OTA INSIGHT

Scraping infrastructure management and solutions

# Continuously yield high quality data while growing from 100 to 100M daily requests

**Glen Henri J. De Cauwsemaecker**
Lead Crawler Engineer @ OTA Insight Ltd.

**100,000,000 requests / day**
~ 300 data sources
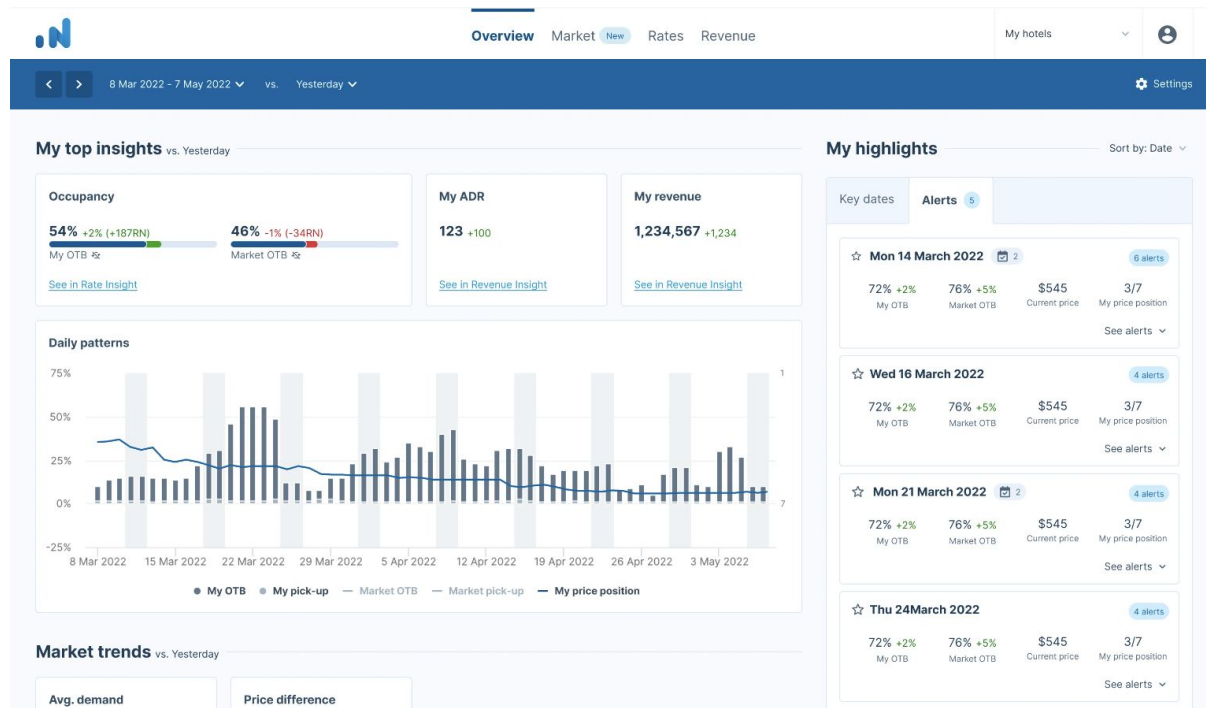9000 GiB (compressed) raw data

(2022)

100 requests / day
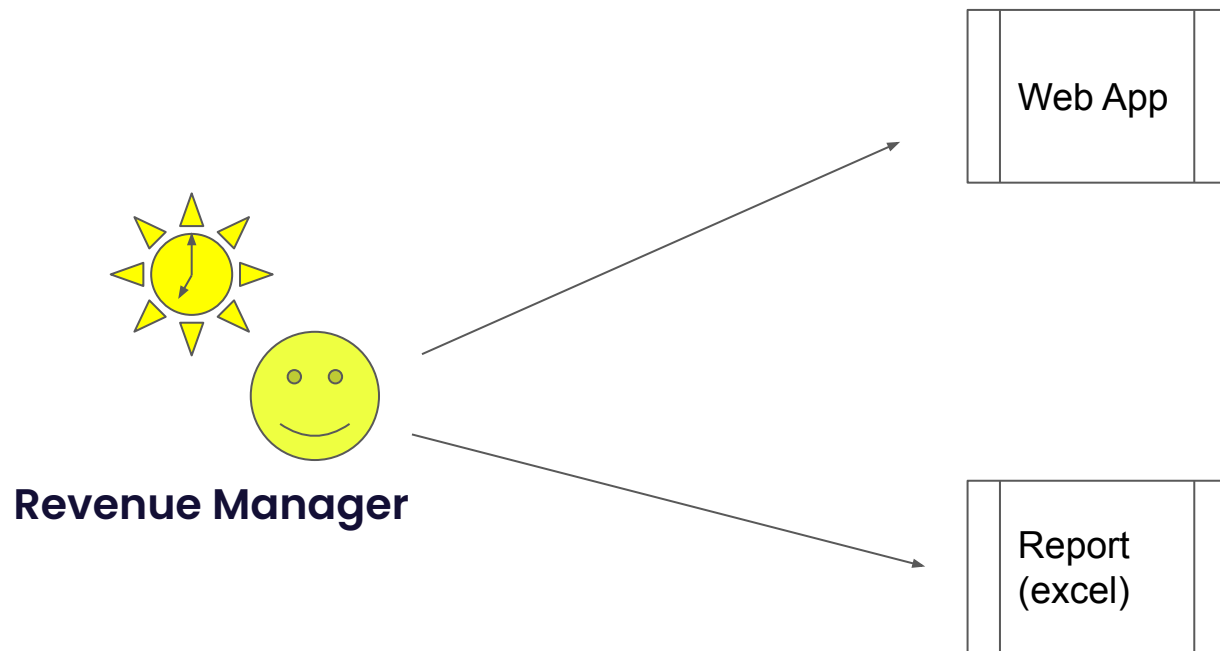(2013)

