



TFG: Aplicación web para clasificación de documentos de texto

Carlos Sanabria Miranda

Índice

- ▷ Alcance y objetivos
- ▷ Aspectos teóricos
- ▷ Planificación y presupuesto iniciales
- ▷ Casos de uso
- ▷ Subsistemas, componentes y despliegue
- ▷ Evolución de la interfaz gráfica
- ▷ Pruebas
- ▷ Demo
- ▷ Conclusiones
- ▷ Ampliaciones
- ▷ Planificación y presupuesto finales

1.

ALCANCE Y OBJETIVOS

Alcance y objetivos

- ▷ **Aplicación que permita la clasificación de documentos de texto**
 - Esta clasificación se realizará en base al contenido de los documentos, sin que exista una organización previa

Alcance y objetivos

- ▷ Para ello, procesa una colección de documentos, para:
 - Identificar las temáticas (**topics**) presentes en la colección
 - Identificar la relación de cada documento con cada topic
 - Asignar cada documento de la colección a un topic
- ▷ Para identificar los topics se utiliza **aprendizaje no supervisado**
 - Los datos no tienen por qué estar etiquetados
 - El número de topics óptimo se determina en base a los documentos de la colección

Alcance y objetivos

- ▷ Funcionalidades adicionales:
 - Identificar los **documentos más representativos** de cada topic
 - Identificar los **documentos** de la colección **más similares** a un nuevo documento
 - Generar un **resumen extractivo** de un documento de texto
- ▷ Implementar una **API REST** que encapsule esta funcionalidad
- ▷ Implementar un **frontend** que consuma la API REST

Alcance y objetivos

▷ **Sus aplicaciones son muy diversas**

▷ **Ejemplos:**

- Organizar en grupos un conjunto de documentos desordenados (noticias, medicina, ...)
- Clasificar un correo electrónico en una determinada carpeta

▷ **En general, cualquier caso de uso que involucre:**

- Documentos de texto
- Su clasificación

2.

ASPECTOS TEÓRICOS

Aspectos teóricos: **LDA**

- ▷ Latent Dirichlet Allocation
- ▷ **Algoritmo** utilizado para la **identificación** de los **topics** y la **clasificación** de los documentos
- ▷ Ideas principales:
 - **Agrupar** las **palabras** presentes en los documentos **en topics**
 - Cada **topic** se describe mediante una **distribución de palabras** (keywords)
 - Cada **documento** se describe mediante una **distribución de topics**
 - Las palabras de cada documento definen pertenencia a los topics
 - No tiene en cuenta el orden de las palabras

Aspectos teóricos: LDA

Documentos

Most Christian denominations view the baptism of Jesus ...



Football and basketball are 2 commonly practiced sports ...

Topics

Religion

Sports

Medicine

Computer science

Palabras

1,7% christianism

1,4% jesus

0,01% basketball

game

...

cancer

football

patient

cancer

database

80%

5%

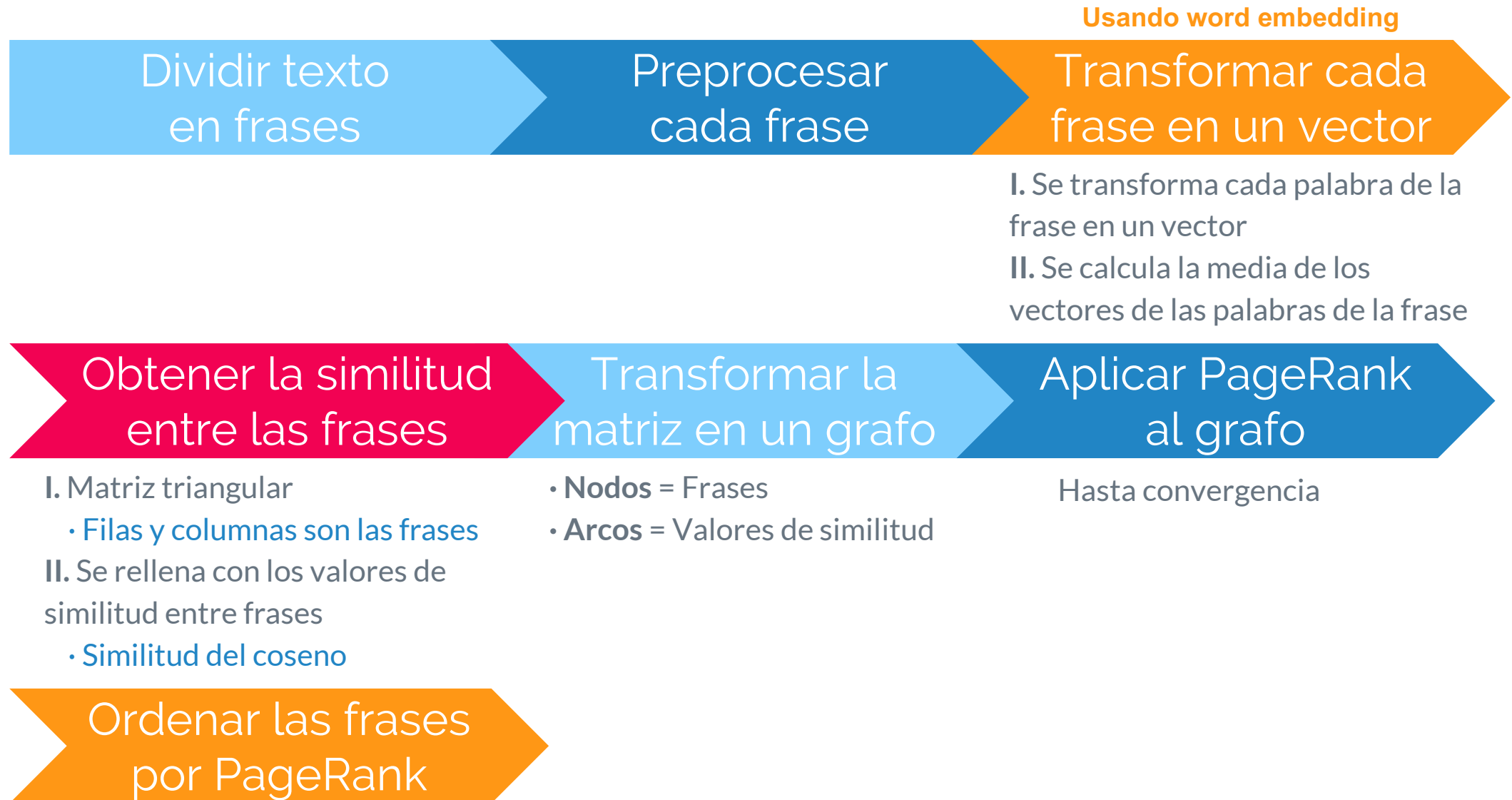
75%

3%

Aspectos teóricos: **TextRank**

- ▷ **Algoritmo** utilizado para la **generación de resúmenes**
- ▷ Basado en el algoritmo PageRank
 - Asigna puntuación numérica a las páginas web / documentos en función de su importancia

