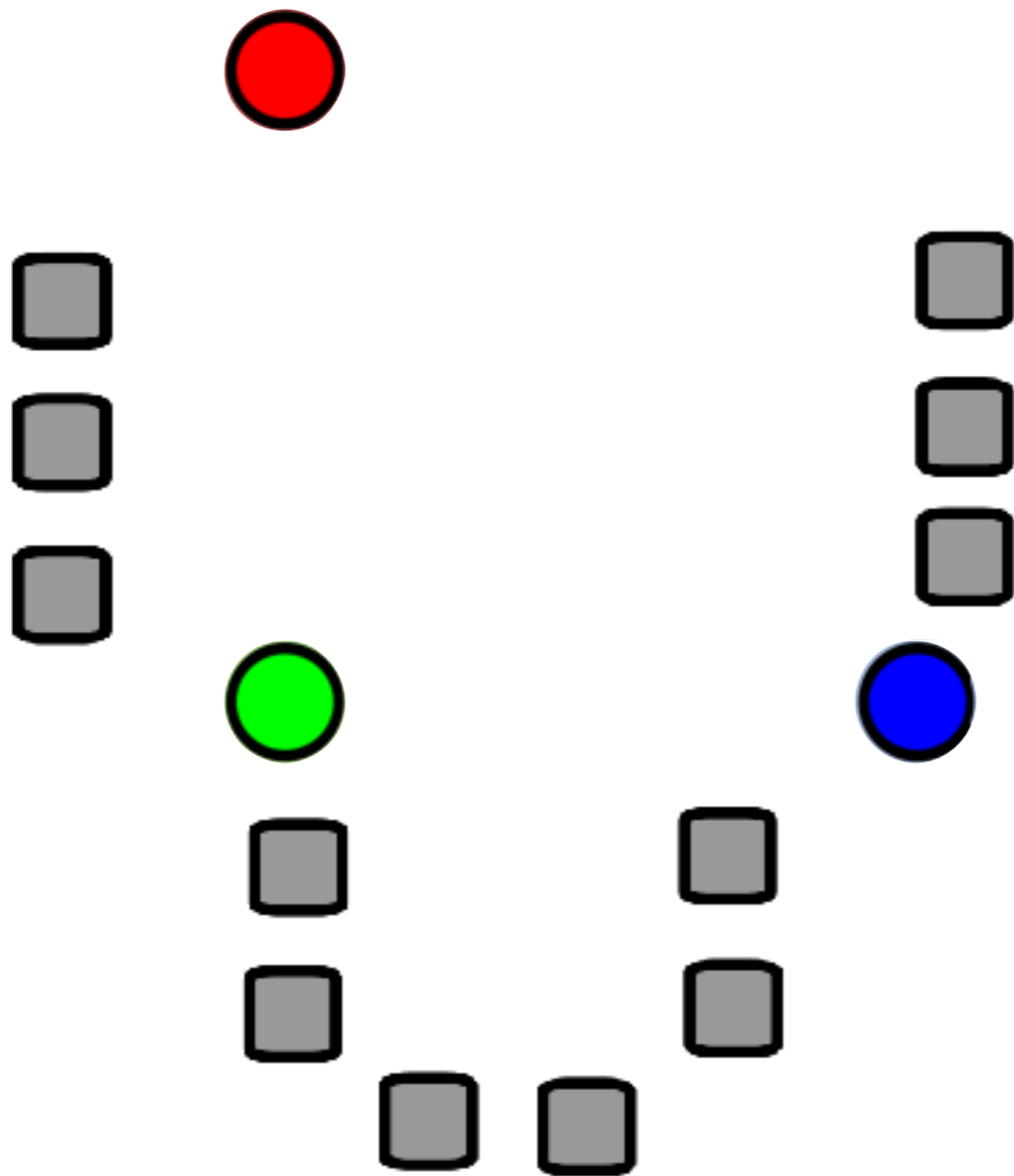
>>> tally.most_common(1)
[((249, 246, 241), 2106)]
```

#f9f6f1

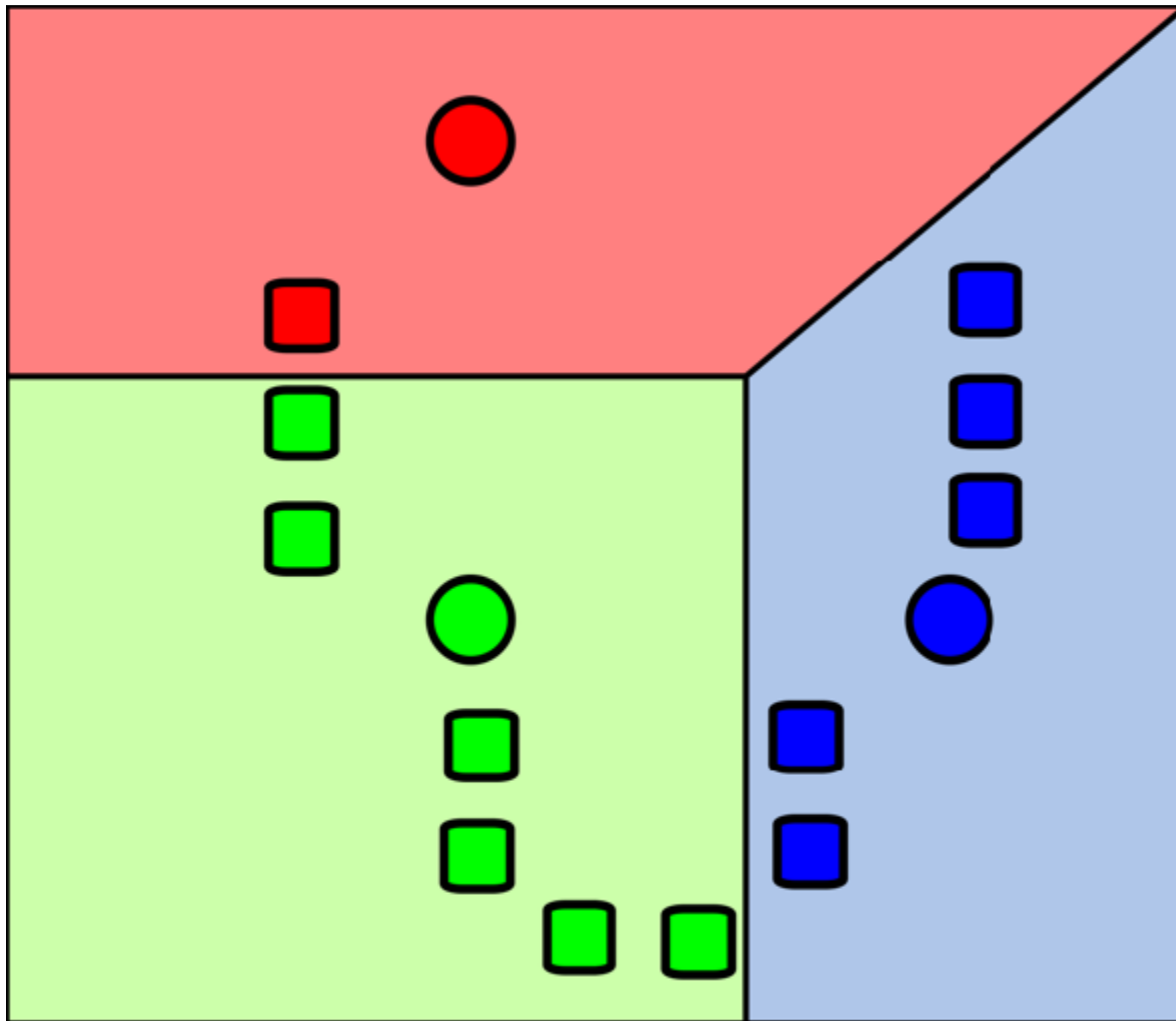
