# COGITO Intelligence API: Main Features

## Release 9.0, March 2015

## CATEGORIZATION

**COGITO Intelligence API** categorizes texts according to different taxonomies. The standard version offers three different **Automatic Categorization** engines.

- **Intelligence** categorization engine: Developed to support the activities of Intelligence Agencies that are required to quickly and accurately assess large amounts of diverse information.
- **Crime & Offense** categorization: Taxonomy focused on the needs of Law Enforcement Agencies.
- **Geotaxonomy** categorization: Categorizes content for geographic locations and information.
- **Computer Crime** categorization: Detailed taxonomy designed for accurate categorization of the Cyber Illegal domain.

With these three taxonomies, **COGITO Intelligence API** uses a comprehensive approach to classify the contents of more than 1000 different items, organized hierarchically. Users can generate results using a specific taxonomy or run all three classifications simultaneously.

### INTELLIGENCE Categorization Engine
With more than 800 entries, the Intelligence categorization engine provides wide coverage of diverse information domains. Categories include:
Campaign Finance
Fuel-Air Explosives
Censorship
Companies Chemicals / Petrochemicals
Oil & Gas Companies / Oil & Gas Drilling
Religious Conflict

A unique feature of this engine is its inclusion of a wide range of diverse topics and issues (chemical weapons, political crisis, wildlife, economic measures, cinema, terrorist groups, automotive, personal investments, etc.) from open domains, such as those found on open sources, social networks, reports and messaging systems (email, sms, chat, etc.).

## CRIME & OFFENSE Categorization Engine

This taxonomy was developed to manage the information typical for the police and law enforcement domain. It includes a variety of categories such as:

Corruption
Murder, grievous bodily injury
Computer related crime

## Computer Crime Categorization Engine

This taxonomy features a selection of specific categories to achieve accurate categorization of the Cyber Illegal domain. To master such a critical domain, we structured the taxonomy to include categories like DoS attack, Intrusion (computer or network), Identity theft, and others.

## GEOTAXONOMY Categorization Engine

This classification engine produces output that highlights the countries and regions of the world in or associated with a text. Cogito uses a semantic approach to process geographic information present in text, deciphering ambiguities and correlating all existing concepts.

## Tagging

Semantic processing of text produces three results that highlight:
- Document summary
- The most semantically relevant words and concepts
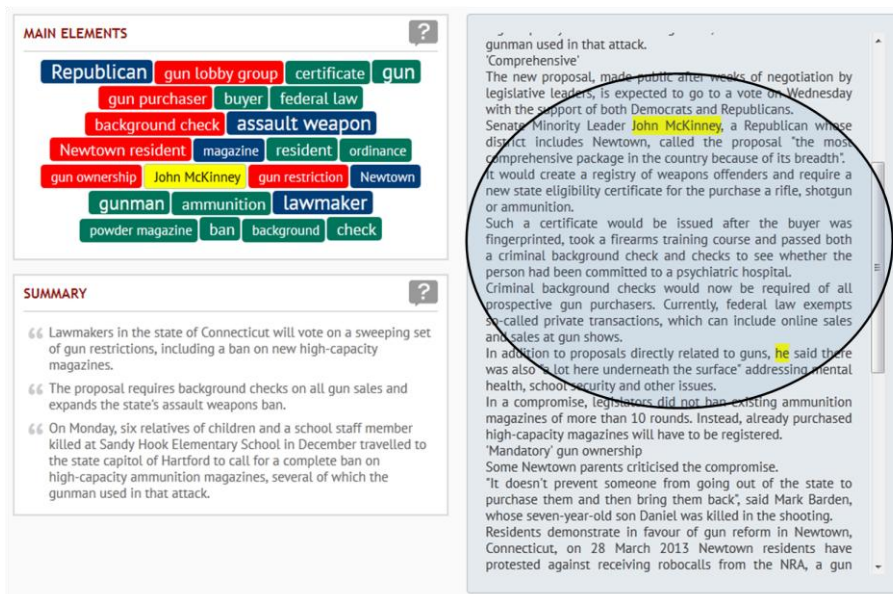- Collocations that characterize the language text

Semantic analysis example:

The results in the circled area above highlight:
- Lemmas or headwords (in *blue*)
- Concepts (in *green*)
- Semantically relevant collocations (in *red*)

Clicking on each element of the tag cloud highlights the element in the corresponding text. For example, clicking on the **concept** "OFFER" highlights the corresponding concept, in this case the synonym "PROPOSAL".

The same result is achieved by selecting any item from the tag cloud and/or phrases that are identified as the most significant (*see the box above the tag cloud*).

Selecting a person mentioned in the tag cloud will prompt the indicated name and personal pronoun, as shown in the following example where the proper name, "John McKinney" and "he" are highlighted.



