



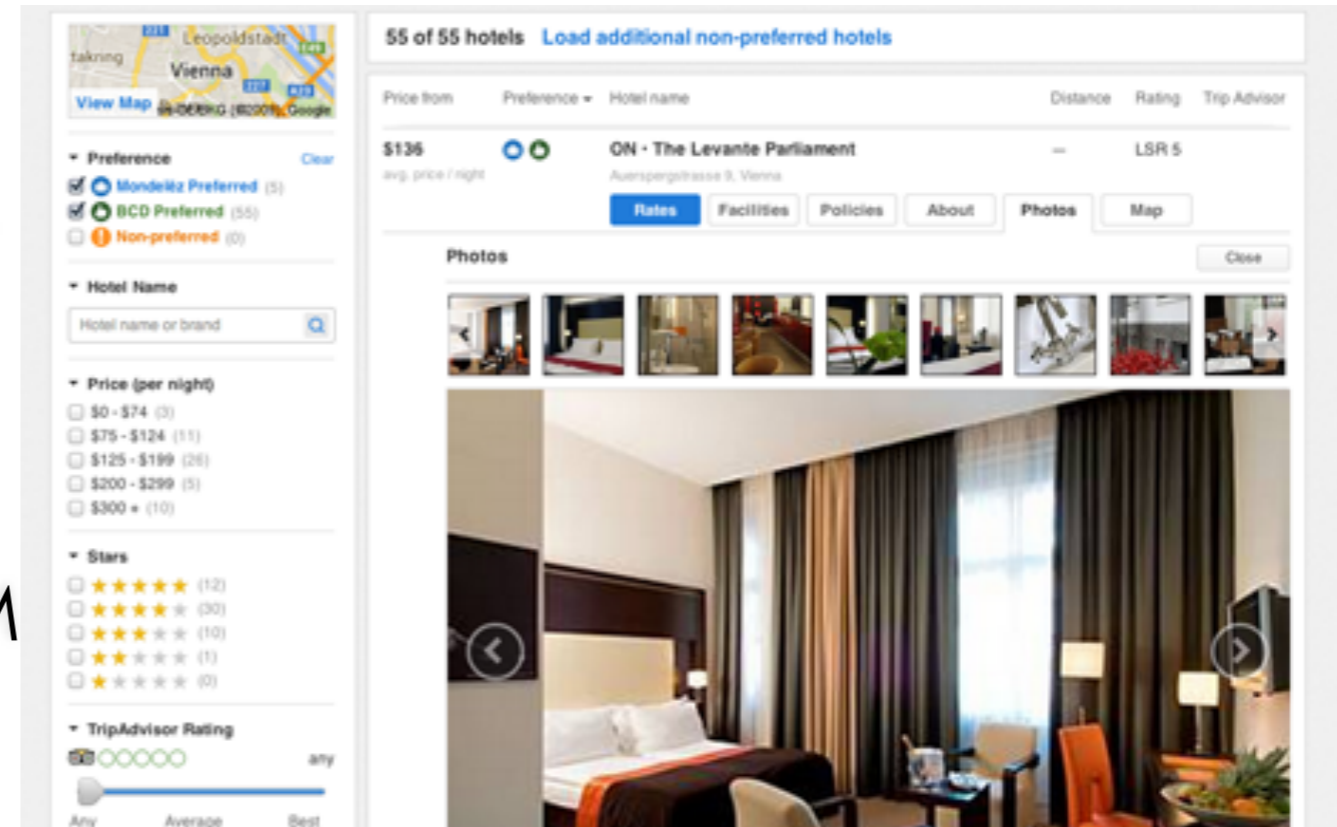
# HOW TO DETECT SIMILAR IMAGES WITH PYTHON

VIKTOR DRACHOV

# ABOUT ME

GETGOING

- [GETGOING.COM](http://GETGOING.COM)
  - B2B, CORPORATE TRAVEL
  - Крупнейший заказчик – BCD TRAVEL
  - 4 провайдера
  - > 100к отелей в каждом
  - ~7М изображений
  - Не хотим показывать одно и то же



# DISCLAIMER

EVERYBODY LIES

- РАССКАЖУ
  - ОБЩИЕ ПОДХОДЫ
  - КАК ЭТО ДЕЛАЛИ МЫ
  - КАК ЭТО НЕ СТОИТ ДЕЛАТЬ
- НЕ РАССКАЖУ
  - РАСПОЗНАВАНИЕ ОБЪЕКТОВ
  - ПОИСК “ПОХОЖИХ” КАРТИНОК
  - SILVER BULLET



# ПОДХОДЫ К СРАВНЕНИЮ ИЗОБРАЖЕНИЙ

- КАЖДЫЙ С КАЖДЫМ
  - ДОЛГО, НЕУДОБНО,  $O(N^2)$
- ХЭШИРОВАНИЕ. PERCEPTUAL HASHING.
  - НЕ ОБЯЗАТЕЛЬНО СРАВНИВАТЬ N С N
  - ЗАНИМАЕТ МЕНЬШЕ МЕСТА
  - МОЖНО ИСКАТЬ (KD-TREE, MVP-TREE, ETC).  $O(N)$  В **ХУДШЕМ** СЛУЧАЕ.
    - Или даже в PostgreSQL
  - По сути мы строим features для классификатора

# OPEN SOURCE

- IMAGEHASH
- PHASH - [HTTP://PHASH.ORG](http://PHASH.ORG)
  - [HTTP://PYPI.PYTHON.ORG/PYPI/PHASH/](http://PYPI.PYTHON.ORG/PYPI/PHASH/)
  - [HTTPS://GITHUB.COM/POLACHOK/PY-PHASH](https://GITHUB.COM/POLACHOK/PY-PHASH)
- IMGSEEK
  - [HTTP://WWW.IMGSEEK.NET/ISK-DAEMON](http://WWW.IMGSEEK.NET/ISK-DAEMON)
- LUCENE
  - [HTTP://WWW.LIRE-PROJECT.NET/](http://WWW.LIRE-PROJECT.NET/)

