

Docode 5: Building a scalable plagiarism detection system

Gaspar Pizarro V.

Web Intelligence Centre

November 22, 2016

Outline

1 Proposed system

2 Architecture

3 Results

4 Conclusions

The core algorithm

- Suspect vs source
- Based on word n-grams
(shingles)

this module works as the user interface

The core algorithm

- Suspect vs source
- Based on word n-grams
(shingles)

this module works as the user interface

The core algorithm

- Suspect vs source
- Based on word n-grams
(shingles)

```
module works user interface
```

- [(module), (works), (user), (interface)]
- [(module works), (works user), (user interface)]
- [(module works user), (works user interface)]

The core algorithm

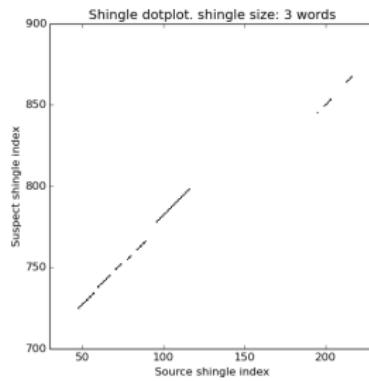
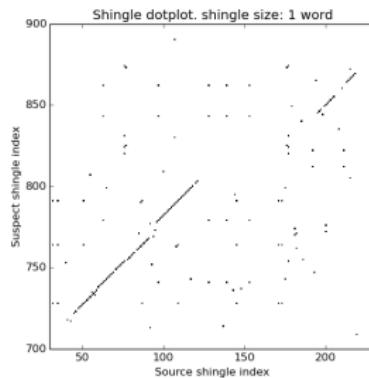
- Suspect vs source
- Based on word n-grams (shingles)

module works user interface

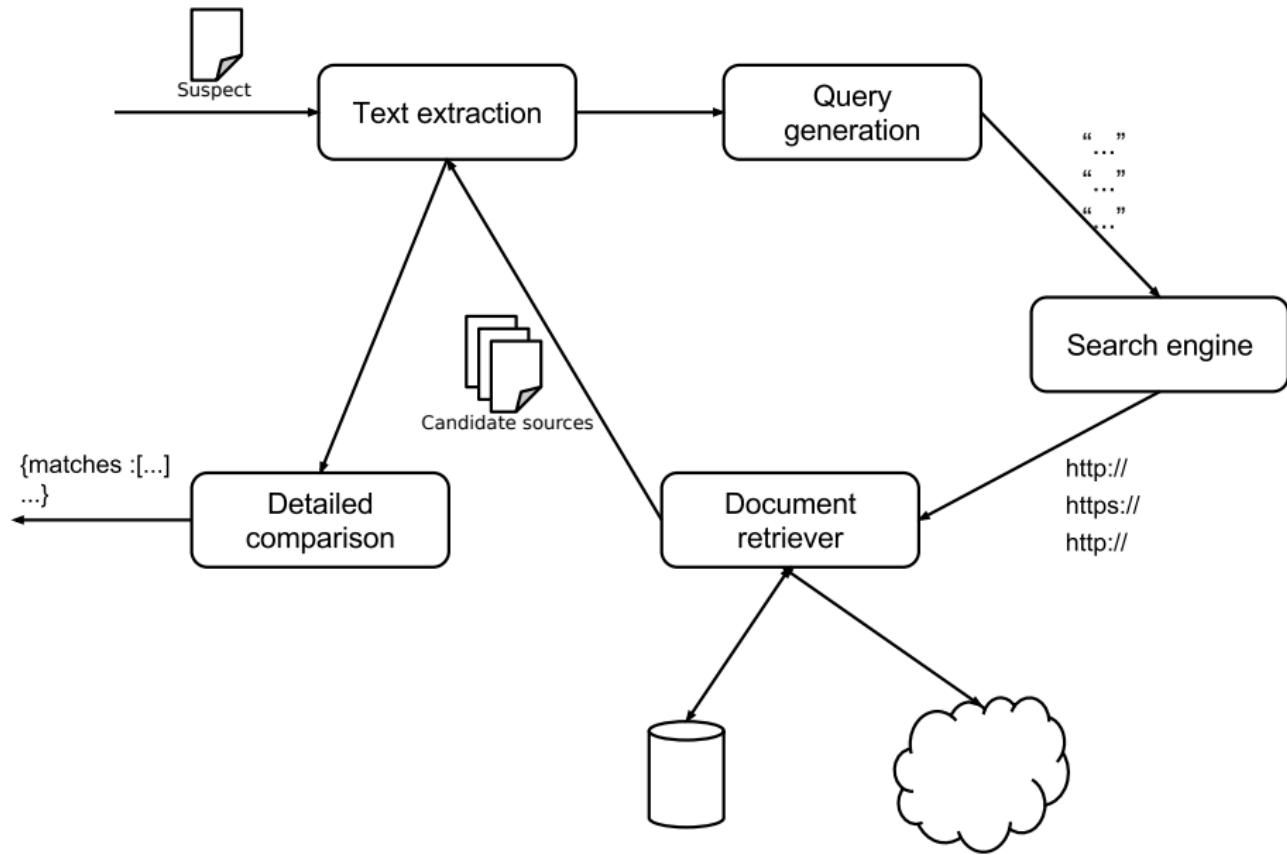
- [(module), (works), (user), (interface)]
- [(module works), (works user), (user interface)]
- [(module works user), (works user interface)]