These elements allow you to mark the context in the document with semantic information which can subsequently be used as follows:
- Highlights the most important sentences in text through a text preview or summary.
- Orders the search results giving a higher ranking to documents with the most semantically relevant keywords or concepts.
- Proposes correlations between different texts on the basis of language collocation.

- Offers synthesized views in order to quickly and accurately grasp the weight and importance of the key information in the text.

## Entity Extraction

The text mining entities are composed of **standard entities** and **domain entities**.

**Cogito Intelligence API** extracts a variety of standard entities including people, organization, places, dates, currencies, addresses, etc., even without the support of lists. The semantic engine is able to identify a proper name in text and always correlate it back to its correct context. For example, "Arthur Andersen" may be categorized as "People" or "Organization" depending on the context.

In addition, to some extent, the anaphora that allows you to extract both the explicit and implicit references is also considered. The following example highlights the proper noun (People -> John McKinney), which relates to the personal pronoun "he" later in the text.



Domain entities may be defined as entities related to a specific context or realm of knowledge (Intelligence and LEA) which is regularly updated with current knowledge and events.

**Cogito Intelligence API** currently manages dozens of domain entities, including: Terrorist Organization, Biological Agents, World leaders, etc.

Thanks to the availability of customization tools (an option which can be acquired together with the training support Expert System provides to clients for special projects), clients can customize the extraction rules, enrich the semantic network or extract new entities based on their needs.

## Relationships Extraction

In addition to the extraction of entities, **Cogito Intelligence API** can extract the *relationships between* semantic entities.

**COGITO Intelligence API** exceeds conventional entity extraction technology by offering coherent suggestions and contextualization of acquired text with reference to specific relations based on a set of over 20 different types of relationships.

Here are some examples:

**COMMUNICATE** ( *report* )

" CNN's Nick Paton Walsh reports on the desperation inside a Syrian town under siege and one doctor trying to make a difference.

**COMMUNICATE** ( *tell* )

" Musharraf's actions came under the purview of high treason, he told parliament.

**LAW** ( *try* )

" In April, the interim government refused to try Musharraf for treason, saying it was beyond its mandate and up to the new government, elected in May.

**BEHAVIOUR** ( *violate* )

" Musharraf violated the constitution twice.

## Fact Mining

This innovative feature allows rapid identification of facts present in text, identifying both the fact, as well as the entities (people, organizations, places) and tags (URLs, phone numbers, emails, etc.) related to it. Both Intelligence and Crime taxonomies vehicle this information to trigger facts of interest and domain entities within.
In the example below, the words "killed" (and "murder", "kidnapping") identified in a fact (Crime) present in text, would generate a set of proposed entities (Bachir, Karachi, Pakistan, etc.) related to the fact.

# COGITO® Intelligence API



In this second example you can see how the concept "blast" determines the fact "Emergency Incident" linked to Karachi, Pakistan, Sindh and Bakir.



With a third "Geography" Fact Mining feature, it is possible to add a further perspective in gathering facts and contextualize them in the exact geographic location as detected in text.

# COGITO® Intelligence API



This change in facts' clusterization provides an innovative outlook from a geographic point of view.

## Emotions

The Emotions feature masters about 80 categories making it possible to detect emotions within the text achieving a better and more focused emotions mining. This innovative feature breaks the boundaries of Sentiment Analysis transforming hidden emotions' content into accessible data and information which can be handled and associated to entities. The Emotions feature provides a new, innovative stage of text analysis.

The example below clearly shows how emotional contents of "fear" were found in the text pointing out all the sentences containing such information. The obtained data brings out a set of key entities (Abubakar Shekau, Boko Haram, Mali, cocaine, etc.) which can be associated to the Emotions' results providing a new analysis perspective.

## Writeprint

The Writeprint feature introduces a whole new level of language analysis providing powerful statistical and semantic text readability indexes to target Biometry and authorship assessment.