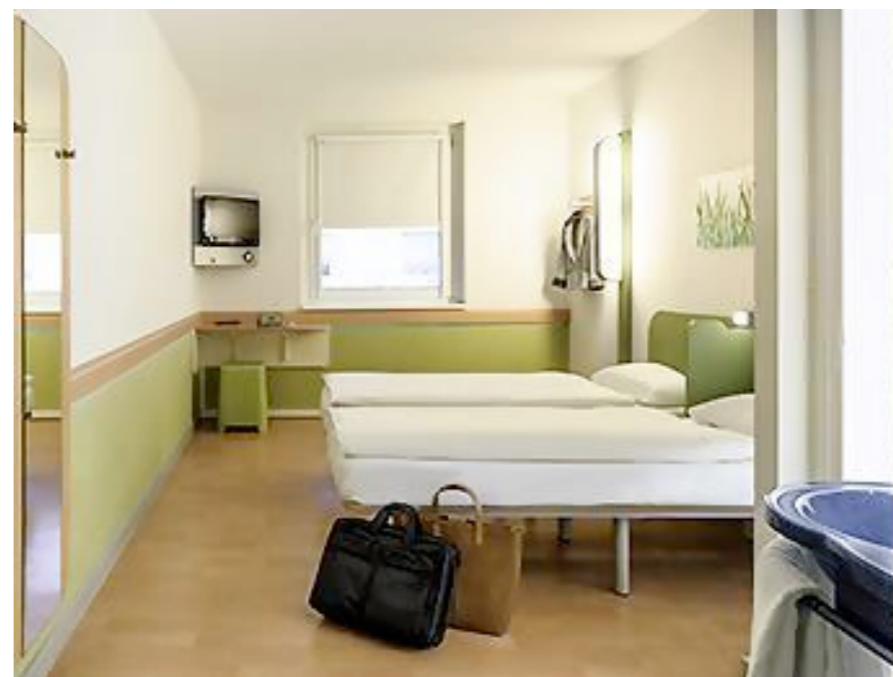
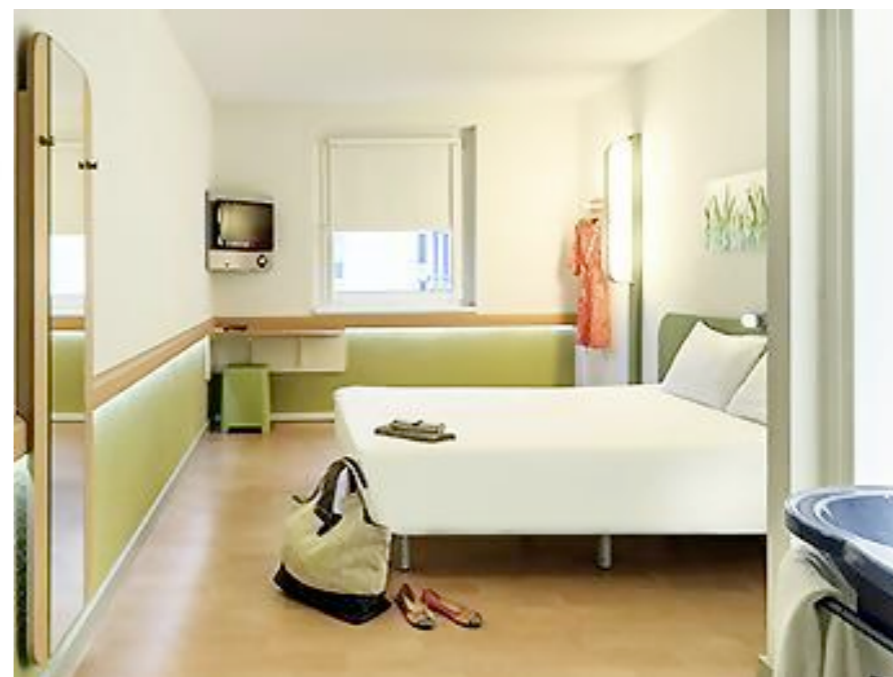
# КАК НАМ ПОРТЯТ ЖИЗНЬ

- ИЗМЕНЕНИЕ ОДНОГО И ТОГО ЖЕ ОРИГИНАЛА
  - РЕСАЙЗ С ИЗМЕНЕНИЕМ ASPECT RATIO
  - ЦВЕТКОРРЕКЦИЯ, БАЛАНС БЕЛОГО, КОНТРАСТ
  - КРОП
  - ПОВОРОТ - НЕТ, НЕ СЛЫШАЛ
- РАЗНЫЕ ФОТО ОДНОГО И ТОГО ЖЕ
  - БОЛЬ!



# ДАННЫЕ - ВСЕ!

- ЧТО ТАКОЕ ПОХОЖИЕ КАРТИНКИ?
- МНОГО ПРИМЕРОВ, ХОРОШИХ И РАЗНЫХ
- МОГУТ БЫТЬ НЕОЖИДАННЫЕ СЮРПРИЗЫ
- ВЫ УЗНАЕТЕ ЧТО БУДЕТ РАБОТАТЬ, А ЧТО НЕТ



