



## Caso de uso - VIGILANCIA SECTORIAL



Infoday sobre Actuaciones del Plan de Impulso de las Tecnologías del Lenguaje:  
oportunidades de la compra pública de innovación – 25/04/2017

# Timeline proyecto vigilancia sectorial

11/2013 - 06/2014

Pilotos de implementación de sistemas Big Data para el apoyo a la gestión de ayudas

- Piloto tecnologías BigData: Hadoop, Hive, Pig, Hue, etc
- Piloto aplicación de grafos
- Piloto correlación oferta y demanda de empleo en el sector TIC

01/2015 - 12/2015

Sistema de soporte a la evaluación de ayudas de la Secretaría de Estado de Telecomunicaciones y para la Sociedad de la Información

20/10/2015

Presentación oficial del Plan de Impulso de las Tecnologías del Lenguaje

05/2016 - 12/2016

Nuevas funcionalidades e integración de técnicas relacionadas con la generación de modelos de semántica latente, semántica basada de ontologías y aprendizaje automático aplicados a la vigilancia sectorial

# Participación y colaboración en el proyecto



- Corpus documental** • Entendido como un conjunto de documentos de una misma tipología, que se agrupan en base a una serie de criterios explícitos.

<b>Ayudas</b>	<b>Ayudas SETSI</b>	<b>Ayudas CDTI</b>	<b>Ayudas SEIDI</b>
<b>Patentes</b>	<b>Patentes ES OEPM</b>	<b>Patentes EEUU TIC USPTO</b>	<b>Patentes EU EPO</b>

En nuestro caso concreto, información textual estructurada y no estructurada para: solicitudes de ayuda, descripciones técnicas de patentes, publicaciones científicas, etc

## ETL

Extracción de datos estructurados y no estructurados



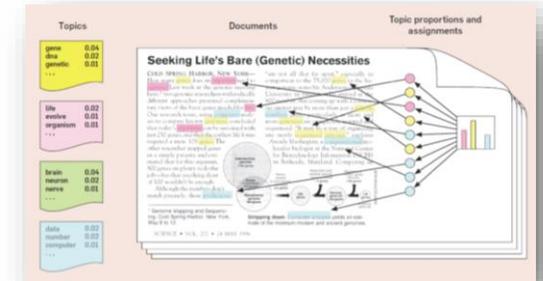
## Técnicas básicas PLN

Tokenizado • Detección de categorías gramaticales (PoS)  
• Lematizado • Wikificado



## Modelado de tópicos

Modelo generativo: estático (LDA, CTM), dinámico (DTM), jerárquico (rLDA) • Vectorización • Huella única del documento



## Técnicas alternativas

Semántica latente avanzada • Semántica basada en ontologías • Machine learning • Word-embeddings • Grafos

# Objetivo: facilitar el trabajo de los evaluadores

## Mejorar el conocimiento general sobre los documentos contenidos en un corpus

Descubrimiento de los temas tratados por cada documento • Vista global de las temáticas tratadas • Vista jerárquica de temas tratados • Obtención de documentos por temática abordada • Identificación de expertos en áreas concretas • Temas que suelen darse de manera conjunta • Evolución en el tiempo de las temáticas

## Herramientas que permiten hacer búsquedas eficientes sobre el corpus documental

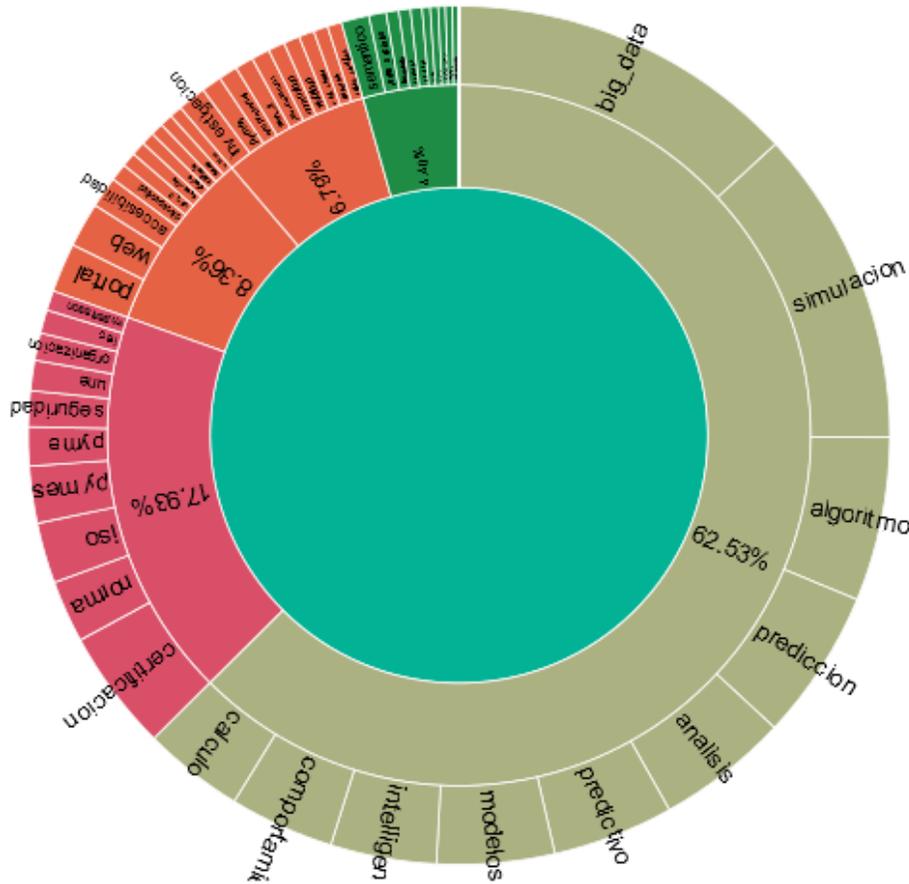
Búsquedas textuales • Búsquedas de documentos parecidos a otro dado • Localización de parejas de documentos temáticamente similares

