# Building an image processing pipeline in Python

Franck Chastagnol, PyCon 2013

# Agenda

- Introduction
- Architecture
- Upload
- Image pre-processing
- OCR
- Structured data extraction
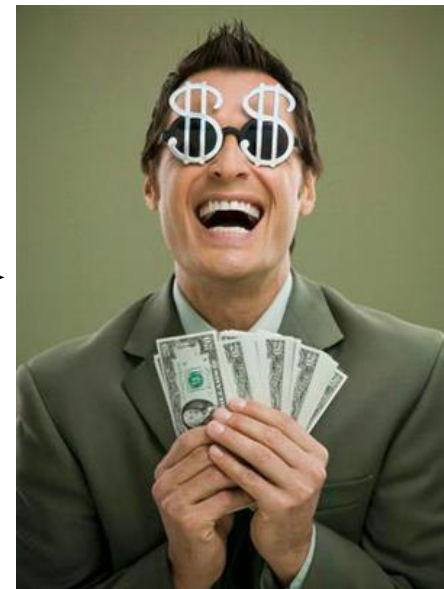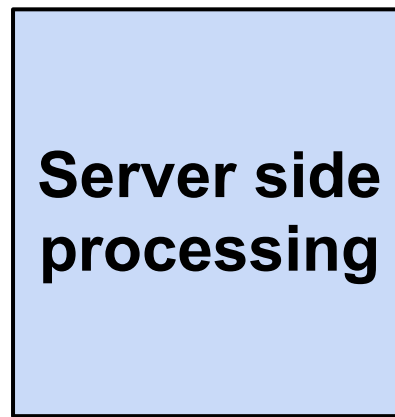- Error handling / re-processing
- Q&A

# **Introduction**

- Background


- Today's case study
  - Image processing pipeline built for Endorse.com

# Endorse.com mobile app

- Reward for buying specific brand products

- Shop anywhere, upload pic of receipt, get $$



