

Desarrollo de métodos de corrección automática de errores de transcripción de imágenes a texto en expedientes digitales

Comisión Nacional de Búsqueda /
CIMAT Monterrey



Alumno: Jorge Luis Vargas Barrera

Encargados:

Dra. Mariana Esther Martinez-Sanchez

Dr. Victor Muñiz



Problema

Descripción de la problemática



Se cuenta con métodos basados en OCR para extracción de texto en documentos digitalizados.

Características de documentos:

- Antiguos
- Con diferentes tipografías
- Censura
- Etc.



627
629
633

[redacted] en ese año de mil, novecientos setenta y cuatro fue muy conocida con el mote que aparecía en el documento de [redacted] y si en la librería que se encontraba o se encuentra en la casa número cuatrocientos diez se decía que había propaganda de tipo Marxista Leninista y de tipo guerrillero; aclarando que en el estado de Puebla no tuvo actividad este tipo de grupos que posteriormente se denominaba la [redacted]

PREGUNTA OCHO.- Que diga el declarante en relación con la pregunta anterior por que motivos los documentos y en especial el que se le ha puesto a la vista no se encuentra rubricado. RESPUESTA.- La información no se rubricaba en virtud de que era transmitida vía telefónica y el mecanógrafo de la base en la ciudad de México en turno era quien asentaba los datos que se le proporcionaba, como ya dije, telefónicamente; y por lo que se refiere al número cuatrocientos seis de la misma calle era ocupado por un grupo de estudiantes quienes se dedicaban a dar "asesoría" a campesinos en ese tiempo de la CCI, a empleados de diferentes empresas que exigían indemnizaciones, y efectuaban reuniones con grupos de izquierda de la propia [redacted]

[redacted] para programar manifestaciones y algunos actos públicos, ya que en esos años la UAP (Universidad Autónoma de Puebla) estuvo infiltrada por el entonces partido comunista mexicano.

PREGUNTA NUEVE.- Que diga el declarante si durante el desarrollo de su trabajo se llegó en alguna ocasión a detener a algún miembro de la [redacted]

RESPUESTA.- Detenciones de ninguna naturaleza en virtud de que el cometido era el cien por ciento informativo, cabe aclarar que a pesar de la efervescencia política interna en la [redacted]

no se detectaron a elementos o acciones de grupos subversivos.

PREGUNTA DIEZ.- Que diga el declarante si recuerda quien era el Director de la Dirección Federal de Seguridad en el año de mil novecientos setenta y cuatro. RESPUESTA.- [redacted]

PREGUNTA ONCE.- Que diga el declarante si recuerda que funciones desarrollaba dentro de la extinta Dirección Federal de Seguridad el [redacted]

[redacted] en el año de mil novecientos setenta y cuatro. RESPUESTA.- El era subdirector DE LA Dirección Federal de Seguridad.



pagina1

lqe

ns

"E A iaa >> ese año af mil novecientos setenta y cuatro (33 152 7. fue muy conocida con el mote que aparecía en el documento de 00 GA y +0 la oros quo so dicotraba o se encuentra en la casa

Y "2220 número cuatrocientos diez se decía que había propaganda de tipo Mandsta De Leninista y de tipo guerrillero; aclarando que el Al estado de Puebla no tuvo actividad este tipo de grupos que posteriormente se denominaba dy

PREGUNTA OCHO.- Que diga plidiblarante, en relación con la pregunta anterior por que motivos los documentos y en especial el que se le ha puesto a la vista no se encuentra SA RESPUESTA La información no se

/% "ubricaba en virtud de que a/a tran pi ida vigelefónica y el mecanágrao dela) : x base en la ciuda :0 en tu 1 era gon asentaba los datos que se le

\$ y Poporcionaba, A soto am por lo que se refiere al número o cuatrocientos sois-88 l risma cali era ocupado por yy río de Syudiantes

NO quienes se dedicalan fiar "asesgrfp" a campesingeén ese tiempo dla COI, a » ES empleados de iares eromfpo oxojari indemnizaciones, y Pfectuaban

3| reuniones con grupps de izquirda de J4 propia

3 GU «8proratar manites alfiapés Y algunos getos públicos, ya que eh |, d30s años I32LAP (Unfersidad AMfónoma:de Puebla) estuva infiltrada por el. : cdo 0 0

