# Information Extraction and Aggregation from Unstructured Web Data for Business Profiling
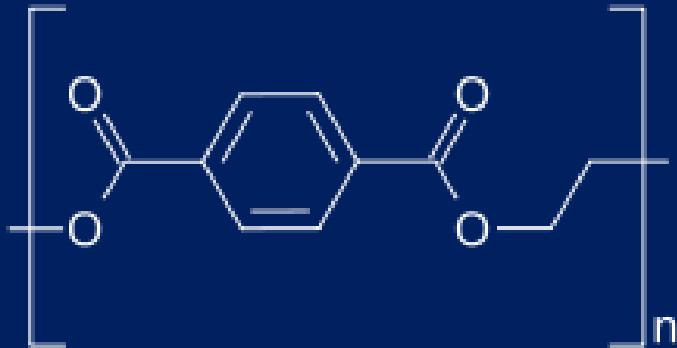
**Student Team : Alexander Michels, Himanshu Ahuja**

**Industry Mentors : Dr. Stephen DeSalvo, Urjit Patel**

# Praedicat: An Insurance Tech Company

## Determining Risk

Ethylene glycol
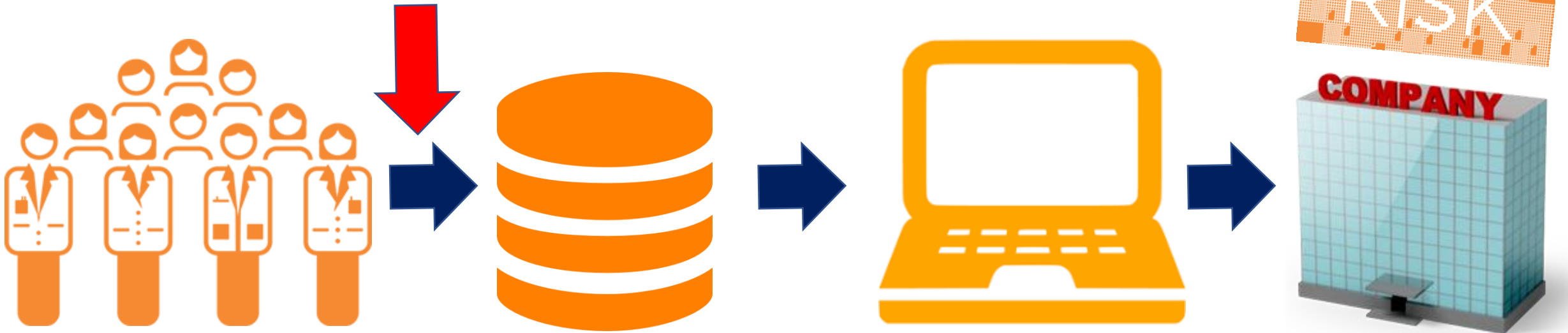- moderately toxic
- generally used in anti-freeze

[1]    https://iaspub.epa.gov/triexplorer/release_fac_profile?TRI=38063MRCNG1236A&TRILIB=TRIQ1&FLD=&FLD=RELLBY&FLD=TSFDSP&OFFDISPD=&OTHDISPD=&ONDISPD=&OTHOFFD=&YEAR=2016

# Praedicat: An Insurance Tech Company

Predicting the likely amount of losses

# Where Do We Fit in?

**RIPS Team**
**Automating**

RISK

COMPANY

| 1. Manual Search | 2. Credible Database | 3. Forward-looking Models | 4. Predict Likely Losses |