# k-means



# k-means

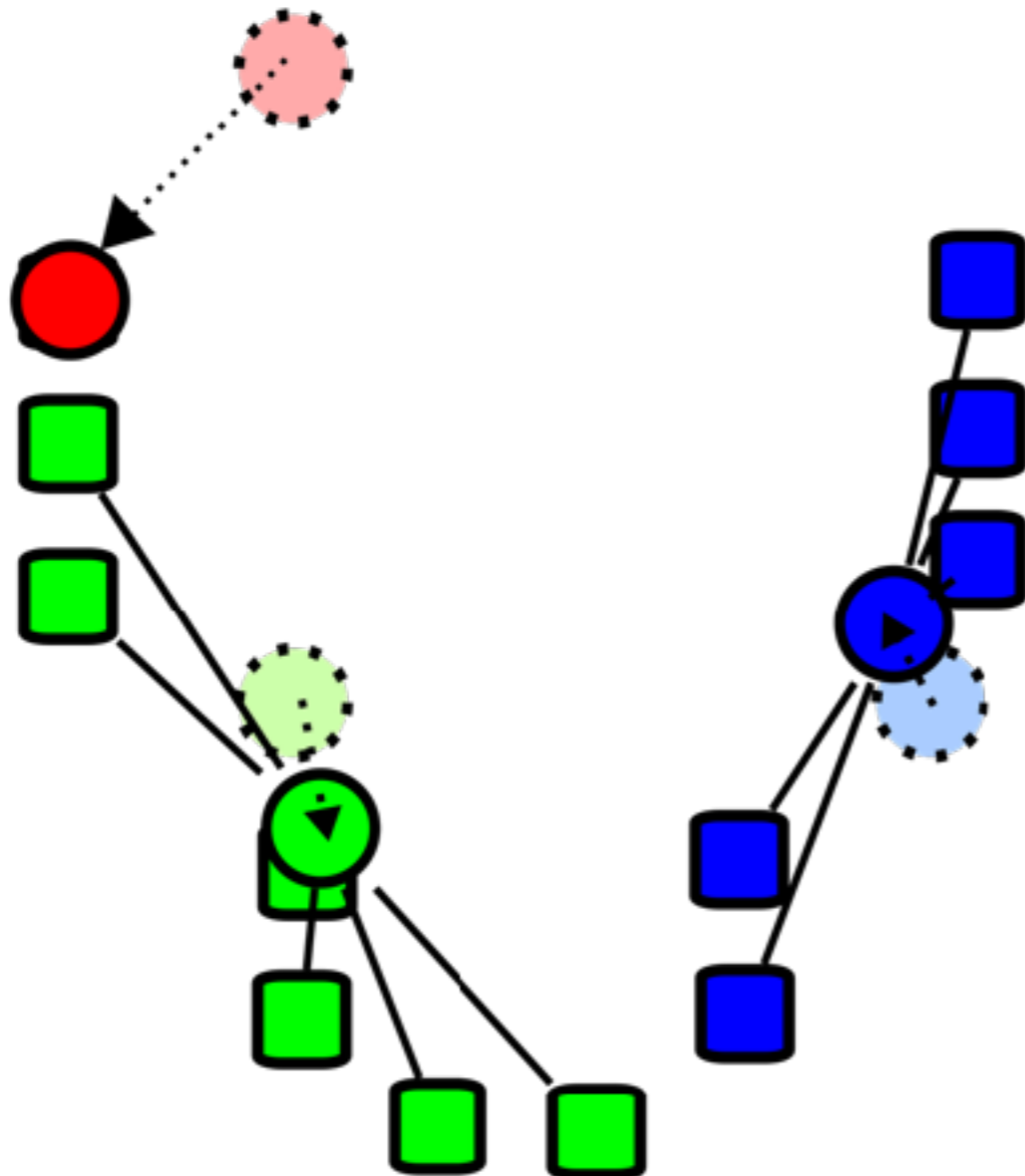


# k-means

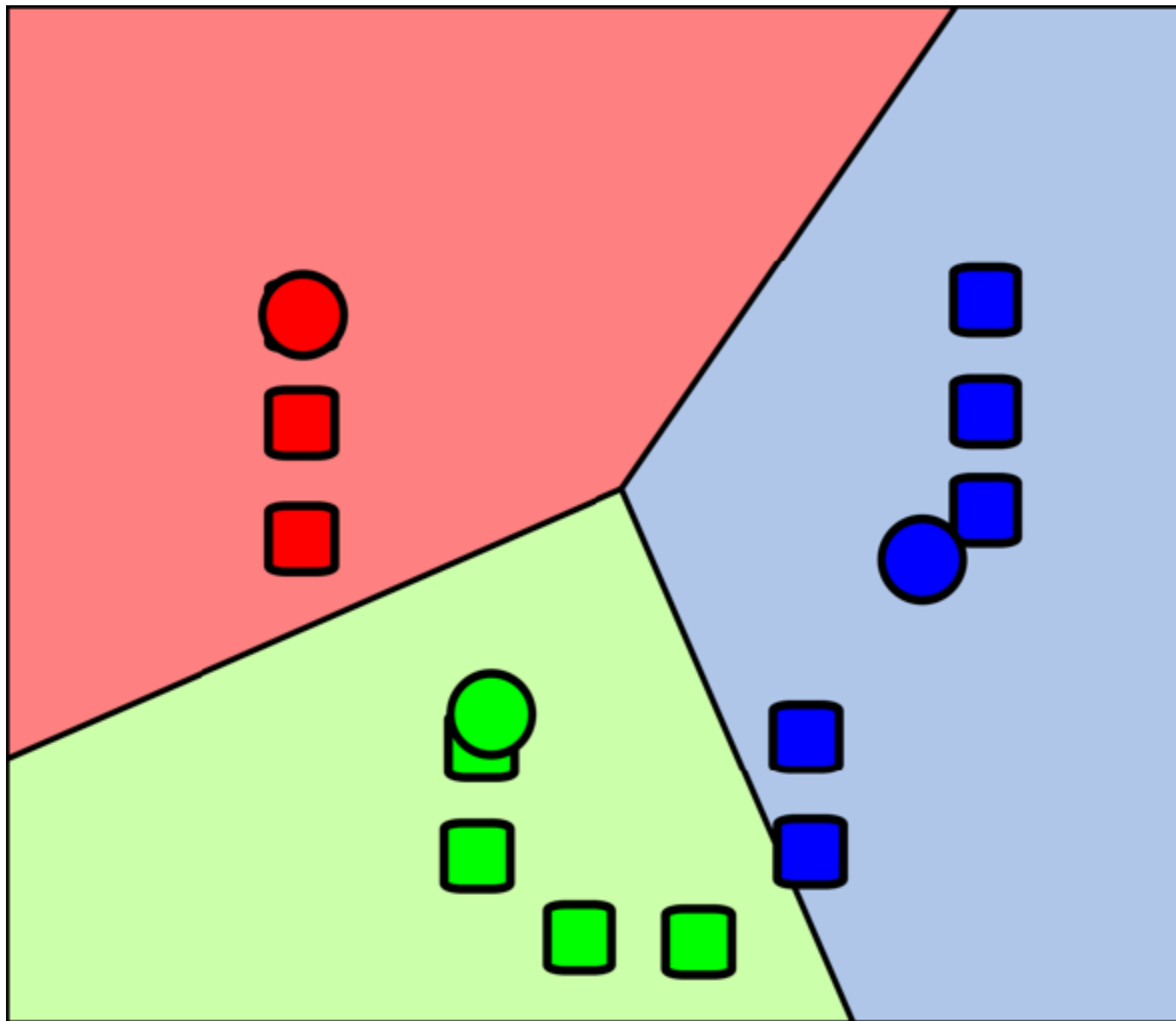




# k-means



# k-means





## Finn Family Moomintroll, by [Tove Jansson](#)

date saved: 2 Oct 2020

tagged with: [fiction](#) [moomins](#) [format:epub](#)