Credits: SpaceX

2

# OTA Insight Ltd.

We deliver smarter revenue, distribution and marketing outcomes.

365+ employees
~ 80 engineers
5 crawler engineers



## Hotel Rates and Ranking Information

**Revenue Manager**

Web App

Report
(excel)

STANDARD TIME ZONES OF THE WORLD

# Hmmm... Data...

# Data Extraction…



Crawler ←——→ Data Source

Extract + Transform

# Data Extraction…

Crawler

Extract

Data Source

Transform

Transformer

*19222#{137#0#166678
90#19267#1#2####}#1
660788000#0*

**1** **Extract**

**0** **Inputs**

👥 + 👤  Double room "Standard"
Rooms left: 5

rates in: **BGN**

**BED&BREAKFAST**
Bed & breakfast
Min stay: 1 night(s)   Sales policy

~~118,00~~
**106,20**
per person

Rooms:
1 ˅

BOOK

2 persons, 1 night(s)
Total: **212,40 BGN**

**HALFBOARD**
Halfboard
Min stay: 1 night(s)   Sales policy

~~145,00~~
130,50
per person

Rooms:
1 ˅

BOOK

2 persons, 1 night(s)
Total: **261,00 BGN**

Capacity: 1-3 persons
Click for info

{
        "max_persons": 2,
        "is_bar": true,
        "is_cancellable": true,
        "price_value": 212.4,
        "currency": "BGN",
        "meal_type_included": 1,
        "room_name": "Double room \"Standard\"",
        "rate_name": "BED\u0026BREAKFAST"
},

{
        "HotelID": "3d51d3127bdde31869c40f4af0aa28af",
        "RoomName": "Double room \"Standard\"",
        "Adults": "2",
        "Currency": "BGN",
        "RateName": "BED&BREAKFAST",
        "Meal": "Bed & breakfast",
        "PricePerNight": "106,20",
        "PriceDescription": "per person",
        "FromDate": "2022-09-08",
        "ToDate": "2022-09-09",
        "ExtractStartTimestampUTC": "2022-08-17 21:04:01"
},

**Transform**

**2**

**3**

9

**0**

Network Storage
( single input / timezone )

CSV
Input

Temporary File Storage
(scraped documents)

Proxy

Data Source

Scrapyd @ VM

**1**

Google
Compute Engine

Temporary File Storage
(fulls backup)

Giant CSV Files

```
message Price {
  // NOTE: These ids have been used in the past, DO NOT REUSE THESE IDS
  reserved 5, 8, 14;

  optional uint32 max_persons = 1 [default = 0];
  optional bool is_bar = 2 [default = false];
  optional bool is_cancellable = 3 [default = false];
  optional int32 cancellation_policy_days = 4 [default = -1];
  optional uint64 cancellation_deadline_date_time_local = 27 [default = 0];
  optional double price_value = 18 [default = 0];
  optional string currency = 6 [default = ""];
  optional string price_id = 7 [default = ""];
  optional double breakfast_cost = 19 [default = 0];
  optional uint32 ranking = 20 [default = 0];
  optional string booking_partner = 21 [default = ""];

  enum MealType {
    NONE = 0;
    BREAKFAST = 1;
    HB = 2;
    FB = 3;
    ALLIN = 4;
```

protobuf
Protocol Buffers

(Json)

Reprocess

Kafka

**3**

BigTable

**2**

Transformer
(Python)

# Error Prone

⚠️ Manual Intensive

⚠️ No Overview

⚠️ Costly Retries

# Integrities (v1)

An integrity is one job assignment, and links to one crawl input.