# 1. Untergeschoss



# 2. Untergeschoss











# TEST IMAGES

KODAK TRUE  
COLOR TEST  
IMAGES





# LOW QUALITY



В 4 РАЗА МЕНЬШЕ  
JPEG QUALITY =  
10

# CONTRAST & BRIGHTNESS





# CROPS





# AFFINE TRANSFORMS



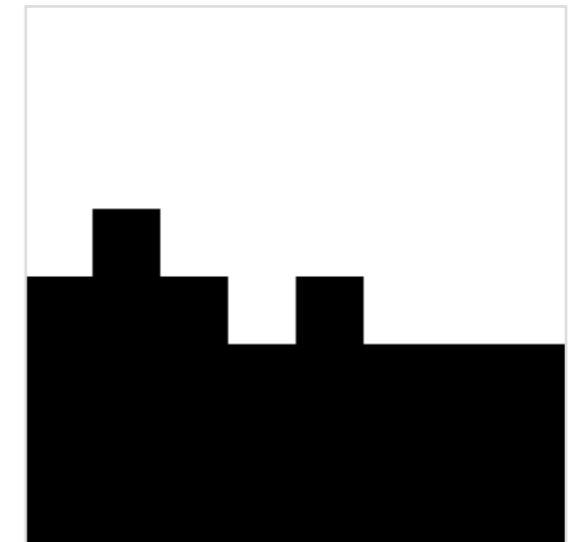
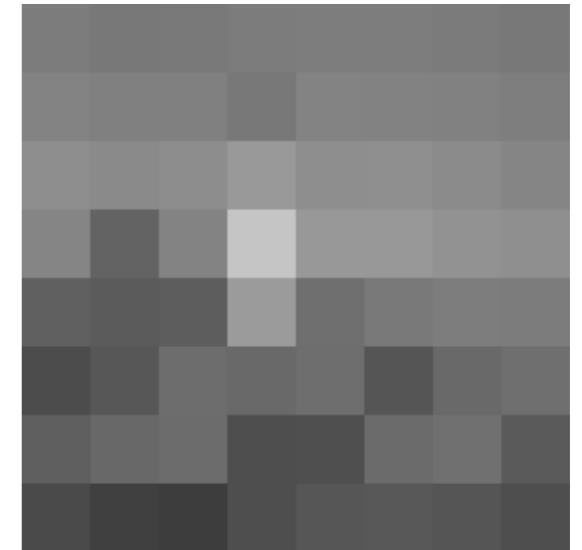
# ПРЕПРОЦЕССИНГ

- РЕСАЙЗИМ ДО УДОБОВАРИМОГО РАЗМЕРА
  - CROP VS ИСКАЖЕНИЕ
  - HINT: БОЛЬШИЕ КАРТИНКИ МОЖНО РЕСАЙЗИТЬ В ДВА ЭТАПА
- НУЖЕН ЛИ ЦВЕТ?
  - ЕСЛИ НУЖЕН, ТО HSV, HSL, YCBCR
- BLUR?

# AVERAGE HASHING

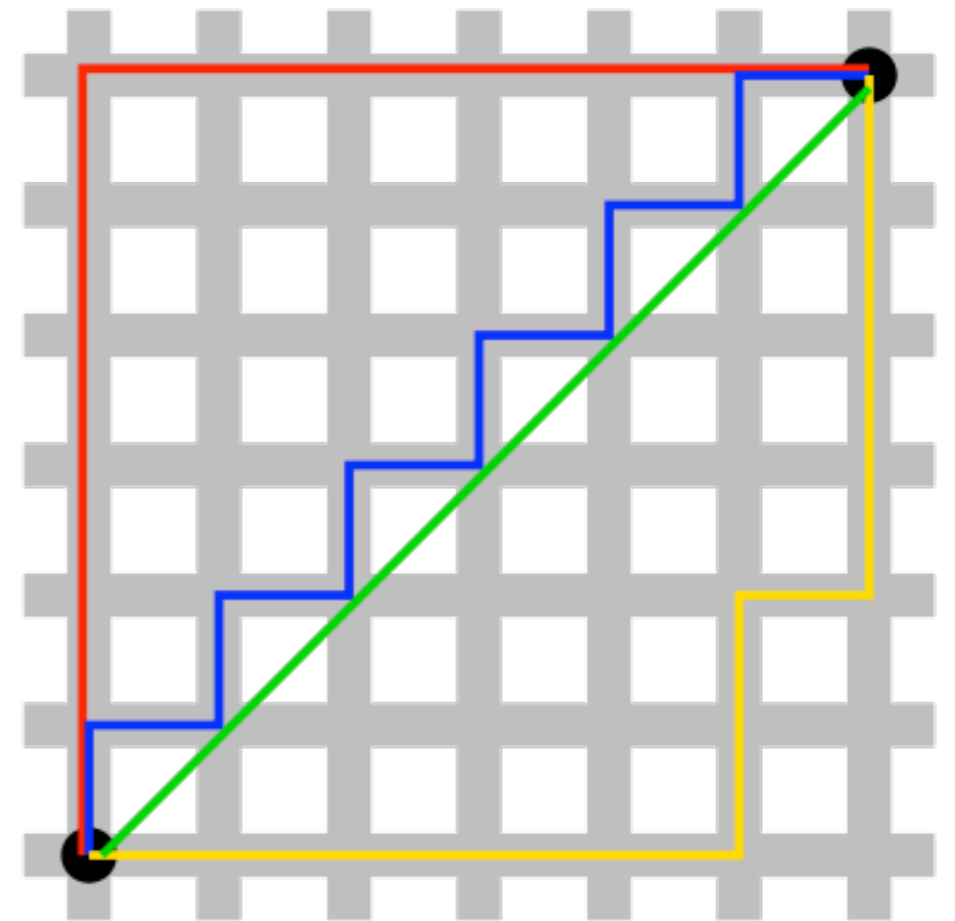
```
image = image.convert('L').resize((8, 8))
pixels = numpy.array(image)
avg = pixels.mean()
diff = pixels > avg
```

```
array([[ True,  True,  True,  True,  True,  True,  True,  True],
       [ True,  True,  True,  True,  True,  True,  True,  True],
       [ True,  True,  True,  True,  True,  True,  True,  True],
       [ True, False,  True,  True,  True,  True,  True,  True],
       [False, False, False,  True, False,  True,  True,  True],
       [False, False, False, False, False, False, False, False],
       [False, False, False, False, False, False, False, False],
       [False, False, False, False, False, False, False, False]], dtype=bool)
```