**Server side processing**

# Pics of receipts are... fun ! (1)

# Pics of receipts are... fun ! (2)

# Pics of shopping receipts are... challenging to process !

- Taken in various environment, lighting

- Resolution varies depending on device

- Quality of receipt printers varies greatly

- It is not english

- Diff. format, no universal UPC / shortnames

## **Agenda**

- Introduction
- <span style="color:red">Architecture</span>
- Upload
- Image pre-processing
- OCR
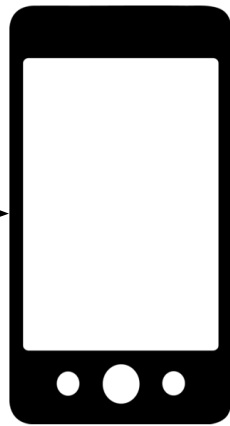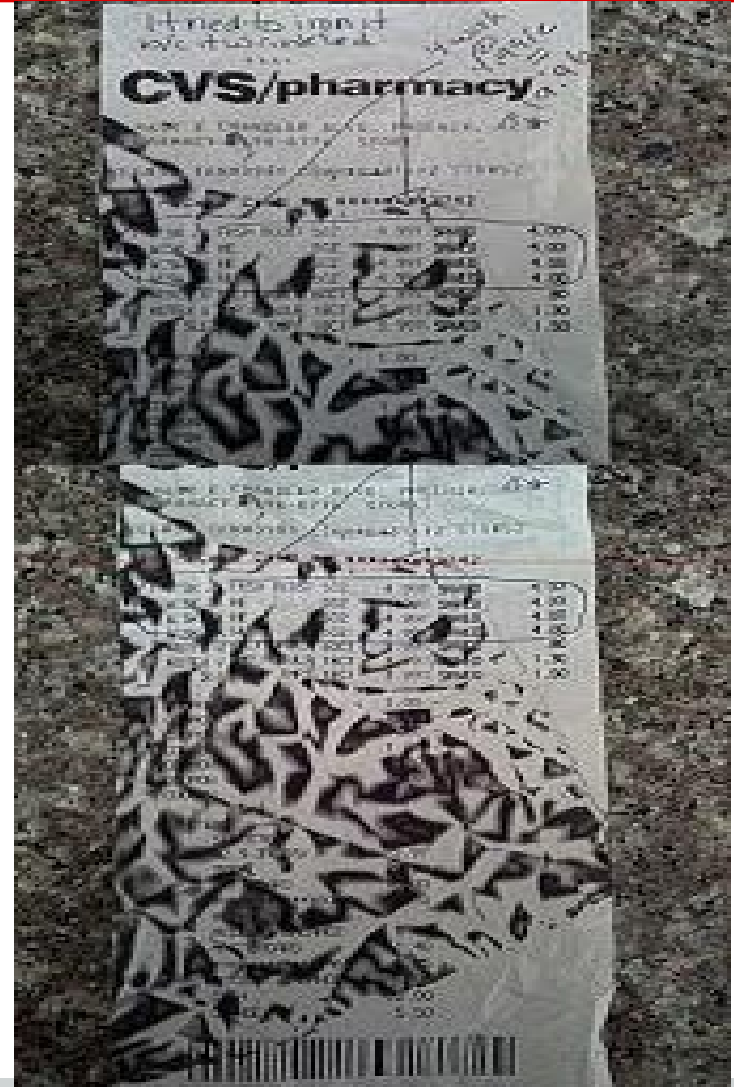- Structured data extraction
- Error handling / re-processing
- Q&A

