



UNIVERSIDAD DE BURGOS
ESCUELA POLITÉCNICA SUPERIOR
Grado en Ingeniería Informática



**TFG del Grado en Ingeniería
Informática**

JIZT

**Generación de resúmenes abstractivos en
la nube mediante Inteligencia Artificial**



Presentado por Diego Miguel Lozano
en la Universidad de Burgos — 1 de febrero de 2021
Tutores: Dr. Carlos López Nozal y
Dr. José Francisco Díez Pastor



UNIVERSIDAD DE BURGOS
ESCUELA POLITÉCNICA SUPERIOR
Grado en Ingeniería Informática



D. nombre tutor, profesor del departamento de nombre departamento, área de nombre área.

Expone:

Que el alumno D. Diego Miguel Lozano, con DNI 71307413-F, ha realizado el Trabajo final de Grado en Ingeniería Informática titulado “JIZT - Generación de resúmenes abstractivos en la nube mediante Inteligencia Artificial.

Y que dicho trabajo ha sido realizado por el alumno bajo la dirección del que suscribe, en virtud de lo cual se autoriza su presentación y defensa.

En Burgos, 1 de febrero de 2021

Vº. Bº. del Tutor:

Vº. Bº. del Tutor:

D. Carlos López Nozal

D. José Francisco Díez Pastor

Resumen

En este primer apartado se hace una **breve** presentación del tema que se aborda en el proyecto.

Descriptores

Palabras separadas por comas que identifiquen el contenido del proyecto Ej: servidor web, buscador de vuelos, android ...

Abstract

A **brief** presentation of the topic addressed in the project.

Keywords

keywords separated by commas.

Índice general

Índice general	III
Índice de figuras	V
Índice de tablas	VI
Introducción	1
Objetivos del proyecto	3
2.1. Objetivos generales	3
2.2. Objetivos técnicos	3
Conceptos teóricos	7
3.1. Pre-procesado del texto	7
3.2. Codificación del texto	10
3.3. Generación del resumen	13
3.4. Post-procesado del texto	16
Técnicas y herramientas	19
4.1. Flask y Flask-RESTful	19
4.2. Docker	19
4.3. Kubernetes	19
Aspectos relevantes del desarrollo del proyecto	21
Trabajos relacionados	23
6.1. Artículos académicos	23
6.2. Proyectos similares	23

Conclusiones y Líneas de trabajo futuras	27
Bibliografía	29

Índice de figuras

3.1. Ejemplo de <i>tokenización</i> con el modelo T5.	11
3.2. Pasaje del libro <i>A Wrinkle in Time</i> . El <i>tóken</i> EOS se ha marcado en rojo.	12
3.3. Proceso de generación de resúmenes, ilustrado con un fragmento del libro <i>The Catcher in the Rye</i>	13
3.4. El <i>framework</i> texto a texto permite emplear el mismo modelo, con los mismos hiperparámetros, función de pérdida, etc., para aplicarlo a diversas tareas de NLP [16].	15

Índice de tablas

6.1. Comparativa de las características ofrecidas por las diferentes alternativas para la generación de resúmenes.	25
---	----

Introducción

El término Inteligencia Artificial (IA) fue acuñado por primera vez en la Conferencia de Dartmouth [1] hace ahora 65 años, esto es, en 1956. Sin embargo, ha sido en los últimos tiempos cuando su presencia e importancia en la sociedad han crecido de manera exponencial.

Uno de los campos históricos dentro de la AI, es el Procesamiento del Lenguaje Natural (NLP, por sus siglas en inglés), cuya significación se hizo patente con la aparición del célebre Test de Turing [2], en el cual un interrogador debe discernir entre un humano y una máquina conversando con ambos por escrito a través de una terminal.

Hasta los años 80, la mayor parte de los sistemas de NLP estaban basados en complejas reglas escritas a mano [3], las cuales conseguían generalmente modelos muy lentos, poco flexibles y con baja precisión. A partir de esta década, como fruto de los avances en Aprendizaje Automático (*Machine Learning*), fueron apareciendo modelos estadísticos, consiguiendo notables avances en campos como el de la traducción automática.

En la última década, el desarrollo ha sido aún mayor debido a factores como el aumento masivo de datos de entrenamiento (principalmente provenientes del contenido generado en la *web*), avances en la capacidad de computación (GPU, TPU, ASIC...) y el progreso dentro del área de la Algoritmia [4].

No obstante, ha sido desde la aparición del concepto de “atención” en 2015 [5, 6, 7] cuando el campo del NLP ha comenzado a lograr resultados cuanto menos sorprendentes [8, 9].

Con todo, la mayor parte de estos avances se han visto limitados al ámbito académico y empresarial. Los modelos cuyo código ha sido publicado, o bien no están entrenados, o bien requieren para ser usados conocimientos

avanzados de matemáticas o programación, o simplemente son demasiado grandes para ser ejecutados en ordenadores convencionales.

