



Ingesting 35M images with Python

In the cloud.

Àlex Vinyals
Software Engineer @ Hotels Data

Unify all the data

Challenges of a metasearch

Sort: Most popular

List

Map

15 results near W Barcelona, sorted by Most popular

Total price for 1 night, 1 guest, 1 room



Reset all filters

Total Price

\$0 - \$55 1

\$110 - \$165 1

\$165 - \$250 1

>\$250 1

Stars

★★★★★ \$195 2

★★★ \$147 4

★★ 4

Any \$14 5

District

Ciutat Vella \$14 15

Meal Plan

Room only 4

Selected hotel



W Barcelona ★★★★★

see large map

Wi-Fi P Parking

Book directly with Starwood for \$359 >



Single room

Room only • Non refundable

\$488 \$356

View deal



Single room

Room only • Non refundable

\$488 \$356

View deal



Single room

Room only • Non refundable

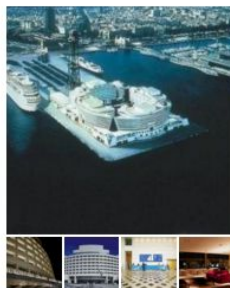
\$488 \$356

View deal

Show more prices >

Very good 8.6
Reviews: 3144

Nearby hotels



Eurostars Grand Marina Hotel GI ★★★★★

0.53 mi. to W Barcelona (map)

Wi-Fi P Parking



Single room

Room only • Non refundable

\$2,146 \$195

View deal



Single room

Room only • Non refundable

\$2,146 \$198

View deal



Single room

Room only • Non refundable

\$2,146 \$210

View deal

Show more prices >

Very good 8.6
Reviews: 3861

Hotel 54 Barceloneta ★★★

0.55 mi. to W Barcelona (map)

Wi-Fi P Parking

Good 7.8
Reviews: 485

Sort: Most popular

List

Map

15 results near W Barcelona, sorted by Most popular

Total price for 1 night, 1 guest, 1 room



Reset all filters

Total Price

\$0 - \$55 1

\$110 - \$165 1

\$165 - \$250 1

>\$250 1

Stars

★★★★★ \$195 2

★★★ \$147 4

★★ 4

Any \$14 5

District

Ciutat Vella \$14 15

Meal Plan

Room only 4

Selected hotel



W Barcelona ★★★★★

see large map

Very good 8.6
Reviews: 3144

W Barcelona

Hotels.com

Single room

Room only • Non refundable

\$488 \$356

View deal

Expedia

Single room

Room only • Non refundable

\$488 \$356

View deal

venere.com

Single room

Room only • Non refundable

\$488 \$356

View deal

Show more prices >

Nearby hotels



Eurostars Grand Marina Hotel GI ★★★★★

0.53 mi. to W Barcelona (map)

Very good 8.6
Reviews: 3861

Wi-Fi P Parking

h. hot.es

Single room

Room only • Non refundable

\$2,146 \$195

View deal

getaroom

Single room

Room only • Non refundable

\$2,146 \$198

View deal

HRS

Single room

Room only • Non refundable

\$2,146 \$210

View deal

Show more prices >

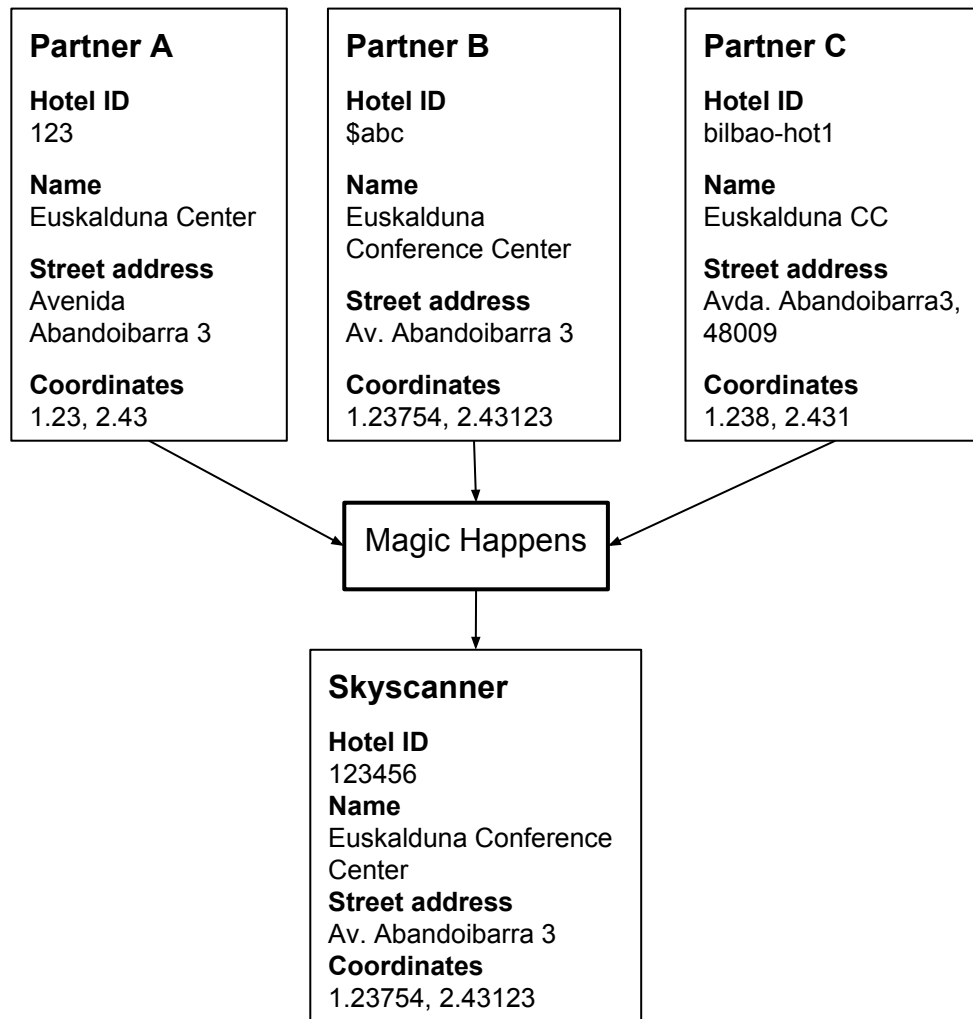


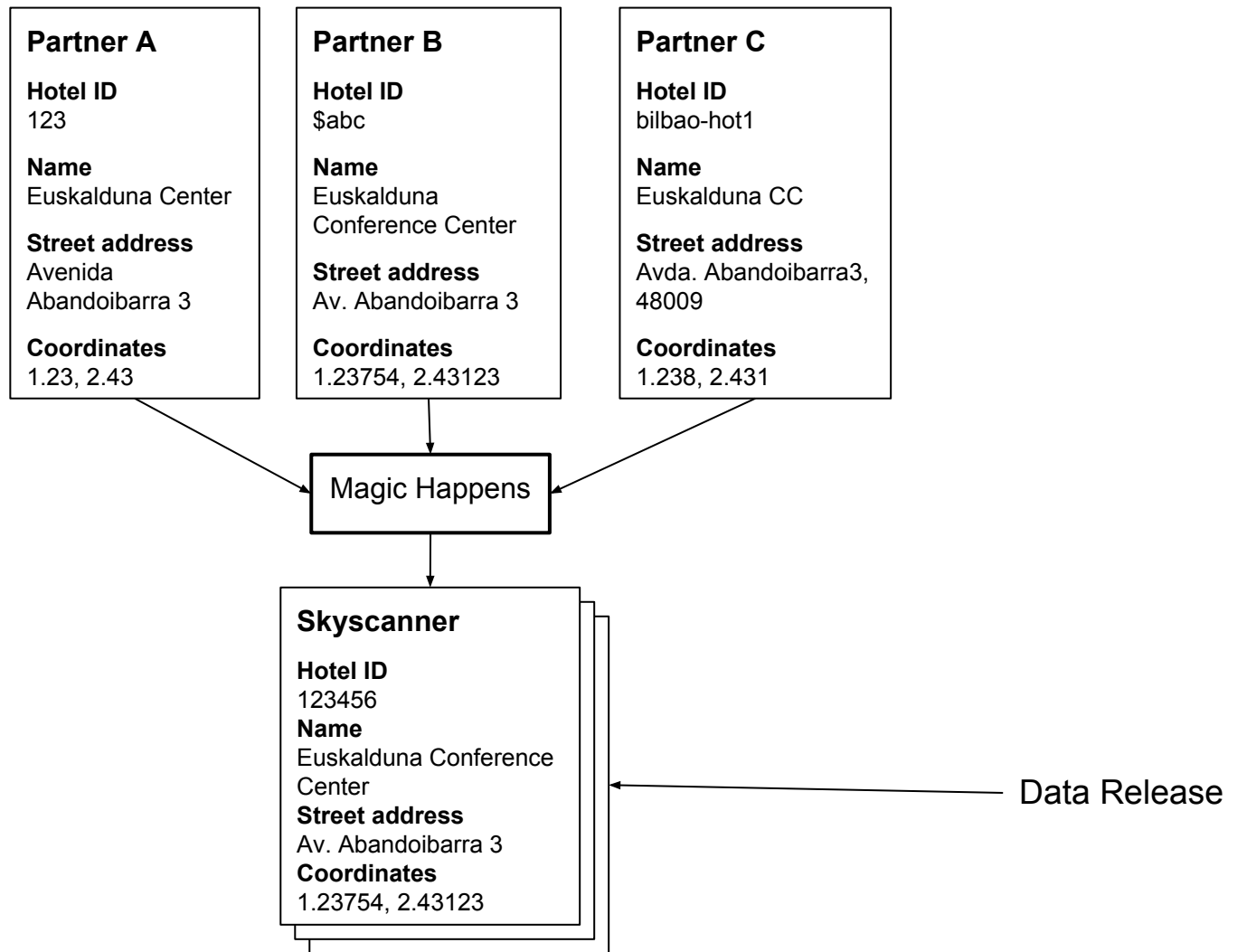
Hotel 54 Barceloneta ★★★

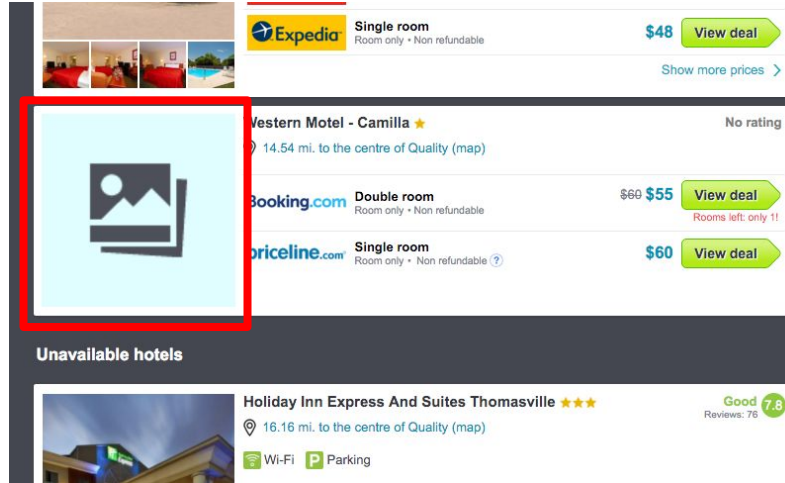
0.55 mi. to W Barcelona (map)

Good 7.8
Reviews: 485

Wi-Fi P Parking







So what about the images?