Aspectos teóricos: TextRank



3. PLANIFICACIÓN Y PRESUPUESTO INICIALES

Planificación inicial

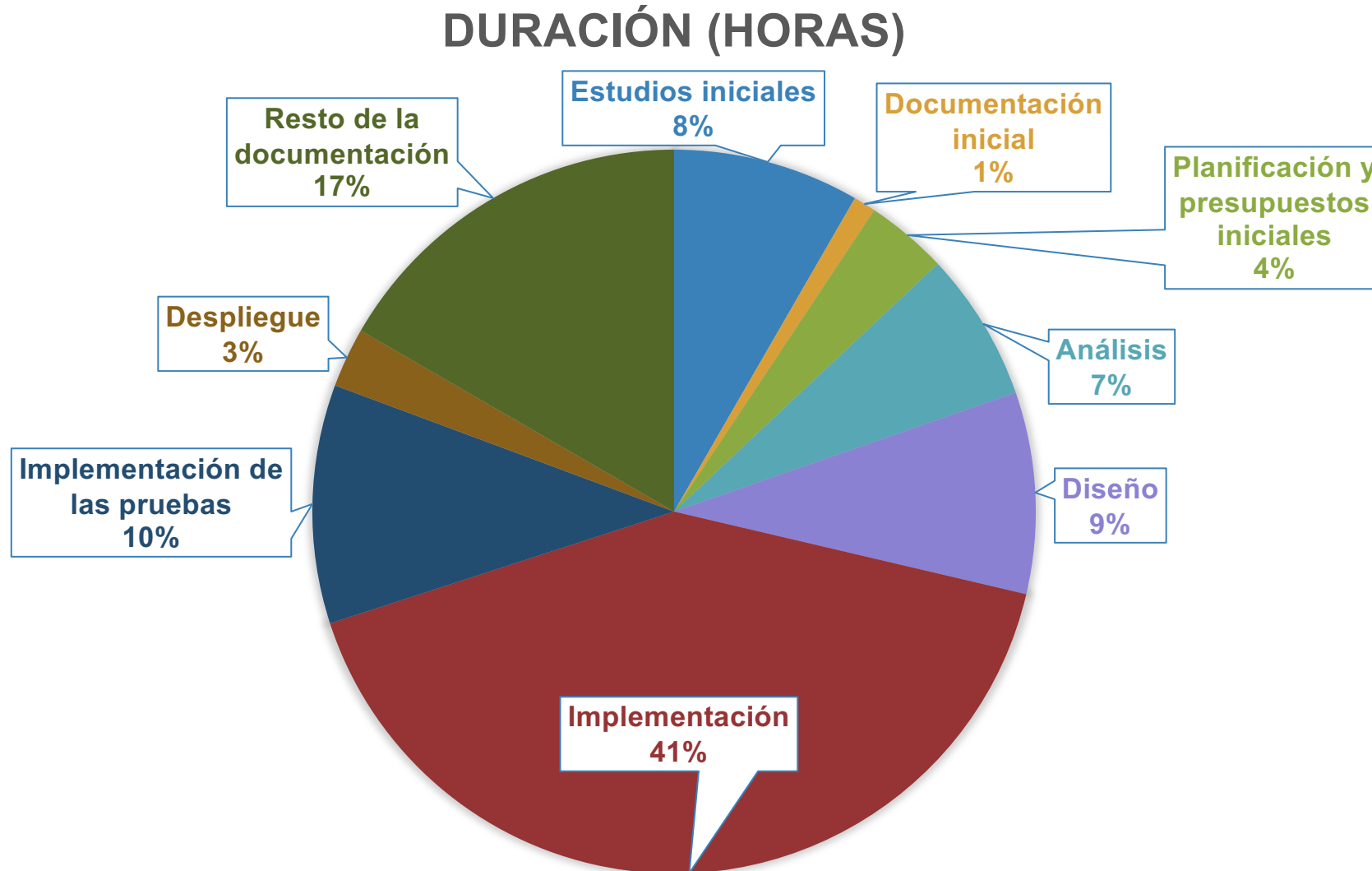
▷ **01 febrero 2019 – 11 junio 2019**

▷ **9 etapas:**

- Estudios iniciales
- Documentación inicial
- Planificación y presupuestos iniciales
- Análisis
- Diseño
- Implementación
- Implementación de las pruebas
- Despliegue
- Resto de la documentación