# Integrities (v1)

Status tracking per integrity, stored in a MySQL Database:

*missing → error or complete*

✅ Auto-Retry on Failure

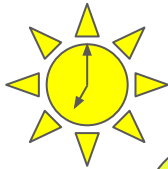| Rates | -8 | -7 | -6 | -5 | -4 | -3 | 0 | 1 | 2 | 3 | 4 | 6 | 7 | 8 | 9 | 10 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Source A | 0.00% | 0.00% | 0.00% | 0.00% | 100.00% | 98.15% | 89.47% | 75.84% | 98.90% | 100.00% | 99.29% | 96.55% | 91.43% | 80.30% | 93.55% | 55.56% | 50.00% |
| Source B | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 12.78% | 93.12% | 7.51% | 59.22% | 81.98% | 48.32% | 54.76% | 71.80% | 94.71% | 96.18% | 90.81% | 98.36% |
| Source C | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 43.45% | 100.00% | 54.62% | 100.00% | 76.24% | 86.93% | 100.00% | 100.00% | 100.00% |
| Source D | 0.00% | 0.00% | 0.00% | 0.00% | 100.00% | 100.00% | 65.50% | 54.25% | 94.64% | 100.00% | 95.00% | 100.00% | 100.00% | 77.78% | 81.82% | 26.67% | 40.00% |
| Source E | | | | | | | 29.50% | 0.00% | | | | | | | | | |
| Source F | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 90.70% | 89.47% | 59.09% | 60.00% |
| Source G | 0.00% | 0.00% | 0.00% | 0.00% | 100.00% | 100.00% | 31.77% | 56.11% | 0.00% | 100.00% | 0.00% | 0.00% | 0.00% | 88.71% | 93.55% | 59.26% | 55.56% |
| Source H | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 4.55% | 33.53% | 97.50% | 100.00% | 98.31% | 100.00% | 100.00% | 85.25% | 93.55% | 59.26% | 50.00% |
| Source I | 0.00% | | | 0.00% | | 100.00% | 98.75% | 100.00% | 100.00% | | 100.00% | | 100.00% | 50.00% | | | |
| Source J | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 19.24% | 47.28% | 100.00% | 100.00% | 100.00% | 100.00% | 72.05% | 97.04% | 99.01% | 66.67% | 42.86% |
| Source K | 0.00% | | 0.00% | 0.00% | 0.00% | | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | | 0.00% | 100.00% | 100.00% | 100.00% | 100.00% |
| Source L | 0.00% | | | 0.00% | | | 100.00% | 100.00% | 100.00% | 100.00% | 100.00% | 100.00% | 60.36% | 95.88% | 97.57% | 100.00% | |
| Source M | | | | | | | | | | | | | | 100.00% | 100.00% | | |
| Source N | | | | | | | | 30.61% | | | | | | | | | |

✅ Overview

# Grouped Integrities

**Revenue Manager**

Report (excel)

Generate Report

Integrities Complete

# Scheduler (v1)

Introduction of a GUI:

✅ High level cross-VM configuration (automate)
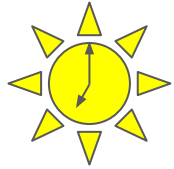
✅ Ability to orchestrate on a higher level

| Spider | Task | Timezone | Batch | Server | Scheduled | Started | Waiting time | Runtime | Pages/min | Request count | Items scraped | Progress | Errors | Exceptions | Actions |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| CtripPrices | Ctrip Prices | 10 | 1 / 1 | crawler6 | 17/10/2016, 14:20 | 14:22 | 0 min | 1 hrs, 4 min | 210 | 15493 | 2581 | 100% | 16 | 3 | Log Retry |
| CtripPrices | Ctrip Prices | 10 | 1 / 1 | crawler3 | 16/10/2016, 14:20 | 14:22 | 0 min | 1 hrs, 41 min | 140 | 15532 | 2679 | 100% | 13 | 0 | Log Retry |