# Technologies

- Common
  - Server Central cloud
  - Linux (ubuntu)
  - Nginx load balancer
  - Tornado app server
  - Python 2.7
  - Redis
  - S3 storage

- Web
  - Mako templates
  - MySQL

- Receipt processing
  - OpenCV
  - NumPy
  - IMagick
  - Tesseract OCR

- Data mining
  - MongoDB
  - Hadoop

# System diagram

# Pipeline

**Receipt Image**



**Structured Doc**

| Pre-Processing | OCR | Parsing | Scoring | Best Result Selection |

Multi-Pass

Retailer = WALMART
Date = 03/11/73 11:00pm
Address: Limoges, FR
Phone #: 650-123-4567

Item1 = 1 x OREO ($1.99)
Item2 = 2 x COKE ($0.99)
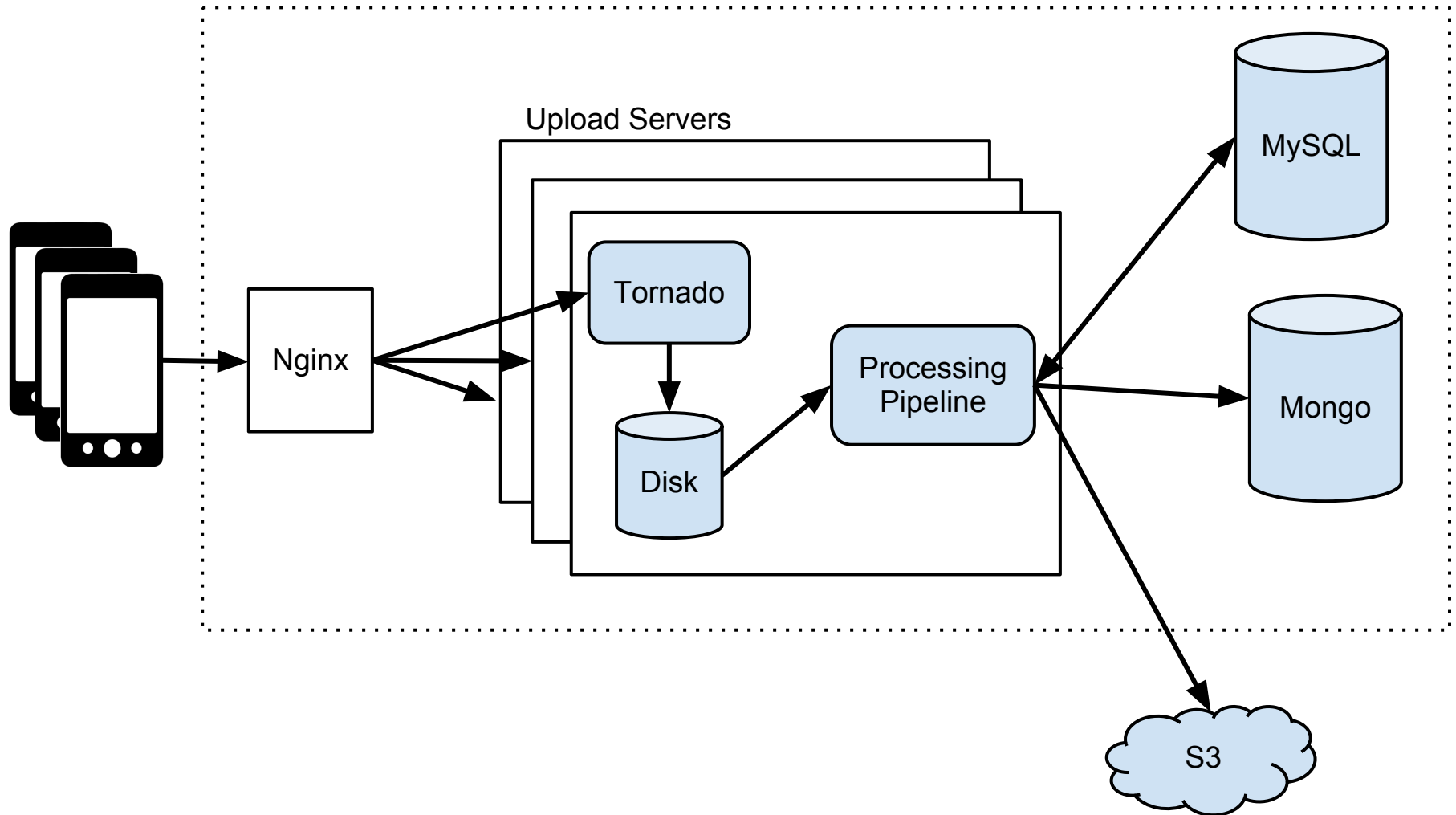Item3 = 1 x MILK  ($3.50)

TAX = $0.87
TOTAL = $10.73

# Agenda

- Introduction
- Architecture
- Upload
- Image pre-processing
- OCR
- Structured data extraction
- Error handling / re-processing
- Q&A

# Mobile uploads

- Images are not small: ~1MB per segment

- Mobile data connection
  - can be spotty
  - upload bandwidth varies

- Ensuring high upload success rate:
  - App capable of re-trying in background
  - Simple and resumable APIs

# Upload workflow



**Server**

**1**

START(nb_segment)

Upload UID

- Insert row in upload table

**2**

UPLOAD(UID, segment_nb, img)

[ segment_received_list ]