Con esta idea en mente, el objetivo de JIZT se centra en acercar los modelos NLP estado del arte tanto a usuarios expertos, como no expertos.

Para ello, JIZT proporciona:

- Una API REST destinada a los usuarios con conocimientos técnicos, a través de la cual se pueden llevar a cabo tareas de NLP.
- Una aplicación multiplataforma que consume dicha API, y que proporciona una interfaz gráfica sencilla e intuitiva. Esta aplicación puede ser utilizada por el público general, aunque no deja de ofrecer opciones avanzadas para aquellos usuarios con mayores conocimientos en la materia.

En un principio, dado el alcance de un Trabajo de Final de Grado, la única tarea de NLP implementada ha sido la de generación de resúmenes. La motivación para esta decisión se ha fundamentado principalmente en la relativa menor popularidad de esta área frente a otras como la traducción automática, el análisis de sentimientos, o los modelos conversacionales. Para estas últimas tareas existe actualmente una amplia oferta de grandes compañías como Google [10], IBM [11], Amazon [12], o Microsoft [13], entre muchas otras. Nuestra mayor limitación reside en que los modelos pre-entrenados que utilizaremos para la generación de los resúmenes funcionan únicamente en inglés. Esperamos que en un futuro próximo aparezcan modelos que admitan otros idiomas.

En un mundo en el que en cinco años se producirán globalmente 463 exabytes de información al día [14], siendo mucha de esa información textual, la generación de resúmenes aliviara en cierto modo el tratamiento de esos datos.

Sin embargo, gran parte del esfuerzo de desarrollo de JIZT se ha centrado en el diseño de su arquitectura, la cual se describirá con detalle en el capítulo de **Conceptos Teóricos**. Por ahora adelantaremos que ha sido concebida con el objetivo de ofrecer la mayor escalabilidad y flexibilidad posible, manteniendo además la capacidad de poder añadir otras tareas de NLP diferentes de la generación de resúmenes en un futuro cercano.

Por todo ello, el presente TFG conforma el punto de partida de un proyecto ambicioso, desafiante, pero con la certeza de que, independientemente de su recorrido, habremos aprendido, disfrutado, y ojalá ayudado a alguien por el camino.

Objetivos del proyecto

2.1. Objetivos generales

- Ofrecer la capacidad de llevar a cabo tareas de NLP tanto al público general, como al especializado. Como se ha mencionado con anterioridad, la única tarea NLP que implementará el presente TFG será la de generación de resúmenes.
- Emplear modelos pre-entrenados estado del arte para la generación de resúmenes abstractivos. Los resúmenes abstractivos se diferencian de los extractivos en que el resumen generado contiene palabras o expresiones que no aparecen en el texto original [15]. Dicho de forma más técnica, existe cierto nivel de paráfrasis.
- Diseñar una arquitectura con aspectos como la flexibilidad, la escalabilidad y la alta disponibilidad como principios fundamentales.
- Poner en práctica lo aprendido a lo largo de la carrera en áreas como Ingeniería del Software, Sistemas Distribuidos, Programación, Minería de Datos, Algoritmia y Bases de Datos.
- Ofrecer la totalidad del proyecto bajo licencias de *Software* Libre.

2.2. Objetivos técnicos

- Los modelos pre-entrenados de generación de texto admiten parámetros específicos para configurar dicha generación, por lo que se deberá

implementar una interfaz que permita a los usuarios establecer dichos parámetros de manera opcional. Por defecto, se proporcionarán los valores que mejores resultados ofrecen, extraídos mayoritariamente de manera experimental.

- Los modelos pre-entrenados de generación estado del arte presentan frecuentemente limitación en la longitud de los textos de entrada que reciben, derivada de la longitud de las secuencias de entrada con las que han sido entrenados. Esta longitud llega a ser tan baja como 512 *tókenes*¹ [16]. Por tanto, se deberá establecer algún mecanismo que permita sortear esta limitación para poder generar resúmenes de textos arbitrariamente largos.
- Gestionar el pre-procesado de los textos a resumir para ajustarlos a la entrada que los modelos pre-entrenados esperan.
- Algunos modelos pre-entrenados generan textos enteramente en minúsculas. Se deberá, por tanto, incluir mecanismos en la etapa de post-procesado que permitan recomponer el correcto uso de las mayúsculas en los resúmenes generados.
- Con el fin de cumplir con el objetivo general referente a la arquitectura, desarrollar una arquitectura de microservicios, basada en la filosofía *Cloud Native* [17, 18]. Este objetivo se divide a su vez en dos puntos:
 - Encapsular cada microservicio en un contenedor Docker.
 - Implementar la orquestación y balanceo de los microservicios a través de Kubernetes.
- Complementariamente al punto anterior, implementar una arquitectura dirigida por eventos [19]. La motivación detrás de la utilización de este patrón arquitectónico se justifica en el capítulo de **Conceptos Teóricos**.
- Implementar una API REST escrita en Python empleando el *framework web* Flask. Dicha API será el punto de conexión con el servicio de generación de resúmenes en la nube.
- Desplegar PostgreSQL como servicio en Kubernetes mediante el Operador PostgreSQL de Crunchy [20]. Esta base de datos cumplirá la doble función de (a) servir como caché para no volver a producir resúmenes ya generados con anterioridad, incrementando la velocidad

¹ Este término se definirá posteriormente. Por ahora, el lector puede considerar que un *tóken* es equivalente a una palabra.

de respuesta, y (b) almacenar los resúmenes generados con fines de evaluación de la calidad de los mismos y extracción de métricas.