# Revenue Manager

Web App

## Scheduled Data Pipelines:

**Morning**

**Afternoon**

**Evening**

No Fresh Data 😢

LIVESHOP

| Date | + OTB ⓘ | Market demand | Hotel A | Hotel B | Hotel C | Hotel D | Hotel E |
|------|---------|---------------|---------|---------|---------|---------|---------|
| Mon 15/08 | -- | 33% | ₹ 7,074 | ₹ 7,500 | ₹ 7,000 | ₹ 6,999 | ₹ 24,400 |
| Tue 16/08 | -- | 36% | ₹ 8,000 | ₹ 7,500 | ₹ 8,000 | ₹ 8,999 | ₹ 13,400 |
| Wed 17/08 | -- | 40% | ₹ 6,321 | ₹ 9,000 | ₹ 8,500 | ₹ 9,449 | ₹ 8,400 |
| Thu 18/08 | -- | 36% | ₹ 6,320 | ₹ 8,000 | ₹ 8,000 | ₹ 8,999 | Sold out |
| Fri 19/08 | -- | 39% | ₹ 6,320 | ₹ 8,000 | ₹ 7,500 | ₹ 8,499 | ₹ 15,400 |
| Sat 20/08 | -- | 41% | ₹ 6,320 | ₹ 8,000 | ₹ 7,500 | ₹ 8,999 | ₹ 7,400 |
| Sun 21/08 | -- | 36% | ₹ 6,320 | ₹ 8,000 | ₹ 7,500 | ₹ 7,999 | ₹ 7,400 |
| Mon 22/08 | -- | 42% | ₹ 7,624 | ₹ 9,000 | ₹ 8,000 | ₹ 7,499 | ₹ 7,400 |
| Tue 23/08 | -- | 41% | ₹ 7,624 | ₹ 14,500 | ₹ 8,000 | ₹ 9,499 | ₹ 5,400 |
| Wed 24/08 | -- | 45% | ₹ 7,624 | ₹ 13,500 | ₹ 10,000 | ₹ 14,999 | ₹ 6,400 |
| Thu 25/08 | -- | 44% | ₹ 7,374 | ₹ 13,500 | ₹ 9,000 | ₹ 10,799 | ₹ 6,400 |
| Fri 26/08 | -- | 43% | ₹ 6,847 | ₹ 9,000 | ₹ 10,000 | ₹ 8,999 | ₹ 6,400 |
| Sat 27/08 | -- | 46% | ₹ 6,847 | ₹ 9,000 | ₹ 8,500 | ₹ 8,999 | ₹ 6,400 |
| Sun 28/08 | -- | 38% | ₹ 6,847 | ₹ 8,000 | ₹ 8,500 | ₹ 8,999 | ₹ 6,400 |
| Mon 29/08 | -- | 40% | ₹ 6,321 | ₹ 9,000 | ₹ 8,000 | ₹ 9,899 | ₹ 6,400 |
| Tue 30/08 | -- | 39% | ₹ 6,321 | ₹ 9,000 | ₹ 8,000 | ₹ 9,899 | ₹ 6,400 |
| Wed 31/08 | -- | 37% | ₹ 6,321 | ₹ 9,000 | ₹ 8,000 | ₹ 8,099 | ₹ 6,400 |

Updated 5 minutes ago

# Liveshops...

## Live

**Always active**

## Scale

**Automatic Resource-based Scaling**

## Inputs

**Stream inputs over Pub/Sub to crawlers on demand.**

**CSV → Protobuf**

Crawler ←——————→ Proxy

SuperProxy
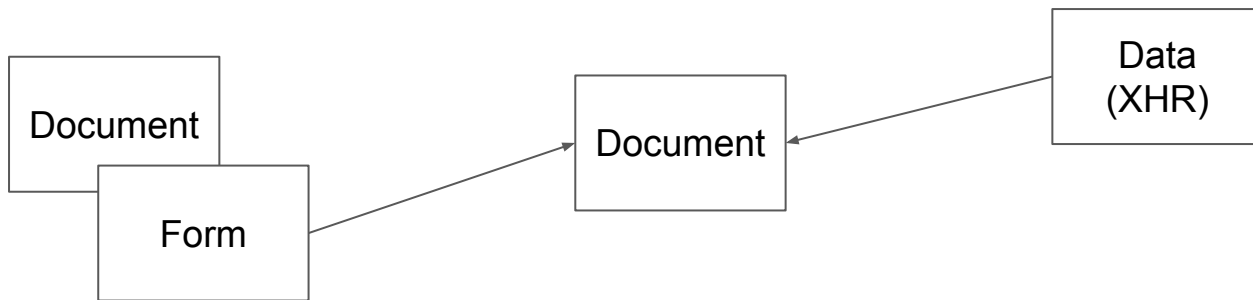
# SuperProxy

A MITM HTTP(S) Proxy written in Golang.

- Designed to handle high throughputs of requests/second;
- Ability to monitor all incoming and outgoing traffic;
- Delegate to proxy providers;