```
>>> from sklearn.cluster import KMeans
```

```
>>> colors = im.getdata()
```

```
>>> KMeans(n_clusters=5).fit(colors).cluster_centers_
```

```
array([[129.4208, 196.3925, 195.2154],  
       [ 32.2470,  98.7502,  73.7802],  
       [236.5401, 170.0094,  60.9854],  
       [ 25.9697, 179.8435, 200.6353],  
       [229.4518, 236.1130, 227.4645]])
```





## Finn Family Moomintroll, by [Tove Jansson](#)

date saved: 2 Oct 2020

tagged with: [fiction](#) [moomins](#) [format:epub](#)

```
>>> import wcag_contrast_ratio as contrast
```

```
>>> for c in candidates:
```

```
...     print(c, contrast.rgb(c, white))
```

1.9778

7.2296

2.0228

2.2523

1.2042

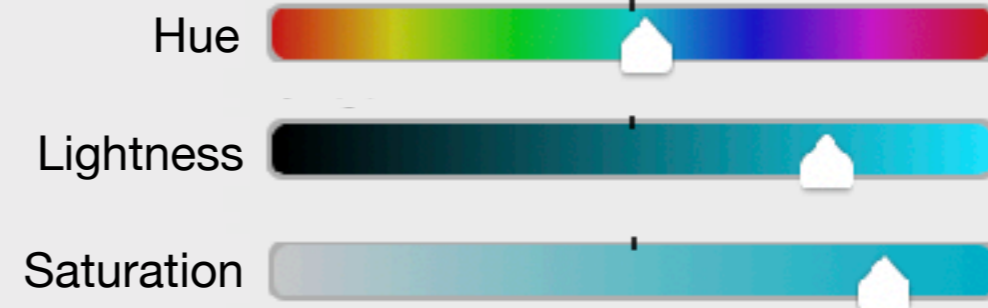
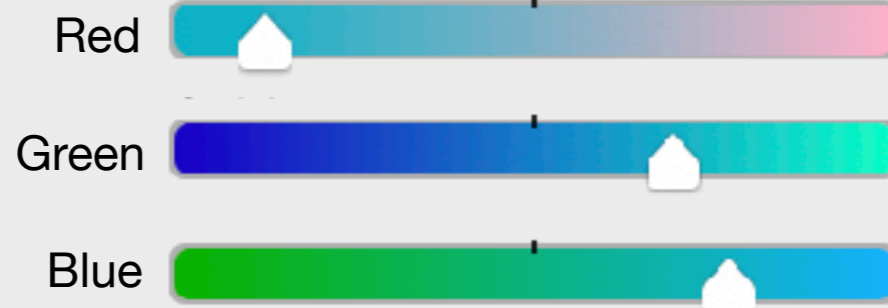




## Finn Family Moomintroll, by [Tove Jansson](#)

date saved: 2 Oct 2020

tagged with: [fiction](#) [moomins](#) [format:epub](#)



