# ICIJ's leak refinery

Mar Cabra - Data editor
International Consortium of Investigative Journalists (ICIJ)

# +190 journalists in more than 65 countries

12 staff members (USA, Costa Rica, Venezuela, Germany, France, Spain)
50% of the team = Data & Research Unit

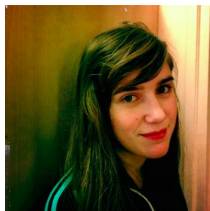ICIJ Data&Research unit

Editor (**Spain**) ⟷ Deputy director (**USA**)

Developer (**Spain**)

Data journalist (**France**)

Research editor (**Venezuela**)

Reporting

**Project**

Data checkers

Data analyst (**Costa Rica**)

Web developer (**Germany**)

Publication

**Project**

# 260 GB - 100,000 companies

SECRECY FOR SALE

PEOPLE'S REPUBLIC OF OFFSHORE

Luxembourg LEAKS

#luxleaks

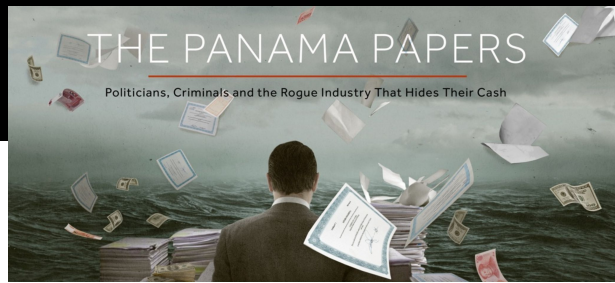More than 100,000 HSBC clients
$100 Billion

SWISS LEAKS

#swissleaks

+500 secret tax agreements (PDFs)

THE PANAMA PAPERS
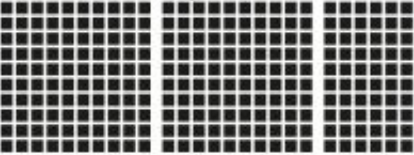Politicians, Criminals and the Rogue Industry That Hides Their Cash

