

Storage for MMIR

Thomas Breuel

UniKL

covering

- ▶ image formats
- ▶ bandwidth / CPU tradeoffs
- ▶ file system
- ▶ relational databases
- ▶ scientific databases
- ▶ common schemas

Image Formats

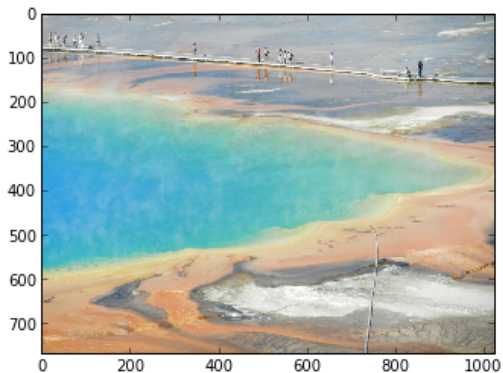
image formats

- ▶ raw (various): $w \times h$, 16 bpp, image as it comes out of the camera
- ▶ uncompressed (PPM): $w \times h \times 3$, usually 8 bpp (sometimes 16 bpp)
- ▶ lossless compression (PNG)
- ▶ lossy compression (JPEG, JPEG2000)

important tools

- ▶ ImageMagick, GNU Magic (convert, etc.)
- ▶ SciPy
- ▶ PIL (Python imaging library)
- ▶ Python exiv2 (EXIF metadata in Python)

```
1 iname = "raw_images/DSCN0243.JPG"  
2 image = imread(iname)  
3 imshow(image)
```



getting metadata

```
1 import pyexiv2
2 metadata = pyexiv2.ImageMetadata(iname)
3 metadata.read()
4 for key in metadata.exif_keys[::10]:
5     print key, "\t", metadata[key]
```

```
Exif.Image.ImageDescription <Exif.Image.ImageDescription [Ascii] =
Exif.Image.ExifTag <Exif.Image.ExifTag [Long] = 252>
Exif.Photo.CompressedBitsPerPixel <Exif.Photo.CompressedBitsPerPixel [Rational] =
Exif.Nikon3.Version <Exif.Nikon3.Version [Undefined] = 0 2 0 0>
Exif.Nikon3.ISOSelection <Exif.Nikon3.ISOSelection [Ascii] = AUTO >
...
Exif.Nikon3.ShotInfo <Exif.Nikon3.ShotInfo [Undefined] = (Binary value suppressed)
Exif.Nikon3.0x00b5 <Exif.Nikon3.0x00b5 [Short] = 4098>
Exif.Photo.SceneType <Exif.Photo.SceneType [Undefined] = 1>
Exif.Photo.Sharpness <Exif.Photo.Sharpness [Short] = 0>
```

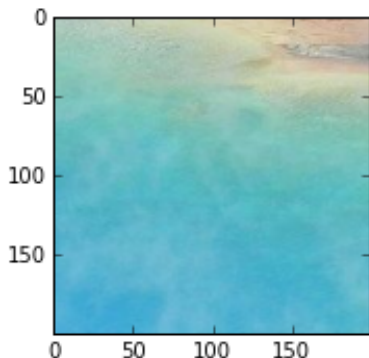
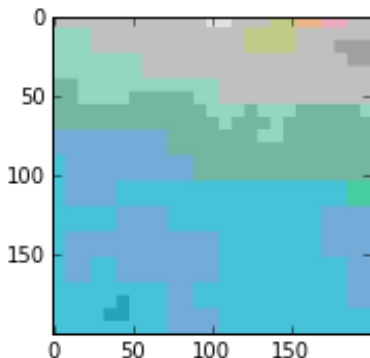
image conversion and size

```
1 !convert raw_images/DSCN0243.JPG DSCN0243.ppm
2 !convert DSCN0243.ppm DSCN0243.png
3 !ls -lh raw_images/DSCN0243.JPG DSCN0243.*
4 !rm DSCN0243.ppm DSCN0243.png
```

```
-rw-r--r-- 1 tmb tmb 916K Oct 31 04:17 DSCN0243.png
-rw-r--r-- 1 tmb tmb 2.3M Oct 31 04:17 DSCN0243.ppm
-rw-r--r-- 1 tmb tmb 201K Oct 30 14:33 raw_images/DSCN0243.JPG
```