```
>>> import colorsys
```

```
>>> max(  
...     rgb_colors,  
...     key=lambda col: colorsys.rgb_to_hls(*col)[2]  
... )
```

# 3/

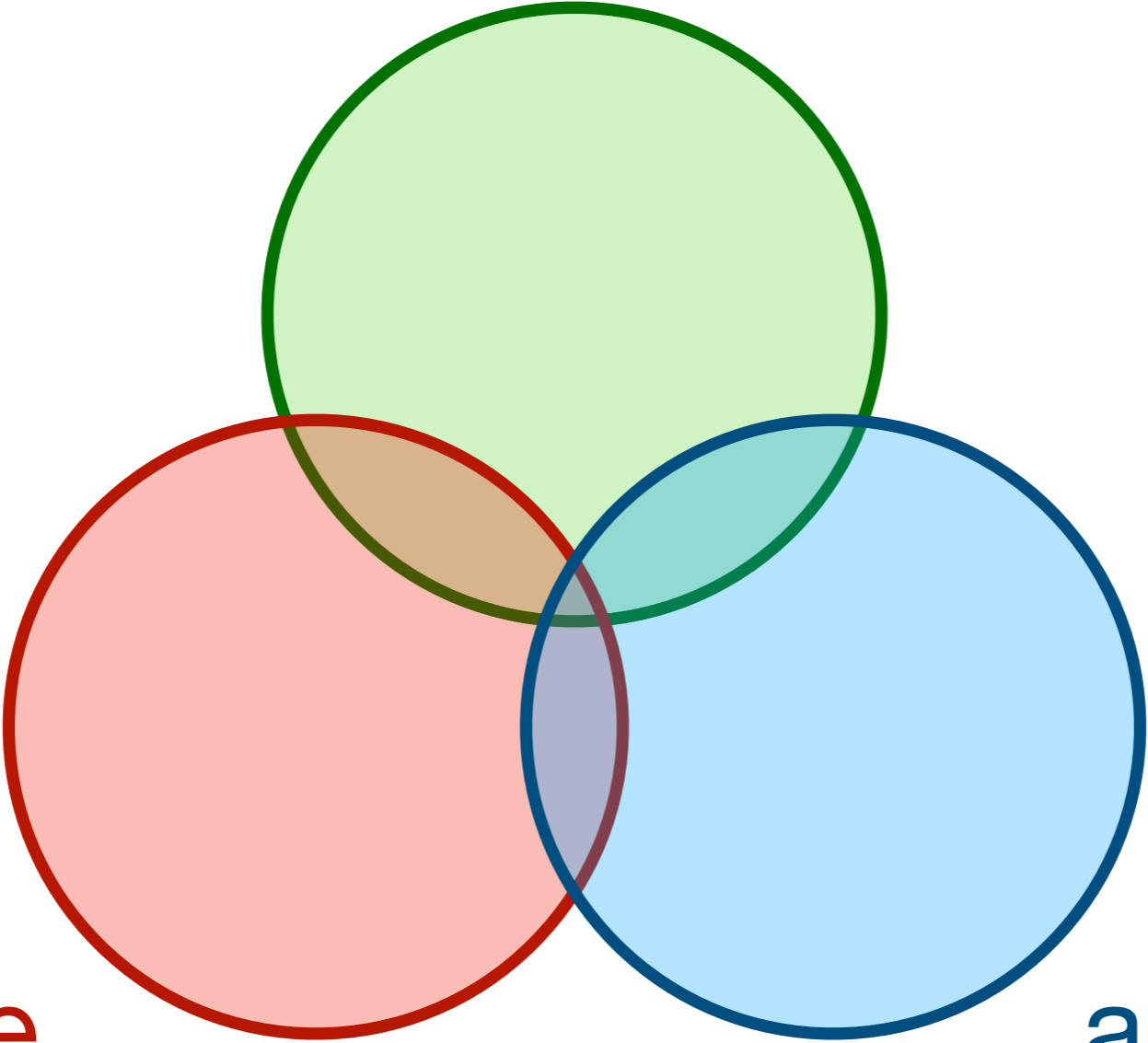
## Previewing the files:

- Use `pdftocairo` to generate thumbnails of PDFs
- Use Quick Look to generate thumbnails of still images
- Create small movies for animated GIFs
- Use Pillow, `k-means` (`scikit-learn`) and WCAG accessibility to choose tint colours

4/

Wrangling the tags

utility  
bills



home

acme  
energy



showing documents 1-54 of 54.  
▶ tag list  
▶ tag cloud



### For All Mankind Lunar Surface Access Module (LSAM) LEGO instructions, by [SpaceXplorer](#)

date saved: 6 Jun 2021  
source: [rebrickable.com](#)  
tagged with: for-all-mankind lego sci-fi



### Journey into our Underworld

date saved: 17 Apr 2021  
source: [www.docsgalore.com](#)  
tagged with: [archives](#)



### Volcano pattern, by [CrossStitchMari](#)

date saved: 26 Feb 2021  
source: [www.fastmail.com](#)  
tagged with: cross-stitch-patterns



### JPEG 2000 spec

date saved: 4 Oct 2019  
tagged with: [file-format-specs](#) [welcome](#)



### Papers from "Making IT Work" conference

date saved: 16 Aug 2016

showing documents 1-54 of 54.

▶ tag list

▼ tag cloud

archives aws by:CrossStitchMari by:Happinest by:James Clear by:Mathysphere by:SpaceXplorer by:The Duckbill Group by:ZAnnaCrossStitch chinese  
 community-safety conferences **cross-stitch-patterns** elasticsearch esu family file-format-specs for-all-mankind  
 from:Gendered Intelligence gsk languages lego **lgbtq-network** music org:mhfa-england org:open-university  
 org:prospect-union personal pride-planets recipes rpg sci-fi self-improvement sketchnotes **stonewall** talks **trans-awareness-training**  
**trans-inclusion-policies** wellcome write-the-docs



### For All Mankind Lunar Surface Access Module (LSAM) LEGO Instructions, by SpaceXplorer

date saved: 6 Jun 2021  
source: [rebrickable.com](http://rebrickable.com)  
tagged with: [for-all-mankind](#) [lego](#) [sci-fi](#)



### Journey Into our Underworld

date saved: 17 Apr 2021  
source: [www.deepstore.com](http://www.deepstore.com)  
tagged with: [archives](#)



### Volcano pattern, by CrossStitchMari

date saved: 26 Feb 2021  
source: [www.fastmail.com](http://www.fastmail.com)  
tagged with: [cross-stitch-patterns](#)



### JPEG 2000 spec

date saved: 4 Oct 2019  
tagged with: [file-format-specs](#) [wellcome](#)

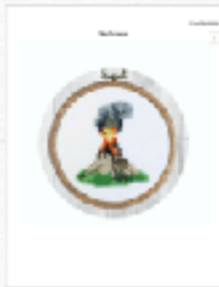
showing documents 1-10 of 10.

---

by: [CrossStitchMari](#) by: [Happines1](#) by: [Mathysphere](#) by: [ZAnnaCrossStitch](#) **cross-stitch-patterns** [music](#) [pride-planets](#)

---

filtering to tag **cross-stitch-patterns** [x]



**Volcano pattern, by [CrossStitchMari](#)**

date saved: 26 Feb 2021  
 source: [www.fastmail.com](http://www.fastmail.com)  
 tagged with: cross-stitch-patterns



**Math Sampler pattern, by [Mathysphere](#)**

date saved: 2 Jan 2020  
 source: [www.fastmail.com](http://www.fastmail.com)  
 tagged with: cross-stitch-patterns



**Solar Eclipse cross-stitch pattern, by [Mathysphere](#)**

date saved: 2 Jan 2020  
 source: [www.fastmail.com](http://www.fastmail.com)  
 tagged with: cross-stitch-patterns



**Bi Planet pattern, by [Mathysphere](#)**

showing documents 1-6 of 6.

by:Mathysphere cross-stitch-patterns [pride-planets](#)

filtering to tags **cross-stitch-patterns** [x] **by:Mathysphere** [x]



**Math Sampler pattern, by Mathysphere**

date saved: 2 Jan 2020  
source: [www.fastmail.com](http://www.fastmail.com)  
tagged with: cross-stitch-patterns



**Solar Eclipse cross-stitch pattern, by Mathysphere**

date saved: 2 Jan 2020  
source: [www.fastmail.com](http://www.fastmail.com)  
tagged with: cross-stitch-patterns



**Trans Pride Planet pattern, by Mathysphere**

date saved: 19 Apr 2019  
tagged with: cross-stitch-patterns [pride-planets](#)



**Bi Planet pattern, by Mathysphere**

# Wrangling the tags:

- How do I find documents that match a set of tags?
- How do I get URLs to filter by tag?
- How do I find mistyped or similar tags?

# find matching documents