Dotplot

- Find lines
- Merge closest ones
- Delete smallest ones



Web/Repository analysis



Web/Repository analysis: Example



the_suspect.pdf

Web/Repository analysis: Example

Chile colonial fue un periodo de la historia ubicada entre Fines del Siglo XVI y Principios del Siglo XIV. Esta época se caracterizó por la creación de una organización institucional muy compleja, el mestizaje y la combinación cultural que darían origen a la sociedad chilena como tal. La colonia Tuvo varios aspectos y características, como la forma de administrar Chile; Los sistemas económicos que eran antiguamente utilizados en España también serían empleados en América, ver qué sería lo que traería más productividad y provecho a nuestra colonia; el sistema educacional que sería empleado, en algunos casos los de mayor situación económica tenían más oportunidades y opciones de estudio que los demás, y como fue la estructura social de esos tiempos, Quienes salían favorecidos por el rango social en el que se encontraban, y quienes eran discriminados y afectados por el bajo rango social; esto produciría desacuerdos en opiniones y complicaciones entre la población que en ese entonces habitaba América, también hubo un factor muy importante que influyó en todos estos aspectos, que fue la religión tenía influencia en todo tipo de decisión En este periodo también se formó una gran combinación racial y cultural, que se originó con la convivencia de españoles e indígenas.

Web/Repository analysis: Example

Chile periodo historia Fines Siglo Principios Siglo XIV época creación

organización mestizaje combinación origen sociedad colonia aspectos características forma
Chile

sistemas España empleados América productividad provecho colonia sistema casos situación

oportunidades opciones estudio estructura tiempos favorecidos rango rango desacuerdos opiniones

complicaciones población América factor aspectos religión influencia tipo decisión periodo

combinación convivencia españoles indígenas

Web/Repository analysis: Example

https://es.wikipedia.org/wiki/Historia_de_Chile

https://es.wikipedia.org/wiki/Mestizaje_en_Am%C3%A9rica

https://es.wikipedia.org/wiki/Econom%C3%ADa_de_Espa%C3%B1a

<https://es.scribd.com/doc/57536530/Estructura-de-Planes-y-Programas.pdf>

<https://es.wikipedia.org/wiki/Gripe>

http://www.pps.k12.or.us/district/depts/edmedia/videoteca/curso3/htmlb/SEC_51.HTM