# А КАК СРАВНИВАТЬ ХЭШИ?

- СОВПАЛ ИЛИ НЕТ
- РАССТОЯНИЕ ХЭММИНГА
  - **00101101** XOR **10011101** == **10110000**
- РАССТОЯНИЕ МАНХЕТТЕНА
- ЕВКЛИДОВО РАССТОЯНИЕ
- КОСИНУСНОЕ РАССТОЯНИЕ
  - $1 - \cos(\Theta)$  МЕЖДУ ВЕКТОРАМИ

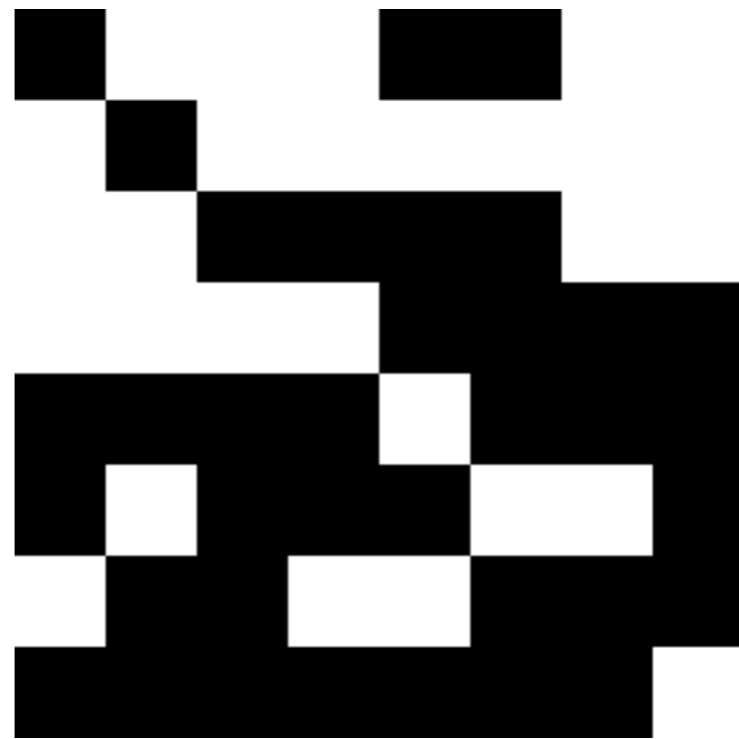
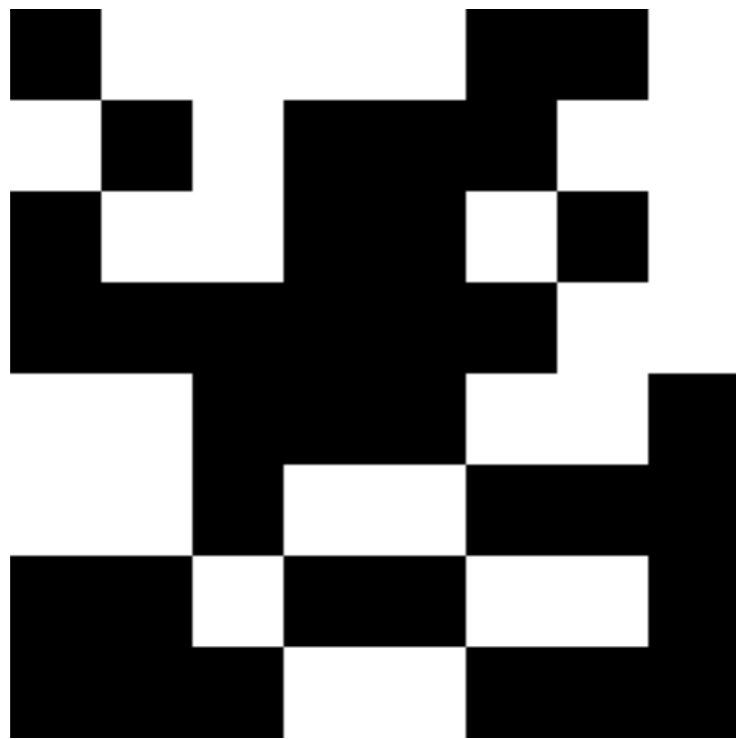


# AVERAGE HASH RESULTS

	4x less	Q=10	Bright	Contrast	Crop small	Crop left	Rotate	Affine	Lenna
64b	1.0	1.0	1.0	1.0	0.94	0.89	0.92	0.97	0.56
256b	1.0	0.99	0.98	1.0	0.89	0.83	0.87	0.81	0.58
1024b	1.0	1.0	0.98	0.99	0.85	0.81	0.82	0.81	0.59

# DIFFERENCE HASHING

```
image = image.convert('L').resize((8, 9))  
pixels = numpy.array(image)  
diff = pixels[1:,:] > pixels[:-1,:]
```





# DIFFERENCE HASH RESULTS

	4x less	Q=10	Bright	Contrast	Crop small	Crop left	Rotate	Affine	Lenna
Vertical	1.0	0.96	0.97	1.0	0.89	0.72	0.83	0.81	0.52
Horizontal	1.0	0.95	0.9	0.95	0.81	0.58	0.8	0.70	0.63