▷ **Duración total: 300h**

Planificación inicial

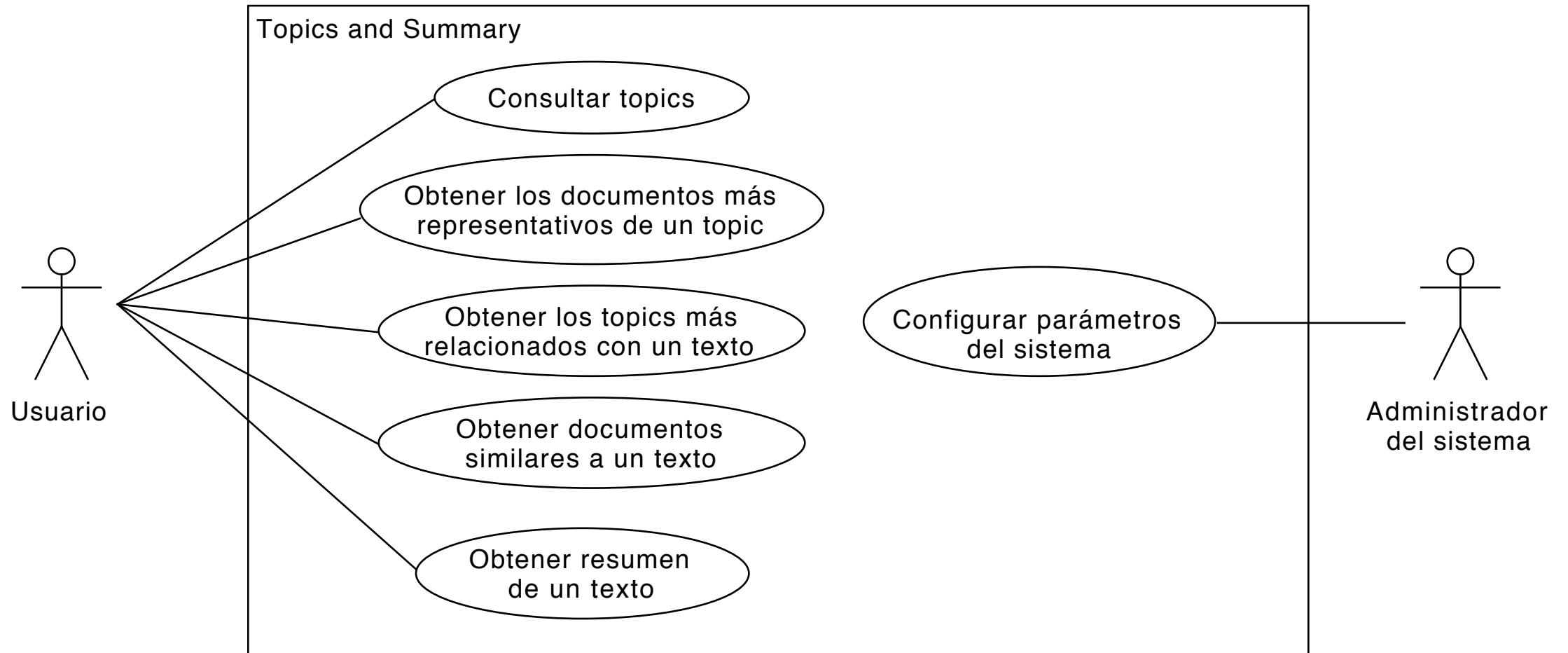


Presupuesto inicial

Presupuesto de costes inicial		
<i>Cód.</i>	<i>Partida</i>	<i>Total</i>
01	Estudios iniciales y planificación inicial	1.350,00 €
02	Análisis y diseño del sistema	1.762,50 €
03	Implementación y despliegue del sistema	6.248,40 €
04	Documentación	1.987,50 €
05	Otros costes	1.134,60 €
Total Coste		12.483,00 €