Web/Repository analysis: Example

https://es.wikipedia.org/wiki/Historia_de_Chile

https://es.wikipedia.org/wiki/Mestizaje_en_Am%C3%A9rica

https://es.wikipedia.org/wiki/Econom%C3%ADa_de_Espa%C3%B1a

<https://es.scribd.com/doc/57536530/Estructura-de-Planes-y-Programas.pdf>

<https://es.wikipedia.org/wiki/Gripe>

<http://www.pps.k12.or.us/district/depts/edmedia/videoteca/curso3/htmlb/SEC-51.HTM>

Web/Repository analysis: Example



Web/Repository analysis: Example

El periodo prehispánico corresponde a la historia de las diferentes etnias amerindias presentes en el territorio, extendiéndose desde alrededor del año 14 800 a. C. hasta la llegada de los españoles. A partir de 1492, se iniciaron las exploraciones europeas en el continente americano. En 1520 Fernando de Magallanes y su expedición fueron los primeros europeos en llegar a Chile por el sur a través del estrecho que hoy lleva su nombre, y en 1536 Diego de Almagro comandó una expedición hasta el Valle del Aconcagua y el norte del actual Chile.

La economía de España es la quinta por tamaño en la Unión Europea y la decimotercera a nivel mundial en términos nominales.¹⁷ En términos relativos o de paridad de poder adquisitivo, se encuentra también entre las mayores del mundo (ver Anexo:Paises por PIB (PPA)). Según un informe de The Economist, España es el 10º país del mundo con mayor calidad de vida.¹⁸ Como en la economía de todos los países europeos, el sector terciario o sector servicios es el que tiene un mayor peso. La moneda de España es, desde 2002, el euro.

La principal estrategia para la consecución de este objetivo en educación básica plantea realizar una reforma integral de la educación básica, centrada en la adopción de un modelo educativo basado en competencias que responda a las necesidades de desarrollo de México en el siglo XXI con miras a lograr mayor articulación y eficiencia entre preescolar, primaria y secundaria

Web/Repository analysis: Example

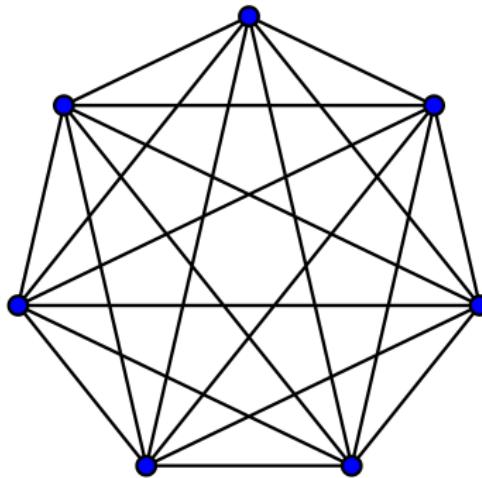
El periodo prehispánico corresponde a la historia de las diferentes etnias amerindias presentes en el territorio, extendiéndose desde alrededor del año 14 800 a. C. hasta la llegada de los españoles. A partir de 1492, se iniciaron las exploraciones europeas en el continente americano. En 1520 Fernando de Magallanes y su expedición fueron los primeros europeos en llegar a Chile por el sur a través del estrecho que hoy lleva su nombre, y en 1536 Diego de Almagro comandó una expedición hasta el Valle del Aconcagua y el norte del actual Chile.

La economía de España es la quinta por tamaño en la Unión Europea y la decimotercera a nivel mundial en términos nominales.¹⁷ En términos relativos o de paridad de poder adquisitivo, se encuentra también entre las mayores del mundo (ver Anexo:Paises por PIB (PPA)). Según un informe de The Economist, España es el 10º país del mundo con mayor calidad de vida.¹⁸ Como en la economía de todos los países europeos, el sector terciario o sector servicios es el que tiene un mayor peso. La moneda de España es, desde 2002, el euro.