Partner A

Hotel ID
123



Partner B

Hotel ID
\$abc



Partner C

Hotel ID
bilbao-hot1



Magic Happens

Skyscanner

Hotel ID
123456



Sort: Most popular

List

Map

15 results near W Barcelona, sorted by Most popular

Total price for 1 night, 1 guest, 1 room



Reset all filters

Total Price

\$0 - \$55 1

\$110 - \$165 1

\$165 - \$250 1

>\$250 1

Stars

★★★★★ 2

★★★ \$147 4

★★ 4

Any \$14 5

District

Ciutat Vella \$14 15

Meal Plan

Room only 4

Selected hotel



W Barcelona ★★★★★

Very good 8.6
Reviews: 3144

see large map

Wi-Fi P Parking

Book directly with Starwood for \$359 >

Hotels.com Single room
Room only • Non refundable \$488 \$356 View dealExpedia Single room
Room only • Non refundable \$488 \$356 View dealvenere.com Single room
Room only • Non refundable \$488 \$356 View deal

Show more prices >

Nearby hotels



Eurostars Grand Marina Hotel GI ★★★★★

Very good 8.6
Reviews: 3861

0.53 mi. to W Barcelona (map)

Wi-Fi P Parking

h. hot.es Single room
Room only • Non refundable \$2,146 \$195 View dealgetaroom Single room
Room only • Non refundable \$2,146 \$198 View dealHRS Single room
Room only • Non refundable \$2,146 \$210 View deal

Show more prices >

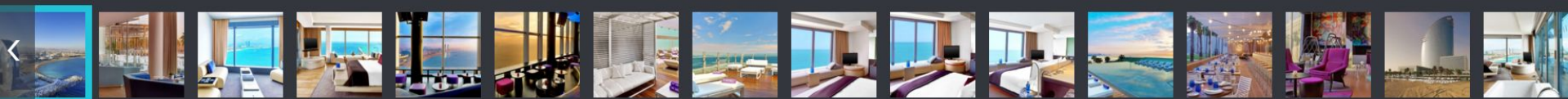
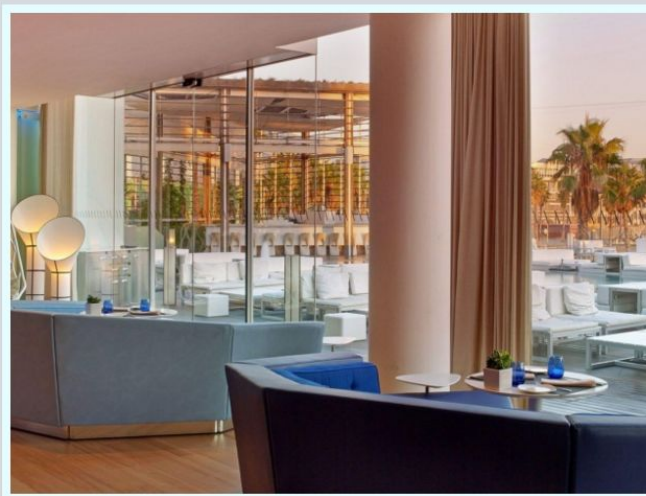


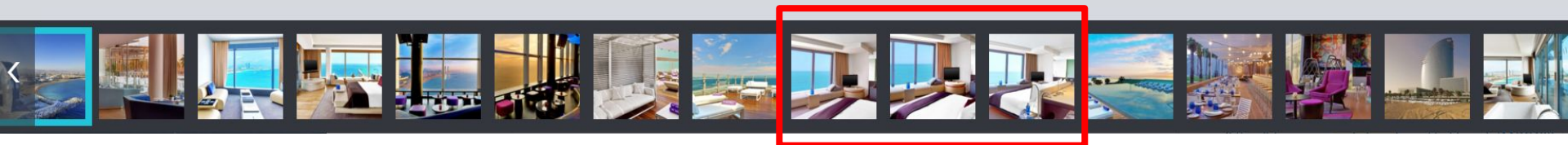
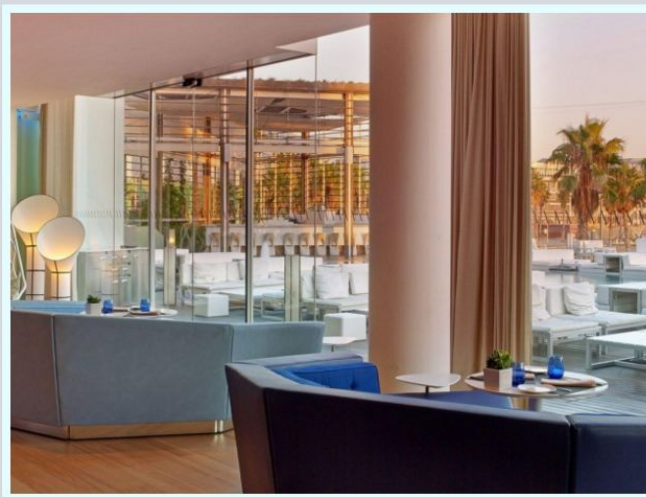
Hotel 54 Barceloneta ★★★

Good 7.8
Reviews: 485

0.55 mi. to W Barcelona (map)

Wi-Fi P Parking







With more than 200 partners
800.000 hotels reach production

Images to process = $K * M * N \sim 35\text{M images}$

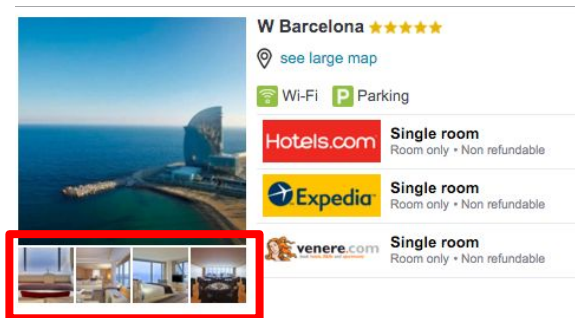
K = number of partners

M = avg number of hotels per partner

N = avg number of images per partner hotel

Resizing is a thing

And we have 14 different configurations



```
[sizes]
628x470, scaledownfit
314x235, scaledownfit
2048x1536, scaledownfit
1024x768, scaledownfit
200x200, scalefill, crop
100x100, scalefill, crop
280x185, 185x280, scaledownfill, crop, contrast 1.05
182x182, scaledownfill, crop, contrast 1.05
102x78, scaledownfill, crop, contrast 1.1
58x58, scaledownfill, crop, contrast 1.1
480x319, 240x319, scaledownfill, crop, contrast 1.05
45x45, scaledownfill, crop, contrast 1.1
627x470, scalefill, crop
313x235, scalefill, crop
```

Tale of an image processing pipeline



Tech Stack

Riding on AWS



Tech Stack

Riding on AWS



SQS
Simple Queue Service



Tech Stack

Riding on AWS



SQS
Simple Queue Service



Compute resources



*with DjangoRestFramework

*without Django ORM

Libraries



*with DjangoRestFramework

*without Django ORM

Libraries





*with DjangoRestFramework

*without Django ORM

Libraries



Kombu

Messaging / queues / amqp



*with DjangoRestFramework

*without Django ORM

Libraries



Kombu

Messaging / queues / amqp

Boto

Amazon stuff



*with DjangoRestFramework

*without Django ORM

Pillow

Image Processing

Libraries



Kombu

Messaging / queues / amqp

Boto

Amazon stuff



*with DjangoRestFramework

*without Django ORM

Pillow

Image Processing

Libraries

Python2.7



Kombu

Messaging / queues / amqp

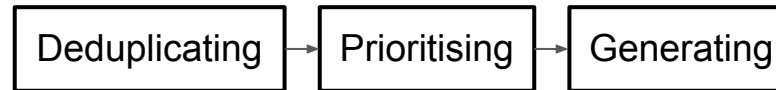
SQLAlchemy

Boto

Amazon stuff



Tale of an image processing pipeline

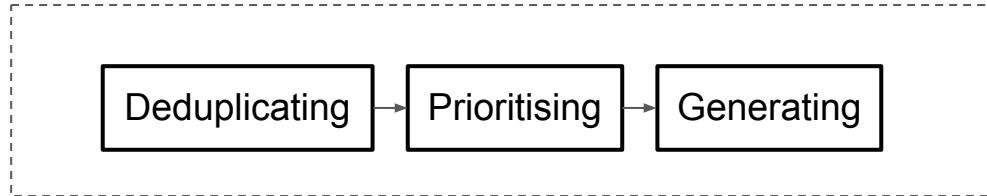


Asynchronous (Always Running)



Tale of an image processing pipeline

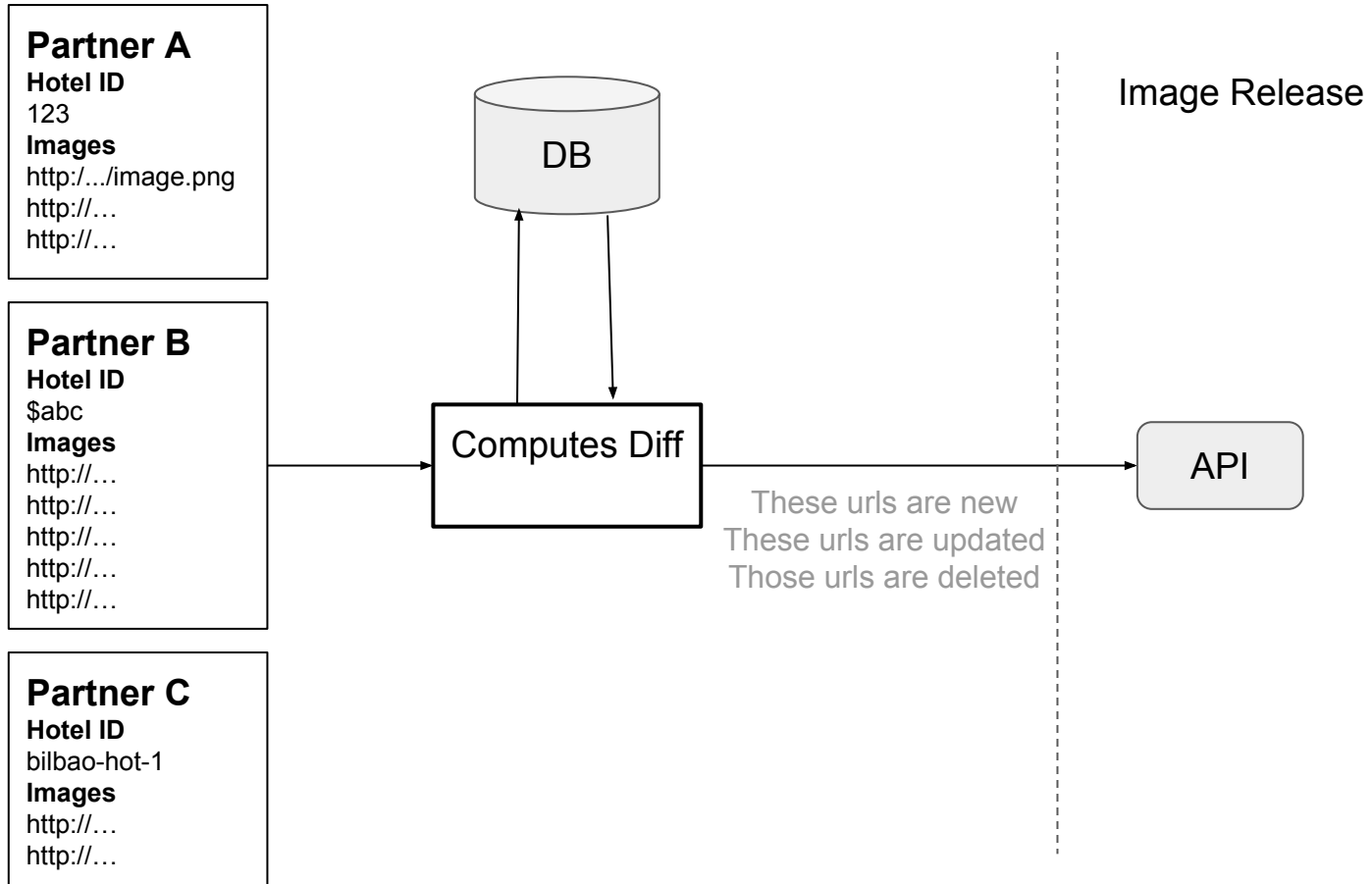
Triggered by the Data Release





Triggering

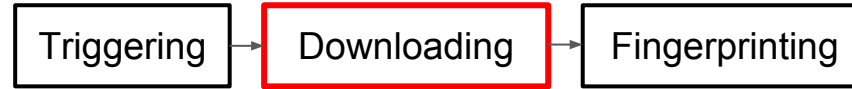




id	url	partner_id	acc_id	status	en
23452927	http://hote...re.com/h	4	3969975	available	t
23452917	http://hote...re.com/h	4	3969975	available	t
23452848	http://img...tels/CAM	5	THN3J4	production	t
23452841	http://img...tels/CAM	5	THN3J4	production	t
23452846	http://img...tels/CAM	5	THN3J4	production	t
23452844	http://img...tels/CAM	5	THN3J4	production	t
23452851	http://img...tels/CAM	5	THN3J4	production	t
23452796	http://d1pa...cloudfro	9	252966	production	t
23452797	http://d1pa...cloudfro	9	252966	production	t
23452794	http://d1pa...cloudfro	9	252966	production	t
23452788	http://www...com/ima	18	931146	filtered	t
23452914	http://hote...re.com/h	4	3969975	available	t
23452798	http://d1pa...cloudfro	9	252966	production	t
23452802	http://d1pa...cloudfro	9	252966	production	t
23452922	http://hote...re.com/h	4	3969975	available	t
23452919	http://hote...re.com/h	4	3969975	available	t
23452795	http://d1pa...cloudfro	9	252966	production	t
23452924	http://hote...re.com/h	4	3969975	available	t
23452799	http://d1pa...cloudfro	9	252966	production	t