With more than 65 core indexes, Writeprint is capable of outlining a document's readability level and the education grade necessary to understand it. The new feature also provides a full set of grammatical and structural analysis items.
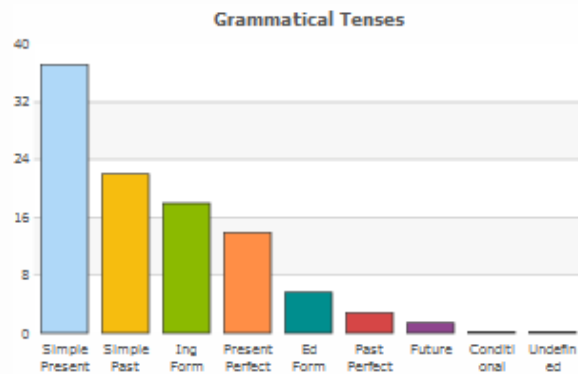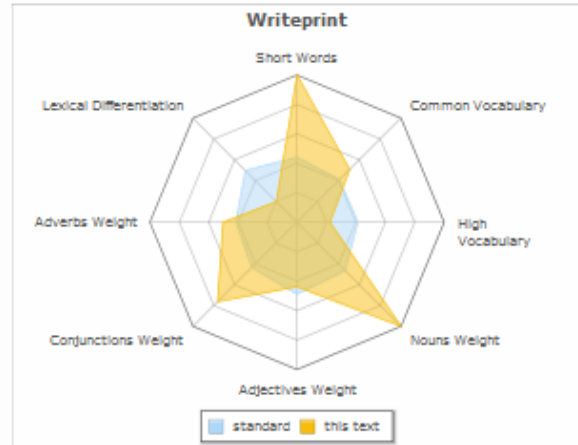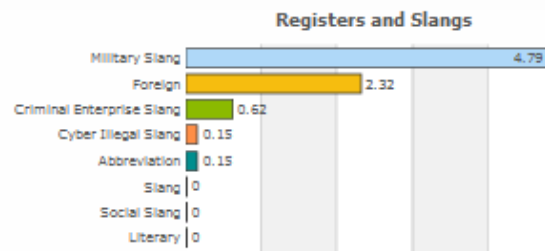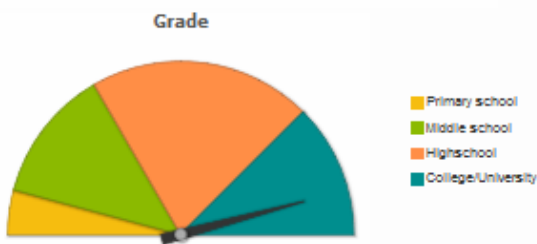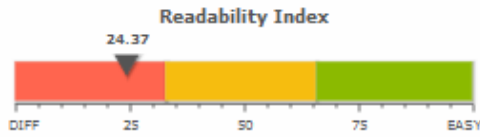
Cogito Intelligence API's readability index is based on standard algorithms for which there's a wide range of scientific literature. We studied and reexamined the Coleman Liau index adding more of our semantic factor greatly increasing our indexes' precision and reliability.

The feature is powered with two sets of statistical and semantic indexes. With more than 5 domain specific slangs and registers, it is possible to achieve better slangs disambiguation, writing styles and topics.

The example below clearly shows how Writeprint is capable of analyzing documents deeply into their textual and grammatical structures working the readability values out. On the right-hand column, the Writeprint graph compares the different scores obtained in the most crucial indexes leaving a "(write-)print" of the author's style and peculiarities.
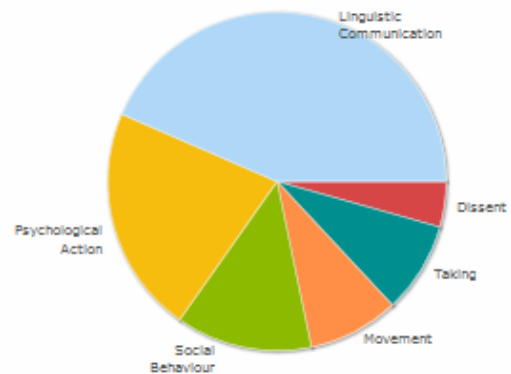
# COGITO® Intelligence API

Home Preview Tagging Categorization Text Mining Semantic Reasoning Fact Mining Emotions People Organizations Places Writeprint

### Readability Index

24.37

DIFF — 25 — 50 — 75 — EASY

### Vocabulary richness

62.68

POOR — 18 — 36 — 54 — 72 — RICH

### Grade

- Primary school
- Middle school
- Highschool
- College/University

### Writeprint

- standard
- this text

### Registers and Slangs

| | |
|---|---|
| Military Slang | 4.79 |
| Foreign | 2.32 |
| Criminal Enterprise Slang | 0.62 |
| Cyber Illegal Slang | 0.15 |
| Abbreviation | 0.15 |
| Slang | 0 |
| Social Slang | 0 |
| Literary | 0 |

### Grammatical Tenses

Simple Present, Simple Past, Ing Form, Present Perfect, Ed Form, Past Perfect, Future, Conditional, Undefined