La principal estrategia para la consecución de este objetivo en educación básica plantea realizar una reforma integral de la educación básica, centrada en la adopción de un modelo educativo basado en competencias que responda a las necesidades de desarrollo de México en el siglo XXI con miras a lograr mayor articulación y eficiencia entre preescolar, primaria y secundaria

Group analysis

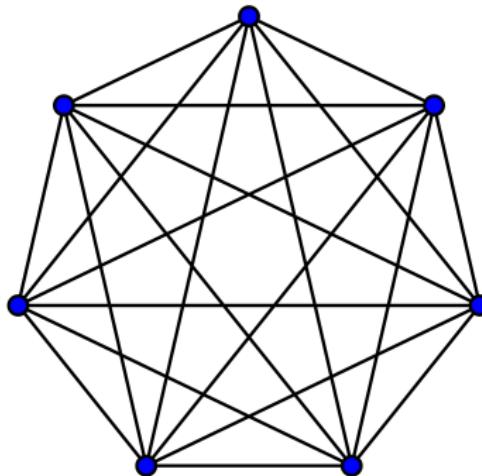
Brute force



- Grows like n^2
- Preprocessing

Group analysis

Brute force



- Grows like n^2
- Preprocessing
- It can be improved

Group analysis

Inverted index

Doc1

El periodo prehispánico corresponde a la historia de las diferentes etnias amerindias presentes en el territorio, extendiéndose desde alrededor del año 14 800 a. C. hasta la llegada de los españoles. A partir de 1492, se iniciaron las exploraciones europeas en el continente americano. En 1520 Fernando de Magallanes y su expedición fueron los primeros europeos en llegar a Chile por el sur a través del estrecho que hoy lleva su nombre, y en 1536 Diego de Almagro comandó una expedición hasta el Valle del Aconcagua y el norte del actual Chile.

Doc2

La economía de **España** es la quinta por tamaño en la Unión Europea y la decimotercera a nivel mundial en términos nominales.¹⁷ En términos relativos o de paridad de poder adquisitivo, se **encuentra** también entre las mayores del mundo (ver Anexo:Paises por PIB (PPA)). Según un informe de The Economist, **España** es el 10º país del mundo con mayor calidad de vida.¹⁸ Como en la economía de todos los países europeos, el sector terciario o sector servicios es el que tiene un mayor peso. La moneda de **España** es, desde 2002, el euro.

Doc3

La principal estrategia para la consecución de este objetivo en educación básica plantea realizar una reforma integral de la educación básica, se **encuentra** centrada en la adopción de un modelo educativo basado en competencias que responda a las necesidades de desarrollo de México en el siglo XXI con miras a lograr mayor articulación y eficiencia entre preescolar, primaria y secundaria

Group analysis

Inverted index

Words to documents

periodo	→	Doc1, 3
España	→	Doc2, 15, 301, 495
Encuentra	→	Doc2, 186 Doc3, 146

Group analysis

Inverted index

Words to documents

periodo	→	Doc1, 3
España	→	Doc2, 15, 301, 495
Encuentra	→	Doc2, 186 Doc3, 146

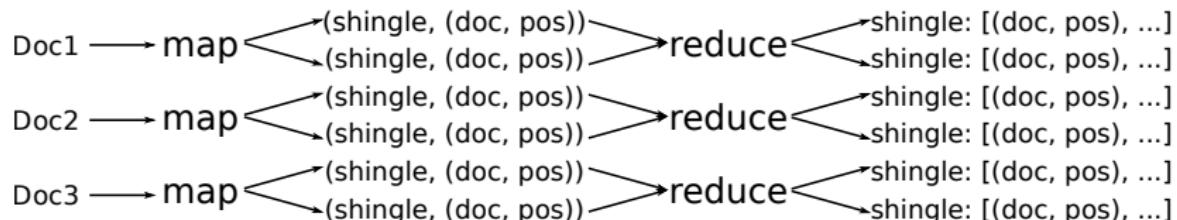
Differences

- Shingles instead of words
- Position in list of shingles instead of position in string