- Desarrollar, con ayuda de Flutter, una aplicación multiplataforma con soporte nativo para Android, iOS, y *web*. Esta aplicación consumirá la API y proporcionará una interfaz gráfica sencilla e intuitiva para que usuarios regulares puedan hacer uso del servicio de generación de resúmenes.
- Seguir el patrón de diseño Clean Architecture [21] y de *offline-first* [22] para la implementación de la aplicación.

Conceptos teóricos

En este capítulo, detallaremos de forma teórica el proceso de generación de resúmenes, desde el momento que recibimos el texto a resumir, hasta que se le entrega al usuario el resumen generado. En el [siguiente capítulo](#), explicaremos las herramientas que hacen posible que todo este proceso se pueda llevar a cabo de forma distribuida «en la nube».

La generación de resúmenes se divide en cuatro etapas fundamentales:

- Pre-procesado.
- Codificación.
- Generación del resumen.
- Post-procesado.

Veamos en detalle en qué consiste cada una de ellas.

3.1. Pre-procesado del texto

El principal objetivo de esta etapa es adecuar el texto de entrada para que se aproxime lo máximo posible a lo que el modelo espera. Adicionalmente, se separa el texto de entrada en frases. Esta separación parecer *a priori* una tarea trivial, pero involucra una serie de dificultades que se detallarán a continuación.

Cabe destacar que, como mencionábamos en la [Introducción](#), los modelos pre-entrenados de los que hacemos uso solo admiten textos en inglés, por

lo que algunas de las consideraciones que tomamos en el pre-procesado del texto solo son aplicables a este idioma.

A grandes rasgos, en la etapa de pre-procesado se divide a su vez en los siguientes pasos:

- Eliminar retornos de carro, tabuladores (`\n`, `\t`) y espacios sobrantes entre palabras (p. ej. "I am" \rightarrow "I am").
- Añadir un espacio al inicio de las frases intermedias (p. ej.: "How's it going?Great!" \rightarrow "How's it going? Great!"). Esto es especialmente relevante en el caso de algunos modelos, como por ejemplo BART [23], los cuales tienen en cuenta ese espacio inicial para distinguir entre frases iniciales y frases intermedias en la generación de resúmenes².
- Establecer un mecanismo que permita llevar a cabo la ya mencionada separación del texto en frases. Esto es importante dado que los modelos tienen un tamaño de entrada máximo. Dos estrategias comunes para eludir esta limitación consisten en (a) truncar el texto de entrada, lo cual puede llevar asociado pérdidas notables de información, o (b) dividir el texto en fragmentos de menor tamaño. En nuestro caso, la primera opción quedó rápidamente descartada ya que los textos que vamos a recibir, por lo general, superarán el tamaño máximo (en caso contrario tendría poco sentido querer generar un resumen). Refiriéndonos, por tanto, a la segunda opción, es frecuente llevar a cabo dicha separación de manera ingenua, únicamente atendiendo al tamaño de entrada máximo. Sin embargo, en nuestro caso decidimos refinar este proceso e implementamos un algoritmo original³ en el que dicha separación se realiza de tal modo que ninguna frase queda dividida. Para garantizar el éxito de este algoritmo, es fundamental que las frases estén correctamente divididas; el porqué se clarificará en la [siguiente sección](#), referente a la codificación del texto.

A continuación, nos centraremos en el proceso de división del texto en frases. A la hora de llevar a cabo este proceso, debemos tener en cuenta que el texto de entrada podría contener errores ortográficos o gramaticales,

² Por el momento, no hacemos uso de este modelo, aunque podría incluirse en el futuro.

³ Utilizamos el término «original» porque no encontramos ningún recurso en el que se tratara este problema, por lo que tuvimos que resolverlo sin apoyos bibliográficos. Esto no quiere decir, sin embargo, que no se hayan implementado estrategias similares en otros problemas diferentes al aquí expuesto.

por lo que debemos tratar de realizar el mínimo número de suposiciones posibles.

No obstante, la siguiente consideración se nos hace necesaria: el punto (.) indica el final de una frase solo si la siguiente palabra empieza con una letra *y* además mayúscula.