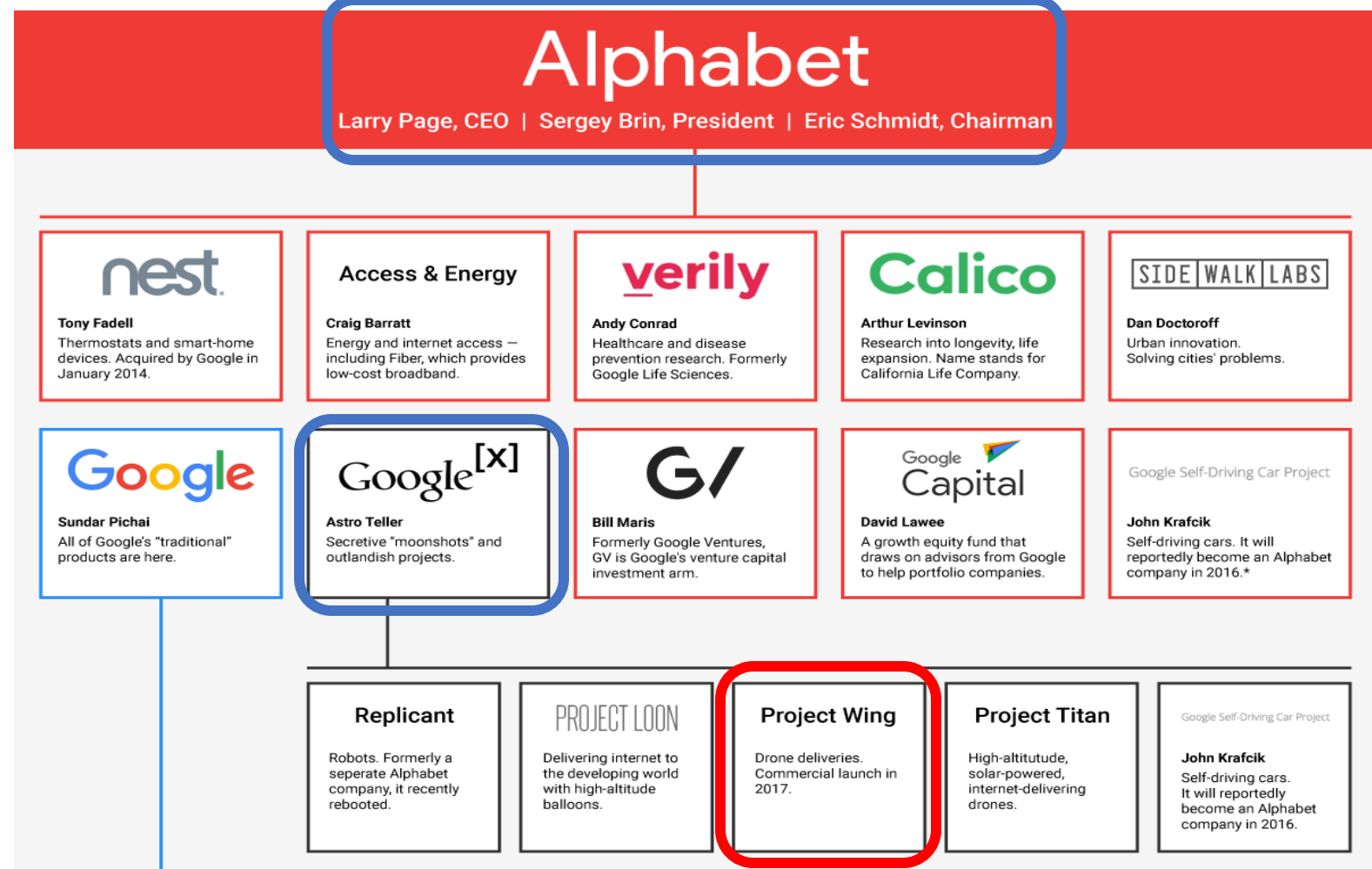# Difficulty of Searching Information

# Difficulty of Allocating Risk

## Example: Allocating Risk for Parent and Subsidiary Companies

Parent Company →
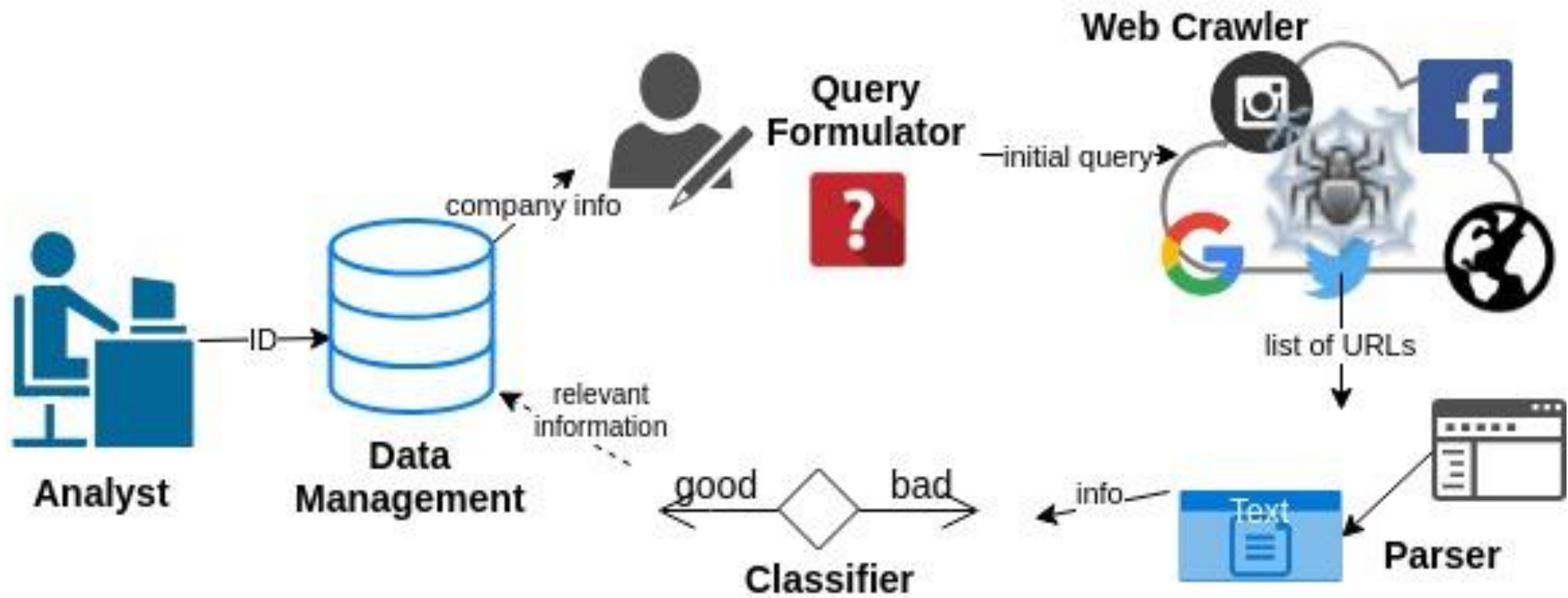