PREGUNTA,NUEVE.- Quydlgglll deci te si duranté el desarrollo de su trabajo se llegó en alyuña ogasión a detoher a algún/miembro de la vB) puestesta. al 'a naturaleza en virtud de

* ques cometio era al cien gor gfonto inforifiavo cabe aclarar que a pesar de 12% WE ofezyscencia política integnafin la

X 5. 2n8stctctaron a elementos o fbclones de tubos subversivos. ----- e o INTA DIEZ.- Que diga elfoleciarante Sfecuerda quien era el Directos de nas a rección Federal de Sdgurblad en el ah de mil novecientos setenta y

y . cuafito, RESPUESTA. Los a

Y PREGUNTA ONCE. Qyo didk el declkrantél si recuerda que funciones Y - desarrollaba dentro de la extintafDirección [Federdl de Seguricad e, - Y A le año de mf novecieftos sotenta y cuatro. RESPUESHA.-

i4 Él era subdirector DE LA Direc: ign Federal'de Seguridad. \$ -----

ee

ELE

RS

mi 300357

0

a Bruno em todo
 por hora e espero que
 me escreva. Salvo
 do Rio de Janeiro. S. de
 D. J. de J. de J. de J.
 Rio de Janeiro 19-2-
 Col. Rep. de J. de J.
 M. de J. de J. de J.

[illegible]
$$A > - \frac{1}{2} - a \quad QA$$



PROCURADURÍA
GENERAL DE
JUSTICIA



Recursos Humanos de la
Coordinación Administrativa
31 de Octubre de 2004
Culiacán, Sinaloa.
Asunto: Se informa sobre
personal.

C. Alfonso Rodríguez Cardenas
Secretario Particular del C. Procurador General de
Justicia del Estado
Presente.

En atención a lo control de acuerdos No. [redacted] 26 de Mayo de 2004, por este
conducido a [redacted] información enviada por la Dirección de Recursos Humanos de
Gobierno del Estado mediante oficio No. [redacted]
de fecha 13 de Mayo y 21 de Septiembre respectivamente del presente año, la cual
consta de [redacted] laborales del Director de Policía Judicial así como del personal
activo con [redacted] al año de 1977.

Sin más por el momento, aprovecho la oportunidad para enviarle un cordial saludo.

Atentamente

Leopoldo Ibarra Medina
Coordinador Administrativo

Cc: C. Juan Fidel González Mendivil - Procurador General de Justicia del Estado - Presidente
C. [redacted]
LIMAS PABLO
POLICIA JUDICIAL

Mérid, Enrique Sánchez Alonzo 1833 Nte., Desarrollo Plan 3 Rios,
C.R. 88030, Culiacán, Sinaloa, México
Tels. (667) 716-11-20 y 716-11-35 Fax: (667) 713-32-60
E-mail: pagustin@ges.gob.mx

Sinaloa
Que nunca dorme



page-367.jpg.txt

pagina1

"Cómo armó la CNDH el rompecabezas", en Proceso No. 1309, 2 de
diciembre de 2001, pp. 34-37.

MORENO BORBOLLA, José Luis. Entrevista a Jaime García Chávez (miembro de la red
urbana del Grupo Popular Guerrillero "Arturo Gámiz" comandado por Oscar González
Equiarte), en Para Romper el Silencio Expediente Abierto, Centro de Investigaciones
Históricas de los Movimientos Armados, A.C., noviembre 1994-enero 1995, pp. 29-37.

NAVARRO ALTAMIRANO, Juan. "Guadalajara: de Maciel a Parrés. La nueva doctrina
universitaria", en Por qué? Revista independiente No. 199, abril 20 de 1972. pp. 20-21.

"Guadalajara ¡A seguir la lucha!", en Por qué? Revista independiente
No. 218, agosto 31 de 1972. pp.12-14.

"Guadalajara. La tarea: fijar a traidores y oportunistas", en Por qué?
Revista independiente No. 183, diciembre 30 de 1971. pp. 25-27.

ORTEGA JUÁREZ, Joel. "El 68: punto final a las fantasías, en Campus Milenio No. 377,
jueves 22 de julio de 2010, pp. 8-9.