# The scale of the leak

Volume of data compared to previous leaks

**1,7 GB**
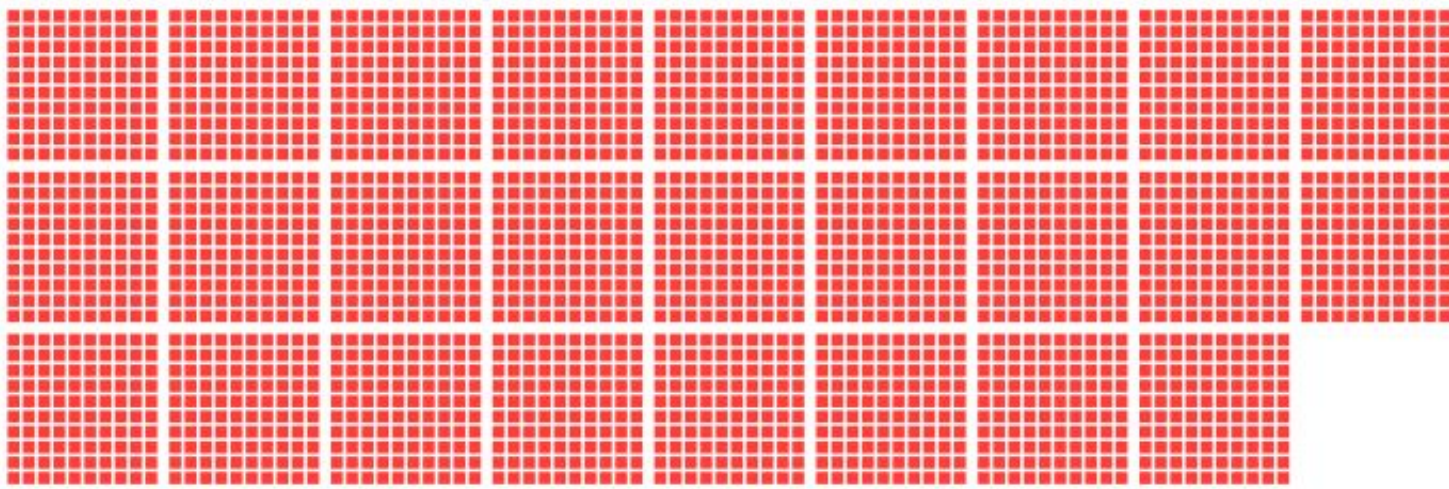Cablegate/Wikileaks (2010)

**260 GB**
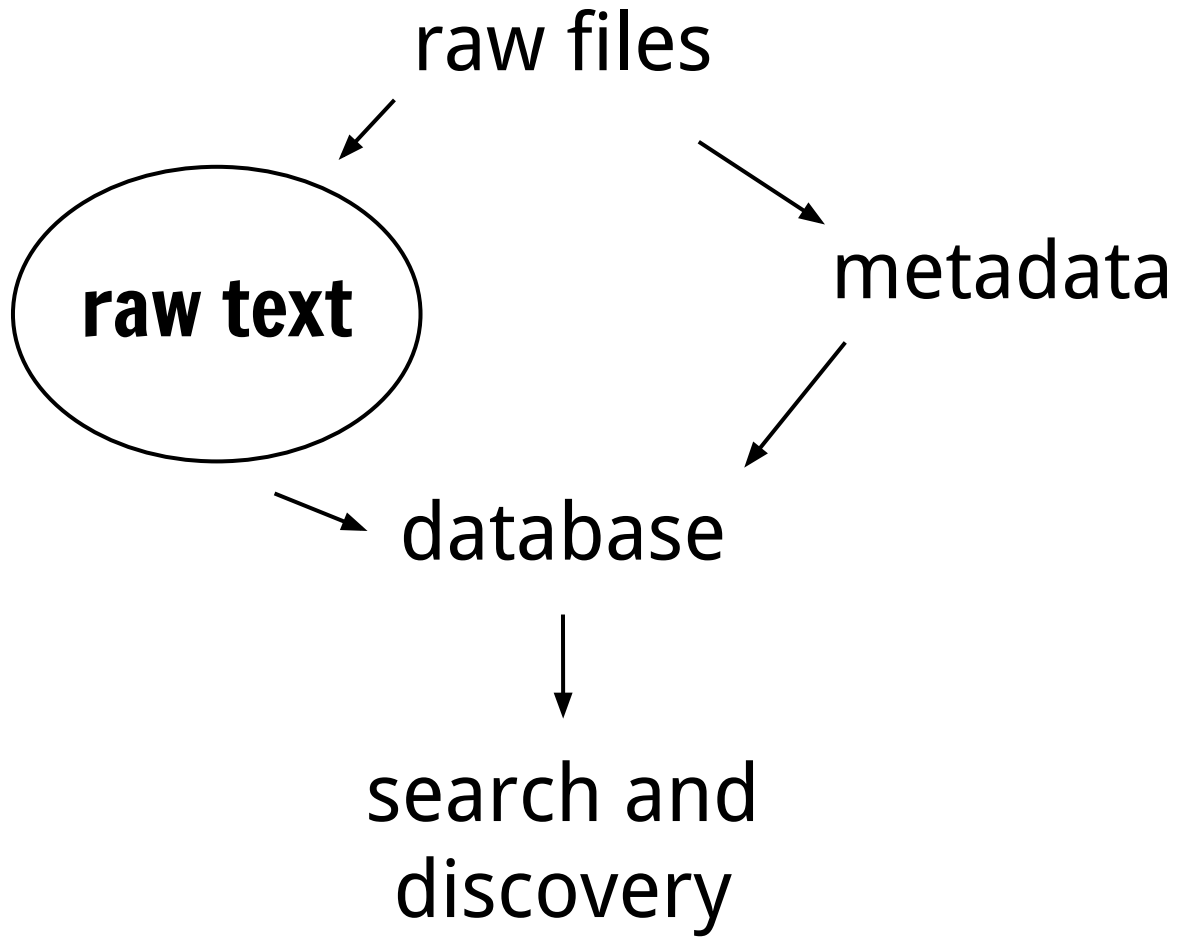Offshore Leaks/ICIJ (2013)

**4 GB**
Luxemburg Leaks/ICIJ (2014)

**3,3 GB**
Swiss Leaks/ICIJ (2015)

**≈ 2,6 TB**
Panama Papers/ICIJ (2016)

■ = 1 GB

raw files

raw text

metadata

database

search and
discovery

file, attachment or embedded object

# detect the type

do we know how to extract the text?

no!
log and tackle later

yes!
extract, OCR and repeat

**3 million files**

**x**

**10 seconds per file**

**=**

**1 year**

Redis queue

35 x g2.xlarge Amazon instances with Ubuntu + Tesseract + Extract

Solr

1 year

÷

35 machines

=

11 days

## Unstructured data extraction

- ICIJ Extract (open source, Java: https://github.com/ICIJ/extract), leverages Apache Tika, Tesseract OCR and JBIG2-ImageIO.