Parent Company →

Subsidiary Company →

# Problem Statement

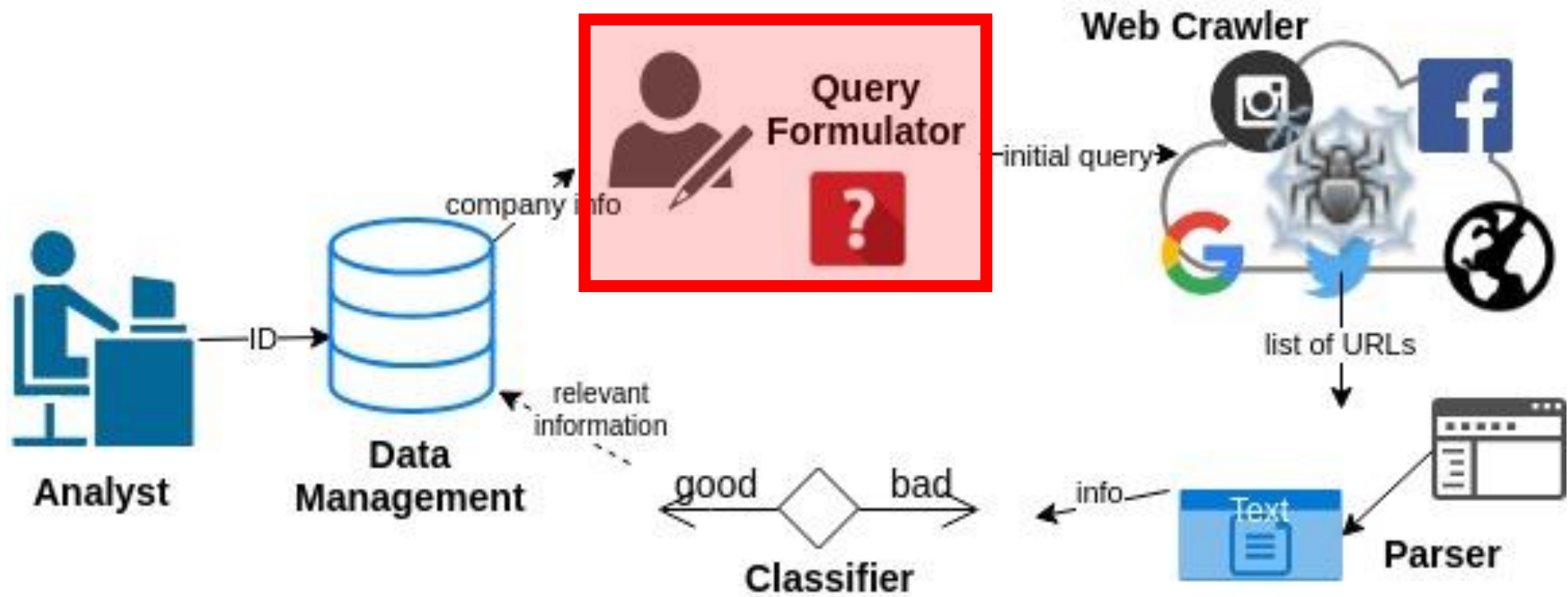How to automate information extraction and aggregation from unstructured data on the Internet for business profiling