PONCE DE LEÓN, Alberto. "El hombre de los 200 asesinatos", en Proceso No. 1361, 1º de
diciembre de 2002, pp. 12.

Por esto! No. 68, 14 de octubre de 1982, "Doloroso alto en el camino ¡Al rescate de Por esto!
Punto Crítico No. 4, abril de 1972, "Chihuahua: la verdad sobre los asaltos" pp. 19-21.

Punto Crítico No. 10, Octubre de 1972, "Tribunal Popular. La Sentencia" pp. 32-33

Punto Crítico No. 69, 31 de Enero de 1977, "Denuncian represión en Chihuahua", pp. 38-40.

Punto Crítico, Revista de Información y Análisis Político, No. 50, primera quincena de abril
de 1976, "Carta Abierta al Pueblo de México de los padres de los Gámiz", p. 4.

Punto Crítico No. 28, mayo de 1974, "Comunicado de prensa, presos políticos de la crujía
"M", p.3,

Punto Crítico, No. 19, agosto de 1973, p. 23.

Punto Crítico, No. 52, primera quincena de mayo de 1976, "Guerrilleros ajustan cuentas con
su pasado", pp. 15-18

Punto Crítico No. 32, Julio de 1975, "En C.U. la Ley del Revolver", p. 5.

Punto Crítico, No. 91, Octubre de 1978, "La amnistia una caricatura", pp.6 -7.

Punto Crítico No. 67, 15 de Noviembre de 1976, "Naucalpan, ¿Estado de México o de sitio?
pp. 14-15.

Punto Crítico, No. 20/21, septiembre-octubre de 1973, "Los Presos Políticos ante los últimos
secuestros. No acepto dicha liberación. Al pueblo de México, pp. 1-2.

Punto Crítico, No. 16, abril de 1973, "La Federación Estudiantil Oaxaqueña denuncia el
cacicazgo de Jamiltepec, p. 1.

RAMÍREZ LADEWIG, Carlos. "Indignidad de los diputados", en Proceso, No. 14, 5 de
febrero de 1977, p. 8.

365

- La calidad de los textos extraídos es deficiente en ocasiones.
- Se cuenta con un método para evaluar la calidad de la extracción.



Objetivo



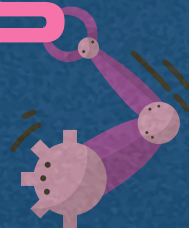
Se pretende desarrollar una serie de métodos para la corrección automática de errores ortográficos/gramaticales.

Basados en:

- Transformaciones morfológicas de palabras basadas en vocabularios.
- Métodos probabilísticos
- Arquitecturas de aprendizaje profundo (redes recurrentes, convolucionales y transformers)



Conjunto de datos



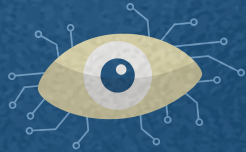


Conjunto de datos

Se tomaron textos con buena calidad de visualización. Los cuales consisten de tesis que