Por ejemplo, en el caso de: "Your idea is interesting. However, I would [...]." se separaría en dos frases, dado que la palabra posterior al punto empieza con una letra mayúscula. Sin embargo: "We already mentioned in Section 1.1 that this example shows [...]." conformaría una única frase, ya que tras el punto no aparece una letra. Procedemos de igual modo en el caso de los signos de interrogación (?) y de exclamación (!). Por ejemplo: "She asked 'How's it going?', and I said 'Great!'." se tomará correctamente como una sola frase; tras la interrogación, la siguiente palabra comienza con una letra *minúscula*.

Con la suposición anterior, también se agruparían correctamente los puntos suspensivos.

Sin embargo, fallaría en situaciones como: "NLP (i.e. Natural Language Processing) is a subfield of Linguistics, Computer Science, and Artificial Intelligence.", en la que la división sería: "NLP (i.e." por un lado, y "Natural Language Processing) is a subfield [...].", por otro, ya que "Natural" empieza con mayúscula y aparece tras un punto.

Asimismo, la razón principal por la que no podemos apoyarnos únicamente en reglas predefinidas, reside en las llamadas Entidades Nombradas (*Named Entities*, en inglés), esto es, palabras que hacen referencia a personas, lugares, instituciones, empresas, etc. Existe toda una disciplina dedicada la identificación de este tipo de palabras, conocida como Reconocimiento de Entidades Nombradas (NER, por sus siglas en inglés), y pese a los buenos resultados conseguidos por algunos de los modelos propuestos, se considera un problema lejos de estar resuelto [24].

En nuestro caso emplearemos un modelo pre-entrenado para solucionar, al menos en parte, el problema de las Entidades Nombradas. Este modelo también solventa situaciones como la descrita anteriormente, en las que las reglas escritas a mano se quedan cortas. En el capítulo de **Técnicas y Herramientas**, hablaremos de dicho modelo y de la implementación concreta en código de los procedimientos expuestos anteriormente.

3.2. Codificación del texto

En esta etapa, se lleva a cabo lo que se conoce en inglés como *word embedding*⁴. Los modelos de IA trabajan, por lo general, con representaciones numéricas. Por ello, las técnicas de *word embedding* se centran en vincular texto (bien sea palabras, frases, etc.), con vectores de números reales [25]. Esto hace posible aplicar a la generación de texto arquitecturas comunes dentro de la IA (y especialmente, del *Deep Learning*), como por ejemplo las Redes Neuronales Convolucionales (CNN) [26].

Esta idea, conceptualmente sencilla, encierra una gran complejidad, dado que los vectores generados deben retener la máxima información posible del texto original, incluyendo aspectos semánticos y gramaticales. Por poner un ejemplo, los vectores correspondientes a las palabras «profesor» y «alumno», deben preservar cierta relación entre ambos, y a su vez con la palabra «educación» o «escuela». Además, su vínculo con las palabras «enseñar» o «aprender» será ligeramente distinto, dado que en este caso se trata de una categoría gramatical diferente (verbos, en vez de sustantivos). A través de este ejemplo, podemos comprender que se trata de un proceso complejo.

Dado que los modelos pre-entrenados se encargan de realizar esta codificación por nosotros, no entraremos en más detalle en los algoritmos concretos empleados, dado que consideramos que se sale del alcance de este trabajo⁵.

Lo que sí que hemos tenido que implementar en esta etapa, ha sido la división del texto en fragmentos a fin de no superar el tamaño máximo de entrada del modelo.

De este modo, podremos realizar resúmenes de textos arbitrariamente largos, siguiendo los siguientes pasos:

1. Dividimos el texto en fragmentos.
2. Generamos un resumen de cada fragmento.
3. Concatenamos los resúmenes generados.

Anteriormente, habíamos mencionado el término *token*. Este concepto se puede traducir al español como «símbolo». En nuestro caso concreto, un *token* es el vector numérico asociado a una palabra al realizar la codificación.

⁴ En el presente documento, hemos traducido este término como «codificación del texto».

⁵ En cualquier caso, el lector curioso puede explorar los algoritmos más populares de codificación, los cuales, ordenados cronológicamente, son: word2vec [27, 28], GloVe [29], y más recientemente, ELMo [30] y BERT [31].

Más concretamente, en modelos más actuales, como el modelo T5 [16], los *tókenes* pueden referirse a palabras completas o a *fragmentos* de las mismas.

Por lo general, las palabras que aparecen en el vocabulario con el que ha sido entrenado el modelo van a generar un único *token*. Sin embargo, las palabras desconocidas, se descompondrán en varios *tókenes*. Lo mismo sucede con palabras compuestas o formadas a partir de prefijación o sufijación. En la **siguiente figura**, podemos ver un ejemplo de ello:

Palabra simple:	lucky	→	[5722]	Un <i>token</i>
Palabra compuesta:	backbone	→	[223, 12269]	Varios <i>tókenes</i>
Palabra con prefijo:	luckily	→	[3, 31299]	Varios <i>tókenes</i>
Palabra no reconocida:	JIZT	→	[446, 20091, 382]	Varios <i>tókenes</i>

Figura 3.1: Ejemplo de *tokenización* con el modelo T5.