# Solution Overview

# Query Formulator: Asking about the right things!



Zero useful results

'Apple Inc.' returns the right results.

PDF result mentions Rentokil Initial PLC involvement in window cleaning.

# Query Formulator: How did we ask the right things?



Mention the file-type

Making some words optional

Making keywords mandatory

Name of the company

Optional alias

# Solution Overview

# Web Crawling: What is web Crawling?

# Web Crawling: Unsupervised machines cannot be trusted



Start with a google search of the company and its business activity.

The business activity appears in the financial report that specifically appears on search services provided by the website.

# Web Crawling: Where and how far?



*The problem:*
We don't know how far to dig, and where to dig?

We don't know the credible sources and where the information lies on the credible sources.

# Web Crawling: Credible data to the rescue



- Interestingly the structured data (available on Federal websites & Wikipedia) is also credible!

- Design of specific crawlers to get data from specific types of data.

- Create the baseline data.

# Parser:
# Getting unstructured data

Use of text abundance to locate meaningful paragraphs.

Filtering out tags containing social media redirects.

Removing graphic contents, advertisements.



Hazardous material spill in front of Dow Chemical facility leads to road closure

By ABC12 News Team | Posted: Mon 2:53 PM, Jul 23, 2018 | Updated: Mon 6:23 PM, Jul 23, 2018

MIDLAND COUNY (WJRT) (7/23/2018) - A stretch of Saginaw Road through the Dow Chemical plant in Midland County was closed after a hazardous material spill.

A truck carrying diethyl phosphorochloridothioate spilled some of the chemical while turning into Gate 17 of the plant around noon. Police say the truck did not crash, but it was unclear what caused it to leak.
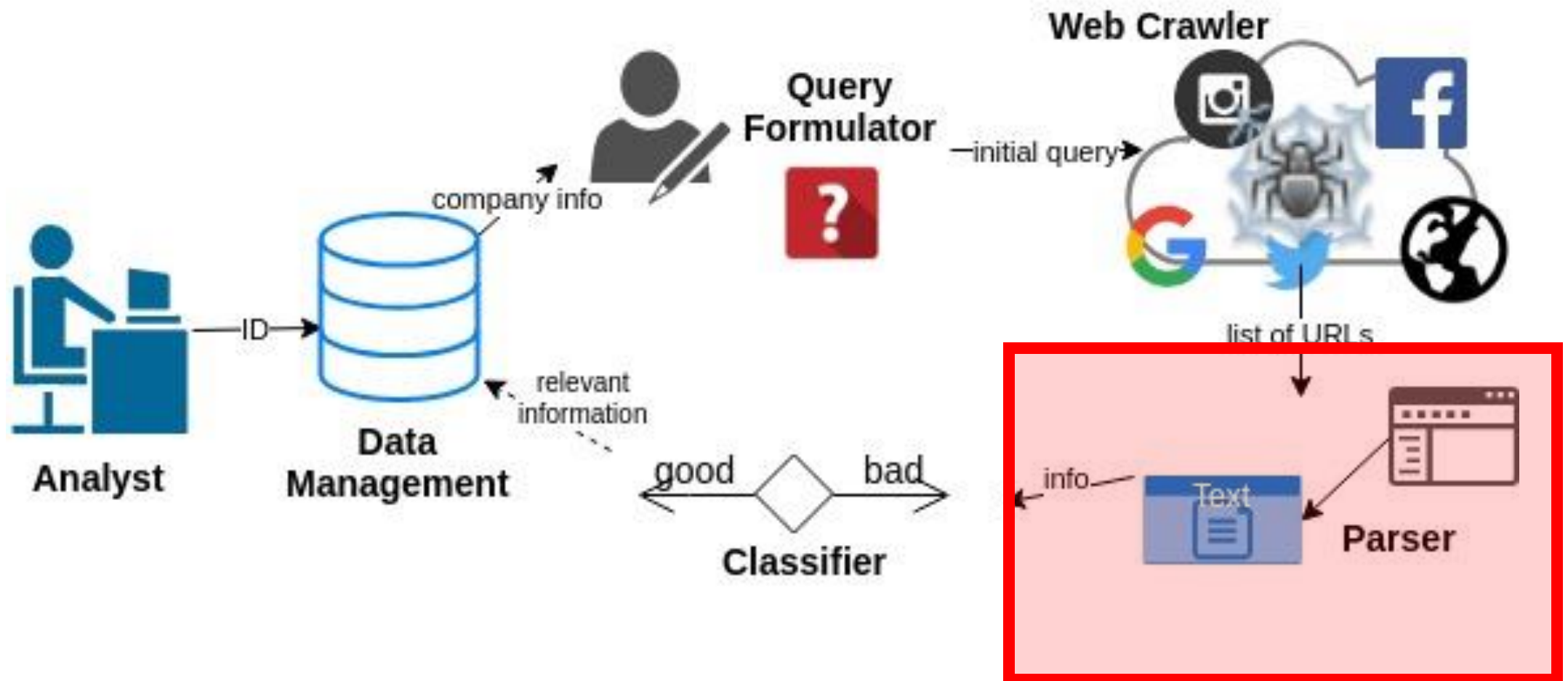
Dow sent a cleanup crew to help clean up the scene, along with Midland County police and firefighters.