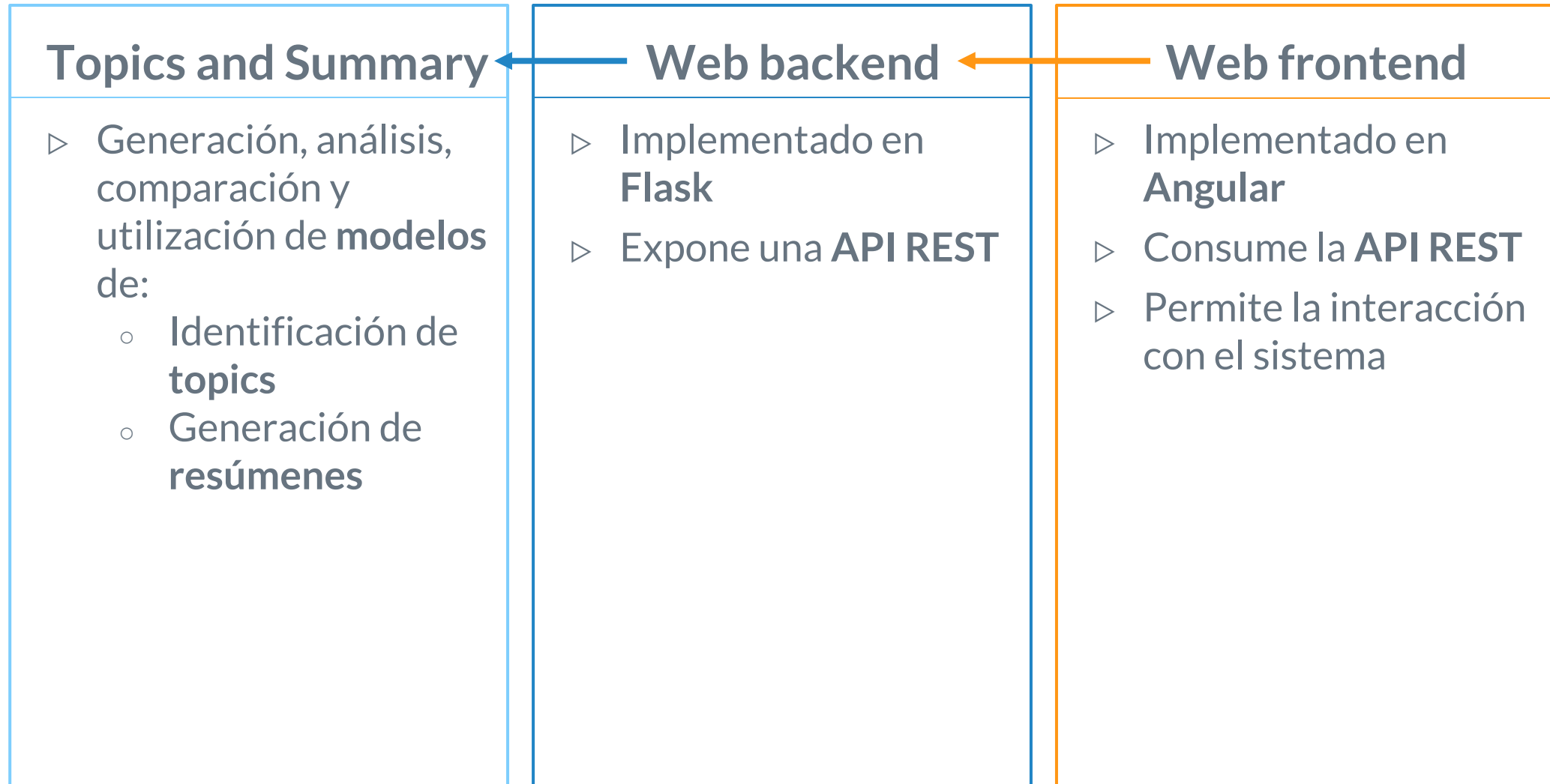
4. CASOS DE USO

Casos de uso

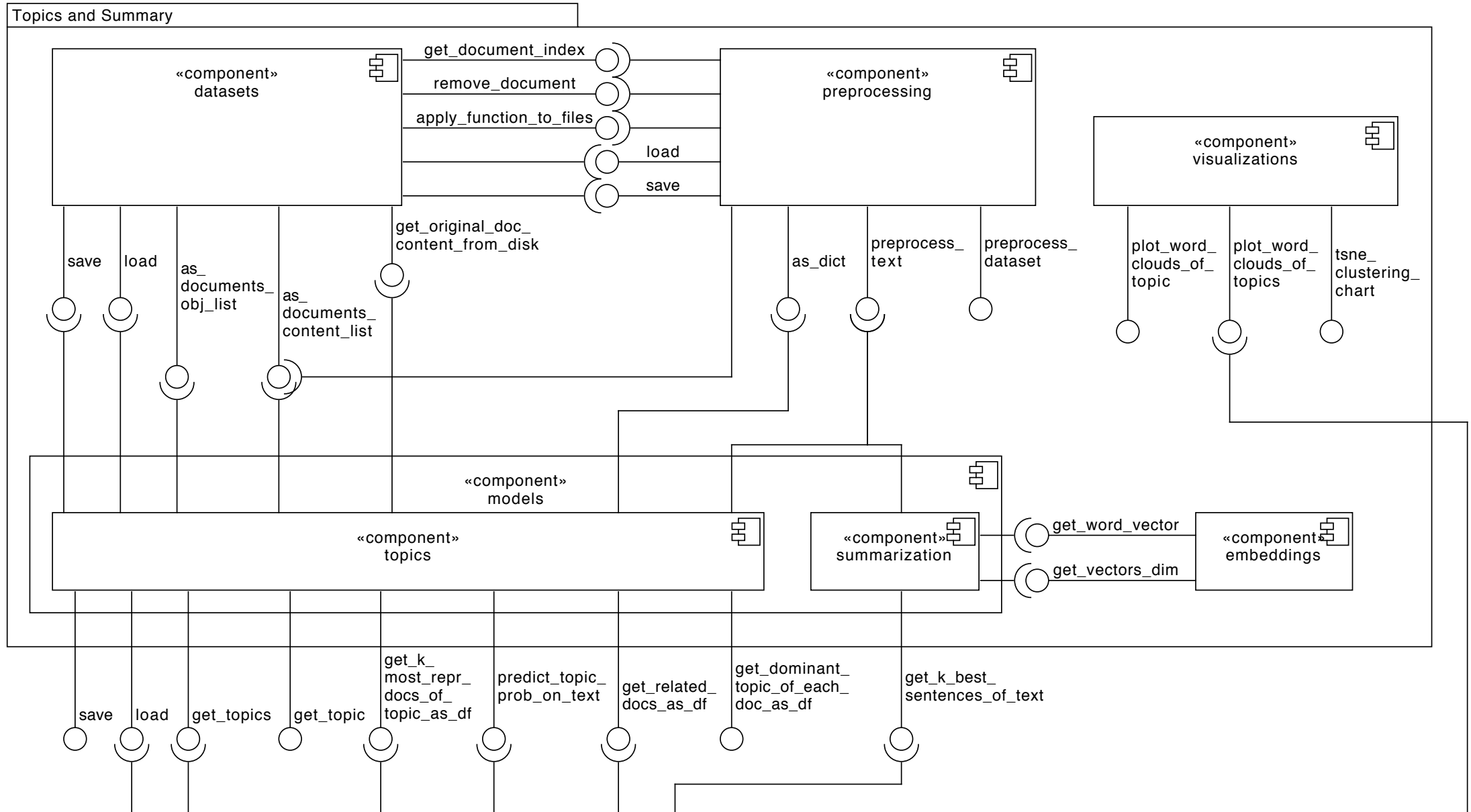


5. SUBSISTEMAS, COMPONENTES Y DESPLIEGUE

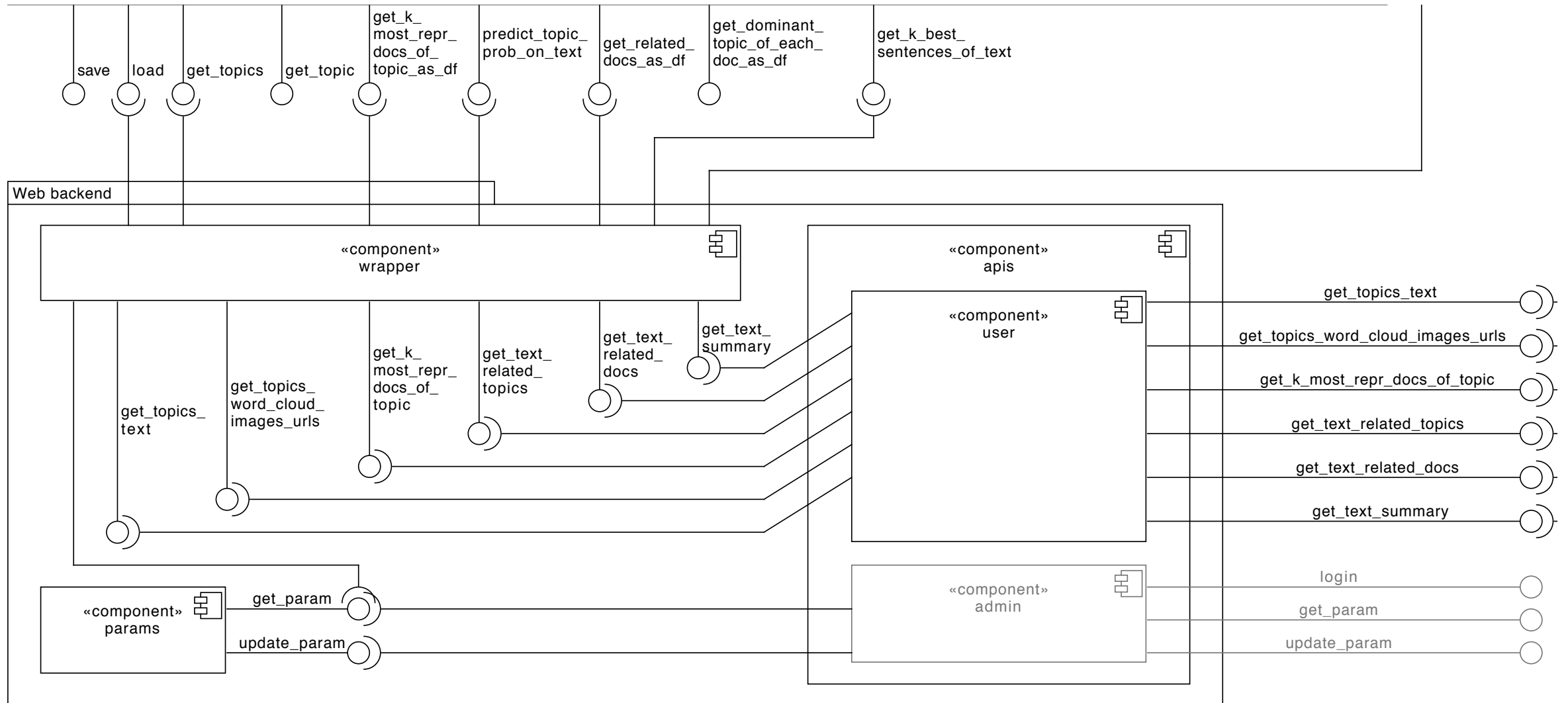
Subsistemas



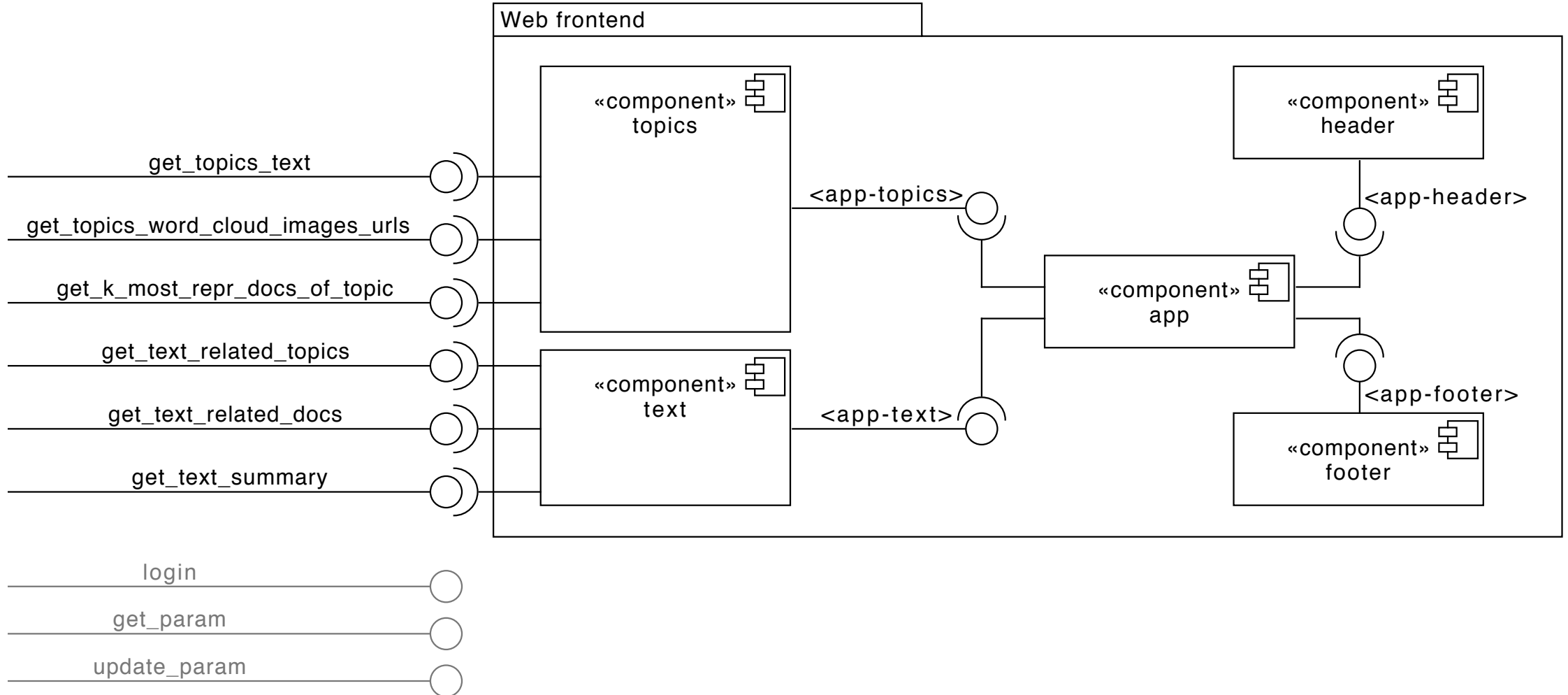
Componentes (Topics and Summary)



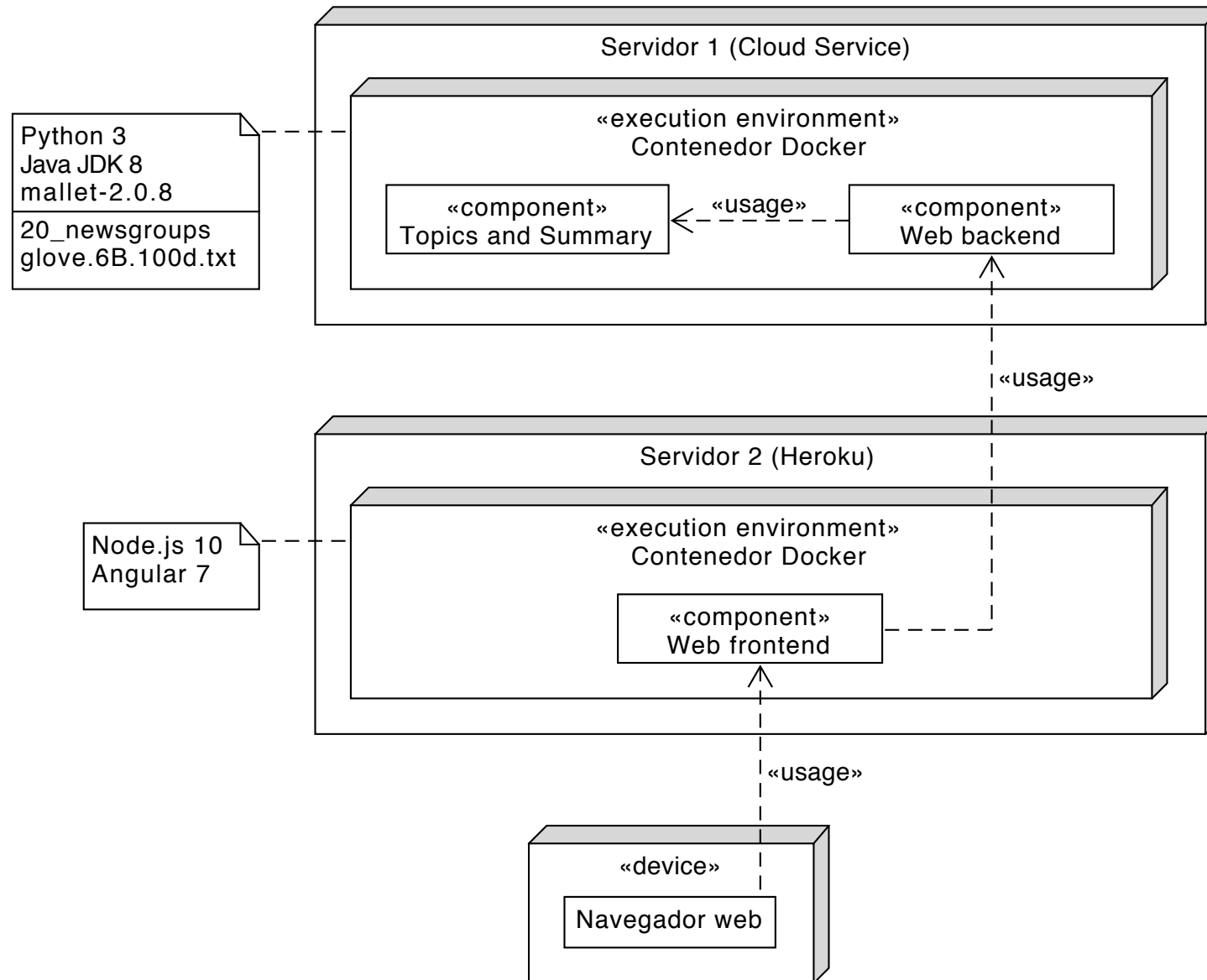
Componentes (Web backend)



Componentes (Web frontend)



Despliegue



6. EVOLUCIÓN DE LA INTERFAZ GRÁFICA

Evolución de la interfaz (fase análisis)

Zona de visualización de topics

Zona de introducción de texto

Evolución de la interfaz (fase análisis)

Topics

Texto	Nube de palabras
Num topic	Keywords
0	study: 0.81% drug: 0.76% disease: 0.62% patient: 0.60%
1	god: 1.51% christian: 1.12% jesus: 0.92% bible: 0.80%
2	windows: 2.81% file: 2.76% program: 2.62% server: 2.60%
..	..