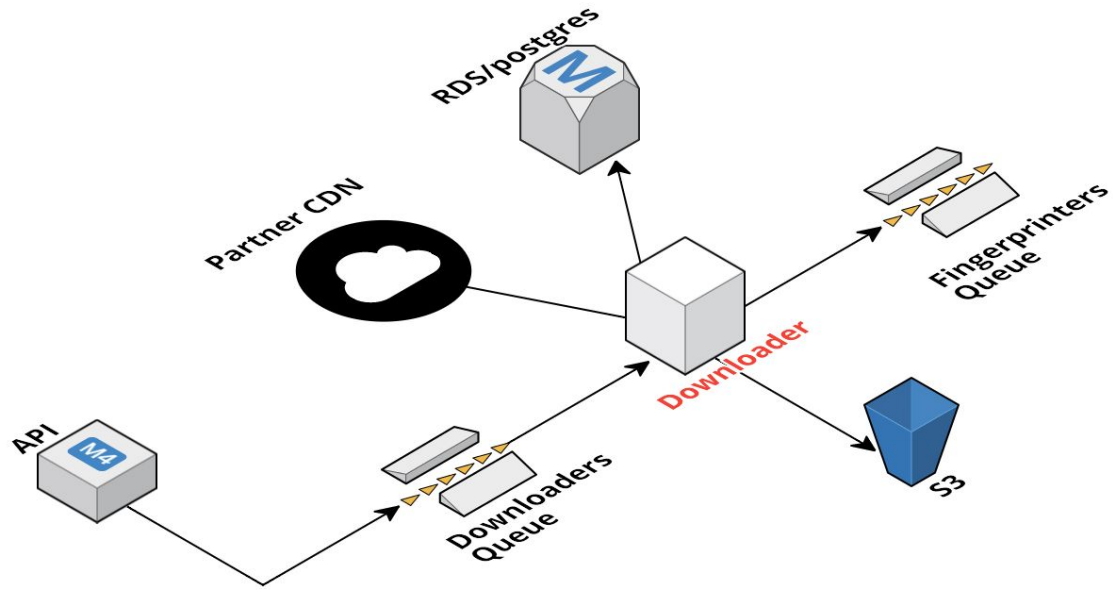
id	url	partner_id	acc_id	status	en
23452927	http://hote...re.com/h	4	3969975	available	t
23452917	http://hote...re.com/h	4	3969975	available	t
23452848	http://img...tels/CAM	5	THN3J4	production	t
23452841	http://img...tels/CAM	5	THN3J4	production	t
23452846	http://img...tels/CAM	5	THN3J4	production	t
23452844	http://img...tels/CAM	5	THN3J4	production	t
23452851	http://img...tels/CAM	5	THN3J4	production	t
23452796	http://d1pa...cloudfro	9	252966	production	t
23452797	http://d1pa...cloudfro	9	252966	production	t
23452794	http://d1pa...cloudfro	9	252966	production	t
23452788	http://www...com/ima	18	931146	filtered	t
23452914	http://hote...re.com/h	4	3969975	available	t
23452798	http://d1pa...cloudfro	9	252966	production	t
23452802	http://d1pa...cloudfro	9	252966	production	t
23452922	http://hote...re.com/h	4	3969975	available	t
23452919	http://hote...re.com/h	4	3969975	available	t
23452795	http://d1pa...cloudfro	9	252966	production	t
23452924	http://hote...re.com/h	4	3969975	available	t
23452799	http://d1pa...cloudfro	9	252966	production	t

id	url	partner_id	acc_id	status	en
23452927	http://hote...re.com/h	4	3969975	available	t
23452917	http://hote...re.com/h	4	3969975	available	t
23452848	http://img...tels/CAM	5	THN3J4	production	t
23452841	http://img...tels/CAM	5	THN3J4	production	t
23452846	http://img...tels/CAM	5	THN3J4	production	t
23452844	http://img...tels/CAM	5	THN3J4	production	t
23452851	http://img...tels/CAM	5	THN3J4	production	t
23452796	http://d1pa...cloudfro	9	252966	production	t
23452797	http://d1pa...cloudfro	9	252966	production	t
23452794	http://d1pa...cloudfro	9	252966	production	t
23452788	http://www...com/ima	18	931146	filtered	t
23452914	http://hote...re.com/h	4	3969975	available	t
23452798	http://d1pa...cloudfro	9	252966	production	t
23452802	http://d1pa...cloudfro	9	252966	production	t
23452922	http://hote...re.com/h	4	3969975	available	t
23452919	http://hote...re.com/h	4	3969975	available	t
23452795	http://d1pa...cloudfro	9	252966	production	t
23452924	http://hote...re.com/h	4	3969975	available	t
23452799	http://d1pa...cloudfro	9	252966	production	t



Downloading





```
import io
import boto
import requests
from PIL import Image

s3 = boto.connect_s3()
bucket = s3.get_bucket('available-images')

@reliable_callback()
def downloader_callback(queued_image):
    """ Overly simplified downloading callback without
    error handling logic """
    response = requests.get(queued_image.url)
    blob = response.content
    key = bucket.new_key(queued_image.basename)
    key.set_contents_from_string(blob)
    image = Image.open(io.BytesIO(blob))
    if should_filter(image):
        return
    fingerprinting_producer.publish(queued_image)
```

```
def should_filter(image):
    height, width = image.size

    short_size = min(width, height)
    if short_size < minimum_short:
        return True

    long_size = max(width, height)
    if long_size < minimum_long:
        return True

    total_pixels = width * height
    if total_pixels > max_pixels:
        return True

    return False
```

```
import io
import boto
import requests
from PIL import Image

s3 = boto.connect_s3()
bucket = s3.get_bucket('available-images')

@reliable_callback()
def downloader_callback(queued_image):
    """ Overly simplified downloading callback without
    error handling logic """
    response = requests.get(queued_image.url)
    blob = response.content
    key = bucket.new_key(queued_image.basename)
    key.set_contents_from_string(blob)
    image = Image.open(io.BytesIO(blob))
    if should_filter(image):
        return
    fingerprinting_producer.publish(queued_image)
```

```
def should_filter(image):
    height, width = image.size

    short_size = min(width, height)
    if short_size < minimum_short:
        return True

    long_size = max(width, height)
    if long_size < minimum_long:
        return True

    total_pixels = width * height
    if total_pixels > max_pixels:
        return True

    return False
```

```
import io
import boto
import requests
from PIL import Image

s3 = boto.connect_s3()
bucket = s3.get_bucket('available-images')

@reliable_callback()
def downloader_callback(queued_image):
    """ Overly simplified downloading callback without
    error handling logic """
    response = requests.get(queued_image.url)
    blob = response.content
    key = bucket.new_key(queued_image.basename)
    key.set_contents_from_string(blob)
    image = Image.open(io.BytesIO(blob))
    if should_filter(image):
        return
    fingerprinting_producer.publish(queued_image)
```

```
def should_filter(image):
    height, width = image.size

    short_size = min(width, height)
    if short_size < minimum_short:
        return True

    long_size = max(width, height)
    if long_size < minimum_long:
        return True

    total_pixels = width * height
    if total_pixels > max_pixels:
        return True

    return False
```

```
import functools
import warnings
from PIL import Image

def reliable_callback():
    def decorator(func):
        @functools.wraps(func)
        def wrapper(*args, **kwargs):
            warnings.simplefilter('error', Image.DecompressionBombWarning)
            try:
                return func(*args, **kwargs)
            except BaseException:
                logger.error("Critical worker error", exc_info=True)
        return wrapper
    return decorator
```

```
import functools
import warnings
from PIL import Image

def reliable_callback():
    def decorator(func):
        @functools.wraps(func)
        def wrapper(*args, **kwargs):
            warnings.simplefilter('error', Image.DecompressionBombWarning)
            try:
                return func(*args, **kwargs)
            except BaseException:
                logger.error("Critical worker error", exc_info=True)
        return wrapper
    return decorator
```

```
import functools
import warnings
from PIL import Image

def reliable_callback():
    def decorator(func):
        @functools.wraps(func)
        def wrapper(*args, **kwargs):
            warnings.simplefilter('error', Image.DecompressionBombWarning)
            try:
                return func(*args, **kwargs)
            except BaseException:
                logger.error("Critical worker error", exc_info=True)
        return wrapper
    return decorator
```

```
import functools
import warnings
from PIL import Image

def reliable_callback():
    def decorator(func):
        @functools.wraps(func)
        def wrapper(*args, **kwargs):
            warnings.simplefilter('error', Image.DecompressionBombWarning)
            try:
                return func(*args, **kwargs)
            except BaseException:
                logger.error("Critical worker error", exc_info=True)
        return wrapper
    return decorator
```

```
from kombu import Connection, Consumer, Exchange, Queue, eventloop

class KombuConsumer(common.BaseConsumer):
    # ... bla bla

    def callback(self, body, message):
        self.handler(body)
        message.ack()

    def listen(self):
        with Connection(self.backend.broker, transport_options={'region': self.backend.region}) as connection:
            with Consumer(connection, self.queue, callbacks=[self.callback], accept=[self.backend.serializer]):
                for _ in eventloop(connection):
                    pass

# What a simplified worker looks like
# Broker URI stored on Backend object, looks like:
#   sqs://{s3_key}:{s3_secret}@

consumer = KombuConsumer(backend, handler=downloader.downloader_callback)
consumer.listen()
```

```
from kombu import Connection, Consumer, Exchange, Queue, eventloop
```

```
class KombuConsumer(common.BaseConsumer):
```

```
    # ... bla bla
```

```
    def callback(self, body, message):
```

```
        self.handler(body)
```

```
        message.ack()
```

```
    def listen(self):
```

```
        with Connection(self.backend.broker, transport_options={'region': self.backend.region}) as connection:
```

```
            with Consumer(connection, self.queue, callbacks=[self.callback], accept=[self.backend.serializer]):
```

```
                for _ in eventloop(connection):
```

```
                    pass
```

```
# What a simplified worker looks like
```

```
# Broker URI stored on Backend object, looks like:
```

```
# sqs://{s3_key}:{s3_secret}@
```

```
consumer = KombuConsumer(backend, handler=downloader.downloader_callback)
```

```
consumer.listen()
```

```
from kombu import Connection, Consumer, Exchange, Queue, eventloop

class KombuConsumer(common.BaseConsumer):
    # ... bla bla

    def callback(self, body, message):
        self.handler(body)
        message.ack()

    def listen(self):
        with Connection(self.backend.broker, transport_options={'region': self.backend.region}) as connection:
            with Consumer(connection, self.queue, callbacks=[self.callback], accept=[self.backend.serializer]):
                for _ in eventloop(connection):
                    pass

# What a simplified worker looks like
# Broker URI stored on Backend object, looks like:
#   sqs://{s3_key}:{s3_secret}@

consumer = KombuConsumer(backend, handler=downloader.downloader_callback)
consumer.listen()
```

```
from kombu import Connection, Consumer, Exchange, Queue, eventloop
```

```
class KombuConsumer(common.BaseConsumer):
```

```
    # ... bla bla
```

```
    def callback(self, body, message):
```

```
        self.handler(body)
```

```
        message.ack()
```

```
    def listen(self):
```

```
        with Connection(self.backend.broker, transport_options={'region': self.backend.region}) as connection:
```

```
            with Consumer(connection, self.queue, callbacks=[self.callback], accept=[self.backend.serializer]):
```

```
                for _ in eventloop(connection):
```

```
                    pass
```

```
# What a simplified worker looks like
```

```
# Broker URI stored on Backend object, looks like:
```