En el anterior ejemplo, si decodificamos los *tókenes* correspondientes a la palabra "backbone", esto es, [223, 12269], obtenemos los fragmentos "back", y "bone", respectivamente.

La idea detrás de esta fragmentación se basa en la composición, uno de los mecanismos morfológicos de formación de palabras más frecuentes [32] en muchos idiomas, como el inglés, español o alemán. Por tanto, presupone que dividiendo las palabras desconocidas en fragmentos menores, podemos facilitar la comprensión de las mismas. Naturalmente, habrá casos en los que esta idea falle; por ejemplo, en la figura anterior, la palabra "JIZT" se descompone en "J", "IZ", "T", lo cual no parece hacerla mucho más comprensible.

Una vez explicado el concepto de *token*, volvamos al problema ya mencionado con anterioridad: los modelos de generación de texto admiten un tamaño de entrada máximo, determinado en función del número de *tókenes*. Debido a que la unidad de medida es el número de *tókenes*, y no el número de palabras, o de caracteres, debemos tener en cuenta algunos detalles, entre ellos el hecho de que los modelos generan *tókenes* especiales para marcar el inicio y/o el final de la secuencia de entrada.

El modelo T5 (el cual como mencionábamos anteriormente, es el único modelo que utilizamos por ahora), genera un único *token* de finalización de

secuencia (EOS, *end-of-sequence*), que se coloca siempre al final del texto de entrada, una vez codificado, y en el caso de este modelo siempre tiene el *id* 1. En la siguiente figura podemos ver un ejemplo con un texto de entrada:

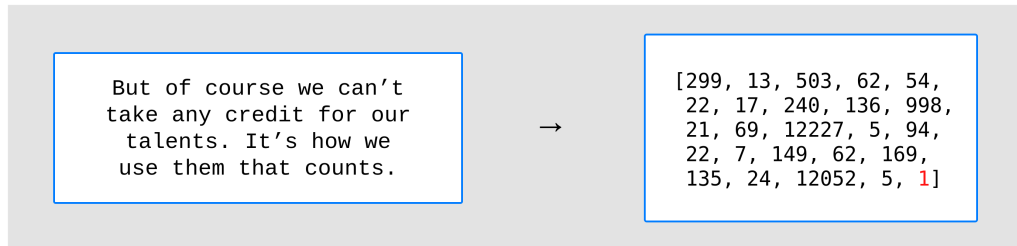


Figura 3.2: Pasaje del libro *A Wrinkle in Time*. El *tóken* EOS se ha marcado en rojo.

Como podemos ver, el *tóken* EOS aparece una única vez por cada texto de entrada, y es independiente de las palabras o frases que este contiene.

Otro aspecto a tener en cuenta, reside en que este modelo no solo es capaz de generar resúmenes, si no que puede ser empleado para otras tareas como la traducción, respuesta de preguntas [16], etc. Para indicarle cuál de estas es la tarea que queremos que desempeñe, curiosamente se lo tenemos que indicar tal y cómo lo haríamos en la vida real; en nuestro caso, simplemente precedemos el texto a resumir con la orden «*resume*» («*summarize*»). Por poner otro ejemplo, si quisiéramos resumir de alemán a español, le señalaríamos: «*traduce de alemán a español*» seguido de nuestro texto («*summarize German to Spanish*»).

Por consiguiente, este prefijo deberá aparecer al principio de cada una de las subdivisiones generadas y, del mismo modo, lo tenemos que tener en cuenta a la hora de calcular el número de *tókenes* de las mismas.

Con las anteriores consideraciones en mente, el objetivo principal será llevar a cabo la división del texto de entrada de forma que el número de *tókenes* varíe lo mínimo posible entre las diferentes subdivisiones, y todo ello sin partir ninguna frase.

Esta es una tarea más compleja de lo que puede parecer. En el capítulo de **Técnicas y Herramientas** se propone un algoritmo que emplea una estrategia voraz para llevar a cabo una primera división del texto; posteriormente procede al *balanceo* de las subdivisiones generadas en el paso anterior, de forma que el número de *tókenes* en cada subdivisión sea lo más parecido posible. Y esto, evidentemente, sin superar el máximo tamaño de entrada del modelo en ninguna de las subdivisiones.

3.3. Generación del resumen

Una vez codificado y dividido el texto apropiadamente, generamos los resúmenes parciales, y posteriormente los unimos, dando lugar a un único resumen del texto completo.