Group analysis

MapReduce, two passes

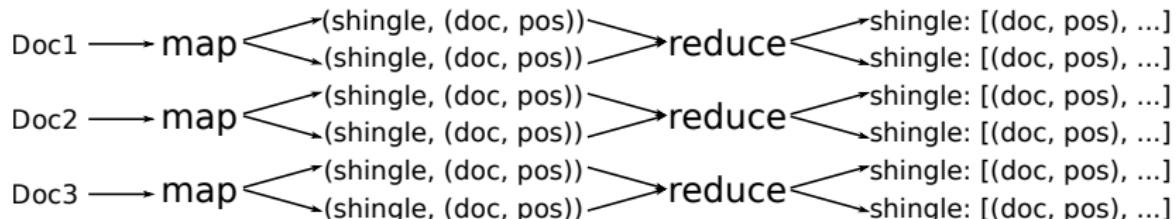
- Indexing job



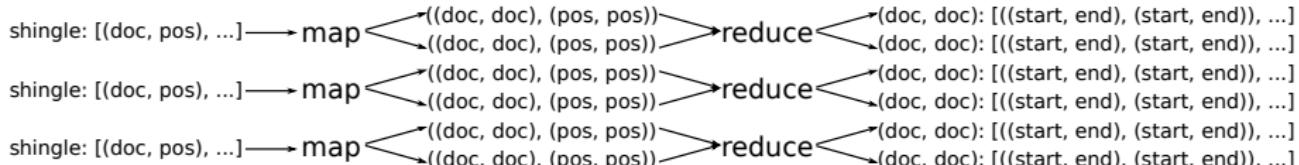
Group analysis

MapReduce, two passes

- Indexing job



- Matching job



Outline

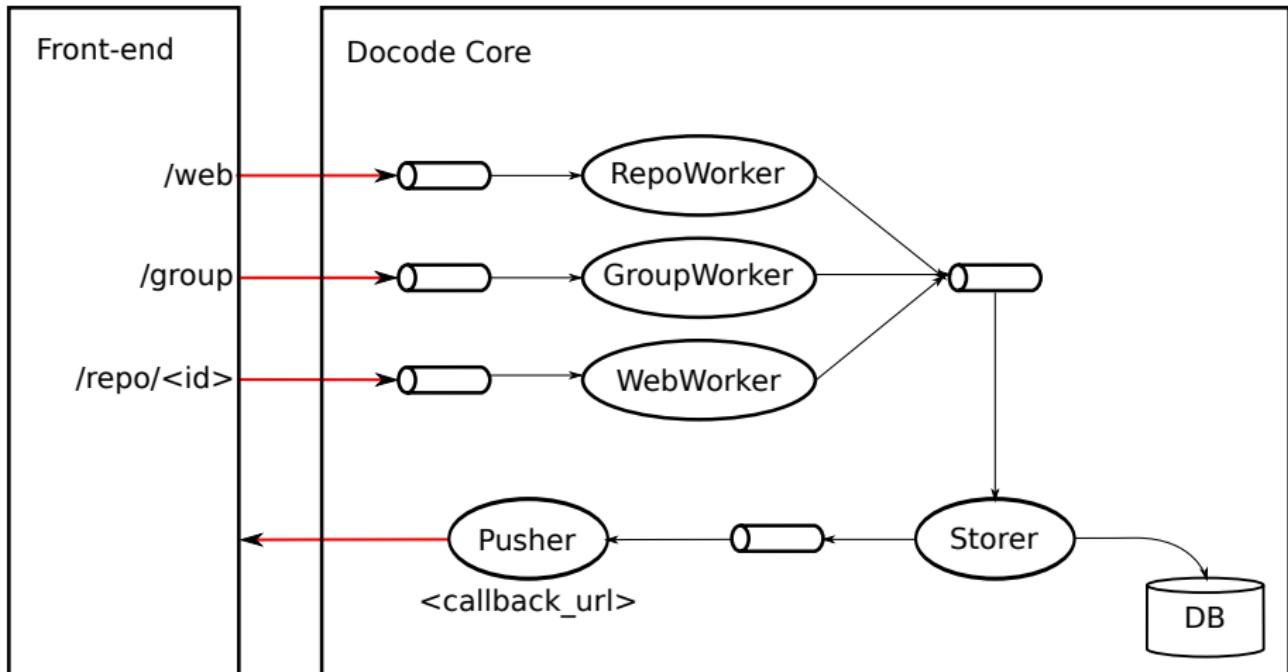
1 Proposed system

2 Architecture

3 Results

4 Conclusions

High-level architecture



Front-end lite

Saldo diario

A 99997249 palabras	B 4998 archivos
---------------------	-----------------

Análisis de originalidad

Por favor, carge el archivo que desea revisar con Docode.**
(Extensiones permitidas: doc, docx, pdf, txt, rtf)