## Comparación entre distintos corpus documentales

Posibilidad de comparar documentos cuyo contenido se desconoce • Compartiendo modelo de tópicos entre organizaciones, pero sin compartir las solicitudes • Detección de solapamientos • Prevención del fraude • Medición de los resultados de la investigación

## ANÁLISIS DETALLADO DE DOCUMENTOS

Corpus: cuestionarios\_2008-2014 Num. de documentos en el corpus: 9786 Algoritmo de perfilado: estatico Num. de perfiles: 20 Entropía media: 0 Fecha: 2013/24/0 (5)

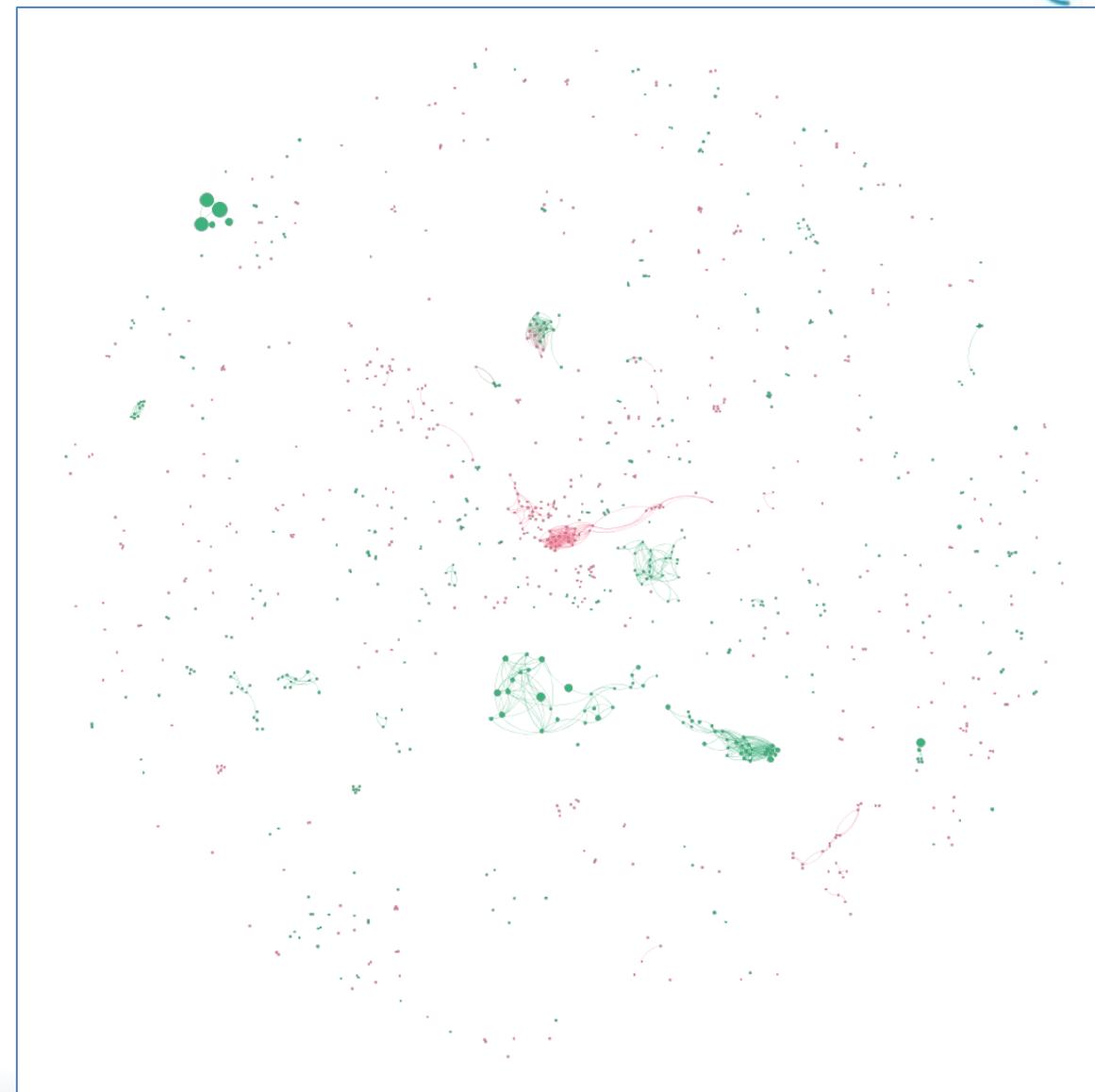
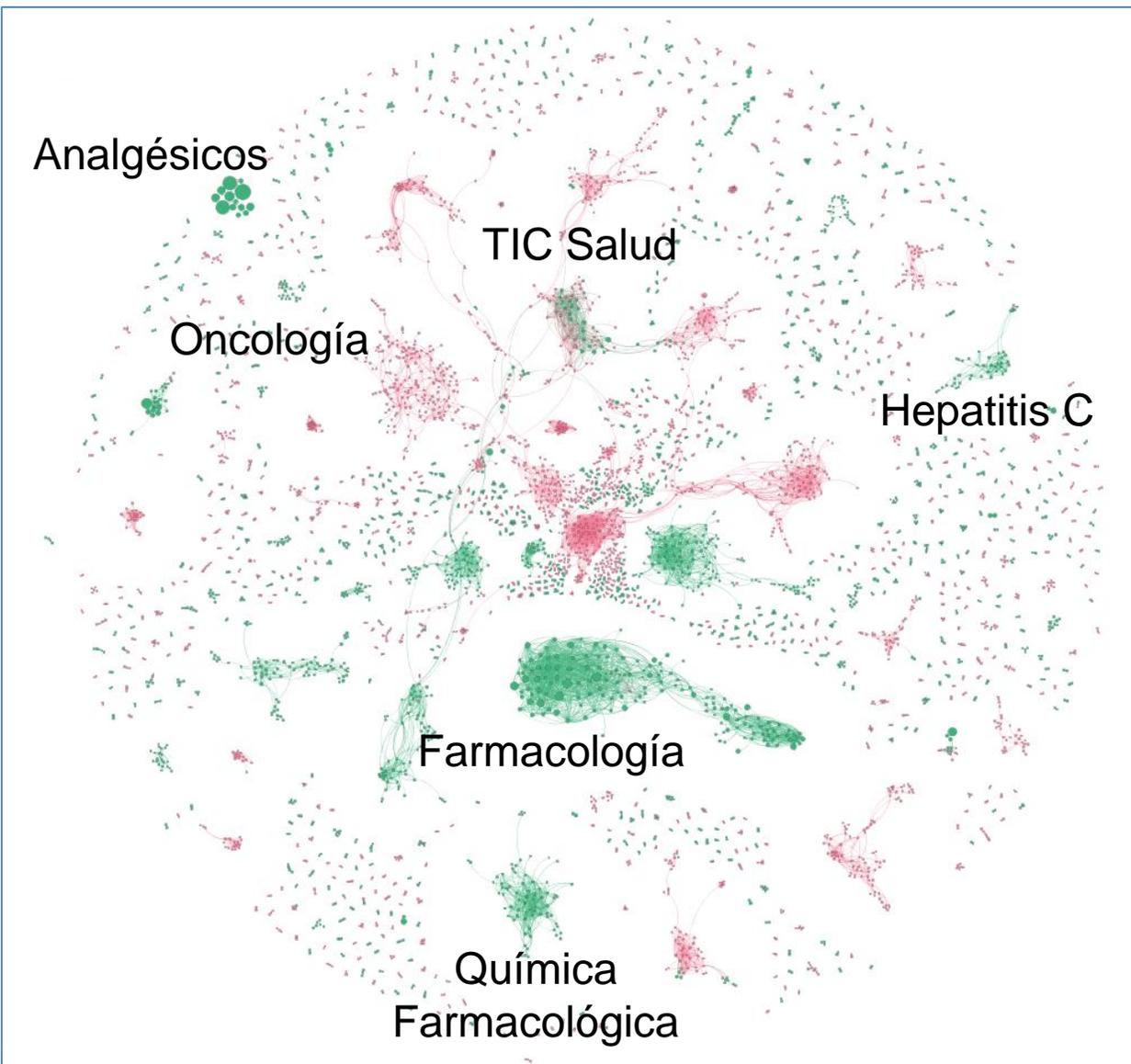