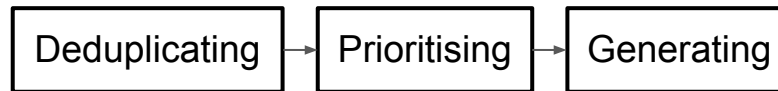
```
# sqs://{s3_key}:{s3_secret}@
```

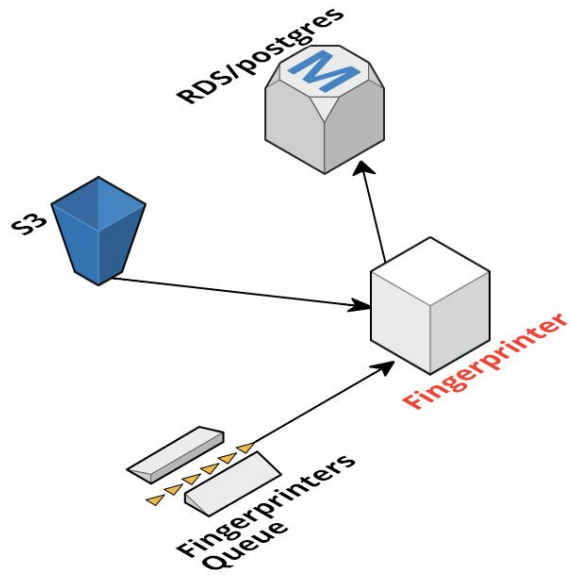
```
consumer = KombuConsumer(backend, handler=downloader.downloader_callback)
```

```
consumer.listen()
```

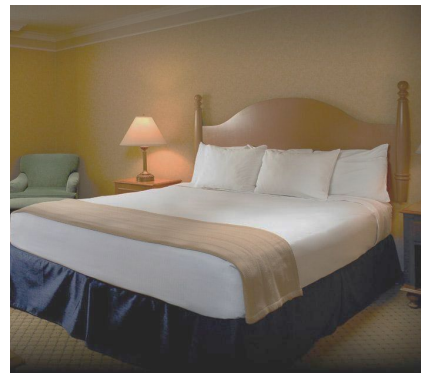
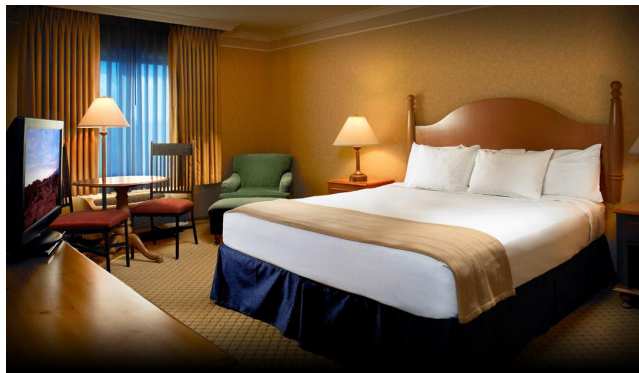


Fingerprinting





Are those images the same?



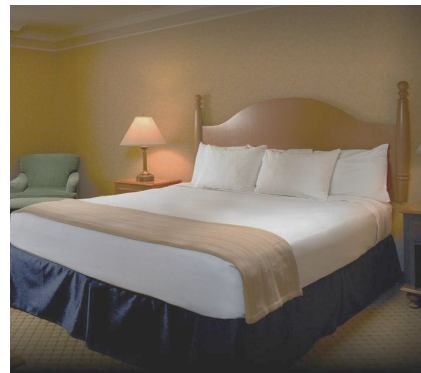
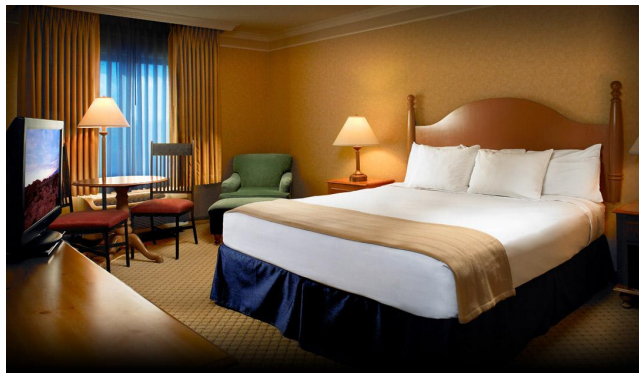


Are those images the same?



Are those images the same?

Yes, they are



```
import imagehash

def fingerprint_callback(queued_image):
    blob = download_image_blob(queued_image.basename)
    image = Image.open(BytesIO(blob))
    result = cropped_hash(image, imagehash.phash)
    store_hashes(queued_image.image_id, result)
```

```
import imagehash

def fingerprint_callback(queued_image):
    blob = download_image_blob(queued_image.basename)
    image = Image.open(BytesIO(blob))
    result = cropped_hash(image, imagehash.dhash)
    store_hashes(queued_image.image_id, result)
```

```

def cropped_hash(image, algorithm, steps=range(0, 51, 10)):
    result = []
    w, h = image.size

    # We want to cut by steps % of the image (default 0%, 10%...50%), which
    # means we need to cut half of that from each side:
    # | N%/2 |           | N%/2 |
    # +-----+-----+-----+
    # |      :           :      | N%/2
    # +- - - +-----+ - - -+-
    # |      |           |      |
    # |      |           |      |
    # +- - - +-----+ - - -+-
    # |      :           :      | N%/2
    # +-----+-----+-----+

    for x in steps:
        x_band = x * w / 200
        for y in steps:
            y_band = y * h / 200
            with image.crop((x_band, y_band, w-x_band, h-y_band)) as sub_image:
                sub_hash = algorithm(sub_image)
                result.append(hash_to_int(sub_hash))

```

```

def cropped_hash(image, algorithm, steps=range(0, 51, 10)):
    result = []
    w, h = image.size

    # We want to cut by steps % of the image (default 0%, 10%...50%), which
    # means we need to cut half of that from each side:
    # | N%/2 |           | N%/2 |
    # +-----+-----+-----+
    # |      :           :      | N%/2
    # +- - - +-----+ - - -+-
    # |      |           |      |
    # |      |           |      |
    # +- - - +-----+ - - -+-
    # |      :           :      | N%/2
    # +-----+-----+-----+

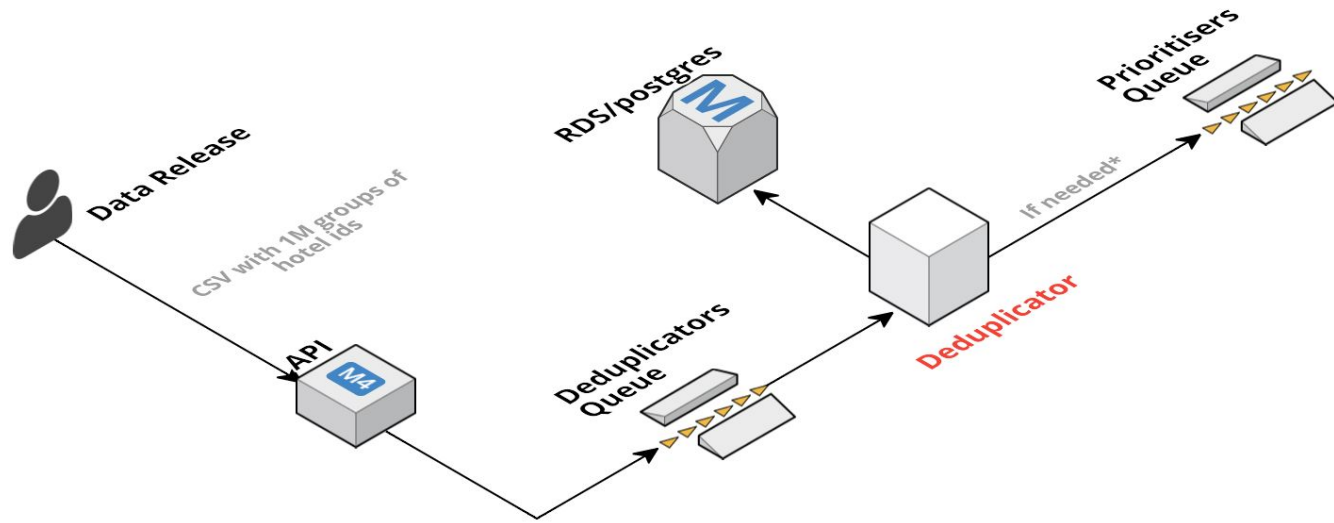
    for x in steps:
        x_band = x * w / 200
        for y in steps:
            y_band = y * h / 200
            with image.crop((x_band, y_band, w-x_band, h-y_band)) as sub_image:
                sub_hash = algorithm(sub_image)
                result.append(hash_to_int(sub_hash))

```



Deduplication Time





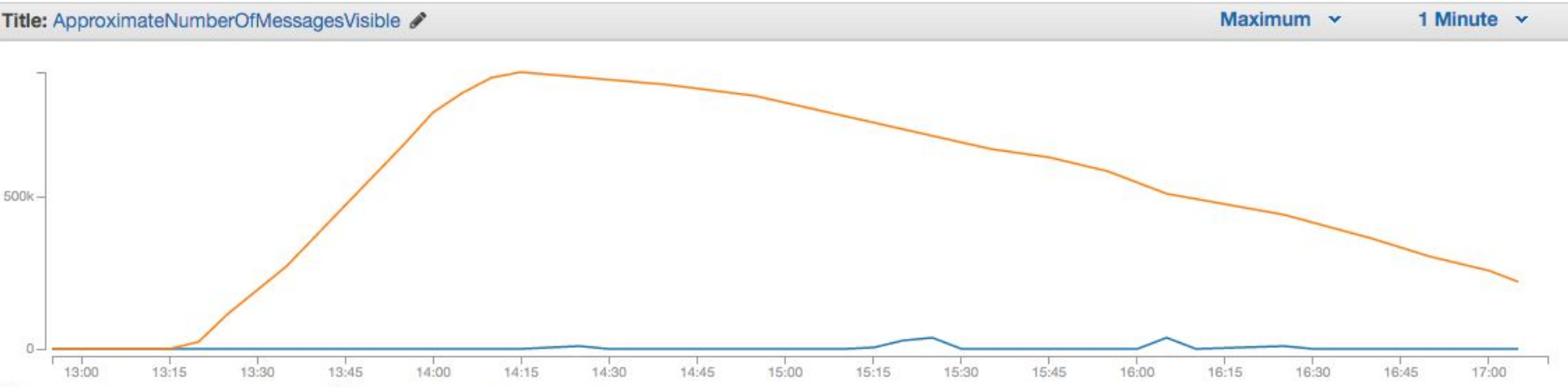
What is a “group”?

