# DOW JONES

# Turning a Billion Articles into Actionable Insights

**David Arnold**
Director of Product Strategy

Dow Jones Professional
Information

# Dow Jones Professional Information Business

**30,000+**
Content Sources

**1.3bn**
Articles

**28**
Languages

**800+**
Regions

# Dow Jones Intelligent Identifiers
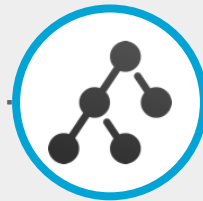
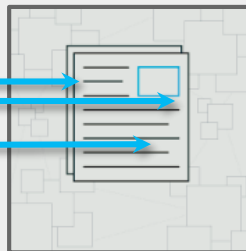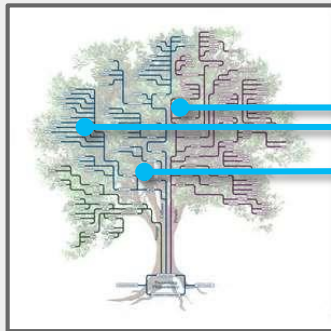Why a Taxonomy?

**1000+** Regions

**1000+** Subjects

**1000+** Industries

## Evolving Taxonomy

- **Not the total number of tags ; We have a dedicated MD team**
- Ensures our taxonomy is **right-sized!**
- Evaluation of new tags (emerging industries, technologies etc…)
- Example: **Drones, Crypto-Currency or AI**
- Evaluation of placement / re-definition
- Evaluation of **tag usage** / removals (VCR technology)
- All part of the 60 day DJID release schedule

**Key Differentiator***: Millions of articles flowing through our pipeline daily, our content is tagged with valuable metadata that underlined defines the article; Harnessing the power of our metadata allows customers to build out models, data analytics, drive visualizations & deep learning.*

**May 2018**

# Keyword Searches vs. Dow Jones Intelligent Identifiers

## Keyword search challenges

Retrieves articles with any mention of the keyword



**Free text = Terroris***

## Solutions using DJID codes

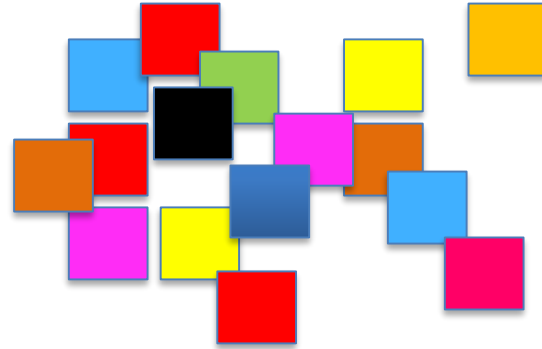Retrieves relevant content



**DJID code = Acts of Terror**

**VS.**

## Structured Data

## Unstructured Data



- Clearly defined
- Expected fields and values
- Easily classified / categorized
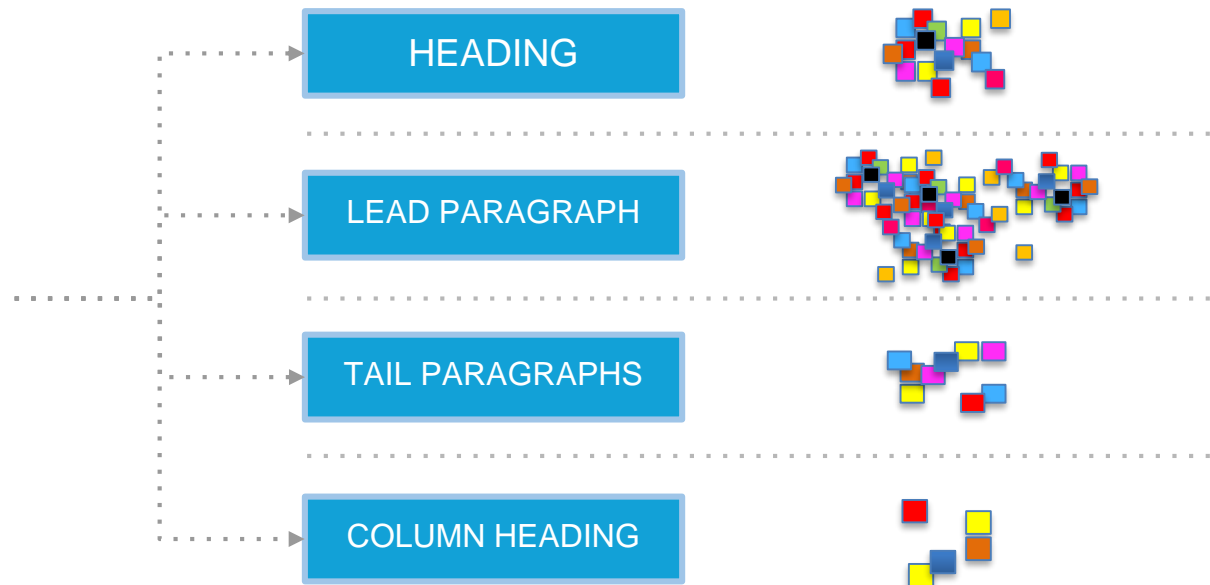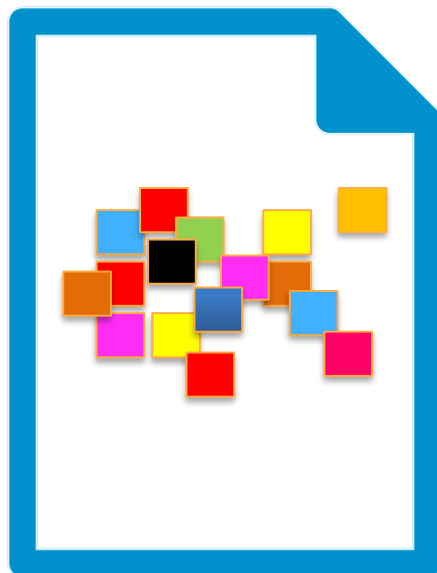
- Data is random / unpredictable
- Data is not defined
- Free form text format
- No expected fields and values