## Representación de la huella única del documento para un modelo de tópicos

Distribución de los tópicos en el documento, porcentaje de texto de la sección técnica de la solicitud:

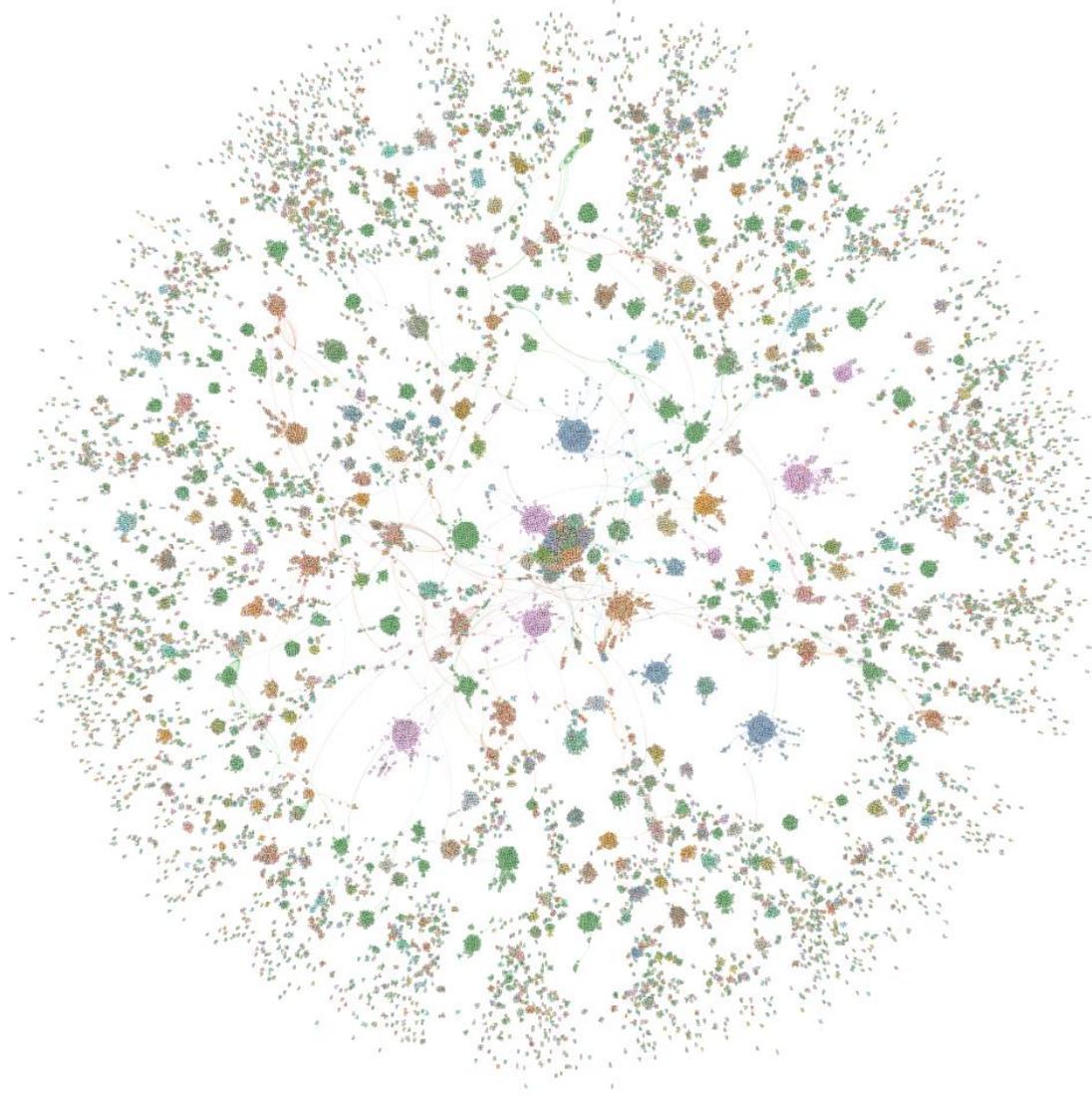
- 62% al tópico “Big data”
- 17% al tópico “Certificaciones”
- 8% al tópico “Portales web”
- 6% al tópico “Investigación”
- 4% al tópico “Semántica”
- ...