id	dr	group_id	status	elements
10220644	70	2053940	done	17,454995,4,501446,13,1003798
10220643	70	2054825	done	34,102912,4,363502,28,37152,46
10220642	70	2054040	done	32,E7D860D5-0576-4BBD-AA26-418
10220641	70	2054394	done	129,243137,131,41113,5,PHAIGV,
10220640	70	2054826	done	33,NOGTIEHTBUMBAVEEBVCV,4,3153

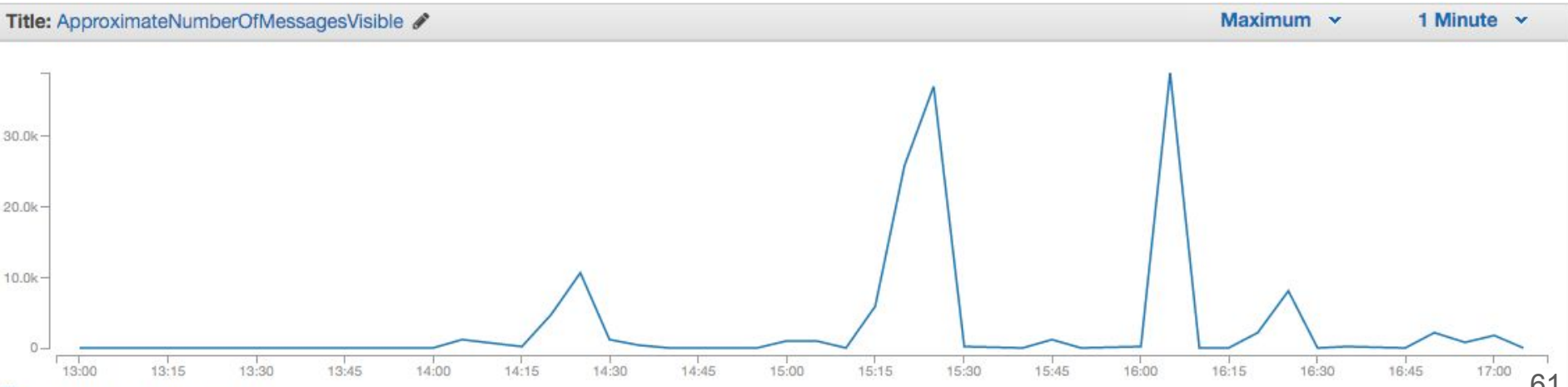
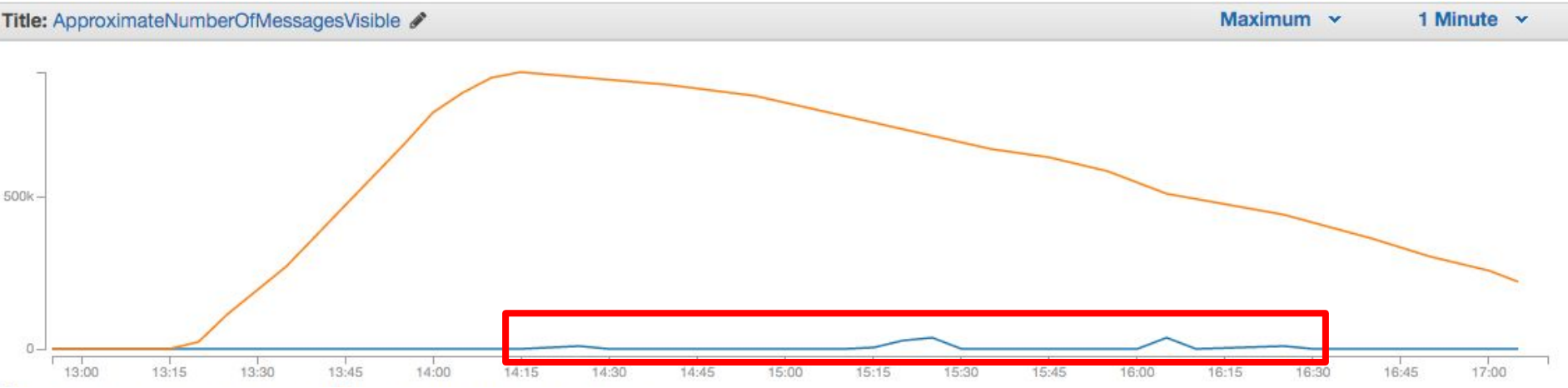
What is a “group”?

id	dr	group_id	status	elements
10220644	70	2053940	done	17,454995,4,501446,13,1003798
10220643	70	2054825	done	34,102912,4,363502,28,37152,46
10220642	70	2054040	done	32,E7D860D5-0576-4BBD-AA26-418
10220641	70	2054394	done	129,243137,131,41113,5,PHAIGV,
10220640	70	2054826	done	33,NOGTIEHTBUMBAVEEBVCV,4,3153

“If needed”



“If needed”



Hotel Group 123

$[(\text{partner_id}_1, \text{accomodation_id}_1), \dots, (\text{partner_id}_n, \text{accomodation_id}_n)]$

Hotel Group 123

$[(\text{partner_id}_1, \text{accomodation_id}_1), \dots, (\text{partner_id}_n, \text{accomodation_id}_n)]$



Hotel Group 123

$[(\text{partner_id}_1, \text{accomodation_id}_1), \dots, (\text{partner_id}_n, \text{accomodation_id}_n)]$

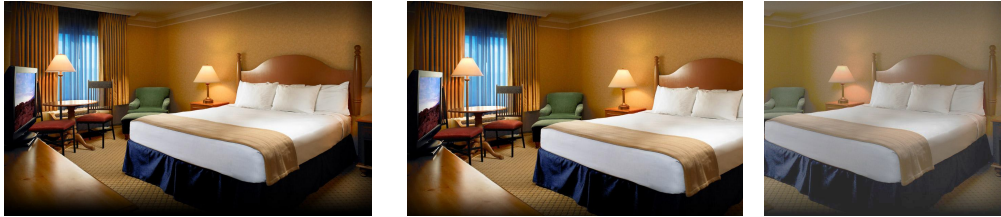


Image Group 456

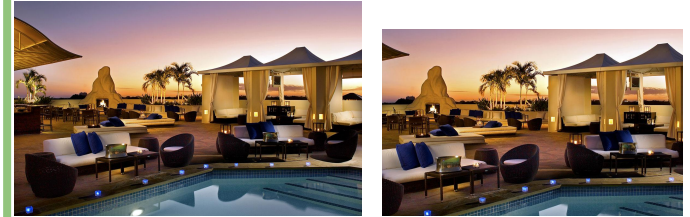


Image Group 203

Hotel Group 123

$[(\text{partner_id}_1, \text{accomodation_id}_1), \dots, (\text{partner_id}_n, \text{accomodation_id}_n)]$

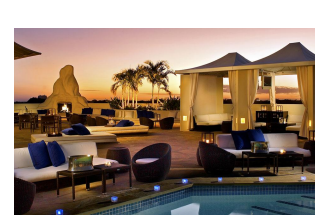


Image Group 456

Image Group 203

Hotel Group 123 has two image groups:
[456, 203]

```
def deduplicate_group(all_hotel_images):
    all_images = set(all_hotel_images)
    groups = []
    while all_images:
        seed_image = all_images.pop()
        group = {seed_image}
        new_additions = {seed_image}
        while new_additions:
            me = new_additions.pop()
            for other in all_images:
                if is_same_picture(me.hashes, other.hashes):
                    group.add(other)
                    new_additions.add(other)
            all_images = all_images - group
        groups.append(group)
    return groups
```

```
def deduplicate_group(all_hotel_images):  
    all_images = set(all_hotel_images)  
    groups = []  
    while all_images:  
        seed_image = all_images.pop()  
        group = {seed_image}  
        new_additions = {seed_image}  
        while new_additions:  
            me = new_additions.pop()  
            for other in all_images:  
                if is_same_picture(me.hashes, other.hashes):  
                    group.add(other)  
                    new_additions.add(other)  
            all_images = all_images - group  
        groups.append(group)  
    return groups
```

```
def deduplicate_group(all_hotel_images):
    all_images = set(all_hotel_images)
    groups = []
    while all_images:
        seed_image = all_images.pop()
        group = {seed_image}
        new_additions = {seed_image}
        while new_additions:
            me = new_additions.pop()
            for other in all_images:
                if is_same_picture(me.hashes, other.hashes):
                    group.add(other)
                    new_additions.add(other)
            all_images = all_images - group
        groups.append(group)
    return groups
```

```
def deduplicate_group(all_hotel_images):
    all_images = set(all_hotel_images)
    groups = []
    while all_images:
        seed_image = all_images.pop()
        group = {seed_image}
        new_additions = {seed_image}
        while new_additions:
            me = new_additions.pop()
            for other in all_images:
                if is_same_picture(me.hashes, other.hashes):
                    group.add(other)
                    new_additions.add(other)
            all_images = all_images - group
        groups.append(group)
    return groups
```

```
def deduplicate_group(all_hotel_images):
    all_images = set(all_hotel_images)
    groups = []
    while all_images:
        seed_image = all_images.pop()
        group = {seed_image}
        new_additions = {seed_image}
        while new_additions:
            me = new_additions.pop()
            for other in all_images:
                if is_same_picture(me.hashes, other.hashes):
                    group.add(other)
                    new_additions.add(other)
            all_images = all_images - group
        groups.append(group)
    return groups
```

```
def is_same_picture(left_hashes, right_hashes):  
    for left_hash in left_hashes:  
        for right_hash in right_hashes:  
            ham_dist = hamdist(left_hash, right_hash)  
            if ham_dist < threshold:  
                return True  
    return False
```

```
def is_same_picture(left_hashes, right_hashes):  
    for left_hash in left_hashes:  
        for right_hash in right_hashes:  
            ham_dist = hamdist(left_hash, right_hash)  
            if ham_dist < threshold:  
                return True  
    return False
```

How do you tune this step?

Guarantees are needed

You build a corpus.

IMPURE
INCOMPLETE



URL MD5: 4fdidd1b909e3555e1af7b0af5a376f5
Group ID: 213



IMPURE
INCOMPLETE



IMPURE
INCOMPLETE



IMPURE
COMPLETE



IMPURE
COMPLETE



IMPURE
COMPLETE



IMPURE
COMPLETE

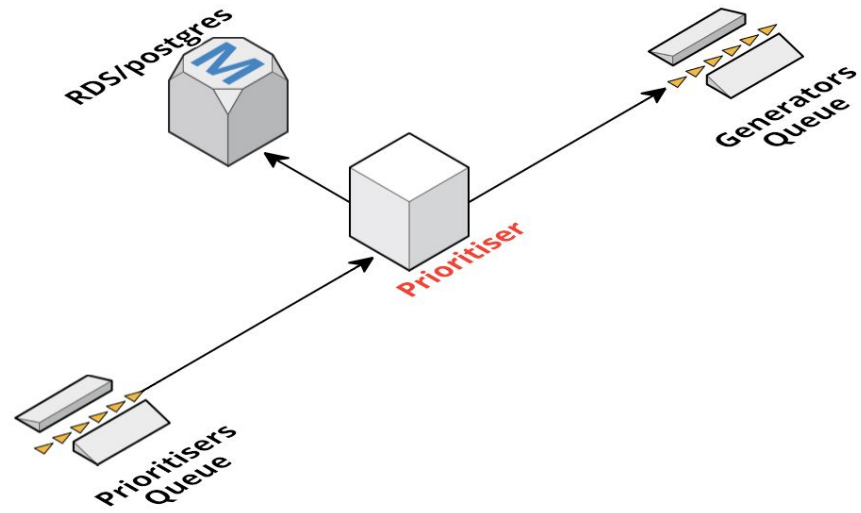


Measure	Run #1	Run #2	Run #3	Run #4	Run #5	Run #6	Run #44	Run #157
General measures								
Homogeneity	83.10%	63.53%	33.76%	99.59%	98.01%	92.20%	99.50%	99.50%
Completeness	99.02%	96.97%	97.74%	99.11%	99.15%	99.27%	99.13%	99.13%
V-measure	90.37%	76.77%	50.18%	99.35%	98.58%	95.60%	99.31%	99.31%
A-score	55.52%	27.61%	10.79%	95.77%	91.06%	72.16%	95.86%	95.86%
Group measures								
Groups complete	88.92%	79.20%	86.03%	91.14%	91.35%	93.12%	91.29%	91.29%
Groups pure	83.81%	97.55%	98.32%	98.23%	93.88%	90.42%	97.82%	97.82%
Image release								
Image loss	29.46%	34.47%	64.13%	-1.80%	5.01%	18.44%	-1.20%	-1.20%
Duplicate probability	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%
Average duplicates	4.14%	10.15%	5.13%	4.46%	3.87%	2.81%	4.26%	4.26%