Método 1

método probabilístico



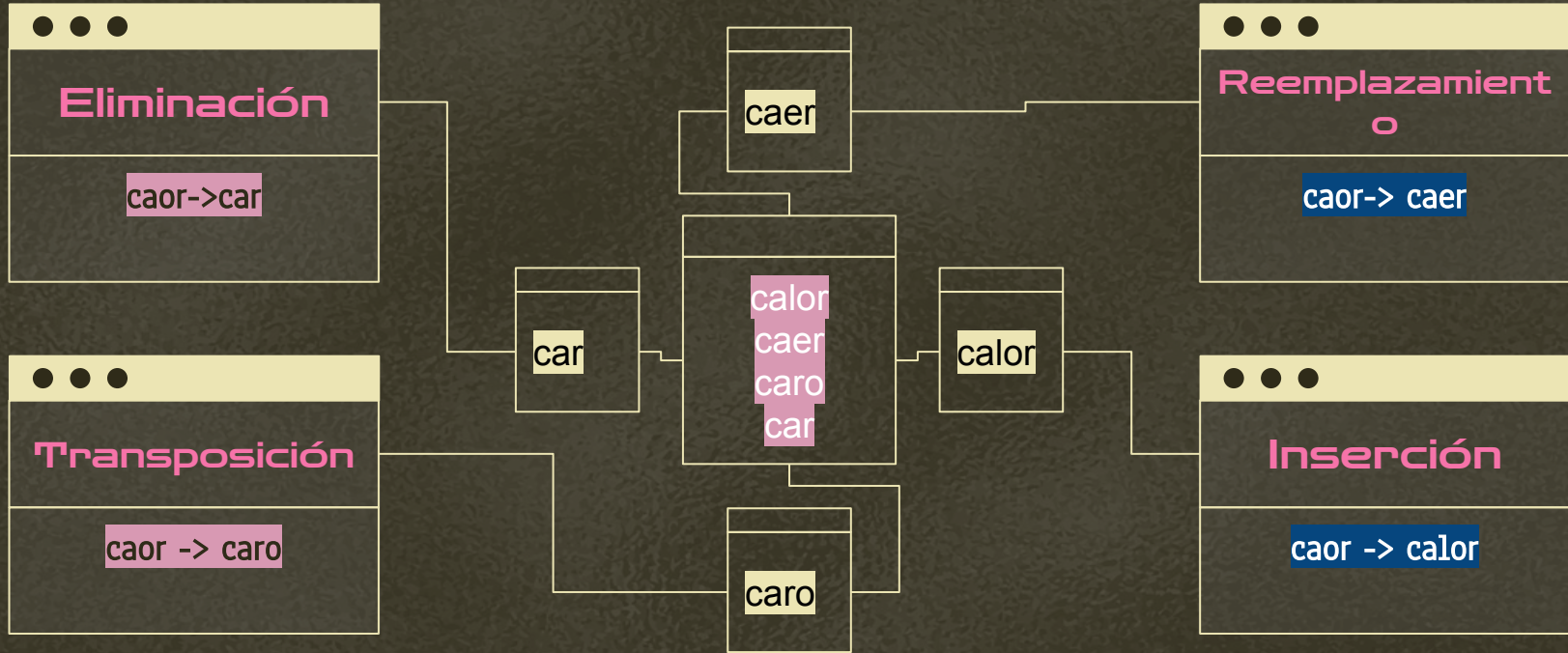
Se utiliza el principio de máxima verosimilitud. Dada una palabra incorrecta s , se busca el mejor candidato w_b que ~~co~~orija de una lista de posibles candidatos $w_i \in C(s)$ con la mayor probabilidad (sin normalizar)

$$w_b = \arg \max_{w_i \in C(s)} P(s|w_i)P(w_i)$$

Componentes

1. Diccionario: Se detectan los errores de escritura y propone candidatos
2. Modelo de error: $P(s|w_i)$ la probabilidad de que s sea escrita en un texto cuando el autor se refería a w_i
3. Modelo de lenguaje: $P(w_i)$ la probabilidad de que w_i aparezca como una palabra en el texto

Generación de candidatos





ORACION ORIGINAL

gobernadores, todos del mismo paktid0, los pkesidente8

ORACION CORREGIDA

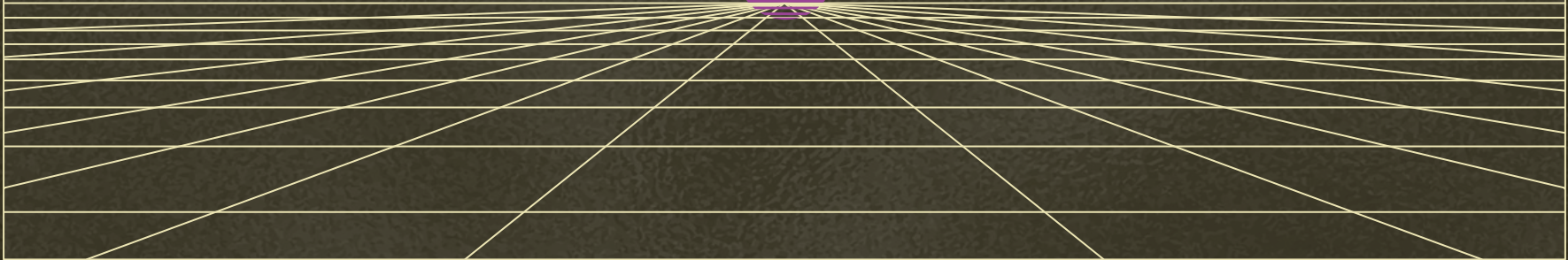
gobernadores todos del mismo partido, los presidente