# Comparativa investigación sanitaria

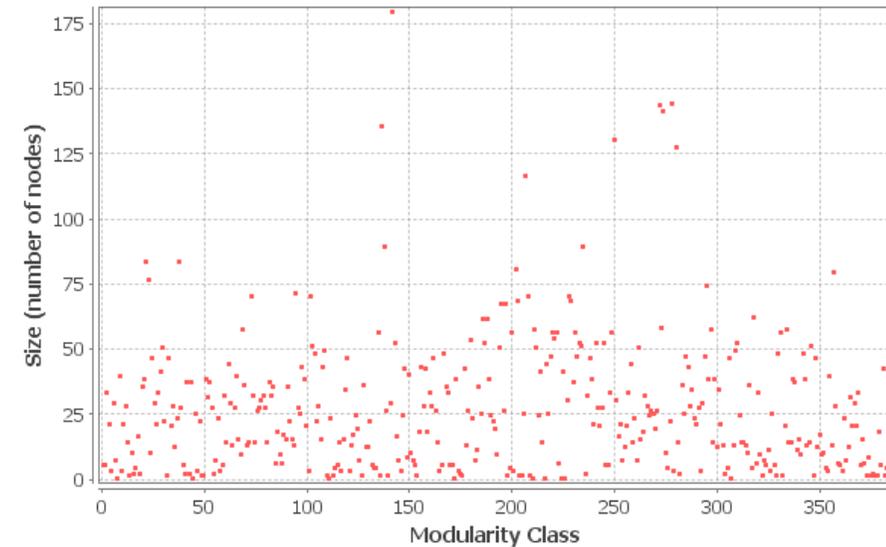


## Results:

Modularity: 0,983  
Modularity with resolution: 49,983  
Number of Communities: 385



## Size Distribution

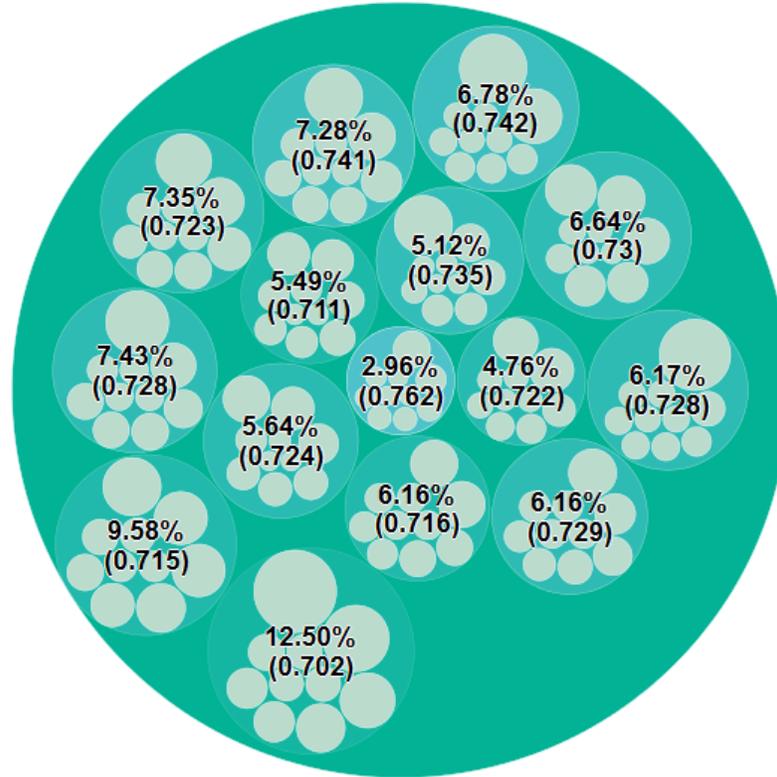


**33,5Kdocs de diferentes programas con el objetivo de identificar temáticas financiadas por varios programas o cambios de nombre del programa a lo largo del tiempo. 385 comunidades.**

# Clasificación de proyectos

## TOPICS OVERVIEW

Corpus: cuestionarios\_2008-2015 Num. de documentos en el corpus: 10220 Algoritmo de perfilado: estatico Num. de perfiles: 15