Prioritisation



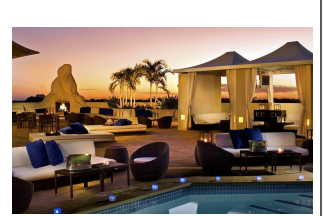
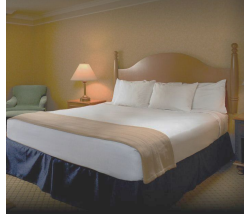


Hotel Group 123

[ImageGroup₄₀₆, ImageGroup₂₀₃]

Hotel Group 123

[ImageGroup₄₀₆, ImageGroup₂₀₃]

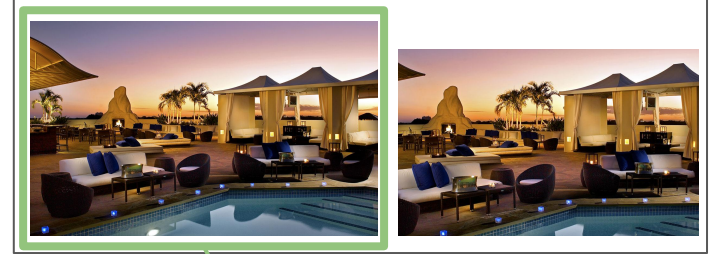


Hotel Group 123

[ImageGroup₄₀₆, ImageGroup₂₀₃]



“Best image”



“Best Image”

Hotel Group 123

[ImageGroup₄₀₆, ImageGroup₂₀₃]



“Best image”



“Best Image”



“Best order”



1

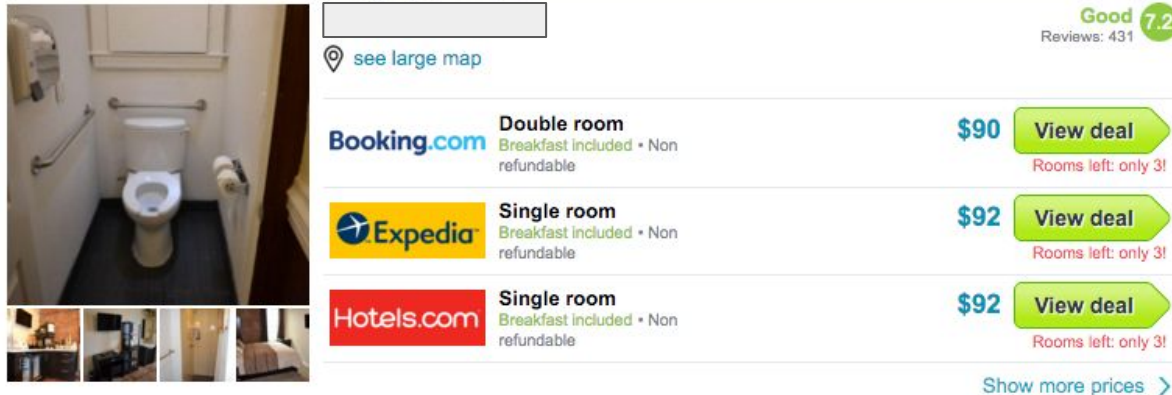


2

reaches production

What could go wrong?

What could go wrong?



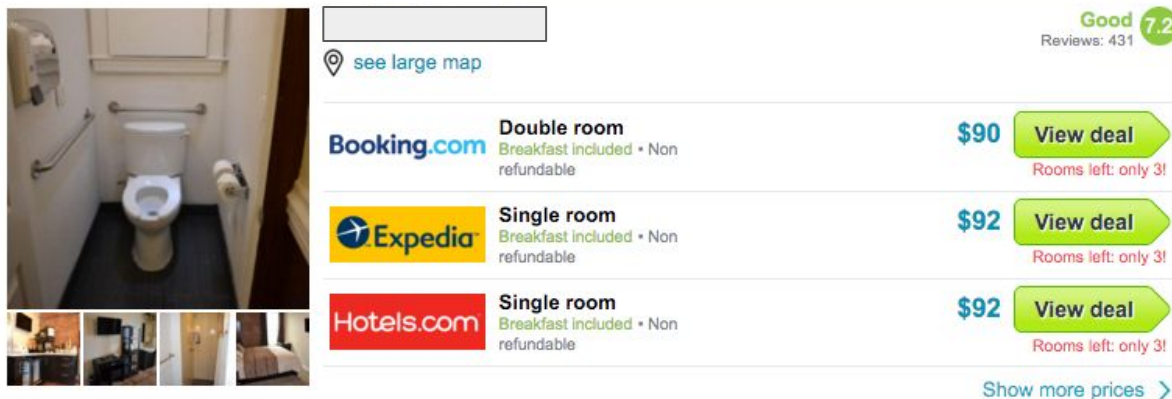
[see large map](#)

Good 7.2
Reviews: 431

Booking.com	Double room Breakfast included • Non refundable	\$90	View deal Rooms left: only 3!
Expedia	Single room Breakfast included • Non refundable	\$92	View deal Rooms left: only 3!
Hotels.com	Single room Breakfast included • Non refundable	\$92	View deal Rooms left: only 3!

[Show more prices >](#)

What could go wrong?



The screenshot displays a hotel booking interface. On the left, there is a large photo of a bathroom with a toilet and a smaller strip of four photos below it showing different room interiors. To the right of the photos is a search bar and a 'see large map' link. Further right, a green circle indicates a 'Good 7.2' rating with 'Reviews: 431'. Below this, a list of room options is shown:

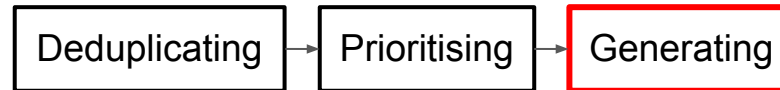
Provider	Room Type	Price	Deal Button	Availability
Booking.com	Double room Breakfast included • Non refundable	\$90	View deal	Rooms left: only 3!
Expedia	Single room Breakfast included • Non refundable	\$92	View deal	Rooms left: only 3!
Hotels.com	Single room Breakfast included • Non refundable	\$92	View deal	Rooms left: only 3!

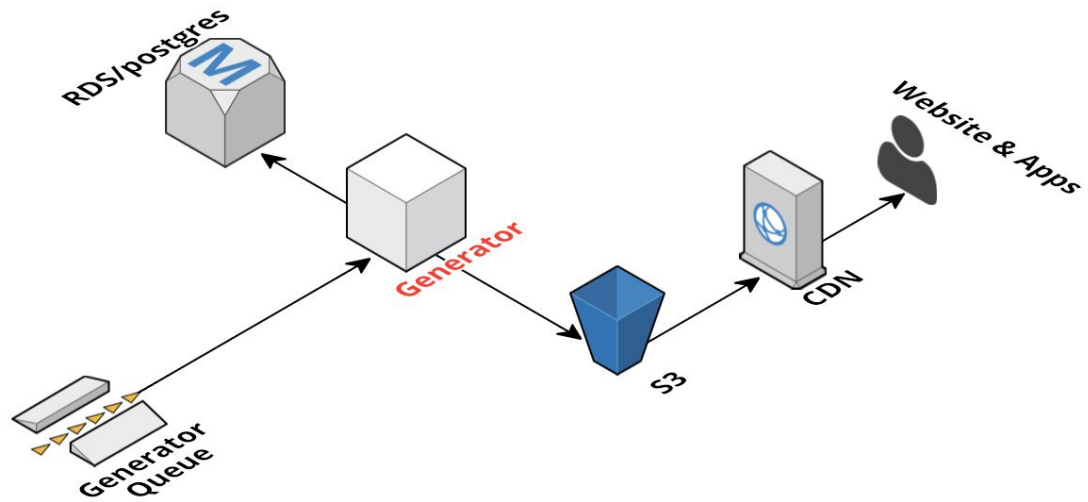
At the bottom right, there is a link 'Show more prices >'.

- * Detect features, prioritise based on that.
- * Tools to manually fix data.



Generation





```

from PIL import ImageEnhance

def scalefit(image):
    sz = image.size
    bs = best_size(image)

    # A bit of math:
    # - if we scale to fit width, the scaling factor is: scale_width = bs[0]/sz[0]
    # - if we scale to fit height, the scaling factor is: scale_height = bs[1]/sz[1]
    # We want to scale to the smaller of them (so that the image fits in both), so we scale to width if:
    # scale_width < scale_height => bs[0]/sz[0] < bs[1]/sz[1] => bs[0]*sz[1] < bs[1]*sz[0]
    # Having it in this form means no floats are needed; all in integers.
    if bs[0] * sz[1] < bs[1] * sz[0]:
        # Scale to width
        w = bs[0]
        h = sz[1] * bs[0] / sz[0]
    else:
        # Scale to height
        w = sz[0] * bs[1] / sz[1]
        h = bs[1]

    return image.resize((w, h), "bilinear")

def contrast(image, value):
    enhancer = ImageEnhance.Contrast(image)
    return enhancer.enhance(float(value))

```

```

from PIL import ImageEnhance

def scalefit(image):
    sz = image.size
    bs = best_size(image)

    # A bit of math:
    # - if we scale to fit width, the scaling factor is: scale_width = bs[0]/sz[0]
    # - if we scale to fit height, the scaling factor is: scale_height = bs[1]/sz[1]
    # We want to scale to the smaller of them (so that the image fits in both), so we scale to width if:
    # scale_width < scale_height => bs[0]/sz[0] < bs[1]/sz[1] => bs[0]*sz[1] < bs[1]*sz[0]
    # Having it in this form means no floats are needed; all in integers.
    if bs[0] * sz[1] < bs[1] * sz[0]:
        # Scale to width
        w = bs[0]
        h = sz[1] * bs[0] / sz[0]
    else:
        # Scale to height
        w = sz[0] * bs[1] / sz[1]
        h = bs[1]

    return image.resize((w, h), "bilinear")

def contrast(image, value):
    enhancer = ImageEnhance.Contrast(image)
    return enhancer.enhance(float(value))

```

```

from PIL import ImageEnhance

def scalefit(image):
    sz = image.size
    bs = best_size(image)

    # A bit of math:
    # - if we scale to fit width, the scaling factor is: scale_width = bs[0]/sz[0]
    # - if we scale to fit height, the scaling factor is: scale_height = bs[1]/sz[1]
    # We want to scale to the smaller of them (so that the image fits in both), so we scale to width if:
    # scale_width < scale_height => bs[0]/sz[0] < bs[1]/sz[1] => bs[0]*sz[1] < bs[1]*sz[0]
    # Having it in this form means no floats are needed; all in integers.
    if bs[0] * sz[1] < bs[1] * sz[0]:
        # Scale to width
        w = bs[0]
        h = sz[1] * bs[0] / sz[0]
    else:
        # Scale to height
        w = sz[0] * bs[1] / sz[1]
        h = bs[1]

    return image.resize((w, h), "bilinear")

def contrast(image, value):
    enhancer = ImageEnhance.Contrast(image)
    return enhancer.enhance(float(value))