ORACION ORIGINAL

pke8idente utilizaba. Bn otro nivel de la pikámide estaban lus

ORACION CORREGIDA

presidente utilizaban en otro nivel de la pirámide estaban los



Resultados

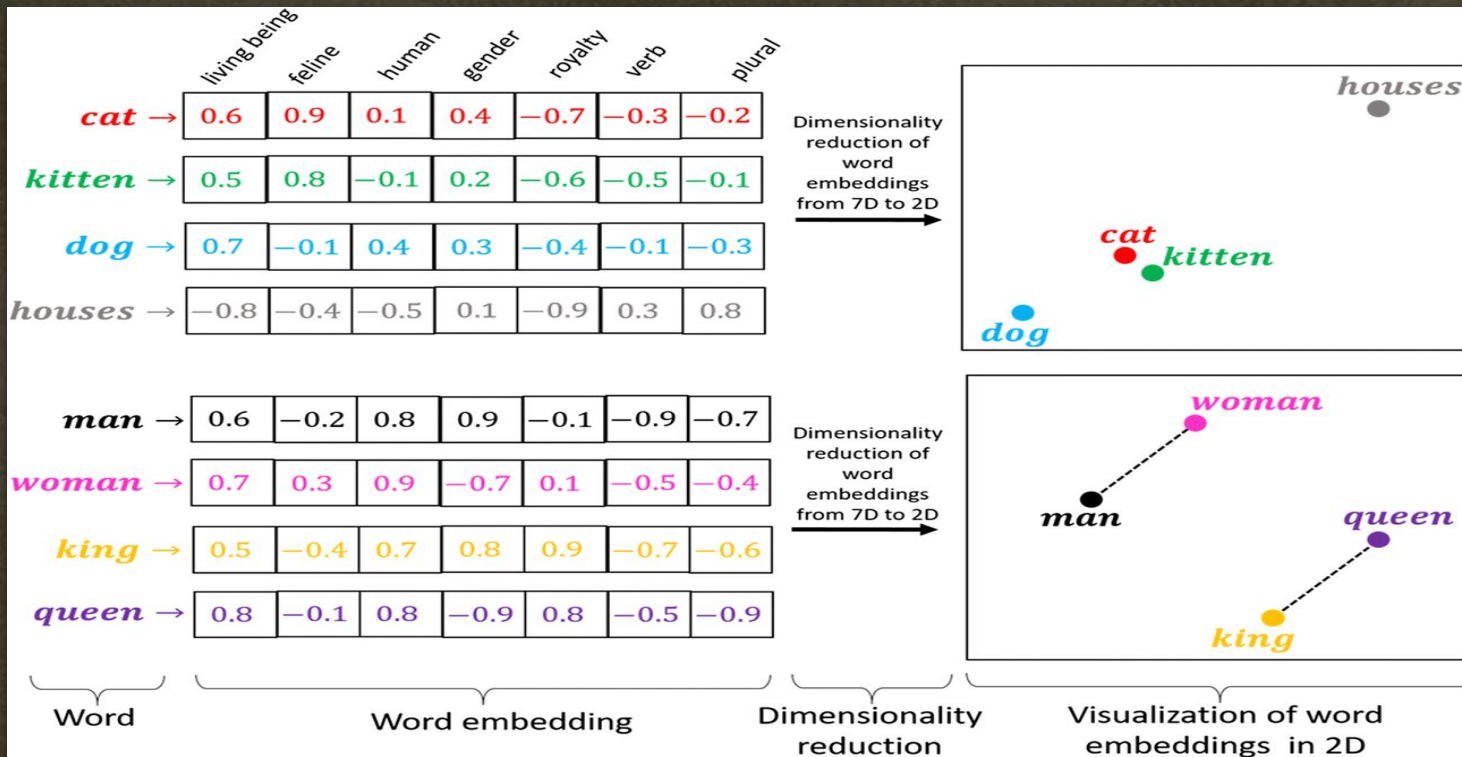
000	
Nivel de ruido Θ_{CR}	Blue score
low	0.5300
medium	0.4355
high	0.3591