## TOPICOS DEL MODELO

### PERFIL 7: 6.17%

SEGURIDAD VULNERABILIDAD SECURITY AISLAMIENTO INFRAESTRUCTURA SW MODELADO PROCESOS\_DE\_NEGOCIO PLANTA CRITICO

### PERFIL 6: 6.64%

JUEGO TURISTICO 3D REALIDAD\_AUMENTADA VIDEOJUEGO JUGADOR REDES\_SOCIALES TURISTA VIRTUAL RED\_SOCIAL

### PERFIL 5: 6.78%

VEHICULO SENSOR ROBOT RUTA URBANO GPS CONDUCTOR CIUDAD DETECCION APARCAMIENTO

### PERFIL 4: 7.28%

ENERGETICO ENERGIA EFICIENCIA\_ENERGETICA AGUA ELECTRICO CONSUMO FABRICACION EDIFICIO CONSUMO\_ENERGETICO PLANTA

### PERFIL 3: 7.35%

DOCUMENTO CIUDADANO FIRMA FACTURA FIRMA\_ELECTRONICA TRAMITACION EXPEDIENTE SEGURIDAD ADMINISTRACION DOCUMENTAL

### PERFIL 2: 7.43%

BIG\_DATA SEMANTICO REDES\_SOCIALES ANALISIS ALGORITMO ONTOLOGIA BUSQUEDA OPINION HADOOP TEXTO

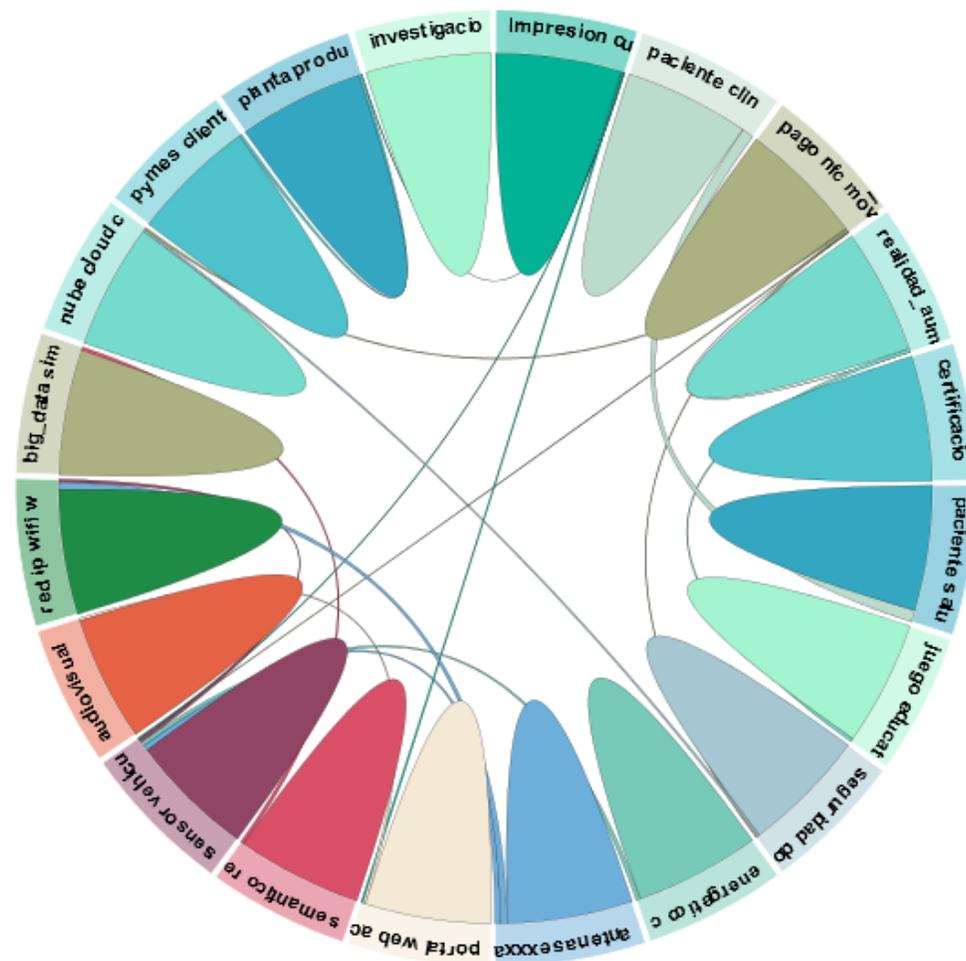
## Visión global del corpus, identificación de temas tratados en los documentos

Visión global del corpus según los tópicos reconocidos, detalle de palabras más frecuentes del tópico y documentos del corpus que mejor se adaptan al tópico

# Reparto de proyectos (correlación)

## ANÁLISIS DE CONEXIONES ENTRE TÓPICOS

Corpus: cuestionarios\_2008-2014 Num. de documentos en el corpus: 9786 Algoritmo de perfilado: estatico Num. de perfiles: 20 Entropía media: 0 Fecha: 20/3/24/0 (5:)



## Reparto de proyectos entre evaluadores

Muestra las relaciones entre tópicos según los documentos analizados.

Permite identificar conocimientos relacionados.