JPEG quality settings

```
1 from PIL import Image
2 im = Image.open(iname)
3 im.save("temp.jpg", quality=1)
4 subplot(121); imshow(imread("temp.jpg")[200:400,200:400])
5 im.save("temp.jpg", quality=80)
6 subplot(122); imshow(imread("temp.jpg")[200:400,200:400])
7
```



Remember JPEG Compression

- ▶ divide image into 8×8 squares
- ▶ perform DCT for each square
- ▶ quantize and select coefficients
- ▶ linearize through scanning
- ▶ compress

in-memory compression requires streams

```
1
2 import StringIO
3 stream = StringIO.StringIO()
4 im.save(stream, 'jpeg')
5 compressed = stream.getvalue()
6 print len(compressed)
```

84770

useful in-memory compression functions

```
1
2 import PIL,StringIO
3
4 def cjpeg(image,quality=98):
5     im = PIL.Image.fromarray(image)
6     stream = StringIO.StringIO()
7     im.save(stream,'jpeg',quality=quality)
8     compressed = stream.getvalue()
9     return compressed
10
11 def djpeg(compressed):
12     stream = StringIO.StringIO(compressed)
13     im = PIL.Image.open(stream)
14     return numpy.array(im)
```

increasing performance

- ▶ compression reduces I/O (less data to move)
- ▶ compression has some overhead
- ▶ overall, storing compressed images is faster
- ▶ devices already use compression for other data types
- ▶ special image compression algorithms are more effective

testing in-memory compression

```
1  
2 compressed = cjpeg(imread(iname))  
3 decompressed = djpeg(compressed)  
4 imshow(decompressed)  
5 print prod(decompressed.shape), len(compressed)
```

2359296 275947



Icons

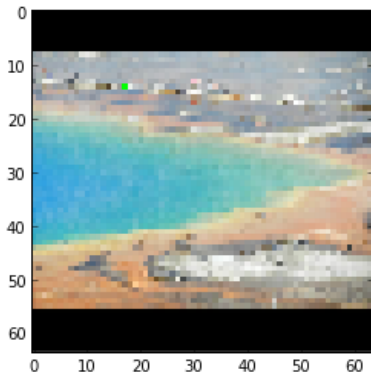
Both for display and for analysis, it's useful to store *icons* of the images being indexed.

- ▶ fixed size and centered (square)
- ▶ may be stored uncompressed (if smaller than the size where compression makes sense)
- ▶ 8 bpp


```
1 from scipy.ndimage import interpolation
2
3 def make_icon(image, size=(256,256)):
4     if amax(image)>1.0: image = image/255.0
5     assert amin(image)>=0.0 and amax(image)<=1.0
6     h,w = size
7     ih,iw,d = image.shape
8     scale = min(h*1.0/ih,w*1.0/iw)
9     scaled = interpolation.zoom(image,(scale,scale,1))
10    sh,sw,_ = scaled.shape
11    result = zeros((h,w,3))
12    dh,dw = h//2-sh//2,w//2-sw//2
13    result[dh:dh+sh,dw:dw+sw,:] = scaled
14    return array(result*255, 'B')
```

```
1 image = imread(iname)
2 icon = make_icon(image,(64,64))
3 print icon.shape,icon.dtype
4 imshow(icon,interpolation='none')
```

(64, 64, 3) uint8



Checksums and Duplicate Detection

```
1 import md5
2
3 def md5file(fname):
4     m = md5.new()
5     with open(fname,"rb") as stream:
6         while 1:
7             s = stream.read(4096)
8             if s=="": break
9             m.update(s)
10    return m.hexdigest()
```

```
1 md5file(iname)
```

```
'5640ebaac01fc22241874189d85dff9a'
```