- Classification difficult
- Prone to error
- Requires sophisticated linguistic tools to categorize

# Creating structure to a News Article

**DISTDOC** - Articles are normalized into a fixed format

HEADING

LEAD PARAGRAPH

TAIL PARAGRAPHS

COLUMN HEADING

# Language is Difficult to Understand

Older systems/machines that only used <span style="color:red">keywords</span> found it very challenging to correctly <u>interpret the context, the way humans can</u>.



let's look at an example

The BMW I bought to replace my Mercedes is a great car.



Which is the "great car"?

The BMW I bought to replace my Mercedes is a great car.

Easy for humans, very difficult for machines... until now !



Keywords can't understand the meaning of this sentence.

The BMW I bought to replace my Mercedes is a great car.

- Keyword rules based engines
- Rely on proximity and distance to other words
- As a way to collect evidence, what the story was about.
- As you can see here, language is "tricky" !

# Language is difficult...

## Lexical ambiguity - Homophones

A word that sounds the same but have different meanings.

- "The **bat** slipped from his hand" or a flying mammal?
- "Cinderella had a **ball**" or is ball an event or physical object?
- "Ron **lies** asleep in his bed" or is he telling a lie in his sleep?
  "I have four **mouths** to feed at home"

"Mouths" (parts of people) =
"People" (the whole)

"The **strings** were praised for their excellent performance."

"The strings" (parts of a violin) =
"violins" or "violinists" (the whole)

"Check out my new **wheels**"

[Translation: Check out my new car!]

"Wheels" (parts of a car)
= "car" (the whole)

Imagine **Regional Dialects** (UK English vs. American English) and **Idioms** ("Raining cats and Dogs" and "Barking up the wrong tree")

# What is Autocoding?

❖ **Autocoding** is the <u>automatic</u> application of Dow Jones Intelligent Identifiers (DJID) codes on all Factiva/DNA content with minimal editorial intervention.  (26 languages)

❖ **Aboutness** is the <u>underlying principle</u> of coding the Dow Jones Intelligent Identifiers taxonomy.

❖ Autocoding specialists **configure** coding systems in an effort to categorize the item based on what it is about, **not** what sector, country, product or group it may be of interest to (also known as "of-interest-to" coding) or for passing mentions.
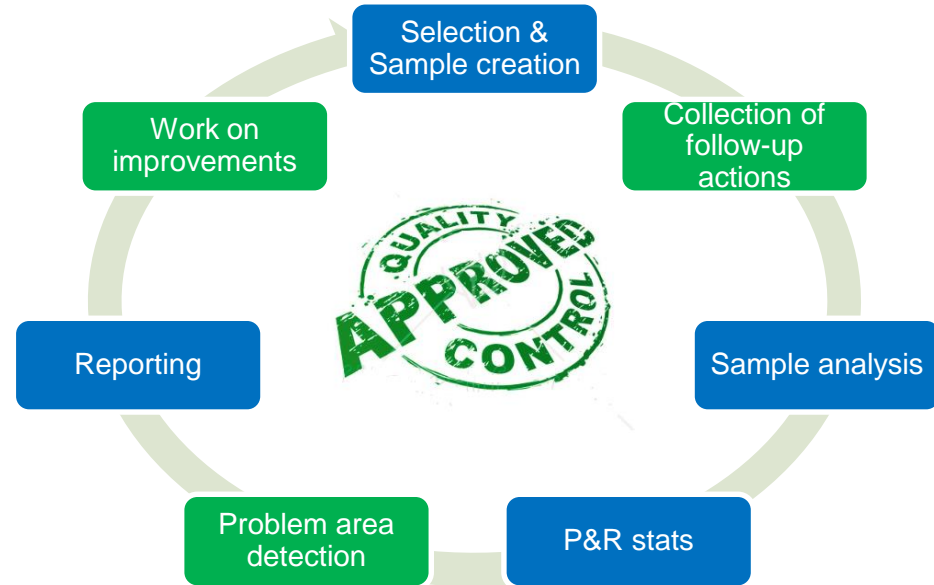
# How do we get around these linguistic challenges?