Gráficamente, los pasos dados tanto en la anterior etapa, la codificación y división del texto, como en esta, la generación del resumen, son:

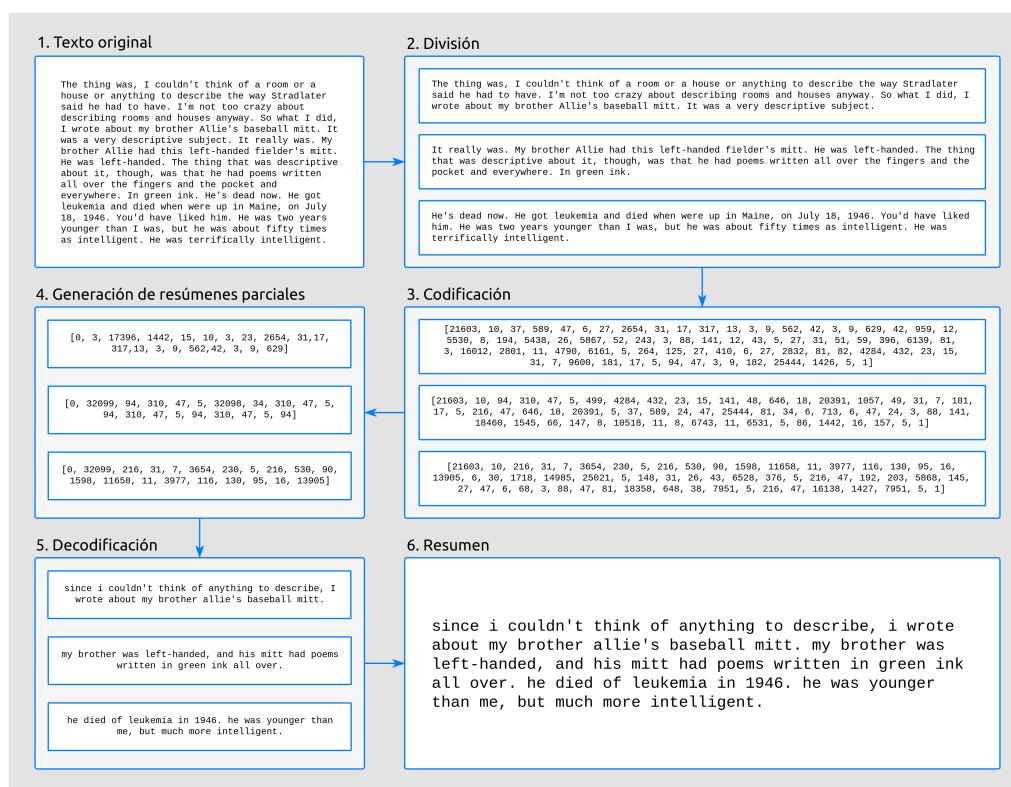


Figura 3.3: Proceso de generación de resúmenes, ilustrado con un fragmento del libro *The Catcher in the Rye*.

Como podemos apreciar en la anterior figura, el modelo generador de resúmenes toma el texto codificado, y devuelve una versión reducida del mismo, también codificado. Por ello, antes de poder unir y devolver el resumen generado, debemos realizar un paso de *decodificación*, que lleva a cabo el proceso contrario a la *codificación*, como veíamos en la [anterior sección](#). Algo con lo que tendremos que lidiar en la siguiente etapa, el post-procesado, será corregir el resumen generado para que se ajuste a las reglas ortográficas vigentes, en especial en lo relativo al uso de mayúsculas.

La ventaja de utilizar modelos pre-entrenados es clara: estos modelos son para nosotros cajas negras, a las que solo tenemos que encargarnos de proporcionarles la entrada en el formato concreto que esperan.

Cabe destacar que, el hecho de realizar la división del texto de esta manera, sin atender a aspectos semánticos, podría partir en dos subdivisiones frases estrechamente relacionadas. Por ejemplo, en la **Figura 3.3**, la frase final de uno de los fragmentos es: «*It was a very descriptive subject*» («Era un tema muy descriptivo»), a la cual le sigue, ya en el siguiente fragmento: «*It really was*» («De veras que lo era»).

Estos casos son difíciles de resolver. Una posible idea sería tratar de determinar si una frase está relacionada con la anterior, quizás mediante el uso de otro modelo, y de ser así tratar de mantenerlas en una misma subdivisión, a fin de que el resumen final mantenga la máxima cohesión y coherencia posibles. Esto incrementaría, no obstante, los tiempos de generación de resúmenes. Por ahora, creemos que los resultados obtenidos son lo suficientemente buenos.

Modelo empleado para la generación de resúmenes: T5

Come hemos mencionado previamente, JIZT hace uso del modelo T5 [16] de Google. Este modelo fue introducido en el artículo *Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer*, presentado en 2019. En él, Colin Raffel *et al.* estudian las ventajas de la técnica del aprendizaje por transferencia (*transfer learning*) al campo del Procesamiento del Lenguaje Natural (NLP).