Relational Databases

opening the database, creating the table

```
1
2 import sqlite3
3 !rm -f sample.db
4 dbname = "sample.db"
5 db = sqlite3.connect(dbname)
6 c = db.cursor()
7 c.execute("create table images (chk_text_primary_key, path_text, w_
           integer, h_integer, jpeg_blob, icon_blob, time_date)")
```

```
<sqlite3.Cursor at 0x5aecf10>
```

what we are storing

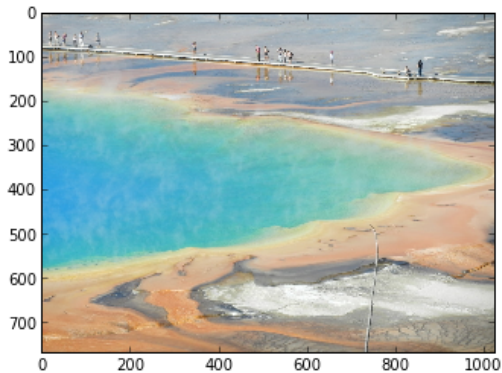
```
1 chk = md5file(iname)      # MD5 checksum of the file
2 image = imread(iname)    # JPEG-compressed image
3 jpeg = cjpeg(image)
4 h,w,d = image.shape      # dimensions of the image
5 icon = make_icon(image)  # an icon (uncompressed)
6 when = datetime.datetime.now()
```


inserting the data into the database

```
1 # note the use of 'buffer' for blobs
2 values = (chk, iname, w, h, buffer(jpeg), buffer(icon), when)
3 c.execute("insert into images values (?, ?, ?, ?, ?, ?, ?)", values)
4 db.commit()
5
```

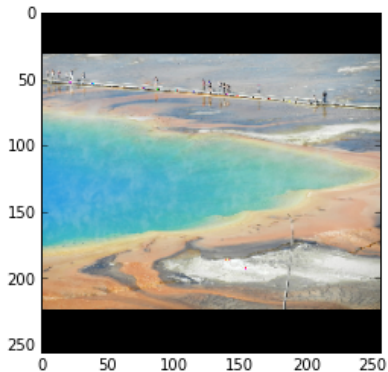
retrieving a jpeg image

```
1 result = list(c.execute("select jpeg from images limit 1"))  
2  
3 imshow(djpeg(result[0][0]))
```



retrieving an icon

```
1 result = list(c.execute("select icon from images limit 1"))  
2  
3 imshow(array(result[0][0], 'B').reshape(256,256,3))
```



getting metadata

```
1 import pyexiv2
2 from datetime import datetime
3
4 def image2db(iname,db,c):
5     chk = md5file(iname)      # MD5 checksum of the file
6     if list(c.execute("select _chk_ from _images_ where _chk=?", (chk,)))
7         != []:
8         c.execute("update _images_ set _path=?_ where _chk=?", (iname, chk
9             ))
10
11         db.commit()
12         return False
13
14 # get metadata
15 metadata = pyexiv2.ImageMetadata(iname)
16 metadata.read()
17 if "Exif.Photo.DateTimeOriginal" in metadata:
18     when = metadata["Exif.Photo.DateTimeOriginal"].value
19 else:
20     when = datetime.now()
21
22 # get content data
23 image = imread(iname)      # JPEG-compressed image
24 jpeg = cjpeg(image)
25 h,w,d = image.shape        # dimensions of the image
```

clear out the database

```
1 c.execute("delete from images")  
2  
3 db.commit()
```

```
1 import glob
2
3 for i,iname in enumerate(sorted(glob.glob("raw_images/*.JPG"))):
4     if i%10==0: print i,iname
5     if not image2db(iname,db,c):
6         print iname,"already in the database"
```

```
0 raw_images/DSCN0149.JPG
10 raw_images/DSCN0204.JPG
20 raw_images/DSCN0254.JPG
30 raw_images/DSCN0286.JPG
40 raw_images/DSCN0320.JPG
...
170 raw_images/P1010409.JPG
180 raw_images/P1010458.JPG
190 raw_images/P1010496.JPG
200 raw_images/P1010536.JPG
```

make sure the checksum-based duplicate detection is actually working

```
1
2 for i,iname in enumerate(sorted(glob.glob("raw_images/*.JPG")[:5]))
3     :
4     if i%10==0: print i,iname
5     if not image2db(iname,db,c):
        print iname,"already in the database"
```

```
0 raw_images/DSCN0225.JPG
raw_images/DSCN0225.JPG already in the database
raw_images/DSCN0557.JPG already in the database
raw_images/DSCN0642.JPG already in the database
raw_images/DSCN0649.JPG already in the database
raw_images/P1010456.JPG already in the database
```


cleaning up and size comparisons

```
1 db.close()
2
3 !du -h raw_images
4 !ls -lh sample.db
```

```
49M raw_images
-rw-r--r-- 1 tmb tmb 106M Oct 31 04:19 sample.db
```

Scientific Databases

HDF5

HDF5 is a database for large scientific datasets in array format.