Anastre el archivo a esta caja o [haga click aquí](#)

100%

Ha cargado el archivo: IntroducciónOscarMolina5.txt.
El archivo contiene: 1669 palabras.

Documentos subidos

Nombre archivo	Fecha	# palabras	Estado	Acciones
Muerte.txt	25/05/2016 11:57	587	Subido	Procesar Eliminar

Documentos en análisis

Nombre archivo	Fecha	# palabras	Acciones
IntroducciónOscarMolina5.txt	25/05/2016 11:57	1669	Actualizar

Historial de documentos analizados

Nombre archivo	Fecha solicitud	Fecha recepción	# palabras	% plagio	Estado	Acciones
transfasisco.txt	25/05/2016 11:40	25/05/2016 11:44	1081	6.9%	Finalizado	Revisar
Letter_JVélezquez.pdf	14/01/2016 16:51	14/01/2016 16:53	914	0%	Finalizado	Revisar
_2_Paulina_Saldívar (copia).docx	13/10/2015 11:26	13/10/2015 11:32	2070	17.1%	Finalizado	Revisar
_2_Paulina_Saldívar.docx	6/10/2015 9:49	6/10/2015 10:07	2070	16.2%	Finalizado	Revisar
Guitarra.txt	5/10/2015 18:03	5/10/2015 18:04	943	26.2%	Finalizado	Revisar
_2_Paulina_Saldívar.docx	4/09/2015 12:51	4/09/2015 12:54	2070	17.2%	Finalizado	Revisar

Front-end lite

Reporte de originalidad

Archivo analizado	<u>Z_Paulina_Galdívar (copia).docx</u>
Peso	12993
Nº de palabras analizadas	2970
Fecha de solicitud	13/10/2015 11:26
Fecha de entrega de reporte	13/10/2015 11:32

Docode encontró coincidencias con 7 páginas de la Web:

Fuente	Índice Plagio (%)	Acciones
La Tormenta - Ensayos para estudiantes https://www.clubensayos.com/Tecnolog%C3%ADa/La-Tormenta/558854.html	0.5%	Ver 1 coincidencia >
Chile colonial - Wikipedia, la enciclopedia libre https://es.wikipedia.org/w/index.php?title=Chile_colonial_(Chile)	11.4%	Ver 11 coincidencias >
Chile colonial - Wikipedia, la enciclopedia libre https://es.wikipedia.org/w/index.php?title=Chile_colonial_(Chile)	11.4%	Ver 11 coincidencias >
Educación en Chile - Ensayos https://www.clubensayos.com/Historia/Educa/VC14/En-Educ-En-Chile/2019080.html	4.7%	Ver 4 coincidencias >
La Colonia (1569-1810) - AP_EducaChile https://sites.google.com/editecnicapdchile/colonias-1	5.1%	Ver 4 coincidencias >
Chile colonial - Wikipedia, la enciclopedia libre https://es.wikipedia.org/w/index.php?title=Chile_colonial	11.4%	Ver 11 coincidencias >
La Colonia http://profesoresenlinea.cl/lecciones/Colonia/VS/	1.7%	Ver 2 coincidencias >

Coincidencias

Percent reported 4 or more hours per day of sedentary behavior

1. 例如，2011年新《选举法》第3条，限制了选民的投票权。

Lia Coloma (1999-2016) · AP_EDUCArte

518

Documento analizado

la otra esencial. El mejor colegio fundado en el país es el Seminario de la Inmaculada. En 1848 la primera escuela de teología se estableció en la ciudad de Mérida. La otra escuela de teología y las Jornadas y Domingos sagrados recibieron grandes aportes de los católicos extranjeros.