Only later we also started using it for emulation of HTTP and TLS Fingerprints.
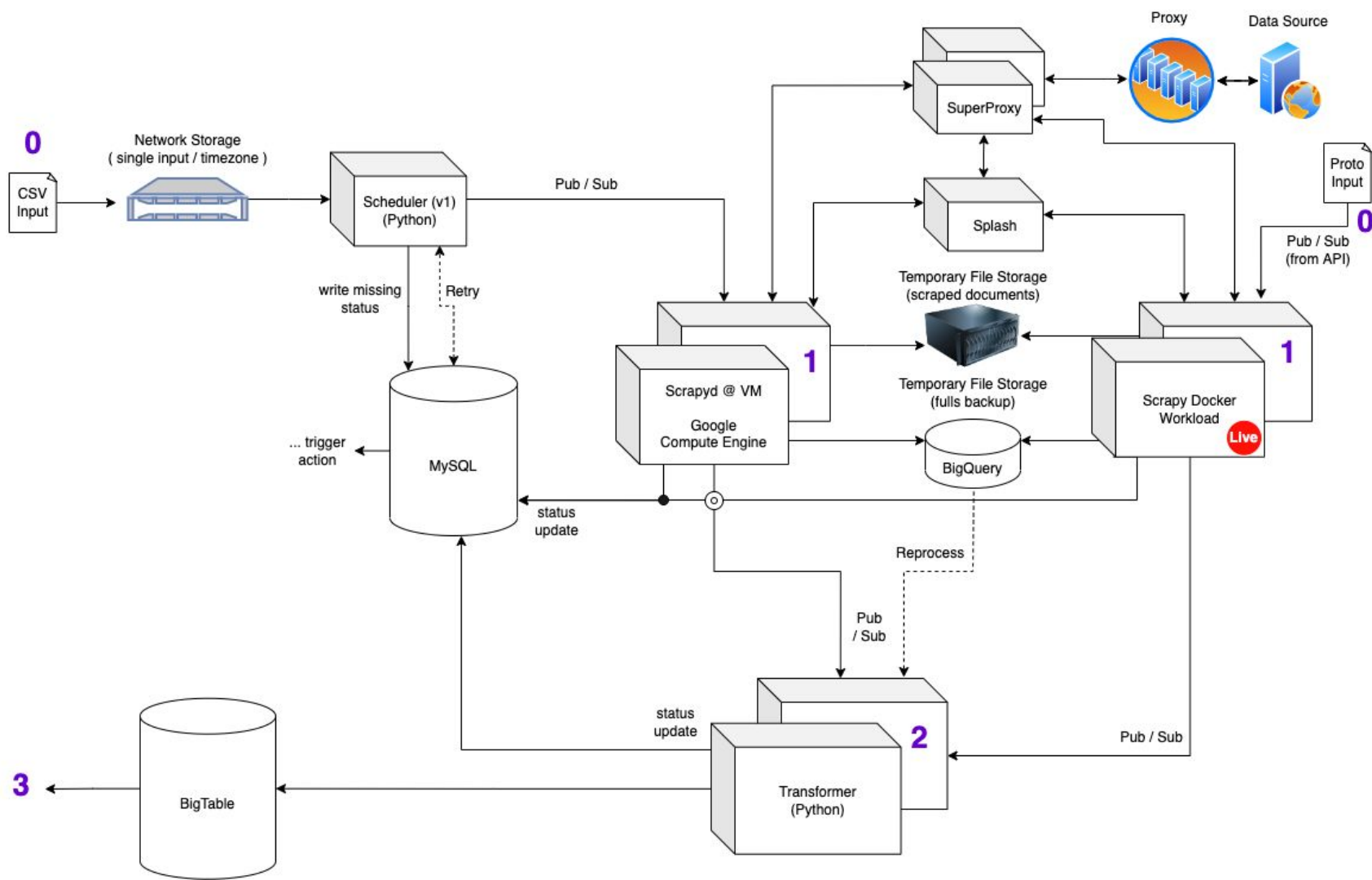
# Internet Browsers

I am a browser.

Document

Form

Document

Data
(XHR)

scrapinghub/**splash**

Lightweight, scriptable browser as a service with an
HTTP API

Proxy

Data Source

**0**

Network Storage
( single input / timezone )

CSV
Input

Scheduler (v1)
(Python)

Pub / Sub

SuperProxy

Proto
Input

**0**

Pub / Sub
(from API)

Splash

write missing
status

Retry

Temporary File Storage
(scraped documents)

**1**

Scrapyd @ VM

Google
Compute Engine

Temporary File Storage
(fulls backup)

**1**

Scrapy Docker
Workload

Live

... trigger
action

MySQL

BigQuery

status
update

Reprocess

Pub
/ Sub

status
update

**2**

Transformer
(Python)

Pub / Sub

**3**

BigTable

22

# 20,000,000

Requests / Day

# Internet Browsers

Replace *Splash*

- Our demands outgrew its capabilities
- Javascript > Lua
- Rendering engine easy to detect

+ Puppeteer

# Anomaly Detection

Based on historical data

- Pattern-based anomaly detection
- Allows us to keep billions of data points
  Accurate at any given point
    - One of many solutions to our "scale" challenges

**1**  **Our client and data needs were growing exponentially**

**2**  **Infrastructure suffered from Legacy Creep**

**3**  **Business needs were harder to integrate into our infrastructure**



Think About It GIF By Louis16art

OTA INSIGHT

# Integrities (v2)

`19214#{184#0#19246900#19231#1#2####}#1660078800#0`

old: `184_1_2_308362_19214_1_rate_19231`

# Integrities (v2)

# Integrities (v2)

Satisfy our needs for atomic writes at scale

# FoundationDB

**The right tool for the job.**

# Integrities (v2)

Apply lessons learned from liveshop crawlers to scheduled crawlers

## k8s

**Orchestrated by scheduler (v2)**

## Protobuf

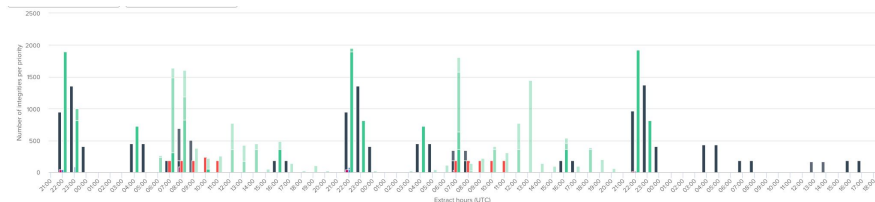**Drop the CSV Inputs in favour of Protobuf-driven Crawl Inputs**