- ▶ no transactions
- ▶ limited concurrent access
- ▶ arrays as primitives
- ▶ very fast random access
- ▶ special compression routines for numerical data, floating point
- ▶ similar to “column databases” and NoSQL databases
- ▶ integrates directly with NumPy/SciPy (Python)
- ▶ supported in a wide range of languages, including C++ and Fortran

```
1 import tables
2 from tables import openFile, Filters, Int32Atom, Float32Atom, Int64Atom
   , UInt8Atom, VLStringAtom
```

Persistent Arrays

HDF5 effectively implements *persistent arrays* or *memory mapped arrays* (the actual implementation may or may not use memory mapping, but it looks that way from Python):

```
1 hdf = tables.openFile("temp.h5", "w")
2 hdf.createEArray("/", "icons", UInt8Atom(), shape=[0, 256, 256, 3])
3 hdf.root.icons.append([icon])
4 hdf.close()
5
6 hdf = tables.openFile("temp.h5", "r+")
7 print len(hdf.root.icons)
8 imshow(hdf.root.icons[0])
9 hdf.close()
```

1



VLArrays, Strings

You can have variable length arrays of variable length strings. Strings are special and a bit different from numerical arrays. You can also use VLArrays with 1D numerical arrays, giving a ragged array. If you want to store variable sized 2D images, you need to store the dimensions in a separate column and reshape on access.

```
1 hdf = tables.openFile("sample6.h5","w")
2 hdf.createVLArray("/", "chk", VLStringAtom())
3 hdf.root.chk.append(u"hello")
4 hdf.root.chk.append("world")
5 hdf.root.chk[0:2]
6 hdf.close()
```

store all the icons from the image database in HDF5

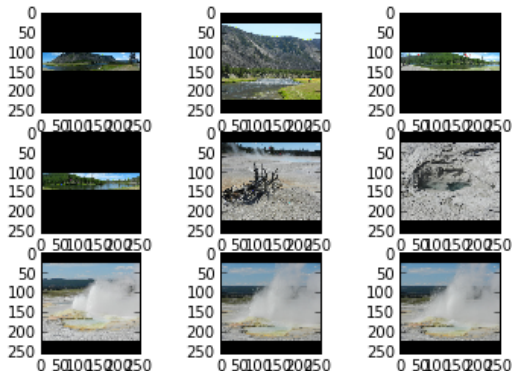
```
1 hdf = tables.openFile("sample.h5","w")
2 hdf.createVLArray("/", "chk", VLStringAtom())
3 hdf.createEArray("/", "icons", UInt8Atom(), shape=[0,256,256,3],
4     filters=tables.Filters(9))
5
6 db = sqlite3.connect(dbname)
7 c = db.cursor()
8 rows = c.execute("select chk, icon from images")
9 for row in rows:
10     chk = str(row[0])
11     icon = array(row[1], 'B').reshape(256,256,3)
12     hdf.root.chk.append(chk)
13     hdf.root.icons.append([icon])
14
15 db.close()
16 hdf.close()
```

```
1 !ls -lh sample.h5
2 !h5ls sample.h5
```

```
-rw-r--r-- 1 tmb tmb 20M Oct 31 04:20 sample.h5
chk          Dataset {207/Inf}
icons        Dataset {207/Inf, 256, 256, 3}
```

quickly accessing HDF5 content

```
1 with tables.openFile("sample.h5","r") as hdf:  
2     for i in range(9):  
3         subplot(3,3,i+1)  
4         imshow(hdf.root.icons[i])
```



Summary - Storage

Storage

- ▶ images can be stored compressed or raw
- ▶ compressed storage increases speed because CPU is cheaper than I/O
- ▶ metadata in EXIF format for images
- ▶ storage on file system, in relational databases, in scientific databases
- ▶ relational databases offer transactions
- ▶ scientific database are optimized for large array storage
- ▶ in-memory compression needed and usually available via “string streams”