En 1860 se creó la Escuela de Medicina y Cirugía, para la formación de los religiosos. También crearon escuelas para indígenas, debido a la necesidad de convertir a los indios a la fe católica. Los jesuitas fundaron la Universidad de San Ignacio de Loyola, transformándose en universidades polifacultades como la Universidad de San Carlos de Guatemala, la Universidad de San Marcos y el Colegio Mayor de Nuestra Señora del Rosario y el Colegio Mayor de San Ildefonso.

En 1776 se fundó la Universidad de San Carlos de Guatemala, que funcionó hasta 1821. En 1798 se creó la Universidad de San Marcos, que funcionó hasta 1821. En 1821 se creó la Universidad de Santiago y Funciones. Y hacia fines de la colonia se fundó la Academia de San Luis Gonzaga, que funcionó hasta 1821. La Universidad de Santiago de Cali es el primer instituto de enseñanza médica de América Latina. La Universidad de Bogotá fue fundada en 1867. La Universidad Japón jugó un papel importante en la colonización americana, especialmente en la formación de profesionales y técnicos. Los profesionales obligados a imponer la enseñanza europea, contribuyeron la fundación real sobre la iglesia, por el que se organizó la preparación en diverso idioma hacia la formación de profesores, catequistas, curas y hospitalarios.

Fuente

35. **El Primer colegio fundado en el país es el Seminario de La Imperial**, en 1658. La primera escuela de la granadera se creó en 1691 en la ciudad de Santiago. La Universidad de Chile es la más antigua escuela de medicina del país. Los Jesuitas y Dominicos impusieron importante grado académico en las universidades.

36. **Los jesuitas fundaron en Chile las primeras escuelas para jóvenes aristócratas:** el Colegios de San Francisco Juan. Allí hizo sus primeros estudios Alonso de Ovalle y también el abuelo don Juan Ignacio Trasierra. Tras la expulsión de los jesuitas regresaron los dominicos, religiosos y laicas, que habilitaba para los estudios superiores.

37. **Cada orden religiosa manejó estudios para la formación de los sacerdotes y se crearon las seminarios de Santiago y de Concepción.**

38. **REGIMEN ESCOLAR EN LA ÉPOCA COLONIAL.**

39. Diferencias religiosas separan las finalidades de la enseñanza colonial de la nuda república. Seis sacerdotes a, por ejemplo, maestros laicos. Los jesuitas creían que el fin de la educación era la salvación; consideraban el primer objetivo de su labor la de formar el espíritu en temas de Dios y la obediencia a la divina voluntad; de acuerdo y del rey. Los maestros laicos creían que el fin de la educación era hacer crecer las diferencias teológicas que atravesaban los jóvenes más apreciados. Los estudios estaban sujetos a la disciplina escolástica y el latín, y terminando hasta las oposiciones sagradas. Es ese tipo de educación que se dio en las universidades y los seminarios. Tal como España, las colonias carecían de un sistema de educación popular. La mayoría de los jefes consideran un deber suyo la educación de sus hijos. Los padres pagaban los estudios de sus hijos en la universidad, por ejemplo, si en las Escuelas Pías, fundadas por el arzobispo José de Gálvez, en 1617. En las colonias, los establecimientos no

Front-end full

The screenshot shows the Docode WIC interface. On the left, a sidebar menu lists 'Negocios', 'Usuarios', 'Analisis Web', 'Repositorios' (selected), 'PAN', 'Papers', and 'Corpus Interno'. The main content area has a title 'Nuevo Repositorio' and displays a table with three rows:

	Nombre corto	Description	Archivos	Institución
X	PAN	Competencia de detección de plagio	1599	WIC
X	Papers	Papers con texto aleatorio	22	WIC
X	Corpus Interno	Corpus Interno	11	WIC

The screenshot shows the Docode WIC interface. On the left, a sidebar menu lists 'Negocios', 'Usuarios', 'Analisis Web', 'Repositorios' (selected), 'PAN', 'Papers', and 'Corpus Interno'. The main content area has a title 'Corpus Interno' and displays a table with 11 rows of document analysis results:

Número	Nombre	Tamaño
1.	source-document00002.txt	2.465
2.	source-document00001.txt	5.932
3.	source-document00003.txt	1.464
4.	source-document00006.txt	3.838
5.	source-document00007.txt	3.062
6.	source-document00009.txt	2.479
7.	source-document00010.txt	3.740
8.	source-document00028.txt	3.113
9.	source-document00027.txt	2.684
10.	source-document00074.txt	3.311
11.	source-document00017.txt	4.011

Outline

1 Proposed system

2 Architecture

3 Results

4 Conclusions

The core algorithm

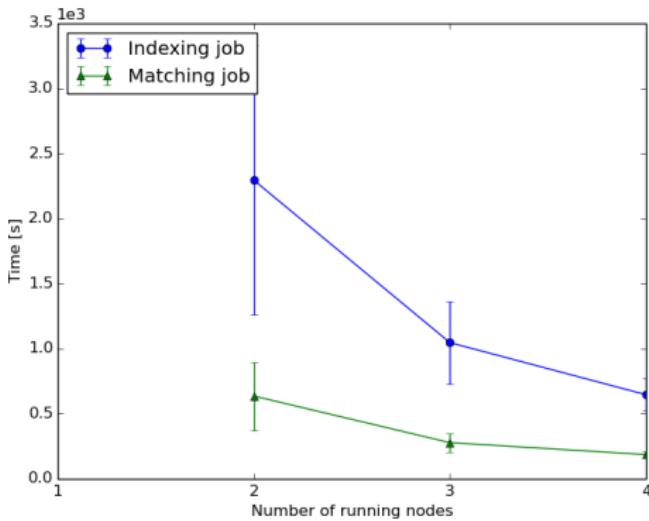
	Precision	Recall	Granularity	Score
Internal corpus	0.99	0.9	1.0	0.95
PAN-PC-13	0.88	0.99	1.04	0.9

Group analysis

- Hadoop cluster, 1 master, 4 slaves, 4GB RAM, 25 GB storage
- 110-MB dataset

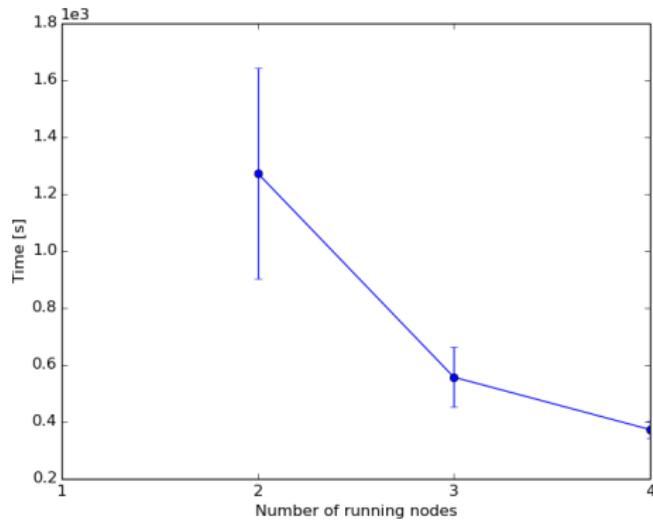
Group analysis

- Hadoop cluster, 1 master, 4 slaves, 4GB RAM, 25 GB storage
- 110-MB dataset



Group analysis

- Hadoop cluster, 1 master, 4 slaves, 4GB RAM, 25 GB storage
- 110-MB dataset



Outline

1 Proposed system

2 Architecture

3 Results

4 Conclusions

Conclusions

- Fully-fledged plagiarism detection engine
- Algorithms
- Deployment at scale

Conclusions

- Fully-fledged plagiarism detection engine
- Algorithms
- Deployment at scale

<http://docode.cl>

Future work

- Text extraction
- User interface
- Search engine

Docode 5: Building a scalable plagiarism detection system

Gaspar Pizarro V.

Web Intelligence Centre

November 22, 2016