Por ejemplo:

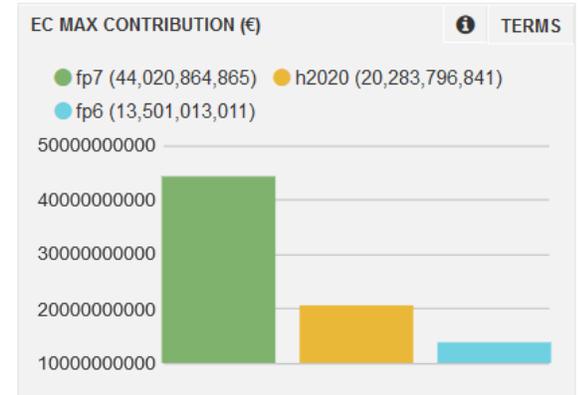
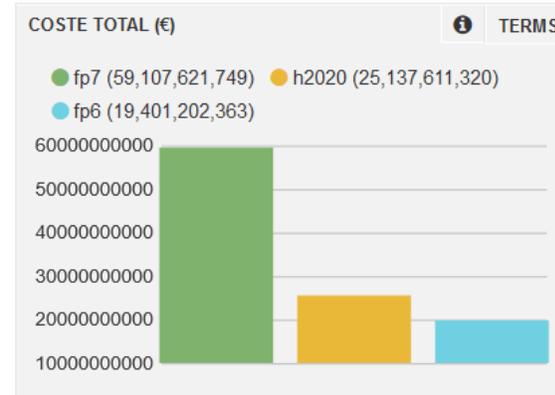
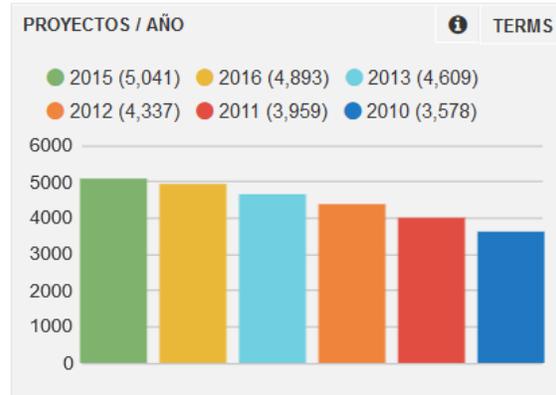
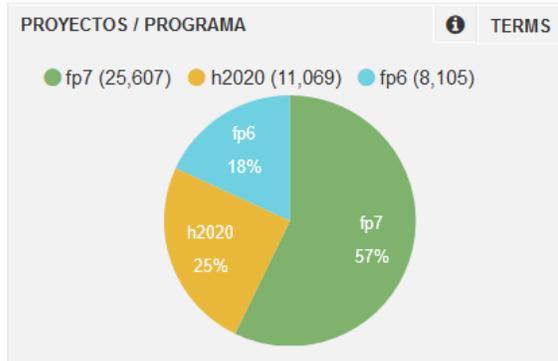
- antena / señal
- red / ip / wifi
- sensor / vehículo

SEARCH i QUERY

Q+

TOTAL HITS i HITS

44,781



FILTERING i FILTERING

**No filters available**

+

EVENTS i TABLE

0 to 10 of 5,000 available for paging

id	programa	title	totalCost	year
100107	fp6	Symbolic Computation Infrastructure for Europe	3351831	2006



# Búsqueda por inferencia temática

Percentil inferior  Percentil superior  Límite de resultados  Método de búsqueda

## Texto

Cyber-physical networks based on embedded systems are part of our society and gain spread and importance. Next generations of aircrafts and cars will be tightly interconnected with each other, with the Internet and other infrastructures. The same holds for many industries and areas of our life. Ubiquitous, highly critical systems go online and create a domain of mixed criticalities, where security and safety requirements of different levels mix. Today, state of the art technology does not provide trustworthiness for such interconnection and mix. The project's cornerstone is MILS (Multiple Independent Levels of Security), a high-assurance security architecture that supports the coexistence of untrusted and trusted components, based on verifiable separation mechanisms and controlled information flow. For the first time in Europe, EURO-MILS does a complete "Common Criteria" security evaluation of a MILS system to its highest levels of assurance, including formal verification engineering. Hardware dependencies are addressed from the beginning by prototype development on two hardware platforms in automotive and avionics. EURO-MILS is strongly market oriented, is carried out in pan-European context, and has an advisory board with government IT security authorities (BSI-Germany and ANSSI-France). In addition, the methodology gained from this high-assurance certification and investigating MILS business, legal, and social acceptance benefit not only the MILS domain but all high-assurance security certifications in Europe. The project brings together eight leading industrial companies, a

Documentos más cercanos

[Generar informe](#)

Identificador 99802

Título Multi-cores Partitioning for Trusted Embedded Systems

[Ver fichero original](#)

## Texto

Growing complexity of applications makes the integration of security and dependability an issue in many domains (e.g. energy supply, transportation, industrial control, aerospace, etc). The engineering of embedded systems needs to take these aspects into account. However, guaranteeing security and dependability in a situation of increasing system complexity is leading to unacceptable development cost and time to market, especially for SMEs, due to the price of tools. The main challenge of this project is supporting mixed criticality embedded systems on multicore open source virtualized platforms in such a way that the development, validation and certification efforts can be lower than the corresponding effort required on independent hardware platforms when using an appropriate methodology. An approach to increasing maintainability and to avoid the growing validation and certification effort is to incorporate mechanisms that establish multiple partitions on the same hardware platform with strict temporal and spatial separation between the individual partitions. In this approach, applications with different levels of dependability can be placed in different partitions and can be validated (and certified if required) in isolation, the MultiPARTES approach. This allows the user to manage complexity while keeping down an escalation of the development effort, but this concept needs to be adapted and applied to multicore and heterogeneous multicore systems. This project aims at developing tools and solutions based on mixed criticality virtualization systems for multicore platforms. The starting point for the virtualization support is