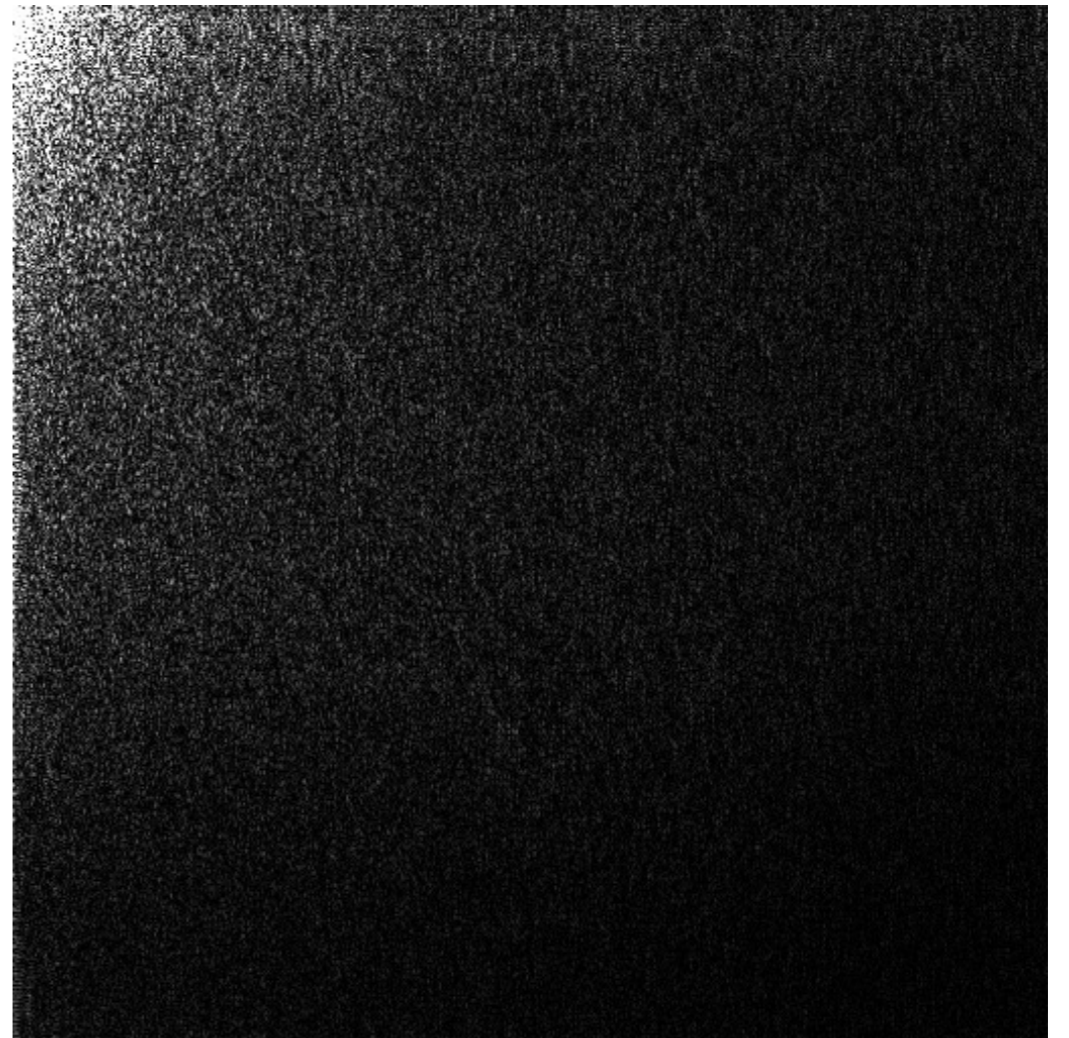
# DCT (DISCRETE COSINE TRANSFORM)

- ЧАСТНЫЙ СЛУЧАЙ DFT (DISCRETE FOURIER TRANSFORM)
- БЕРЕМ ВЕЩЕСТВЕННУЮ СОСТАВЛЯЮЩУЮ, ОТБРАСЫВАЕМ КОМПЛЕКСНУЮ
- МЕНЬШЕ ВЫЧИСЛЕНИЙ, РЕЗУЛЬТАТ ПРАКТИЧЕСКИ ТОТ ЖЕ

$$\begin{aligned} X_{k_1, k_2} &= \sum_{n_1=0}^{N_1-1} \left( \sum_{n_2=0}^{N_2-1} x_{n_1, n_2} \cos \left[ \frac{\pi}{N_2} \left( n_2 + \frac{1}{2} \right) k_2 \right] \right) \cos \left[ \frac{\pi}{N_1} \left( n_1 + \frac{1}{2} \right) k_1 \right] \\ &= \sum_{n_1=0}^{N_1-1} \sum_{n_2=0}^{N_2-1} x_{n_1, n_2} \cos \left[ \frac{\pi}{N_1} \left( n_1 + \frac{1}{2} \right) k_1 \right] \cos \left[ \frac{\pi}{N_2} \left( n_2 + \frac{1}{2} \right) k_2 \right]. \end{aligned}$$

# СТРОИМ DCT ХЭШ

```
pixels = np.array(image, dtype=np.float)
pixels_dct = dct(dct(pixels, axis=1), axis=0)
low_freq = pixels_dct[:hash_size, :hash_size]
median = np.median(low_freq)
diff = low_freq > median
```