Documentos más representativos del topic 2

> Resumen del documento más representativo del topic 2

Probabilidad: 65%

Contenido del documento más representativo del topic 2

> ...

Texto

Introduce un texto...

> Configuración topics relacionados

> Configuración documentos similares

> Configuración resumir texto

<Resultados de la opción seleccionada en la Zona Texto>

Topics

Texto	Nube de palabras		
Topic 0	Topic 1	Topic 2	Topic 3
study patient drug effect disease	god christian bible moises jesus	windows server file image program	space nasa system earth mission

Documentos más representativos del topic 2

> Resumen del documento más representativo del topic 2

Probabilidad: 65%

Contenido del documento más representativo del topic 2

> ...

Texto

Introduce un texto...

> Configuración topics relacionados

> Configuración documentos similares

> Configuración resumir texto

<Resultados de la opción seleccionada en la Zona Texto>

Evolución de la interfaz (fase diseño)

https://www.topics-and-summary-models.com

Topics

Texto

Nube de palabras

Núm. Topic	Keywords	
0	study: 0.81% drug: 0.76% disease: 0.62% patient: 0.60% effect: 0.53%	
1	god: 1.51% christian: 1.12% jesus: 0.92% bible: 0.80% moises: 0.79%	
2	windows: 2.81% file: 2.76% program: 2.62% server: 2.60% image: 2.53%	
...

Documentos más representativos del topic 2

Esta es la frase más importante del documento más representativo del topic 2.

Probabilidad: 65%
Esta es una frase del documento más representativo del topic 2.
Esta es otra frase del documento más representativo del topic 2. ...

Esta es la frase más importante del segundo documento más representativo del topic 2.

Texto

Introduce un texto

Topics relacionados

Número máx. de topics devueltos

1620

Mostrar visualización gráfica ☒

Obtener

> Documentos similares

> Resumir texto

https://www.topics-and-summary-models.com

Topics

Texto

Nube de palabras

Topic 0

study patient
drug effect
disease

Documentos más representativos

Topic 1

god bible
christian moises
jesus

Documentos más representativos

Topic 2

windows server
file image
program

Documentos más representativos

Topic 3

space nasa
system mission
earth

Documentos más representativos

Documentos más representativos del topic 2

Esta es la frase más importante del documento más representativo del topic 2.

Probabilidad: 65%
Esta es una frase del documento más representativo del topic 2.
Esta es otra frase del documento más representativo del topic 2. ...

Esta es la frase más importante del segundo documento más representativo del topic 2.