### Text Statistics

| Index | Value | | |
|---|---|---|---|
| Sentences Count | 28.00 | | - |
| Words per Sentence | 23.11 | * | 10.06-16.92 |
| Characters per Sentence | 134.50 | * | 55.34-109.92 |
| Short Words Count | 13.92 | * | 5.20-9.50 % |
| Different Words Count | 59.87 | * | 66.10-81.26 % |
| Uncommon Words Count | 5.73 | | 5.22-9.62 % |
| Common Vocabulary Weight | 66.24 | | 42.04-100.00 % |
| High Vocabulary Weight | 2.97 | * | 0.96-1.15 % |
| Technical Vocabulary Weight | 1.91 | * | 58.70-68.44 % |
| Nouns Count | 22.1 | * | 8.32-11.67 % |
| Verbs Count | 11.44 | * | 2.94-4.96 % |
| Adjectives Count | 11.28 | | 9.92-12.97 % |
| Conjunctions Count | 4.79 | * | 2.36-4.58 % |
| Adverbs Count | 3.71 | | 2.06-4.68 % |

### Verb Classes

- Linguistic Communication
- Dissent
- Taking
- Movement
- Social Behaviour
- Psychological Action

## Semantic Reasoning

Semantic reasoning brings a new, innovative function which extends the Text mining feature. In fact, for entities related to the Intelligence and Security domain, Semantic Reasoning is able to automatically infer information NOT present within the text, thus providing consistent information about domain entities. As shown in the example below, the Reasoning triggers a process of consistent information enrichment which supports inferences formulation and data analysis.



## Inferential Entities

Inferential Entities are unveiled by assumption of their strong connection with the entities in the text, as highlighted by Semantic Reasoning. So the inferred information is not taken from text but is strongly consistent and related to the detected entities.

# COGITO® Intelligence API

Inferential Entities apply to all People, Organizations and Places entities providing a higher abstraction level to extract further information from inferred metadata and picture new semantic overviews.

## Social Tags Normalization

This innovative feature allows immediate detection of social and web related data, and processes the information to retrieve its meaning. Our Social Tags Normalization feature is the first tool to actually gather genuine information from social networks by detecting crucial entities in hashtags, profile names and URLs.

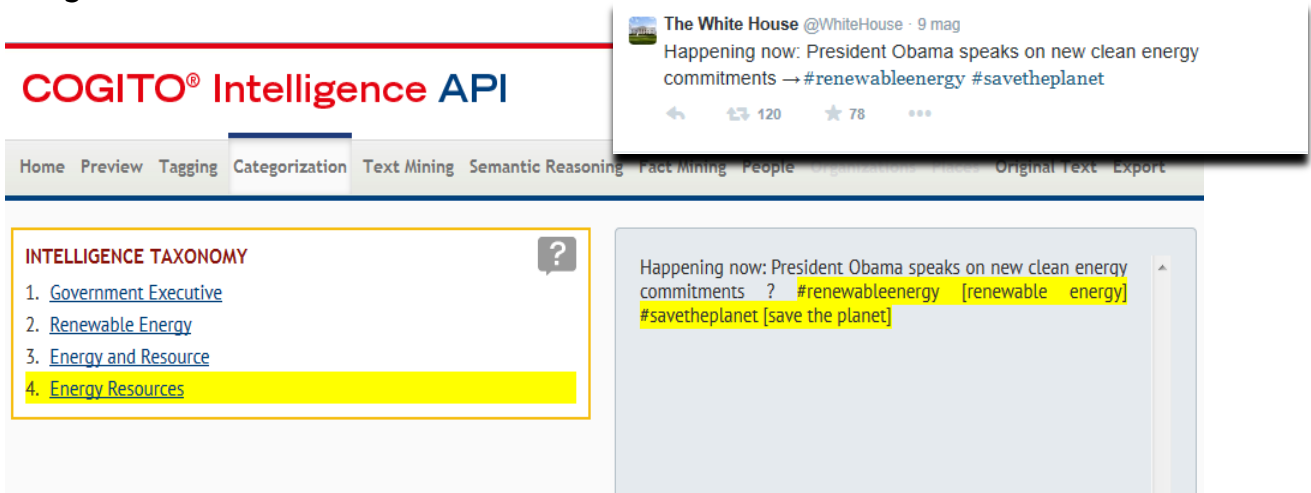| @syrianelectronicarmy | → **Cogito turns it into** → | Syrian electronic army |
| @_cypherpunks_ | → **Cogito turns it into** → | cypher punks |
| @alqaeda | → **Cogito turns it into** → | al qaeda |
| #bringbackourgirls | → **Cogito turns it into** → | bring back our girls |
| #narcoterrorism | → **Cogito turns it into** → | narco terrorism |
| #elections2014 | → **Cogito turns it into** → | elections 2014 |

All of the highlighted texts are core entities which are recognized and extracted from social data. That would not be possible without our normalization feature.

# COGITO® Intelligence API

Normalization also provides crucial information to achieve precise semantic categorization.