26 Languages

❖ Hybrid strategic approach - Requiring multiple tools

❖ **No one tool / system** today can do it all

❖ Multiple tools fill tackle these linguistic challenges

❖ Mission: is to hit extremely high **quality** metrics / KPI – No 100%

# Quality – A "never-ending" process

- Customers want to hear we "*stay on top*" of our taxonomy and tagging.  Remember "*evolving*" DJID?

- We also stay on top of Quality.

- Dedicated Quality team

- Monthly monitoring

- Monitored by staff in Quality, Metadata, Autocoding & external

- Gaps identified assigned for follow-up (by Autocoding)

- Future? More about machine learning and Auto-Feedback loop – *More to come…*



Selection & Sample creation

Collection of follow-up actions

Work on improvements

QUALITY APPROVED CONTROL

Reporting

Sample analysis

Problem area detection

P&R stats

# Recent GDPR Business Case

- ❖ **GDPR** - General Data Protection Regulation

- ❖ EU regulation on **data protection and privacy** for all individuals within the European Union.

- ❖ It aims primarily to **give control** to citizens and residents over their personal data

- ❖ It becomes enforceable on **25 May 2018**

# What is a SIP?



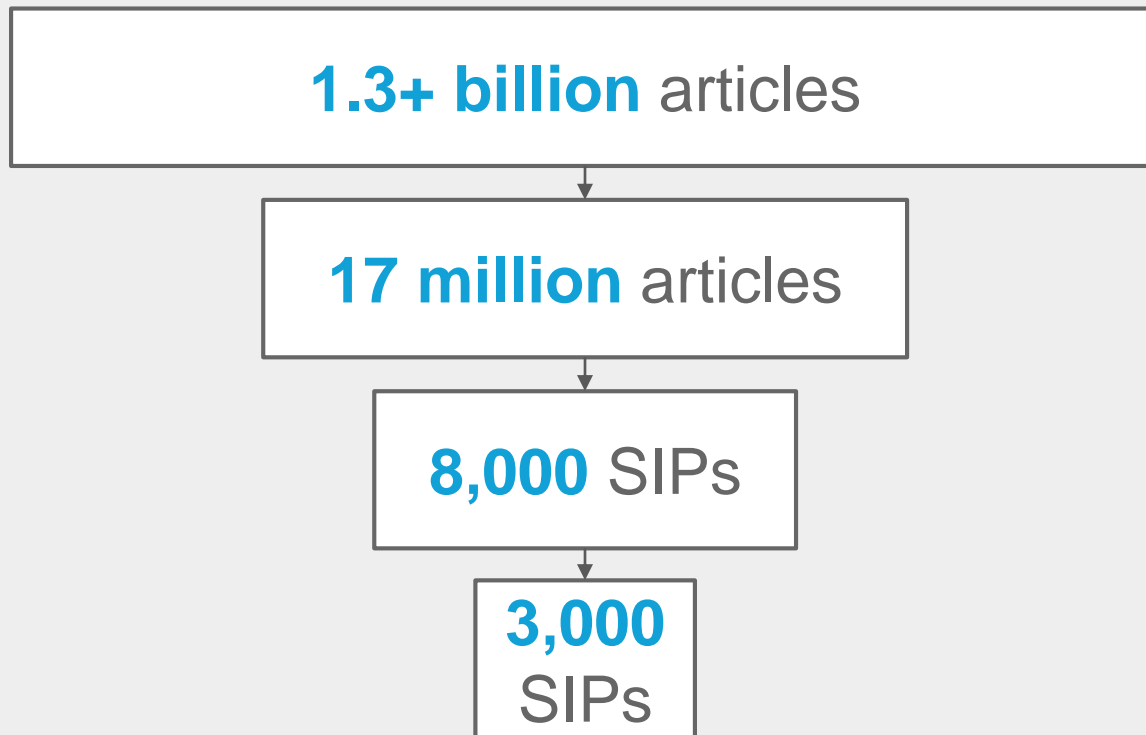Financial

Terrorism

Organized Crime

Trafficking

# New GDPR Regulation

❖ Many SIPs were **AQUITTED** of the crime

❖ **200k SIPs** needed to be reviewed

# Solution to the Business Challenge

1.3+ billion articles

17 million articles

8,000 SIPs

3,000 SIPs

# Benefit

- ❖ Saved **11 research years**

- ❖ **GDPR Compliant** Data

# Conclusion

- ❖ AI powers Dow Jones **Professional Information Business**

- ❖ **Autocoding** uses AI to apply codes to articles

- ❖ AI utilized to read over **1.3 billion** articles