# Integrities (v2)

Automate workload definitions

# Planner

**Forecast and simulate**

# Scheduler

**Orchestrate workloads
and deliver crawl inputs**
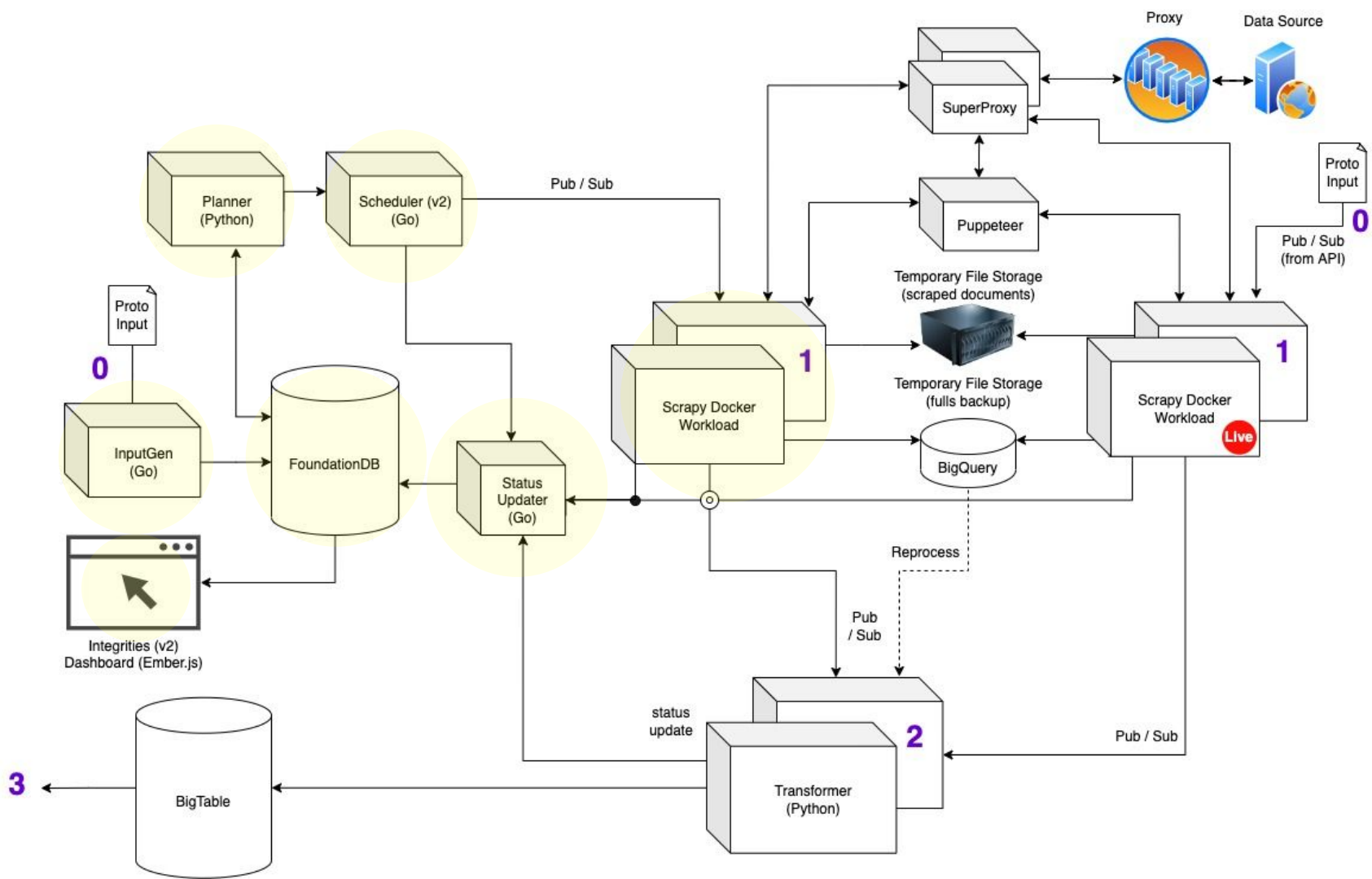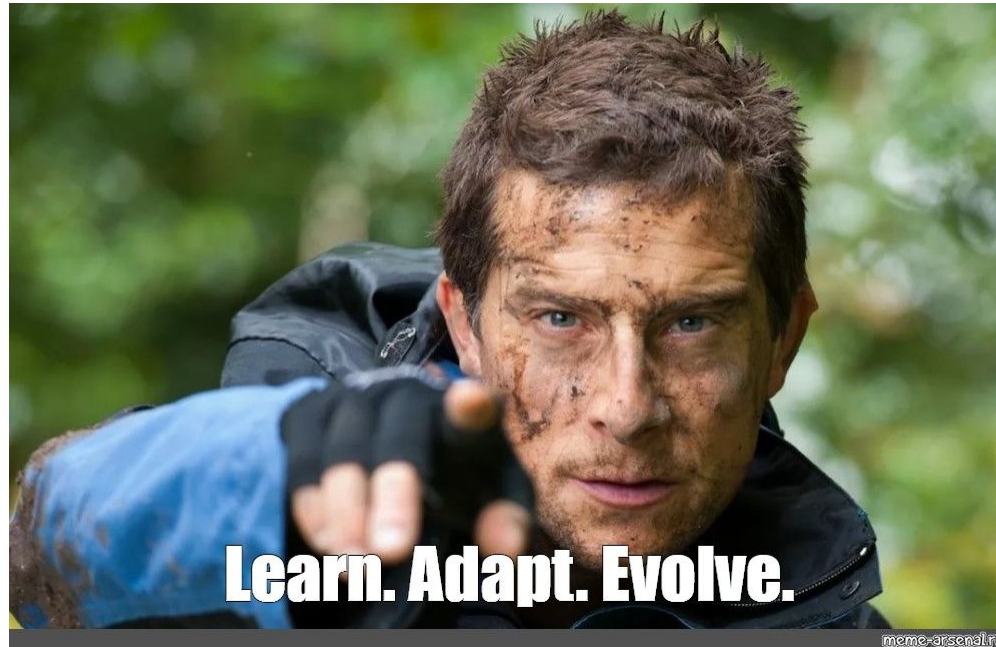
# Integrities (v2)