Tradicionalmente, cada nuevo modelo se entrenaba desde cero. Esto ha cambiado con la inclusión del aprendizaje por transferencia; actualmente, la tendencia es emplear modelos pre-entrenados como punto de partida para la construcción de nuevos modelos.

Las tres principales ventajas del empleo del aprendizaje por transferencia son [33]:

- Mejora del rendimiento de partida. El hecho de comenzar con un modelo pre-entrenado en vez de un modelo ignorante (*ignorant learner*), proporciona un rendimiento base desde el primer momento.
- Disminución del tiempo de desarrollo del modelo, consecuencia del punto anterior.

- Mejora del rendimiento final. Esta mejora ha sido estudiada tanto en el caso del NLP [34], como de otros ámbitos, como la visión artificial [35], o el campo de la medicina [36].

La principal novedad de este artículo se encuentra en su propuesta de tratar todos los problemas de procesamiento de texto como problemas texto a texto (*text-to-text*), es decir, tomar un texto como entrada, y producir un nuevo texto como salida. Esto permite crear un modelo general, al que han bautizado como T5, capaz de llevar a cabo diversas tareas de NLP, como muestra el siguiente diagrama:

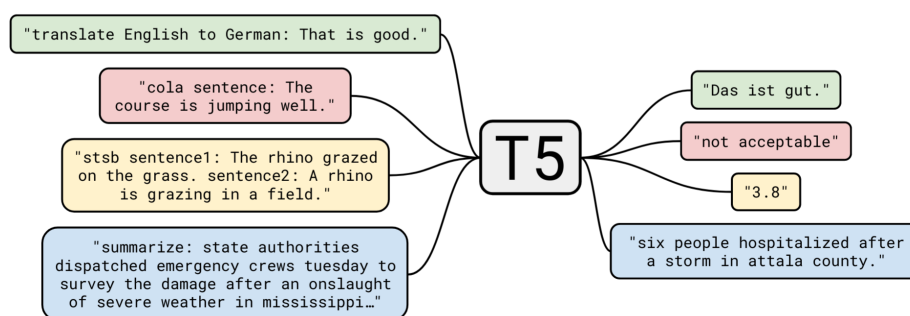


Figura 3.4: El *framework* texto a texto permite emplear el mismo modelo, con los mismos hiperparámetros, función de pérdida, etc., para aplicarlo a diversas tareas de NLP [16].

En cualquier caso, se puede realizar un ajuste fino del modelo para una de las tareas, a fin de mejorar su rendimiento en dicha tarea específica.

Las posibilidades que este modelo nos ofrece son muy interesantes, dado que en un futuro, nuestro proyecto podría incluir otras tareas de Procesamiento de Lenguaje Natural, haciendo uso de un solo modelo.

Principales estrategias de generación de resúmenes

Una de las técnicas más extendidas de generación de lenguaje es la auto-regresión, la cual se basa en el supuesto de que la distribución de probabilidad de una secuencia de palabras puede descomponerse en el producto de las distribuciones condicionales de las siguientes palabras [37]. Expresado matemáticamente:

$$P(w_{1:t}|W_0) = \prod_{t=1}^T P(w_t|w_{1:T-1}, W_0), \text{ siendo } w_{1:0} = \emptyset$$

donde W_0 es la secuencia inicial de *contexto*. En nuestro caso, esa secuencia inicial va a ser el propio texto de entrada. La longitud de T no se puede conocer de antemano, dado que se corresponde con el momento $t = T$ en el que el modelo genera el *token* de finalización de secuencia (EOS), mencionado anteriormente.

Una vez introducido el concepto de auto-regresión, podemos explicar brevemente las cuatro principales estrategias de generación de lenguaje, las cuales se pueden aplicar todas ellas a la generación de resúmenes: búsqueda voraz, *beam search*, muestreo *top-k*, y muestreo *top-p*.

Búsqueda voraz

La búsqueda voraz, en cada paso, simplemente selecciona la palabra con mayor probabilidad de ser la siguiente, es decir, $w_t = \operatorname{argmax}_w P(w|w_{t-1})$ para cada paso t .

Por ejemplo, dada la palabra "El", la siguiente palabra elegida sería "cielo", por ser la palabra con mayor probabilidad, y a continuación "está", y así sucesivamente.

Este tipo de generación tiene dos problemas principales:

- Los modelos, llegados a cierto punto, comienzan a repetir las mismas palabras una y otra vez. En realidad, esto es un problema que afecta a todos los modelos de generación, pero especialmente a los que emplean búsqueda voraz y *beam search* [38, 39].
- Palabras con probabilidades altas pueden quedar enmascaradas tras otras con probabilidades bajas. Por ejemplo, en el anterior ejemplo, la secuencia "El cielo es azul" nunca se dará, porque a pesar de que "azul" presenta una probabilidad alta, está precedida por "es", la cual no será escogida por tener una probabilidad baja.