```
>>> class Document:
...     ...
...     def matches_tags(self, query):
...         return query.issubset(self.tags)

>>> matching_documents = [
...     doc
...     for doc in documents
...     if doc.matches_tags(query)
... ]
```

# URLs to tag filters

`docs.python.org/3/search.html?q=subset`

`/files?tag=books`

`wellcomecollection.org/images?`

`query=snake&locations.license=pdm&color=e02020`

`/files?tag=books&tag=programming`

# hyperlink

```
>>> import hyperlink
>>> url = hyperlink.URL.from_text("/files")

>>> url = url.add("tag", "books")
"/files?tag=books"

>>> url = url.add("tag", "programming")
"/files?tag=books&tag=programming"

>>> url = url.remove("tag", "books")
"/files?tag=programming"
```



council-tax

counciltax

concihtax



*Suspiciously similar*

car-tax

cambridge-council

train-tickets



*Definitely different*

# rapidfuzz

```
>>> from rapidfuzz import fuzz
```

```
>>> fuzz.ratio("council-tax", "counciltax")  
95.23809523809524
```

```
>>> fuzz.ratio("council-tax", "conciiltax")  
90.0
```

```
>>> fuzz.ratio("council-tax", "car-tax")  
55.55555555555556
```

```
>>> fuzz.ratio("council-tax", "train-tickets")  
25.0
```

https://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/



# rapidfuzz