- Store image file
- Update upload row

Repeat for each segment

# Upload - scalability

- Nginx
  - sticky session module

- Tornado writes img files to local disk

- Job picks up img files once upload finished
  - Store originals in S3
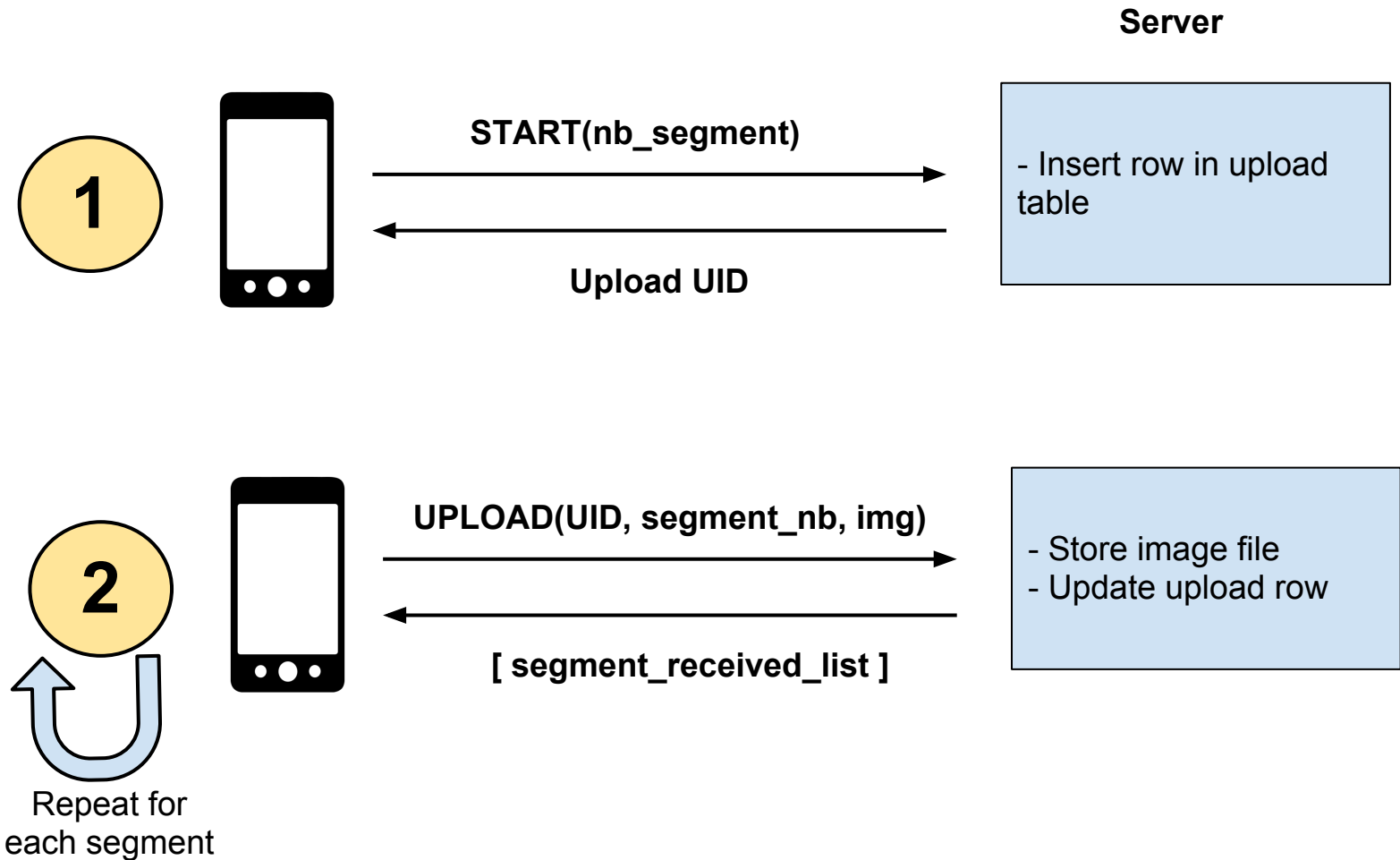  - Run pipeline

# **Agenda**

- Introduction
- Architecture
- Upload
- <span style="color:red">Image pre-processing</span>
- OCR
- Structured data extraction
- Error handling / re-processing
- Q&A

# But why ??

- OCR is a solved problem... for book scans

- Clean b&w 300 dpi images of book pages scanned under perfect conditions
  => recognition rate = 95% to 99%

- Wrinkled paper, bad quality print, inconsistent lighting, noise, angle, etc...
  => recognition rate = ~25% or less

# Pre-processing steps

- From color to b&w
  - unblur / sharpen filters
  - un-highlight color regions
  - adaptive thresholding