## Clasificación IPC

Etiqueta	Confianza
H	100%
H03	0.98782206%
H04	0.9999916%

## Tic

✘ 0.9985097 %

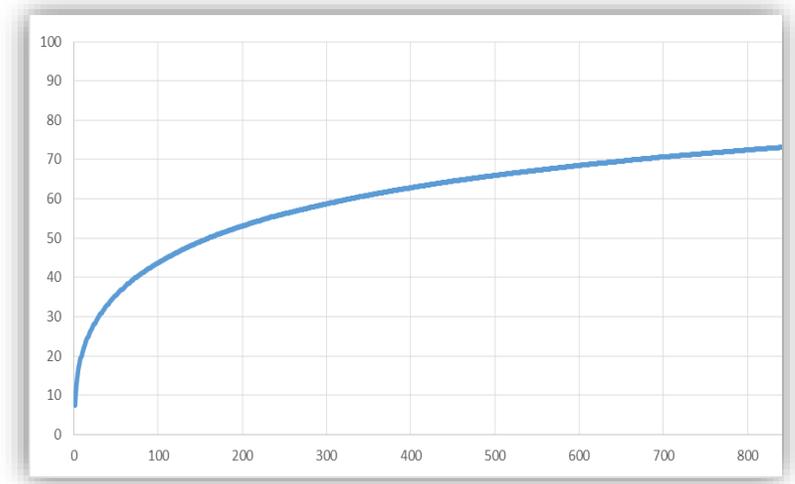
## E-commerce

✔ 0.92791665 %

Permite la búsqueda de documentos del corpus cuyo contenido es más parecido a un texto dado

Se emplea una búsqueda no textual basada en la distancia de sus proyecciones sobre el modelo de tópicos

- **Un 7.5% de las citas buscadas aparecieron como primer resultado de la búsqueda automática.**
- Resultados detallados de la búsqueda de citas:
  - El 20.8% aparecieron entre los 10 primeros resultados
  - El 35.5% entre los 50 primeros
  - El 43.7% entre los 100 primeros
  - **El 50% aparecieron entre los 161 primeros resultados**
- El 87.7% de todas las citas buscadas aparecieron entre las primeras 2.700 patentes seleccionadas por el sistema. **Esto es, para una amplia mayoría de los casos, se redujo la búsqueda al 1% del total de patentes analizadas.**
- Estos resultados se obtuvieron sin aplicar ninguna restricción a la búsqueda. **Aplicando filtros adicionales por metadatos (IPC u otros) se podría reducir más aún el ámbito de búsqueda.**

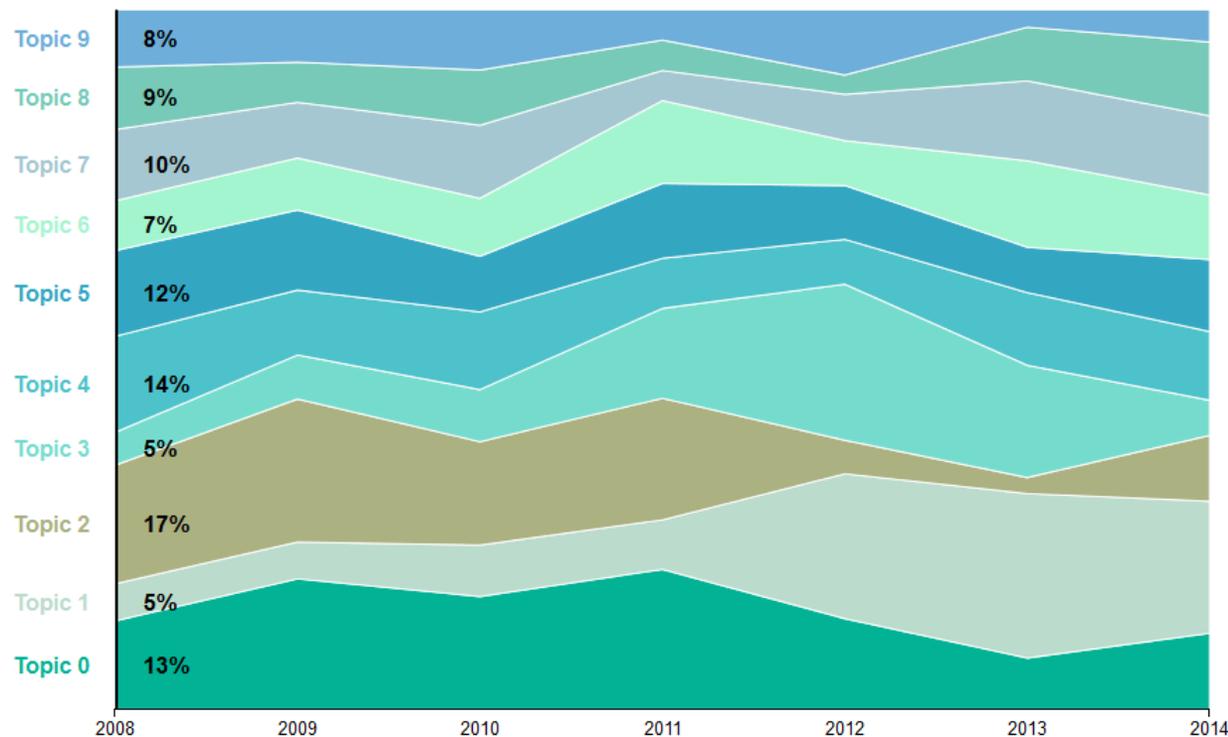


# Evolución temporal

## ANALYSIS OF DYNAMIC PROFILES

Ayudas Concedidas 2008-2014

Corpus: concedidas\_2008-2014 Num. de documentos en el corpus: 2811 Algoritmo de perfilado: dinámico Num. de perfiles: 10 Fecha: 20/6/21/0 (5:)