## Structured data extraction

- A bunch of Python

## Database

- Apache Solr (open source, Java)
- Redis (open source, C)
- Neo4j (open source, Java)

## App

- Blacklight (open source, Rails)
- Linkurious (closed source, JS)

# Stack

1 SELECTED NODES

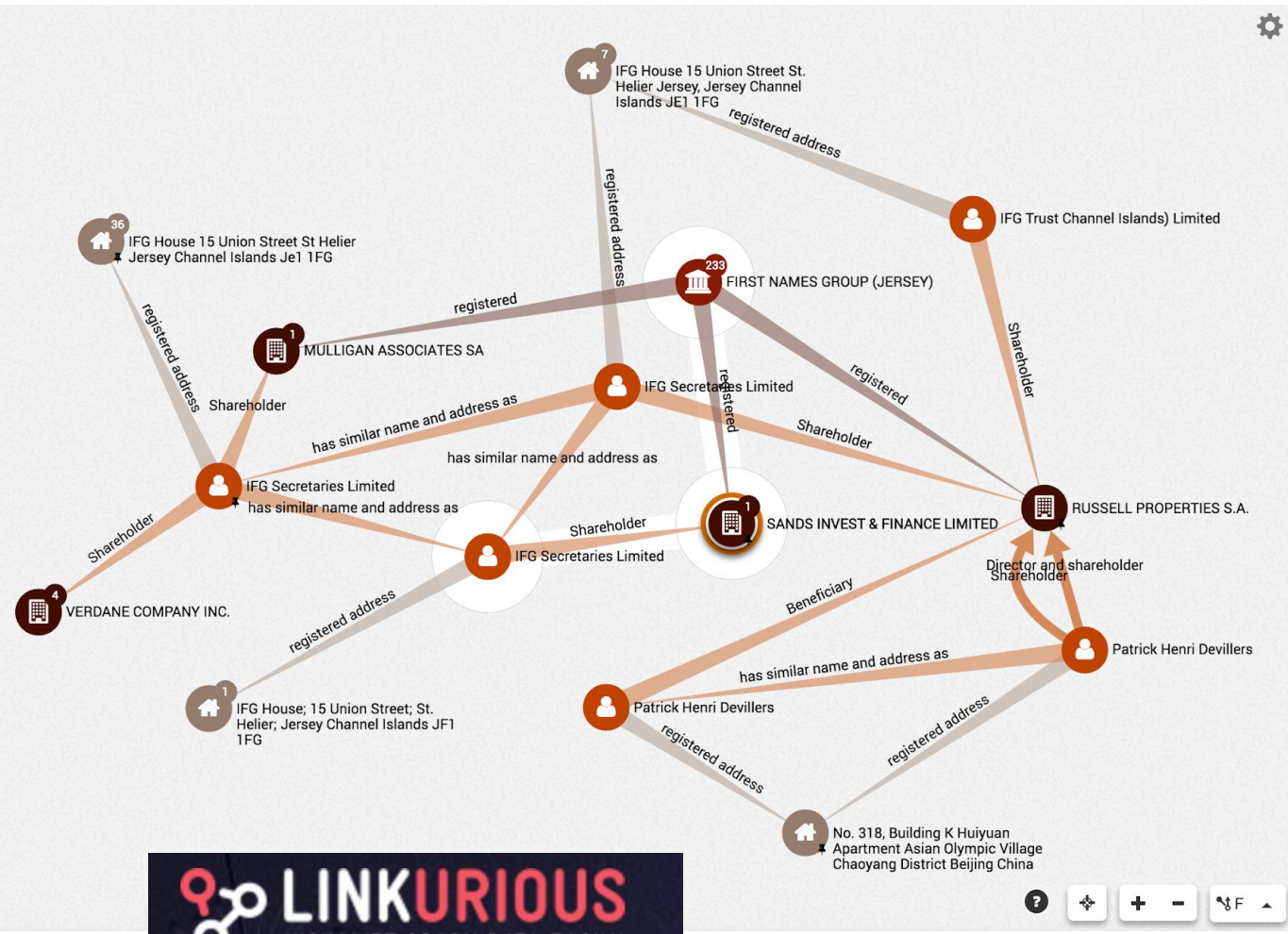**SANDS INVEST & FINANCE LIMITED**
#246682

pinned   Company

Select...   Deselect   Hide   Expand...