512x512 DCT

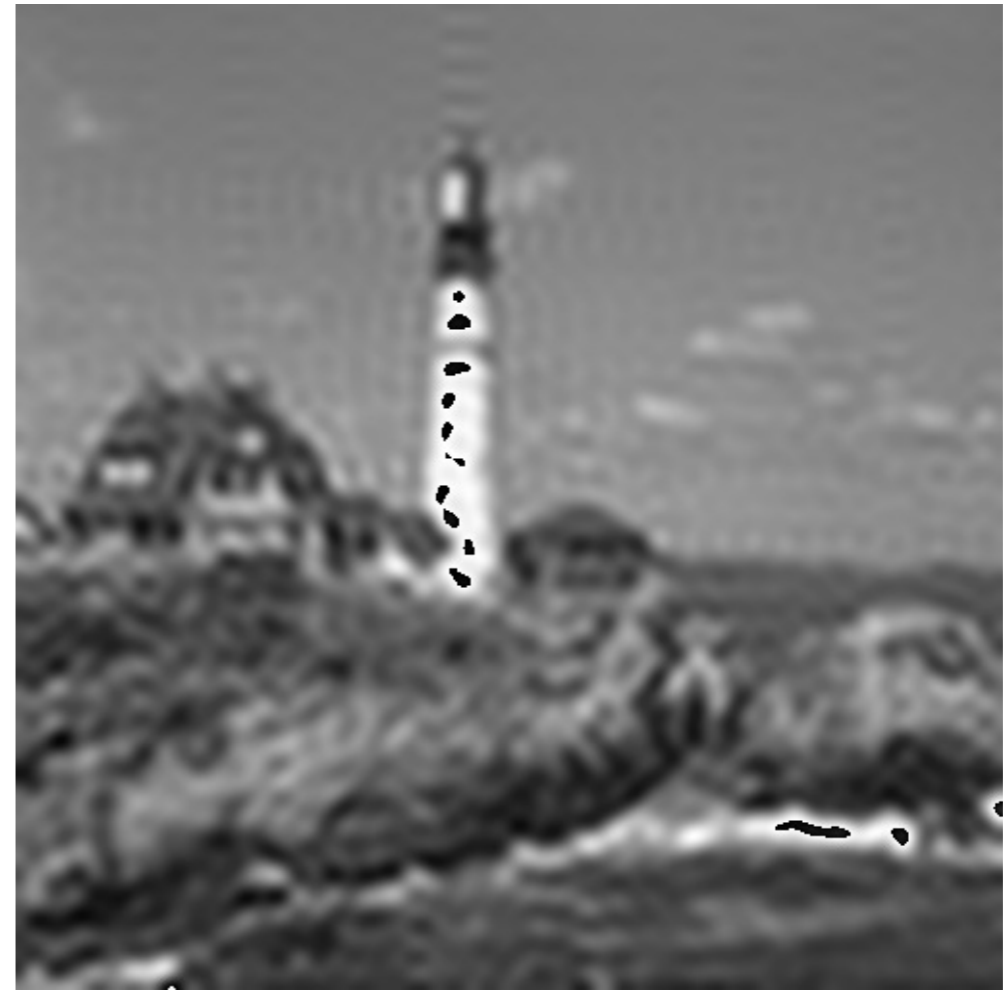


256x256 DCT

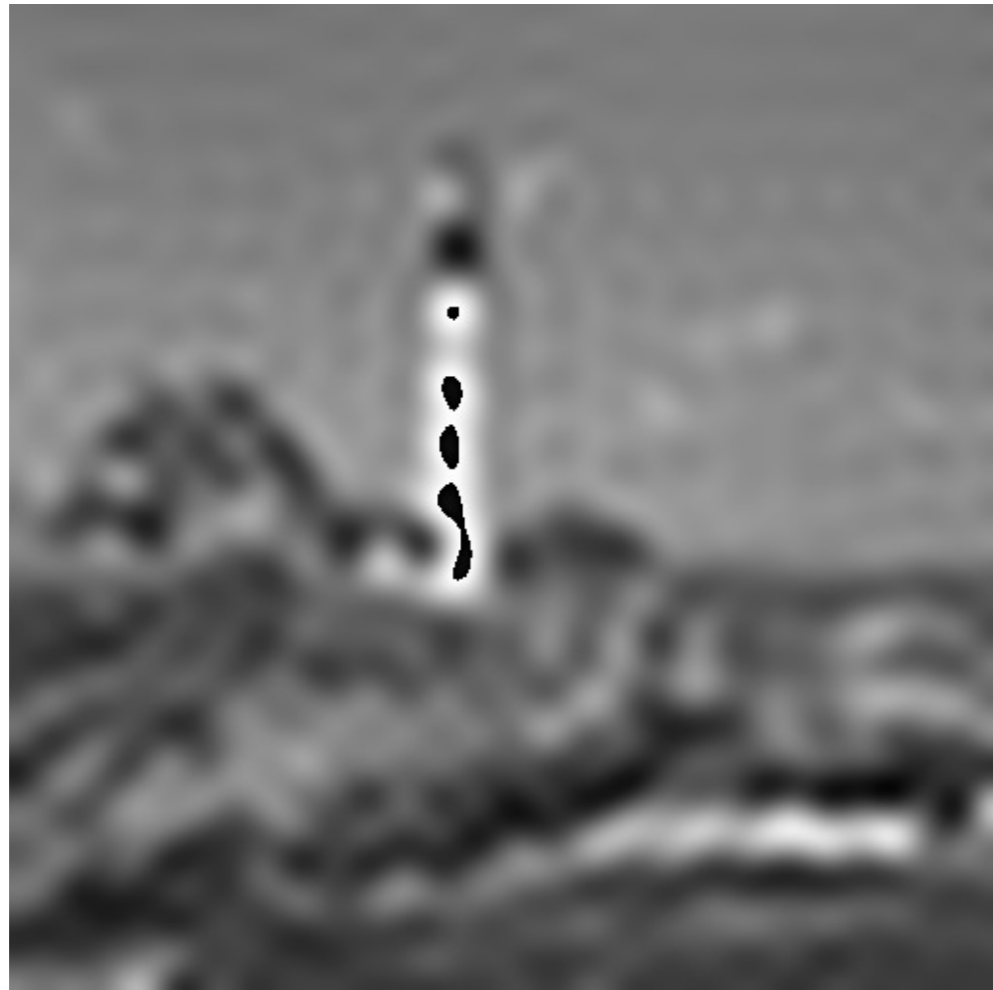




128x128 DCT



64x64 DCT



32x32 DCT



16x16 DCT

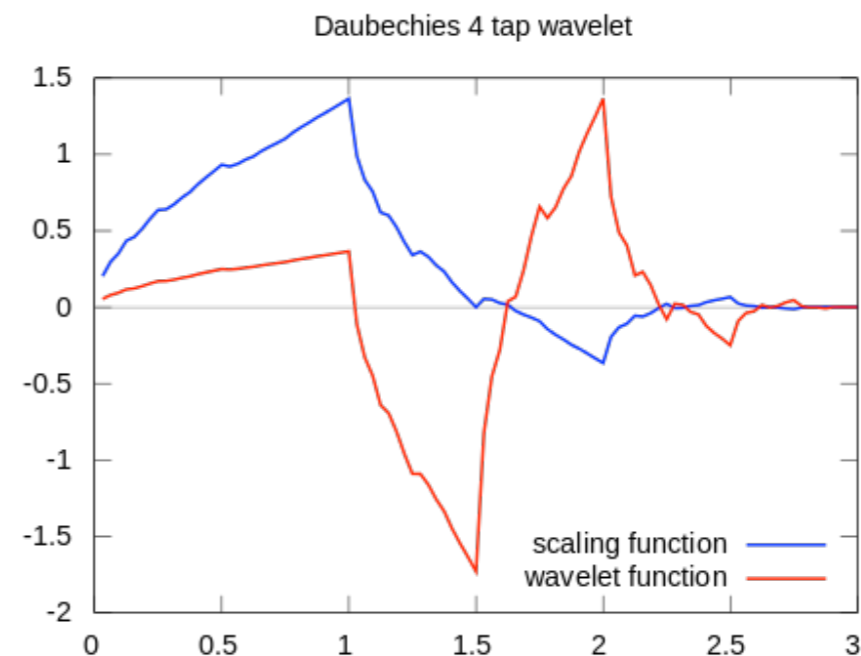
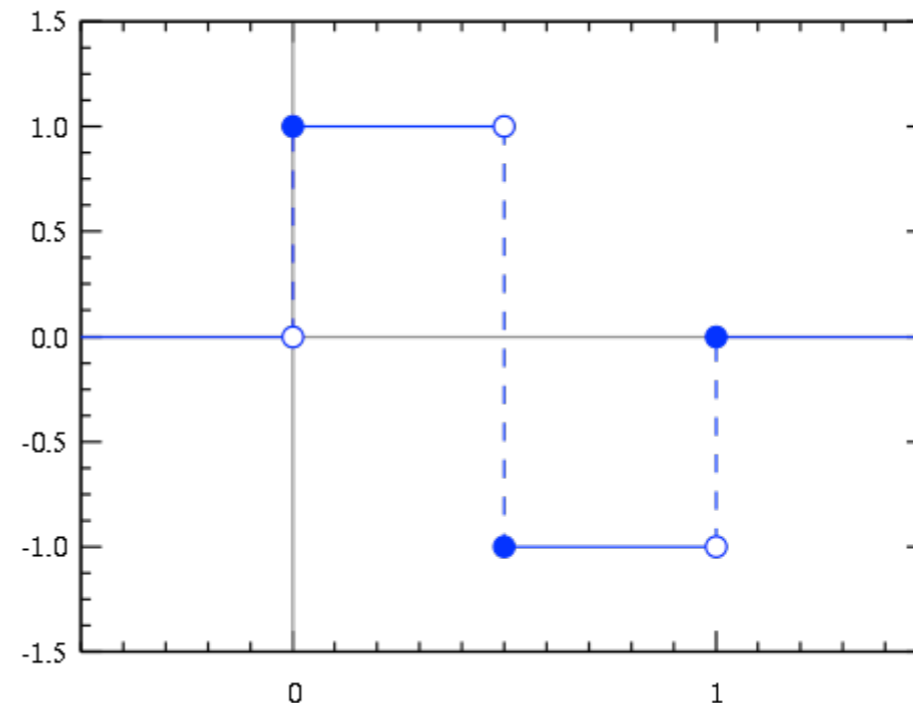


# DCT HASH RESULTS

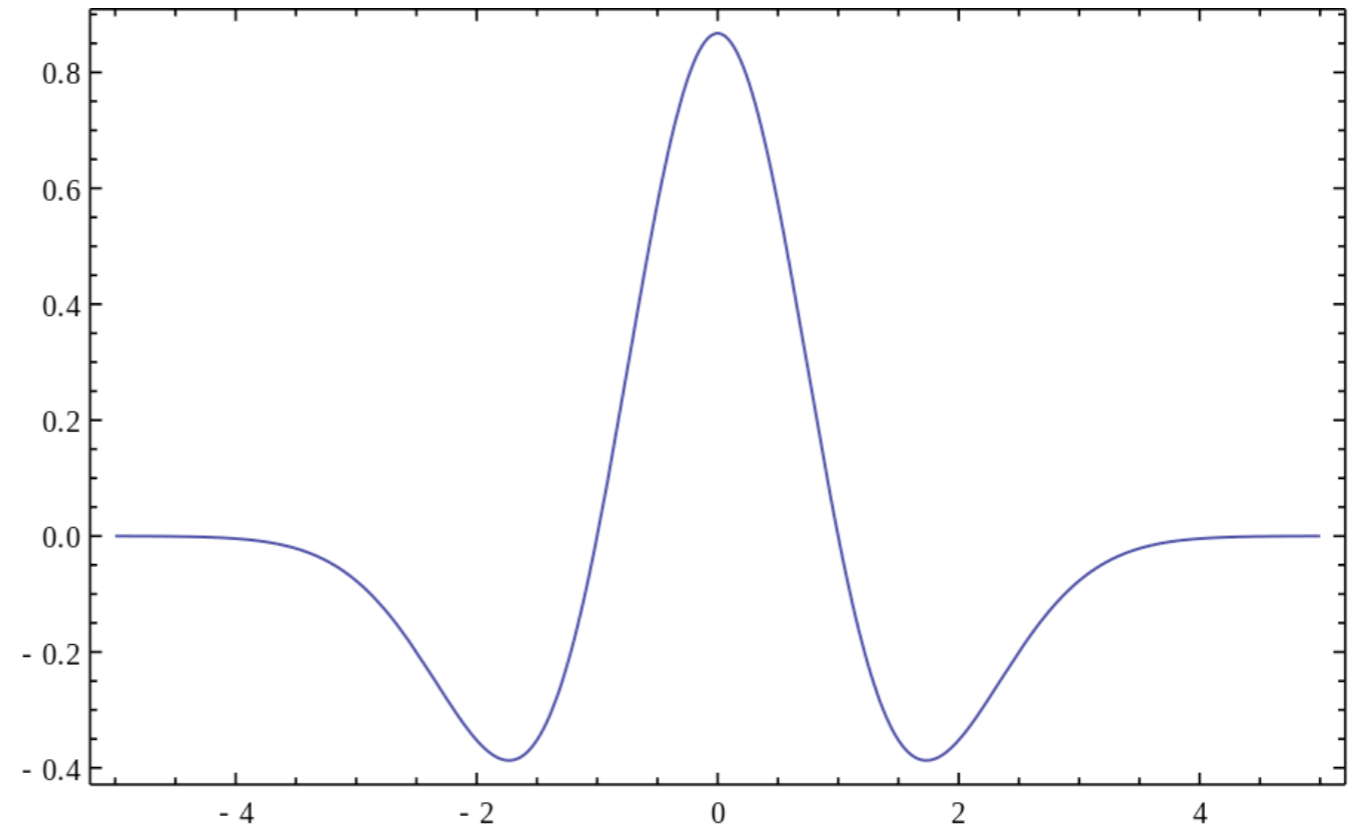
	4x less	Q=10	Bright	Contrast	Crop small	Crop left	Rotate	Affine	Lenna
Our	1.0	0.97	0.84	0.94	0.84	0.63	0.78	0.75	0.59
pHash	0.88	0.97	0.88	0.94	0.84	0.47	0.75	0.67	0.56
Our 256b	0.99	0.97	0.89	0.96	0.77	0.48	0.64	0.61	0.53
Our 1024b	0.998	0.97	0.9	0.95	0.58	0.49	0.56	0.52	0.5