Migration to our new infrastructure

## Test

**Tests and Traffic Emulators**

## Parallel

**Run both versions of Integrities in parallel for an extended period to compare workloads.**

Proxy

Data Source

SuperProxy

Puppeteer

Proto Input

0

Pub / Sub (from API)

Planner (Python)

Scheduler (v2) (Go)

Pub / Sub

Proto Input

0

InputGen (Go)

FoundationDB

Status Updater (Go)

Scrapy Docker Workload

1

Temporary File Storage (scraped documents)

Temporary File Storage (fulls backup)

Scrapy Docker Workload

Live

1

BigQuery

Reprocess

Integrities (v2) Dashboard (Ember.js)

Pub / Sub

status update

Transformer (Python)

2

Pub / Sub

BigTable
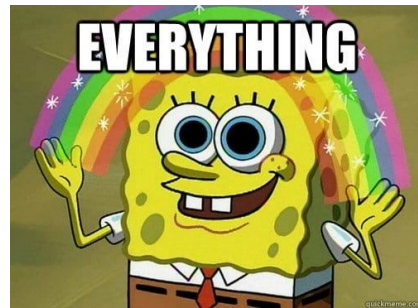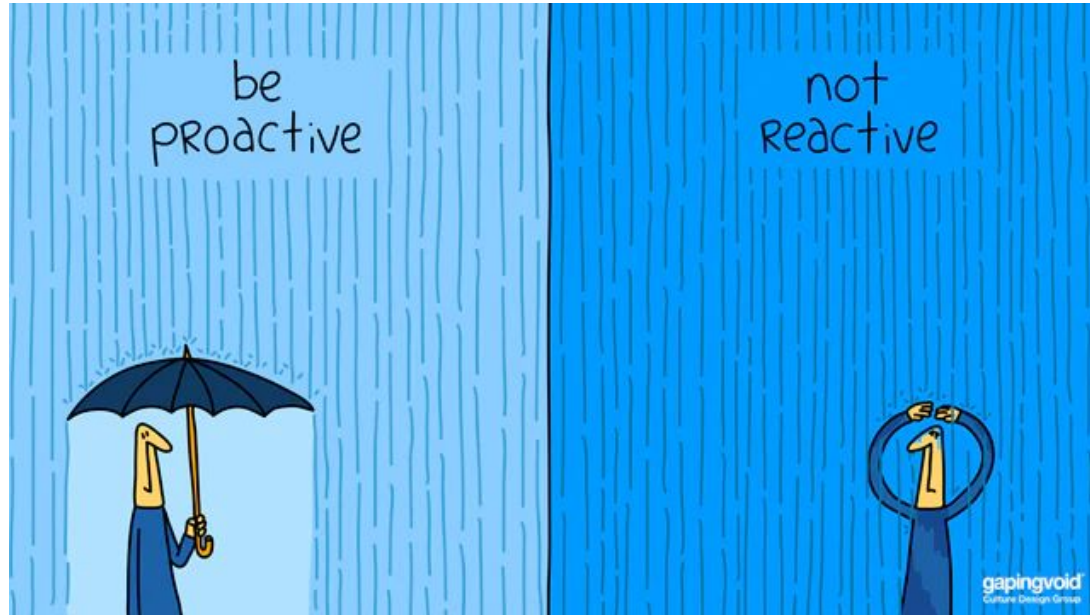
3