Método 2

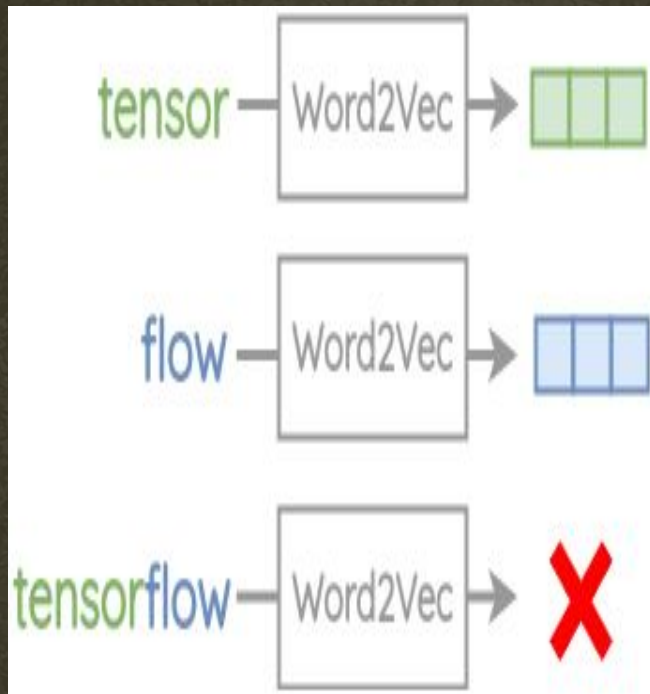
método basado en word embeddings



Embeddings



Problemas wor2vec



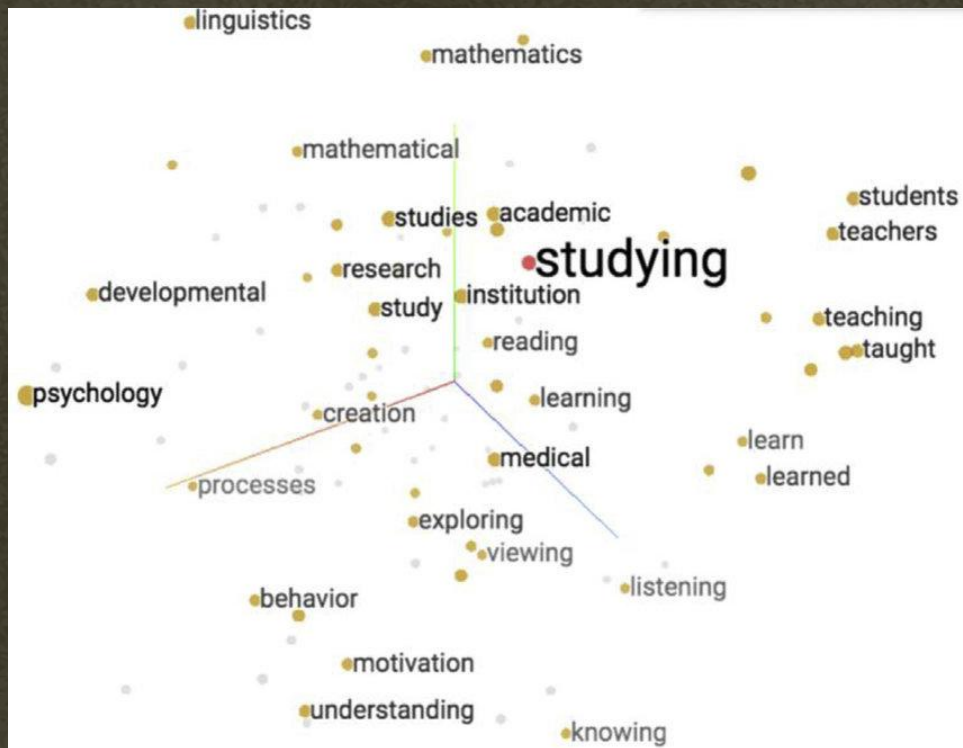
Shared radical

eat eats eaten eater eating

Fasttext

3-grams <eating>
<ea eat ati tin ing ng>







ORACION ORIGINAL

gobernadores, todos del mismo paktid0, los pkesidente8

ORACION CORREGIDA

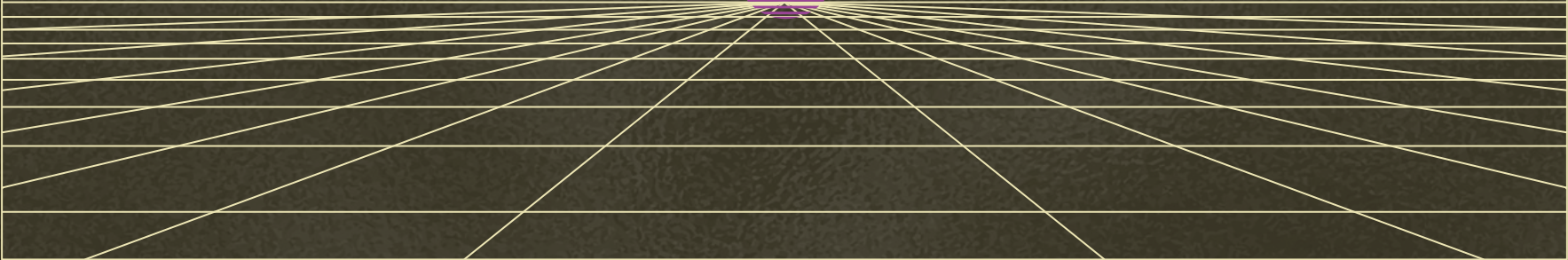
gobernadores todos del mismo paktiyā los presidente de

ORACION ORIGINAL

pke8idente utilizaba. Bn otro nivel de la pikámide estaban lus

ORACION CORREGIDA

idente utilizaba plutonio otro nivel de la pirámide estaban



Resultados con Fasttext

000	
Nivel de ruido OCR	Blue score
low	0.5025
medium	0.4162
high	0.3034

Método 3

método basado en transformers

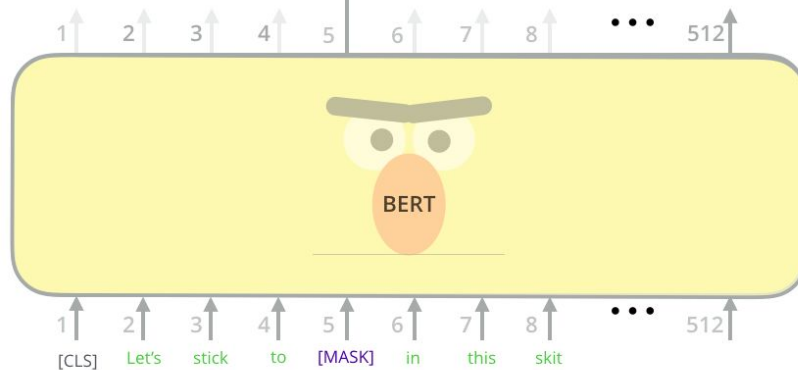


Use the output of the masked word's position to predict the masked word

Possible classes:
All English words

0.1%	Aardvark
...	...
10%	Improvisation
...	...
0%	Zyzzzyva

FFNN + Softmax



Randomly mask
15% of tokens

Input



ORACION ORIGINAL

setenta han pu6licadu obras testimoniales en on e8foekzo por recuperar la memoria

ORACION CORREGIDA

setenta han publicado obras testimoniales en un desarrollo por recuperar la memoria

ORACION ORIGINAL

hist6rica de aquellos a6os. Sin em6ar90, esos testimonios abarcan en 8u mayoría a la

ORACION CORREGIDA

hist6rica de aquellos a6os . sin embargo , esos testimonios abarca en su mayoría a la



Resultados con BET θ

000	
Nivel de ruido θ CR	Blue score
low	0.67899
medium	0.6421
high	0.5092

Siguientes pasos

Vocabulario

Formar un vocabulario más
grande

Preprocesamiento

Mejorar el preprocesamiento
del texto

Neuspell

Hacer funcionar neuspell
para nuestro propósito



Gracias