# WAVELETS

- ОНИ ЕСТЬ РАЗНЫЕ
- ОПИСЫВАЮТ ТЕ ЖЕ ВОЛНЫ, ТОЛЬКО В ПРИВЯЗКЕ К ПОЛОЖЕНИЮ
- JPEG2000
- ИСПОЛЬЗУЮТСЯ В IMGSEEK
- ЕСТЬ РЕАЛИЗАЦИЯ В PHASH



# MEXICAN HAT



4x less	Q=10	Bright	Contrast	Crop small	Crop left	Rotate	Affine	Lenna
0.77	0.82	0.88	0.9	0.51	0.53	0.5	0.5	0.5

# STATISTICAL MOMENTS

- СТАТИСТИЧЕСКИЙ МОМЕНТ

$$M_{ij} = \sum_x \sum_y x^i y^j I(x, y)$$

- CENTROID:  $\{ X, Y \} = \{ M_{10}/M_{00}, M_{01}/M_{00} \}$

- CENTRAL MOMENT

- TRANSLATION INVARIANT

$$\mu_{pq} = \sum_m^p \sum_n^q \binom{p}{m} \binom{q}{n} (-\bar{x})^{(p-m)} (-\bar{y})^{(q-n)} M_{mn}$$

- ORIENTATION

- SCALE INVARIANT MOMENT

$$\eta_{ij} = \frac{\mu_{ij}}{\mu_{00}^{\left(1 + \frac{i+j}{2}\right)}}$$

# HU INVARIANT MOMENTS

- Их 7 (8)
- ИНВАРИАНТНЫ К
  - ТРАНСЛЯЦИИ
  - МАСШТАБИРОВАНИЮ
  - ПОВОРОТУ
- SKIMAGE.MEASURE.MOMENTS\_HU
- [HTTPS://EN.WIKIPEDIA.ORG/WIKI/  
IMAGE\\_MOMENT](https://en.wikipedia.org/wiki/Image_Moment)

	Euclidian distance	Equalized
4x less	0.039	0.088
Q=10	0.0079	0.054
Bright	1.71	0.27
Contrast	0.52	0.0077
Crop	0.37	0.41
Crop	0.12	0.048
Rotate	0.24	0.044
Affine	0.22	0.092
Lenna	1.18	1.43

# RADIAL VARIANCE HASH

- ДИСПЕРСИЯ ИНТЕНСИВНОСТИ ВДОЛЬ ОСЕЙ
- ОСИ ПОД УГЛОМ, ПРОХОДЯТ ЧЕРЕЗ ЦЕНТР
- ЕСТЬ РЕАЛИЗАЦИЯ В RHASH

4x less	Q=10	Bright	Contrast	Crop small	Crop left	Rotate	Affine	Lenna
0.99	1.0	0.77	0.98	0.98	0.34	0.94	0.73	0.46



# ИНТЕРЕСНЫЕ ТОЧКИ

- КРАЯ (EDGES)
- УГЛЫ (CORNERS)
- BLOBS
- RIDGES (КОНТУР?)
- РЕГИОНЫ
- ТОЧКИ - ХОРОШО
- ДЕСКРИПТОРЫ - ЛУЧШЕ
- OPENCV - FEATURES2D
- SCIKIT IMAGE -  
SCKIMAGE.FEATURE

# ИТОГИ

- ПРОБЛЕМА СЛОЖНАЯ
- РЕШЕНИЯ ЕСТЬ
  - ПРОСТЫХ НЕТ
- КОМБИНИРУЙТЕ РАЗНЫЕ ПРИЗНАКИ
- ИЗ 7М КАРТИНОК МЫ УБРАЛИ 600К ДУБЛИКАТОВ
  - 20 C3.LARGE EC2 МАШИН ЗА СУТКИ

# КОНТАКТЫ

- [HTTP://GITHUB.COM/VENTURA/TALKS](http://github.com/ventura/talks)
- [VIKTOR.DRACHOV@GMAIL.COM](mailto:VIKTOR.DRACHOV@GMAIL.COM)
- ВНИМАНИЕ, ГОЛОСОВАНИЕ
  - [HTTPS://GITHUB.COM/VENTURA/IMAGE\\_COMPARE](https://github.com/ventura/image_compare)