Texto

Introduce un texto

Topics relacionados

Número máx. de topics devueltos

1620

Mostrar visualización gráfica ☒

Obtener

> Documentos similares

> Resumir texto

7. PRUEBAS

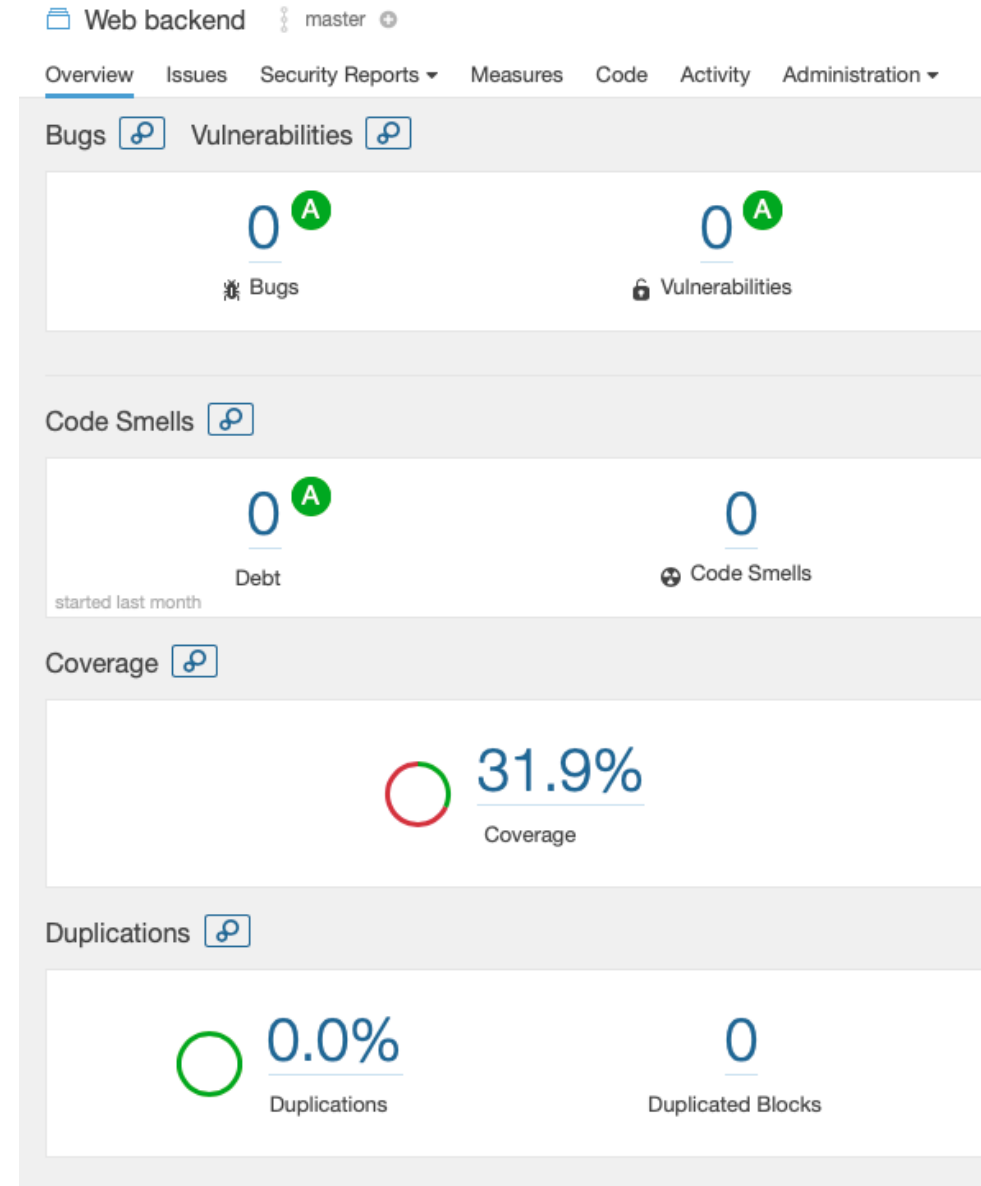
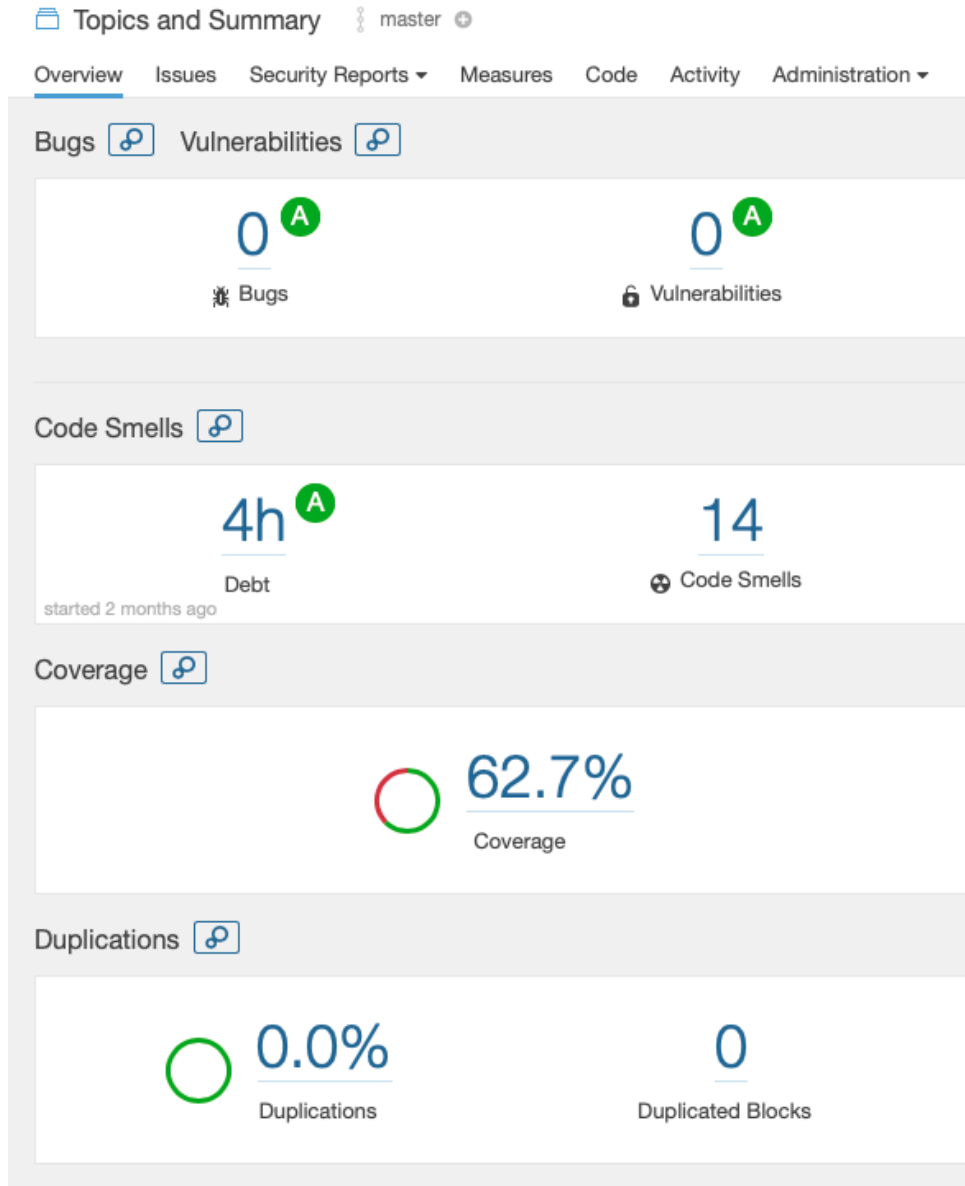
Pruebas unitarias automatizadas

▷ Subsistema **Topics and Summary**

▷ Subsistema **Web backend:**

- Componente params
- Resumen de un texto cuando TextRank:
 - Converge
 - No converge

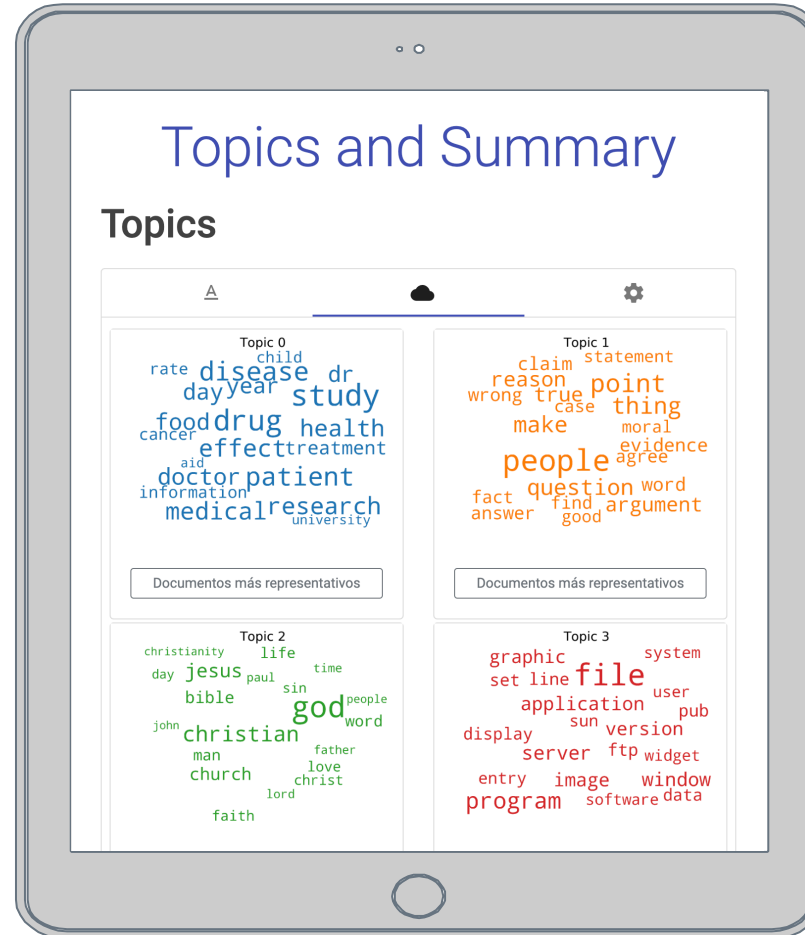
Análisis estático de código



Pruebas de usabilidad

- ▷ Actividades guiadas y cuestionarios a **3 usuarios**
- ▷ **Resultados:**
 - Número errores y tiempo tareas **bajo**
 - **Más solicitado:** Notificación cuando aparece la información dinámicamente
- ▷ **Posibles cambios:**
 - **Aumentar** número mensajes **ayuda**
 - **Disminuir tiempo** generación **resúmenes**
 - Permitir **datasets español**
 - **Añadir scroll automático** que desplace al usuario a la información que aparece dinámicamente