PROPERTIES

Find a property...

| | |
|---|---|
| file_number | 6015821 |
| inactivationDate | 18-FEB-2013 |
| jurisdiction | BVI |
| name | SANDS INVEST & FINANCE LIMITED |
| registrationDate | 20-APR-1998 |
| status | DIS |
| struck_off_date | 31-OCT-2013 |

2 EDGES / 3 IN DATABASE

IFG House 15 Union Street St. Helier Jersey, Jersey Channel Islands JE1 1FG

registered address

IFG Trust Channel Islands) Limited

registered address

IFG House 15 Union Street St Helier Jersey Channel Islands Je1 1FG

FIRST NAMES GROUP (JERSEY)

registered

MULLIGAN ASSOCIATES SA

Shareholder

Shareholder

IFG Secretaries Limited

registered

Shareholder

has similar name and address as

has similar name and address as

IFG Secretaries Limited

has similar name and address as

RUSSELL PROPERTIES S.A.

Shareholder

IFG Secretaries Limited

SANDS INVEST & FINANCE LIMITED

Shareholder

VERDANE COMPANY INC.

Director and shareholder
Shareholder

registered address

Beneficiary

IFG House; 15 Union Street; St. Helier; Jersey Channel Islands JF1 1FG

has similar name and address as

Patrick Henri Devillers

Patrick Henri Devillers

registered address

registered address

No. 318, Building K Huiyuan Apartment Asian Olympic Village Chaoyang District Beijing China

Neo4j

# Our platforms by the numbers

- + 14 million documents, 20 formats
- 500 users* - 100 active each week, close to 200 per month
- 1 full-time programmer for improvements

**Open-source** + exclusive expertise**

*not unique users

*except graph database

# Not just a network…

## … but a community

Search people, places, and more...

# ICIJ
THE INTERNATIONAL CONSORTIUM
OF INVESTIGATIVE JOURNALISTS

MAIN    FORUM    GROUPS    **FILES**    LINKS    PHOTO    MEMBERS    SEARCH

## RECENT FILES

UPLOAD FILE ⊕

🕐 RECENT FILES    ❤ FEATURED FILES    ★ TOP RATED ITEMS    ✅ CATEGORIES    🏷 BROWSE BY TAGS

**PDF**
**Bolton Group -
Greek...**
By Harry Karanikas

**PDF**
**Bombardier...**
By Kristof Clerix

**PDF**
**Bombardier...**
By Kristof Clerix

**PDF**
**Hutchison
Whampoa...**
By Colm Keena

**PDF**
**Hutchison
Whampoa...**
By Lars Bové

**PDF**
**kastra investments
ltd**
By Colm Keena

**PDF**
**TMT II Luxco Sarl**
By Colm Keena

**PDF**
**Black & Decker...**
By Colm Keena

**PDF**
**skype technologies**
By Colm Keena

**PDF**
**Disney // Disney
Stores...**
By Jan Kleinnijenhuis

Pages:  «  1  2  3  4  5  6  7  8  9  10  ...  »  »»

Oxwall

THE INTERNATIONAL CONSORTIUM
OF INVESTIGATIVE JOURNALISTS

ICIJ

KF Knight Foundation

# Global I-HUB

Please read ou...

### SEARCH

Search people, plac...

### USERS

### FORUM TOPICS

Global I-Hub ...
General » Gener...

## PLEASE SIGN IN

Username/Email

••••••••

Your Google Authenticator current valid code:

••••

☐ Remember me

Forgot passphrase

SIGN IN ➔

...sting any material.

Collaboration space - Mar 12

# Our community by the numbers

- +600 users*
- Shared status in +21,000 occasions**
- 1,500 forum topics with 5,400 posts**
- Uploaded +600 files**

*not unique users

**October 2015

# the next steps

- entity extraction
- making our data silos talk

# Next: ICIJ Knowledge Center

# Our current challenge *(global data sharing)*

# Encrypted email (easy)







**Review of tools**

# Encrypted email: PGP

# Questions?

# Thanks!

mcabra@icij.org