```

And that's it for the pipeline

```
image_db=> select status, count(*) from image group by status;
```

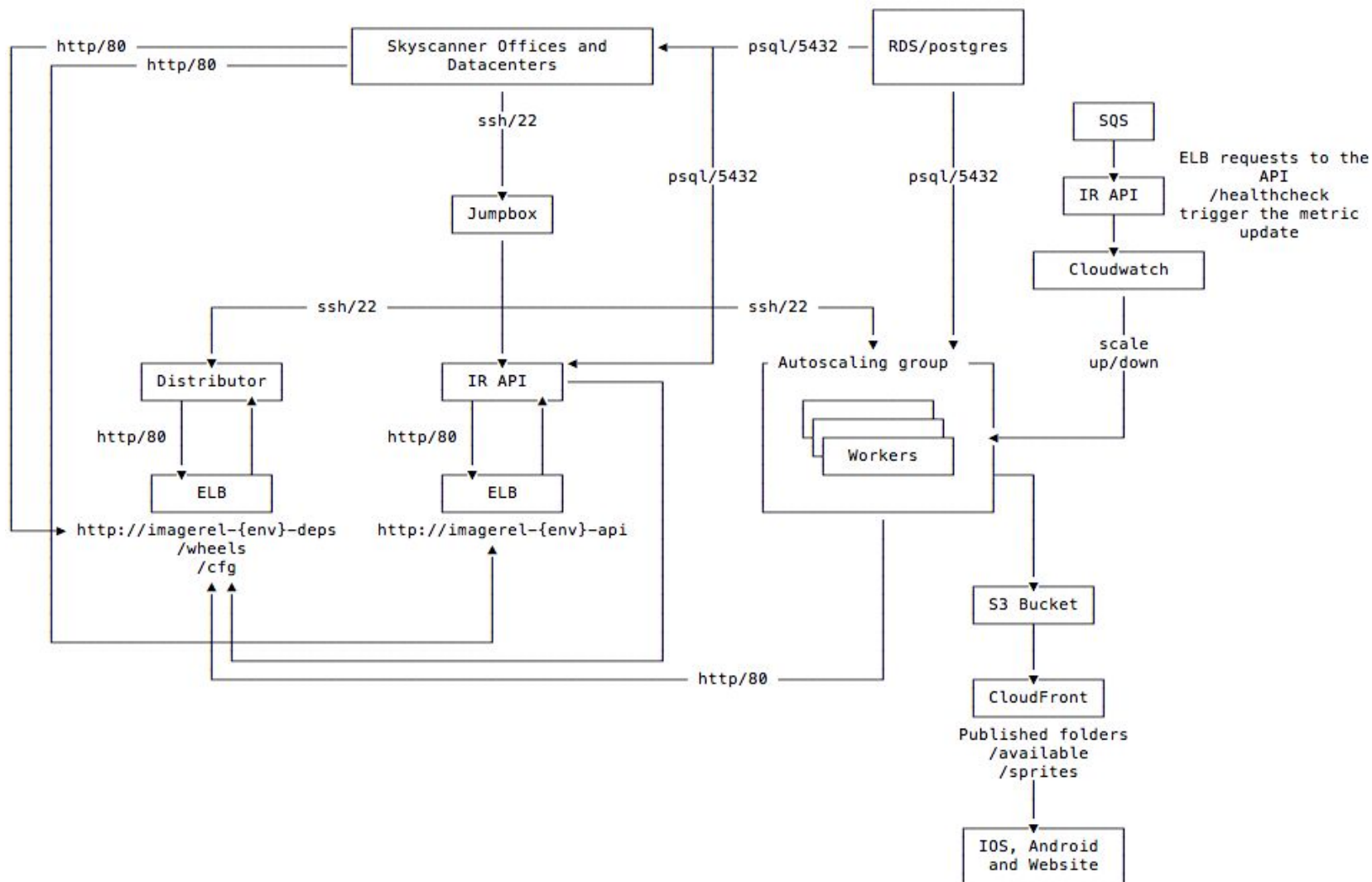
status	count
production	7260896
downloaded	48729
filtered	3560235
error_download	1332522
available	16665606
not_found	2248861
new	10387

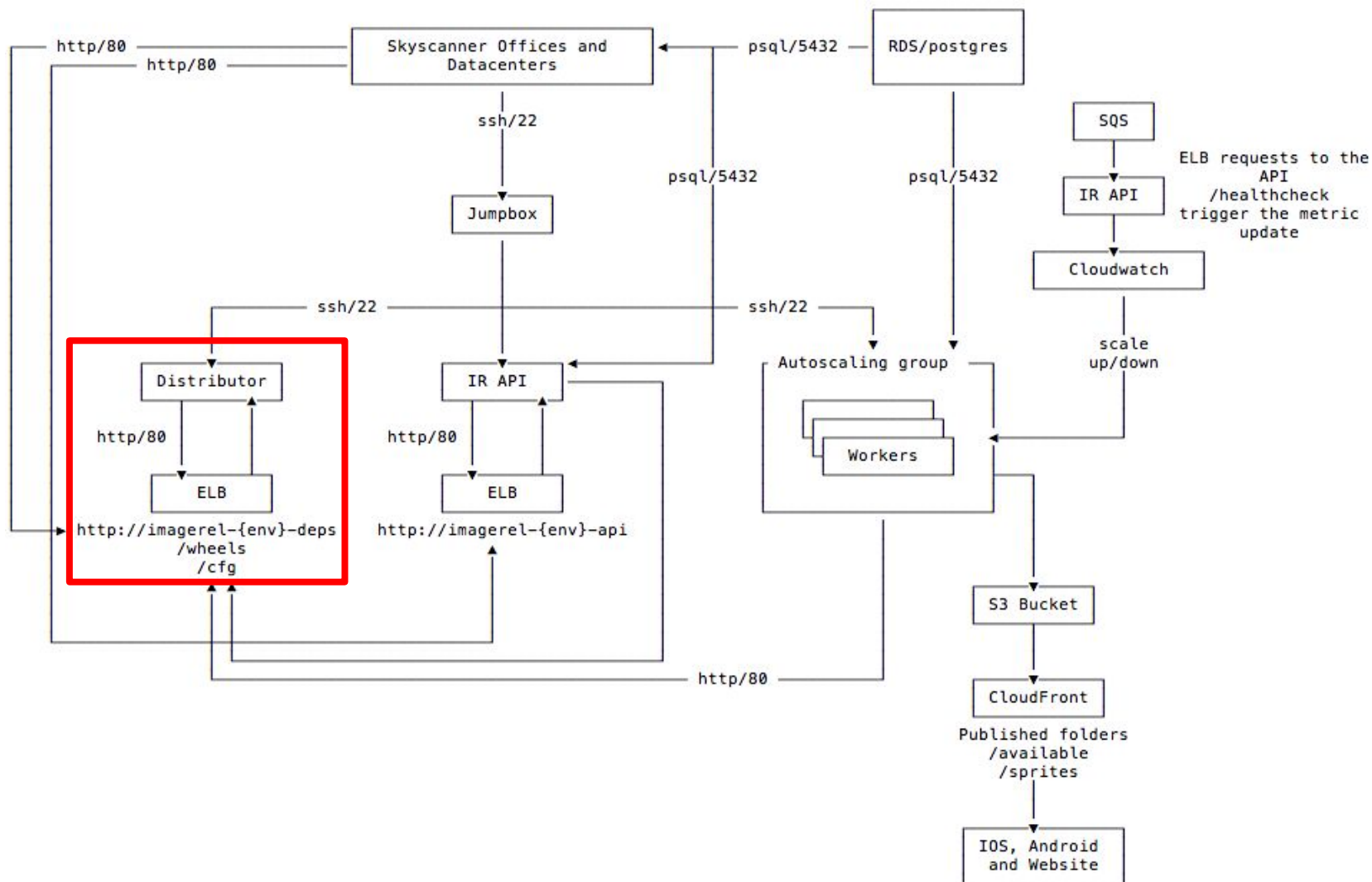
(7 rows)

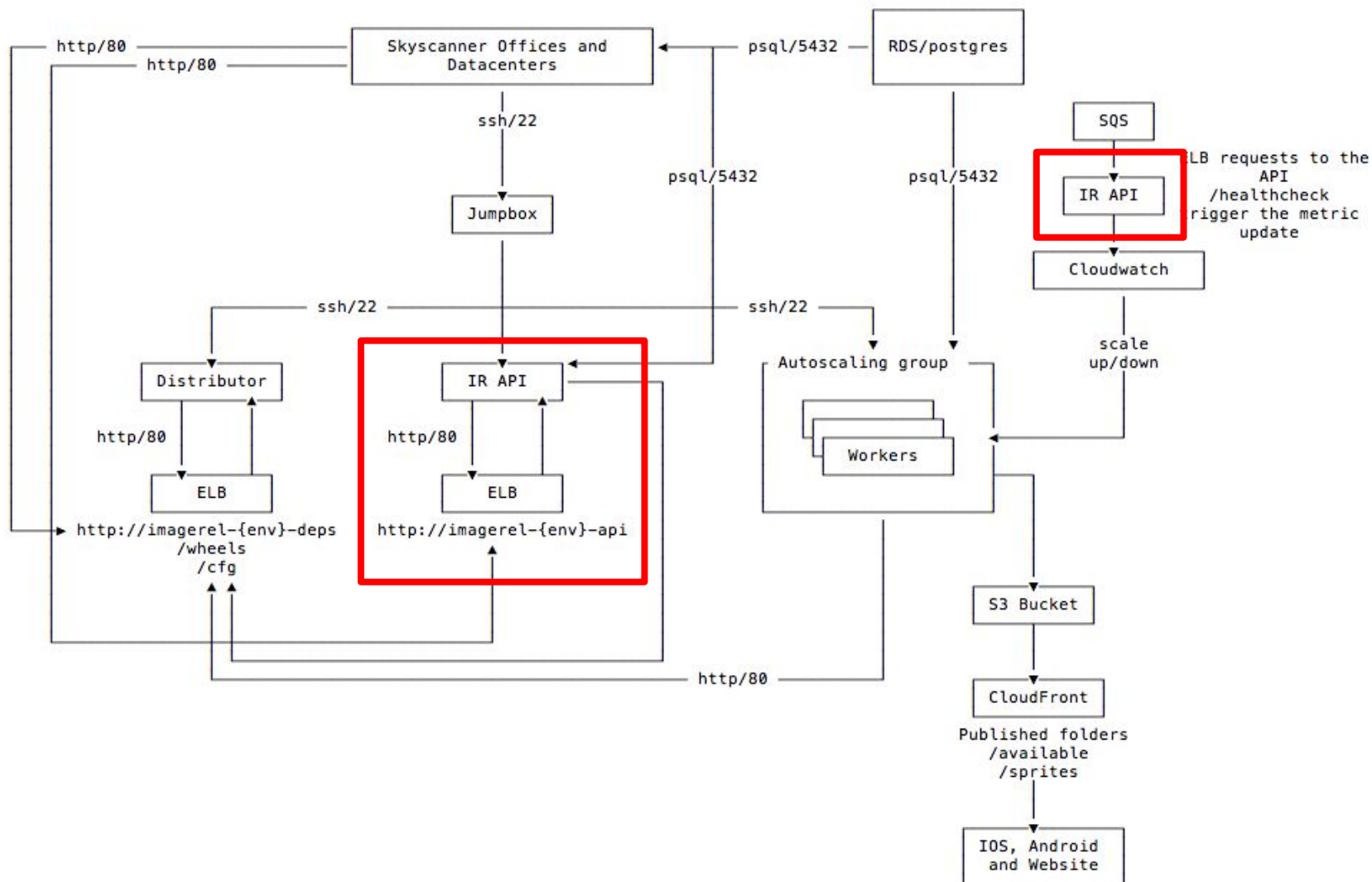
```
image_db=> select status, count(*) from image group by status;
```