Pruebas de adaptabilidad



8. DEMO

9. CONCLUSIONES

Conclusiones

- ▷ El sistema **cumple los objetivos**
- ▷ Los **modelos** dan **buenos resultados**
 - Aunque el tiempo de generación de resúmenes es en ocasiones elevado con textos grandes
- ▷ El **frontend** sigue los **prototipos** diseñados
- ▷ La **arquitectura** en 3 subsistemas fue un **acierto**
- ▷ Ha supuesto un **gran aprendizaje**, ya que se ha:
 - Profundizado en Procesamiento de Lenguaje Natural
 - Pasado por todas las fases del desarrollo de software

10. AMPLIACIONES

Ampliaciones



Implementar las funcionalidades del administrador



Probar los modelos de extracción de topics con más datasets



Evaluar los resultados del algoritmo TextRank



Probar otros algoritmos de generación de resúmenes



Aumentar el número de mensajes de ayuda en la GUI



Generar los resúmenes de los documentos del dataset en batch



Permitir documentos en español



Añadir i18n al frontend



Hacer la interfaz web accesible

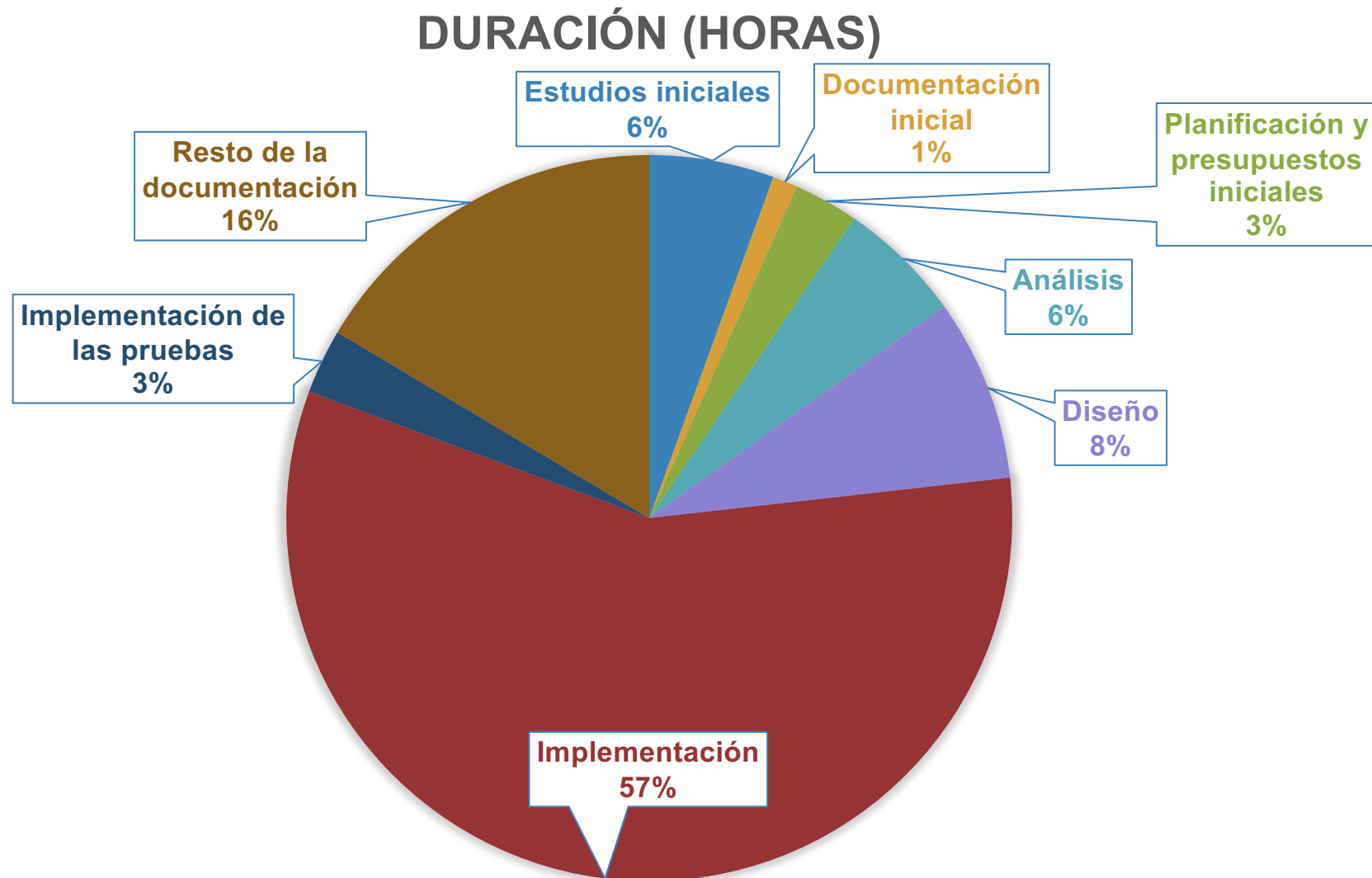
11.

PLANIFICACIÓN Y PRESUPUESTO FINALES

Planificación final

- ▷ **01 febrero 2019 – 07 julio 2019**
- ▷ Mismas etapas, salvo despliegue
- ▷ Duración total: ~~300h~~ 450h

Planificación final



Presupuesto final

Presupuesto de costes final		
<i>Cód.</i>	<i>Partida</i>	<i>Total</i>
01	Estudios iniciales y planificación inicial	1.425,00 €
02	Análisis y diseño del sistema	2.306,25 €
03	Implementación y despliegue del sistema	10.312,45 €
04	Documentación	2.962,50 €
05	Otros costes	1.354,20 €
Total Coste		18.360,40 €

Total Coste Presupuesto Inicial: 12.483,00 €

Incremento en el coste: 5.877,40€

¡Gracias!
¿Preguntas? 

Créditos

Un agradecimiento especial a todas las personas que hicieron y lanzaron estos increíbles recursos de forma gratuita:

▷ Plantilla de la presentación de [SlidesCarnival](#)