- Cropping
  - The carpet problem

- Extracting lines
  - OCR does poorly on non-straight lines
  - Lines recognition

=> OpenCV + Numpy is great

# Image pre-processing example

**Original**

**Cropping**

**Lines extract.**

# **Agenda**

- Introduction
- Architecture
- Upload
- Image pre-processing
- OCR
- Structured data extraction
- Error handling / re-processing
- Q&A

# Tesseract

- Tesseract
  - Open source
  - Started at HP in the 90s
  - Google uses it for Book scan project
  - C++ core engine, APIs
  - Python bindings

# OCR Training

- Shopping receipt fonts are not standard !
  - Training process is no fun
    - scanned various receipt types
    - extracted each letter from alphabet
    - generated synthetic receipts used for training

- Shopping receipts are not english !
  - OCR uses dictionaries to improve its output quality:
    - words dictionary with frequency in language
    - word pairs probability
    - punctuations / non alpha character rules

# **Agenda**

---

- Introduction
- Architecture
- Upload
- Image pre-processing
- OCR
- <span style="color:red">Structured data extraction</span>
- Error handling / re-processing
- Q&A

# You got text, now what ?

```
( 903 ) 657 - 5707
MANAGER R0BERT JACKSON
2121 US HIGHWAY 79 S
HENDERSON TX 75654
ST# 0165 DP# 00000018 TE# 08 TR# 06834
ELECTROLYTE 007874206418 F 3.14 X
GATORADE 005200032016 F 1.00 X
YOGURT MELT 001500004730 F 2.48 N
RTD APPLE 002800098443 F 2.38 N
BREAD 007874298114 F 1.50 0
FFBRFZE 003700025221 4.97 X
2PK BK SLP B 004721365070 5.00 T
SVBT0TAL 38. 16
TAX1 8.250 X 1.24
TOTAL 39 .40
CASH TEND 100.40
CH8NGE DVE 61.00
TC# 3312 2198 4945 1493 8462
03/05/13 16:47.18
```

- Parser
  - In: Text
  - Out: Structured doc

- Receipt
  - Store
  - List
    - Items (UPC, price)
  - SubTotal
  - Taxes
  - Total

# Regex = headache

- Wide variety of mistakes in OCR output makes using regex hard / impossible

- Levenshtein distance is your friend
  - Similarity score between 2 strings (e.g. nb edits)
  - Pure Python implementation is slow. C lib + Python bindings faster

- "fuzzy matcher"
  - **Pattern**:     "%s TAX  (%d.d%%) = $%d.%d ON $%d.%d"
  - **Input**:     "CA T8X (8.0%) = $4.00 ON $50.00
  - **Output**:     Score = 1 (e.g. 1 edit)

# Extracting + storing structured data

- Shopping receipts come in a variety of format
  - Specific parsers for most common formats
  - Generic parser for others
  - Store document in Mongo

- Mongo DB benefits
  - schemaless
  - map-reduce capabilities makes it a scalable data-mining solution