### TOPICOS DEL MODELO

RED COMUNICACION ACCESO MOVIL  
INTERNET SOLUCION

### TOPIC 3

DIGITAL JUEGO CREACION MOVIL 3D  
SOLUCION CLIENTE VIDEOJUEGO

### TOPIC 2

INVESTIGACION EXPERIMENTAL VIABILIDAD  
DIVULGACION NOVEDAD ESPAXXXOL

### TOPIC 1

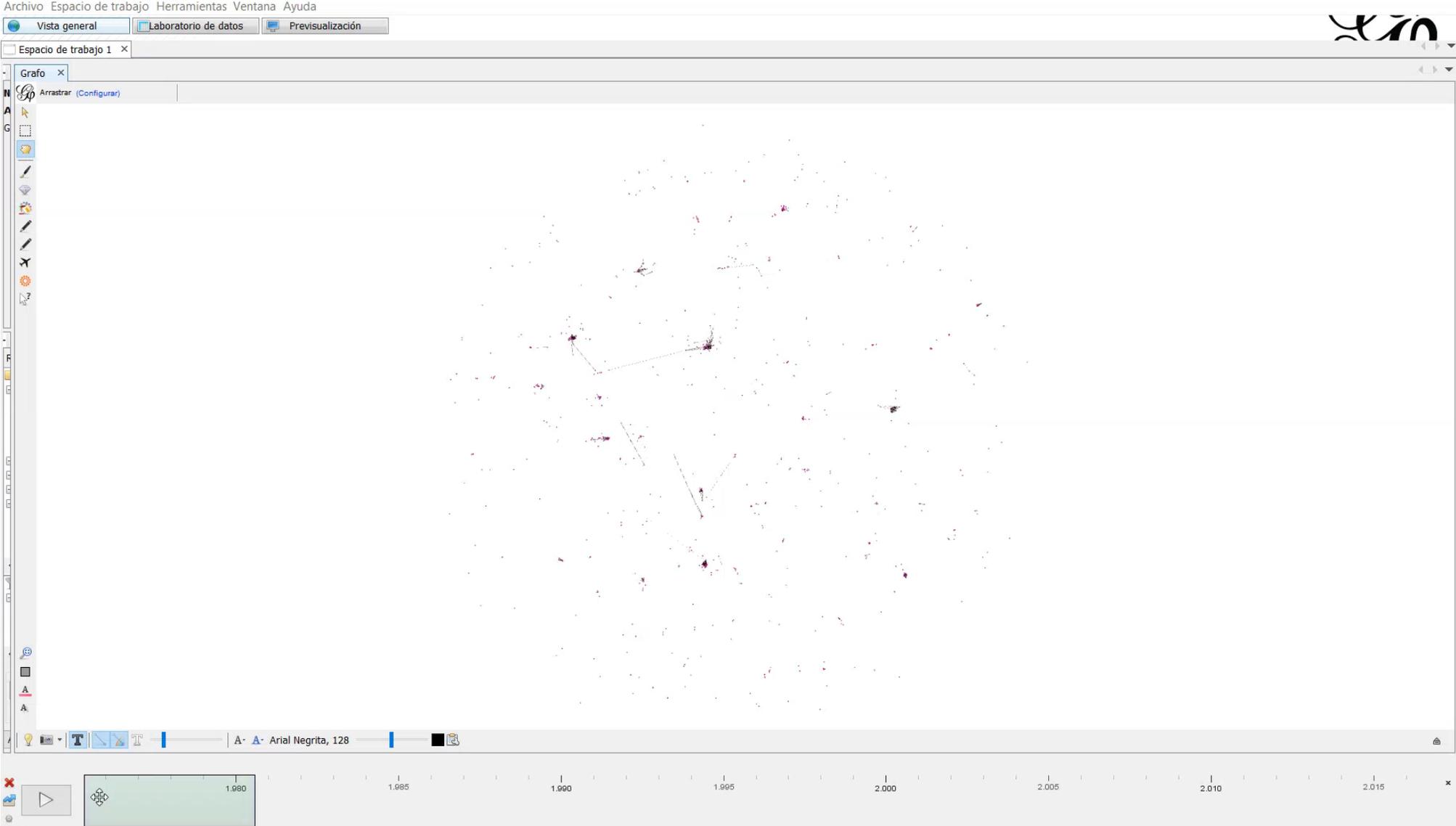
SOLUCION ANALISIS NUBE CLOUD  
ALGORITMO BIG\_DATA CLOUD\_COMPUTING  
MODELOS

### TOPIC 0

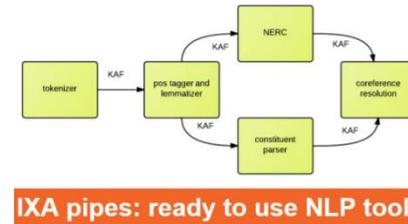
GESTION CLIENTE SOLUCION SEGURIDAD  
PYMES MEJORA TIC ORGANIZACION

## Evolución de la temática contenida en los documentos del corpus

# Evolución temporal National Science Foundation



# Herramientas opensource



- Uso de técnicas **semánticas** empleando un conocimiento base (ontologías específicas dominio)
- Profundizar en técnicas para **independencia idioma**
- Análisis de **redes semánticas**
- **Análisis temporal** sobre modelos de tópicos
- **Reducción del peso** de los modelos
- Aplicación de nuevas técnicas en cálculo de distancia
- Incorporación de **nuevos corpus documentales** y compartición de modelos
- Publicación REST APIs - PaaS / Interoperabilidad (IaaS)
- Avanzar hacia un posible **análisis de impacto** (utopía o realidad??)

# Muchas gracias!!

[david.nieto@satec.es](mailto:david.nieto@satec.es)



satec 

[www.satecgroup.com](http://www.satecgroup.com)