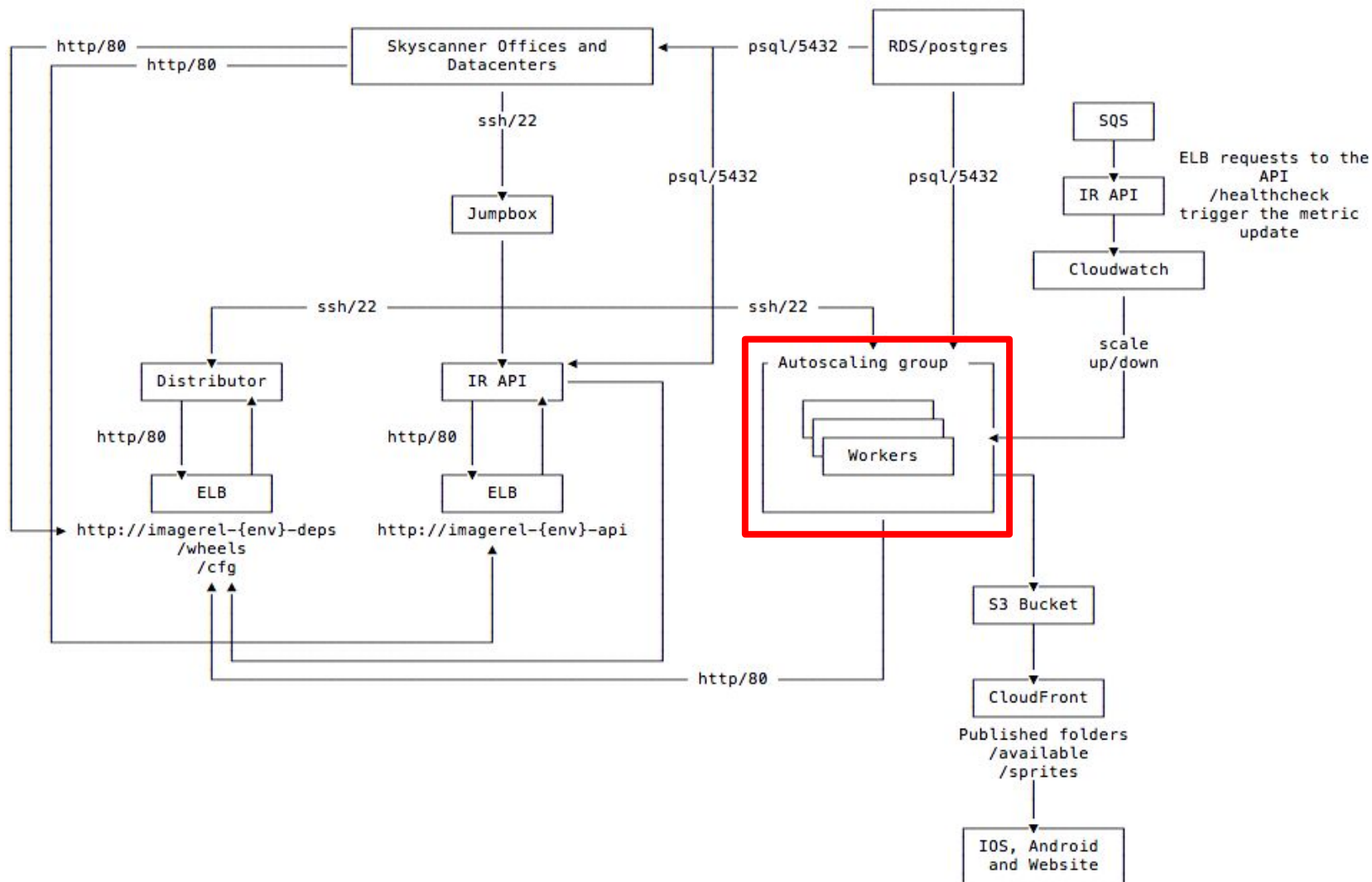
status	count
production	7260896
downloaded	48729
filtered	3560235
error_download	1332522
available	16665606
not_found	2248861
new	10387

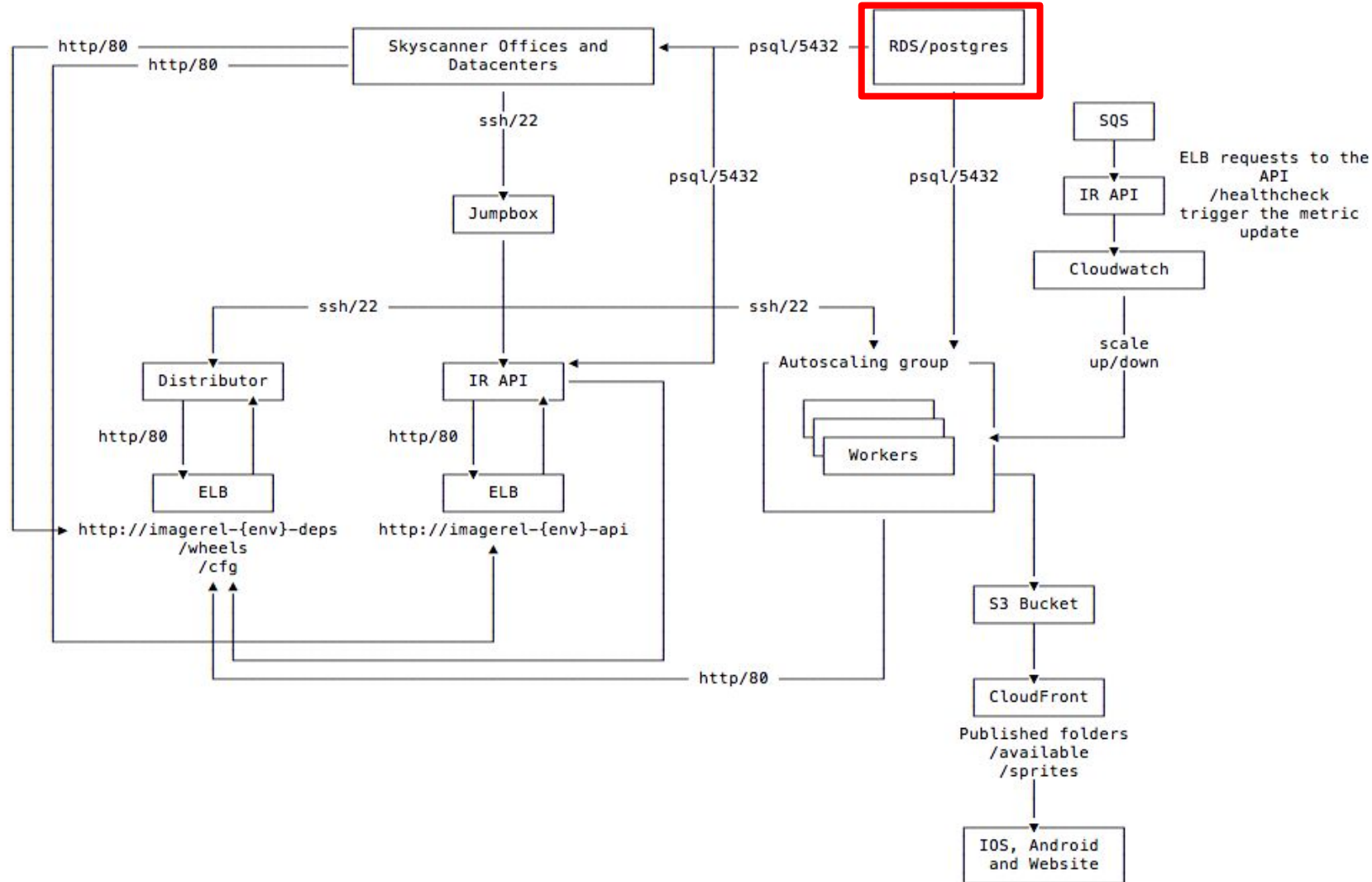
(7 rows)

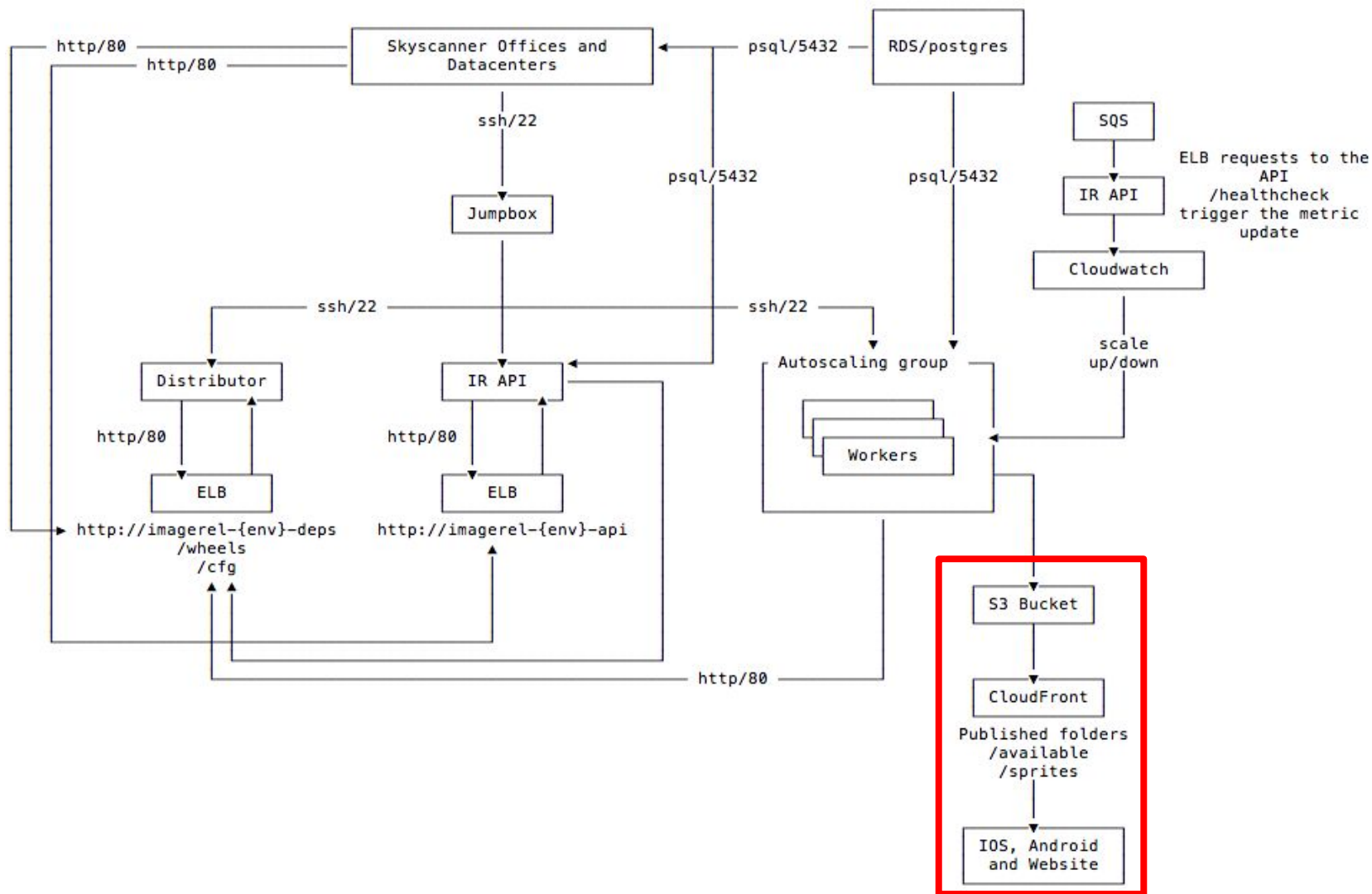














Thanks for listening

Any questions?