# **Agenda**

- Introduction
- Workflow
- Upload
- Image pre-processing
- OCR
- Structured data extraction
- <span style="color:red">Error handling/re-processing</span>
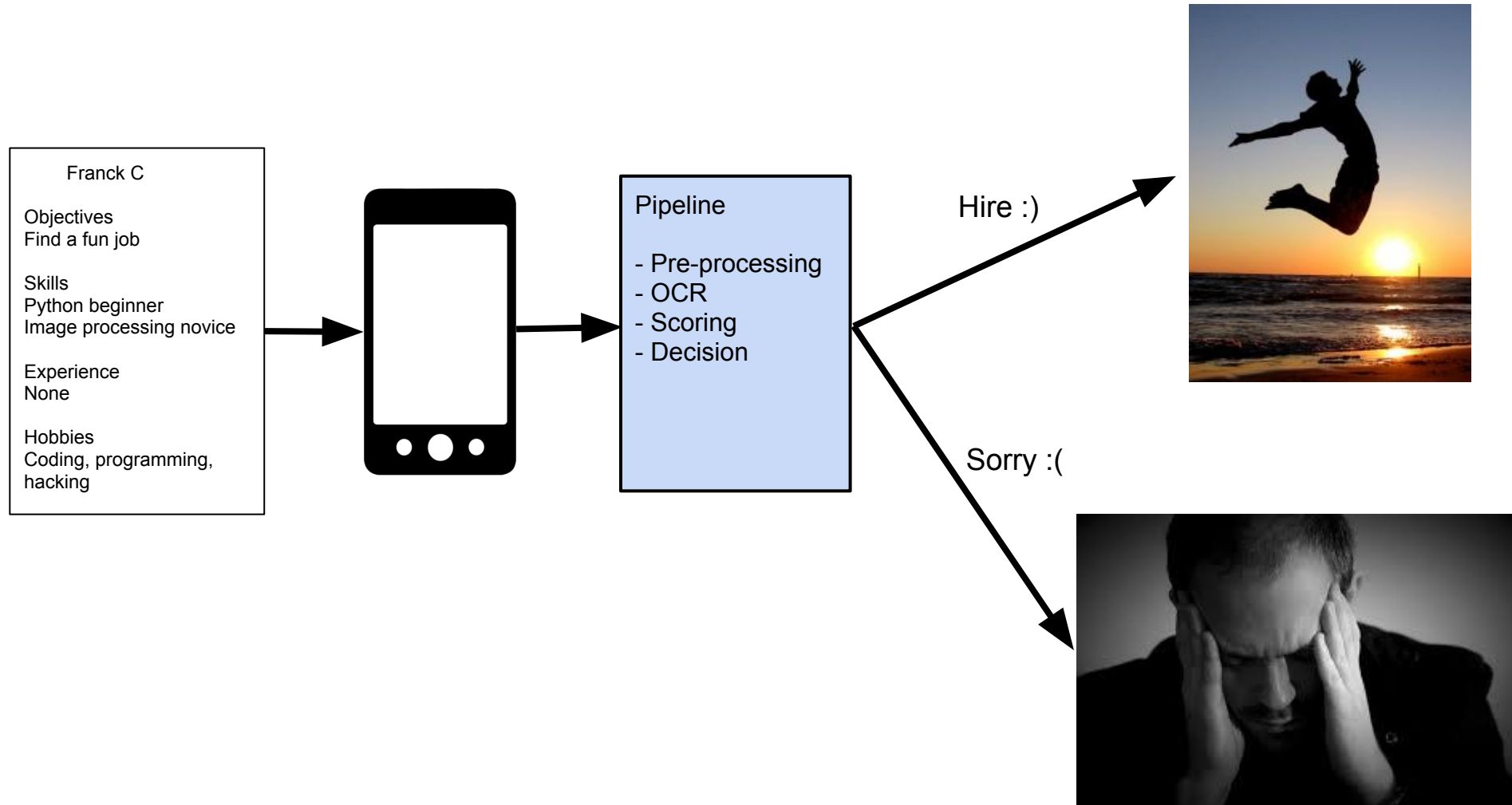- Q&A

# **Breakage will happen**



- You are a great coder, but...
  - Your co-workers ? interns ?
  - Pipeline will crash, servers will die

- How to get some good sleep at night ?
  - Good strategy for storing originals
  - Support re-runs

# **Agenda**

- Introduction
- Workflow
- Upload
- Image pre-processing
- OCR
- Structured data extraction
- Error handling/re-processing
- <span style="color:red">Q&A</span>

# Hiring pipeline (in Python)

Franck C

Objectives
Find a fun job

Skills
Python beginner
Image processing novice

Experience
None

Hobbies
Coding, programming,
hacking

Pipeline

- Pre-processing
- OCR
- Scoring
- Decision

Hire :)

Sorry :(

# Questions &
# (hopefully some) Answers