# Automate

Priority #1

- <u>Error Monitoring</u>: Sentry
- <u>Alerting</u>: slack + email
- <u>Incident Response</u>: PagerDuty → OpsGenie
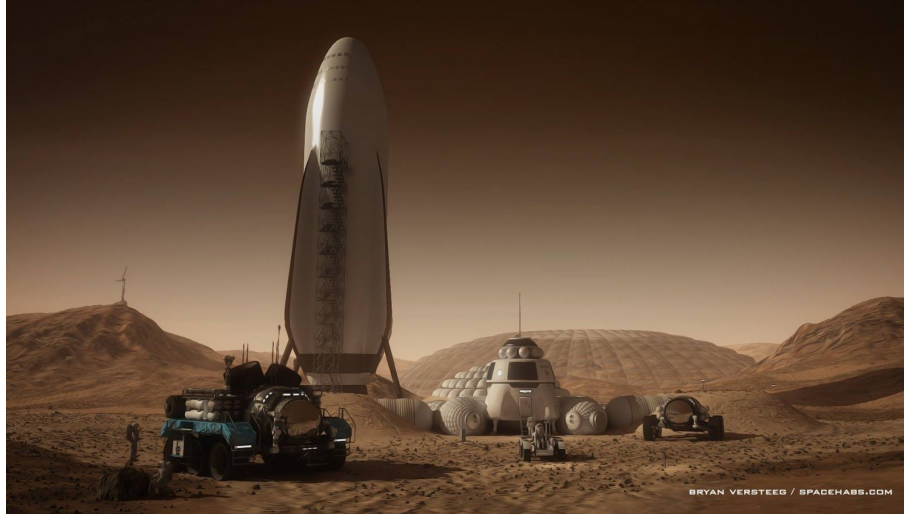- <u>Tooling</u>: Metric exports & Dashboards

# Be Pragmatic

Cost-Reward Balance

- Proxy?
- Browser?
- Automate?


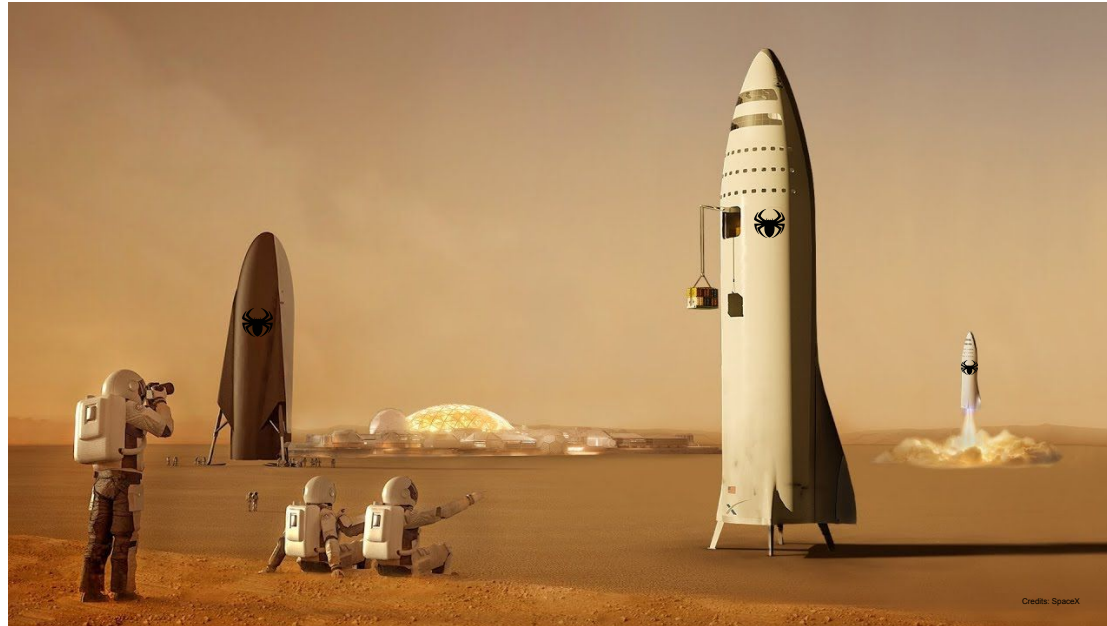AT WHAT COST?

# Present to Future Roadmap

- Hire & Educate
- Improve UX
- R&D: Anti-Fingerprint

careers.otainsight.com

# Give a crawler engineer scrapy and you can extract data for a day, teach him how to scale and you will go to space.



Credits: SpaceX

# Thank You

**Web Data Extraction**
**Summit 2022**

Powered by Zyte

**extract.summit@glendc.com**