The chemical is a form of acid, which can be flammable.

In a statement posted on Facebook, Dow says continuous air monitoring is ongoing around the scene and has shown no danger to the general public. No evacuations have been ordered.

However, Saginaw Road was closed between Bay City and Salzburg roads while cleanup continued. The road reopened around 4:20 p.m.

# Parser: Why go towards unstructured data?



More credible information

Structured Data

Lower volume of data

More indicative information

Unstructured data

Higher volume of data

# Classifiers

- Our web crawlers and parsers can give us a huge quantity of data, but not all sources are created the same.
- Our web crawler needs to distinguish between high quality data and YouTube comments.

# Term Frequency-Inverse Document Frequency (TF-IDF)

- Each term (word, *i* ) in a document (*j*) has its own *term frequency (tf)*

$$tf_{i,j} = \frac{n_{i,j}}{\sum_k n_{i,j}}$$

- Each word has an *inverse-document frequency (idf)*, *which*

$$idf(w) = log(\frac{N}{df_t})$$

$$w_{i,j} = tf_{i,j} \times \log\left(\frac{N}{df_i}\right)$$

$tf_{i,j}$ = number of occurrences of $i$ in $j$
$df_i$ = number of documents containing $i$
$N$ = total number of documents



Magnitude of frequency

The  a  an  of  in     leather  cosmetics     pneumonoultramicroscopicsilicovolcanoconiosis

Common words     Less frequent words     Rare words

Keyword Generation with word2vec

# doc2vec



- Doc2vec uses a *Distributed Memory version of Paragraph Vector* (PV-DM)
- You can "tag" documents with topics.
- We can attempt to cluster or classify documents using tags.

# Data Management

- Many companies can't make use of their data
- Our data needs to be:
  - Traceable back to the source
  - Queryable
- Our solution: ProfileManager and WebResourceManager.

# The Story So Far...

- Our web crawling and parsing capabilities are bottlenecked by our ability to classify information.

- We currently have a database of high-credibility information about 52,629 companies and corporate entities including:

  - 10Ks
  - 8Ks
  - EX21s
  - Wikipedia

# Going Forward

# Knowledge Graph

- Knowledge graphs represent **statements of facts** as subject-predicate-object triples

  For example ("Barack Obama", "is a", "Muslim")

- G=(V,E) where V is a set of concept nodes and E is a set of predicate edges

- To determine the degree of truth of a statement of fact, the shortest path from subject to object is traversed and paths that pass through high-degree nodes are assigned a low truth value

# How we are hoping to use the knowledge graph

## Credibility Checking and Classification

- A document that contains a statement of fact **P** that is contradictory to the set of possible-worlds as described by our knowledge graph is rejected.

- We can also use competitive learning to prioritize trustworthy sources.

## Aggregates Information in a Queryable Format

- This makes the information much more accessible to the analysts at Praedicat, Inc. and makes the job of summarization easier

**Profile Manager**

company info →

**Document Summarizer and Content Extractor**

readable summary of findings →

**Master Document (Summarized Company Profile)**

# Milestone 2

# Questions?