```
from itertools import combinations

from rapidfuzz import fuzz

def find_similar_tags(tags, *, threshold=80):
    for t1, t2 in combinations(sorted(tags), 2):
        if fuzz.ratio(t1, t2) > threshold:
            yield t1, t2
```

# Wrangling the tags:

- How do I find documents that match a set of tags?  
**List comprehensions are fine for small N.**
- How do I get URLs to filter by tag?  
**Use the hyperlink library.**
- How do I find mistyped or similar tags?  
**Use `fuzzywuzzy` or `rapidfuzz`.**

```
docstore add 'BillJul21.pdf' \  
  --source_url='https://acme.co/gas' \  
  --title='2021-07: Gas bill for 2021' \  
  --tags='home,utility-bills,acme-energy'
```

showing documents 1-54 of 54.  
▶ tag list  
▶ tag cloud



### For All Mankind Lunar Surface Access Module (LSAM) LEGO instructions, by [SpaceXplorer](#)

date saved: 6 Jun 2021  
source: [rebrickable.com](#)  
tagged with: for-all-mankind lego sci-fi



### Journey into our Underworld

date saved: 17 Apr 2021  
source: [www.docsplore.com](#)  
tagged with: [archives](#)



### Volcano pattern, by [CrossStitchMari](#)

date saved: 26 Feb 2021  
source: [www.fastmail.com](#)  
tagged with: cross-stitch-patterns



### JPEG 2000 spec

date saved: 4 Oct 2019  
tagged with: [file-format-specs](#) [welcome](#)



### Papers from "Making IT Work" conference

date saved: 16 Aug 2016

showing documents 1-54 of 54.

▶ tag list

▼ tag cloud

archives aws by:CrossStitchMari by:Happinest by:James Clear by:Mathysphere by:SpaceXplorer by:The Duckbill Group by:ZAnnaCrossStitch chinese  
 community-safety conferences **cross-stitch-patterns** elasticsearch esu family file-format-specs for-all-mankind  
 from:Gendered Intelligence gsk languages lego **lgbtq-network** music org:mhfa-england org:open-university  
 org:prospect-union personal pride-planets recipes rpg sci-fi self-improvement sketchnotes **stonewall** talks **trans-awareness-training**  
**trans-inclusion-policies** wellcome write-the-docs



### For All Mankind Lunar Surface Access Module (LSAM) LEGO Instructions, by SpaceXplorer

date saved: 6 Jun 2021  
source: [rebrickable.com](https://rebrickable.com)  
tagged with: [for-all-mankind](#) [lego](#) [sci-fi](#)



### Journey Into our Underworld

date saved: 17 Apr 2021  
source: [www.deepstore.com](http://www.deepstore.com)  
tagged with: [archives](#)



### Volcano pattern, by CrossStitchMari

date saved: 26 Feb 2021  
source: [www.fastmail.com](http://www.fastmail.com)  
tagged with: [cross-stitch-patterns](#)



### JPEG 2000 spec

date saved: 4 Oct 2019  
tagged with: [file-format-specs](#) [wellcome](#)



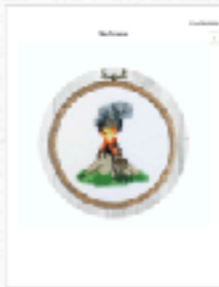
showing documents 1-10 of 10.

---

by: [CrossStitchMari](#) by: [Happines1](#) by: [Mathysphere](#) by: [ZAnnaCrossStitch](#) **cross-stitch-patterns** [music](#) [pride-planets](#)

---

filtering to tag **cross-stitch-patterns** [x]



**Volcano pattern, by [CrossStitchMari](#)**

date saved: 26 Feb 2021  
 source: [www.fastmail.com](http://www.fastmail.com)  
 tagged with: cross-stitch-patterns



**Math Sampler pattern, by [Mathysphere](#)**

date saved: 2 Jan 2020  
 source: [www.fastmail.com](http://www.fastmail.com)  
 tagged with: cross-stitch-patterns



**Solar Eclipse cross-stitch pattern, by [Mathysphere](#)**

date saved: 2 Jan 2020  
 source: [www.fastmail.com](http://www.fastmail.com)  
 tagged with: cross-stitch-patterns



**Bi Planet pattern, by [Mathysphere](#)**

showing documents 1-6 of 6.

by:Mathysphere cross-stitch-patterns [pride-planets](#)

filtering to tags **cross-stitch-patterns** [x] **by:Mathysphere** [x]



**Math Sampler pattern, by Mathysphere**

date saved: 2 Jan 2020  
source: [www.fastmail.com](http://www.fastmail.com)  
tagged with: cross-stitch-patterns



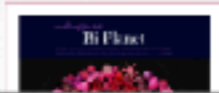
**Solar Eclipse cross-stitch pattern, by Mathysphere**

date saved: 2 Jan 2020  
source: [www.fastmail.com](http://www.fastmail.com)  
tagged with: cross-stitch-patterns

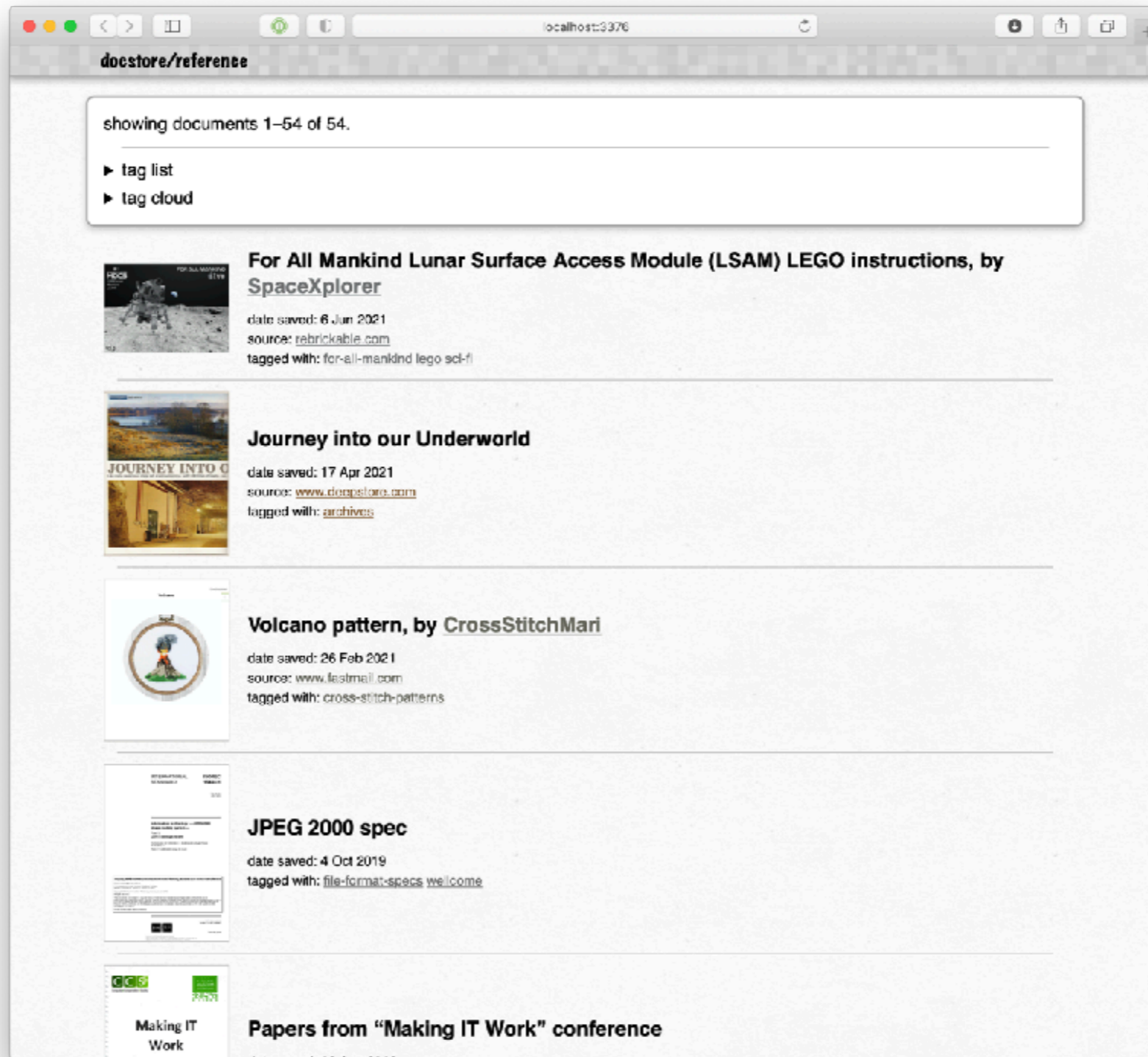


**Trans Pride Planet pattern, by Mathysphere**

date saved: 19 Apr 2019  
tagged with: cross-stitch-patterns [pride-planets](#)



**Bi Planet pattern, by Mathysphere**



alexwlchan/docstore





Storing the files

Storing the metadata

Previewing the files

Wrangling the tags

Takeaways:

# Takeaways:

- You're all going to go away and download docstore!!



# Takeaways:

- You're all going to go away and write your own file management software!!

# Takeaways:

- Writing tools that are just for you is easier than writing tools for everyone
- Writing tools that are just for you can have interesting lessons and ideas for writing “real” software
- Source code is a terrible way to share ideas

# Using Python to organise my physical paper

Alex Chan · [alexwlchan.net](http://alexwlchan.net) · [they/she](https://they/she)

Slides: [alexwlchan.net/files/2021/campug.pdf](http://alexwlchan.net/files/2021/campug.pdf)