3.4. Post-procesado del texto

Como veíamos en la Figura 3.3, el resumen producido por el modelo T5, una vez decodificado, se encuentra todo en minúsculas. Por lo demás, el

modelo parece hacer un buen trabajo a la hora de generar el texto en lo que a colocación de puntuación y espacios se refiere, luego la principal labor de esta etapa será poner mayúsculas allí donde sean necesarias, lo que en inglés se denomina *truecasing* [40].

Las mayúsculas, tanto en inglés como español, se emplean principalmente en dos ocasiones:

- Al inicio de cada frase. Como veíamos en la sección referente al **pre-procesado** del texto, la separación de un texto en frases no es, por lo general, una tarea trivial. En este caso, podemos reutilizar lo aplicado en dicha etapa. Teniendo el resumen generado dividido en frases, podemos fácilmente poner la primera letra de cada una de ellas en mayúsculas.
- En los nombres propios. En este aspecto, de nuevo vuelve a aparecer el problema del Reconocimiento de Entidades Nombradas (NER). De modo similar a como procedíamos en el pre-procesado, emplearemos un modelo estadístico que realiza la labor de *truecasing*.

Tras esta etapa, el resumen está listo para ser entregado al usuario.

Técnicas y herramientas

4.1. Flask y Flask-RESTful

Flask es uno de los *frameworks* más populares para la creación de aplicaciones *web* en Python[41]. Está concebido para ser lo más simple posible. En nuestro caso, la hemos empleado para implementar la lógica de la API REST. Además, hemos utilizado una conocida extensión de Flask, Flask-RESTful [42], que facilita aún más dicha implementación.

4.2. Docker

Se trata de una serie de servicios como plataforma (PaaS), que proporcionan virtualización a nivel de sistema operativo, permitiendo ejecutar *software* en paquetes llamados *contenedores* [43].

A diferencia de las máquinas virtuales, en las cuales el sistema operativo subyacente se comparte a través del hipervisor, cada contenedor Docker ejecuta su propio sistema operativo.

Docker nos va a permitir encapsular cada servicio en un contenedor, posibilitando la implementación de la arquitectura de microservicios.

4.3. Kubernetes

Aspectos relevantes del desarrollo del proyecto

Mencionar Cloud Native.

Trabajos relacionados

6.1. Artículos académicos

Las publicaciones más relevantes para nuestro proyecto están relacionadas con los modelos de los que hacemos uso para la generación de resúmenes: T5 y Truecase.

tRuEcasIng

En este artículo, fruto de la colaboración entre la universidad Carnegie Mellon (Pensilvania, EE. UU.) e IBM, los autores Lucian Vlad Lita *et al.*, exploran los problemas del *truecasing*, es decir, el proceso de recomponer las mayúsculas de un texto, y proponen un *truecaser* estadístico que alcanza una precisión del 98 % en artículos de noticias [40].

6.2. Proyectos similares

A continuación, enumeraremos algunos proyectos relacionados con nuestro trabajo.

Bert Extractive Summarizer

Este proyecto *open-source* implementa un generador de resúmenes extractivos haciendo uso del modelo BERT [44] de Google para la codificación de palabras, y aplicando *clustering* por *k-means* para determinar las frases que se incluirán en el resumen. Este proceso se detalla en [45].

El generador de resúmenes puede ser *dockerizado*, pudiéndose ejecutar como servicio, proporcionando una REST API para solicitar los resúmenes. El autor ofrece *endpoints* gratuitos con limitaciones a la hora de realizar peticiones, y *endpoints* privados de pago para aquellos particulares o empresas que requieran de mayores prestaciones.

Se puede acceder al proyecto a través del siguiente enlace:

<https://github.com/dmmiller612/bert-extractive-summarizer>.

ExplainToMe

ExplainToMe es un proyecto también *open-source* centrado en la generación de resúmenes extractivos de páginas *web*, permitiendo cómodamente pegar y copiar el *link* de la *web* que se quiere resumir.

Emplea el algoritmo de TextRank [46], el cual a su vez está inspirado en el conocido PageRank [47], el algoritmo basado en grafos que empleaba originalmente Google en su motor de búsqueda. En su caso, TextRank aplica los principios del algoritmo de Google a la extracción de las frases más importantes de un texto.

Como en el caso anterior, también implementa una API REST.

El proyecto no ha sido actualizado desde finales de 2018. Se puede visitar a través de: <https://github.com/jjangsangy/ExplainToMe/tree/master>.

Smmry

Se trata de uno de las primeras opciones que aparecen en los motores de búsqueda a la hora de buscar «*summarizers*». También genera resúmenes extractivos, aunque a diferencia de los anteriores, no es un proyecto *open-source*.

Destacan su velocidad (*cachea* los textos resumidos recientemente), y sus múltiples opciones de resumen, como por ejemplo: ignorar preguntas, exclamaciones o frases entrecomilladas en el texto original, o la generación de mapas de calor en función de la importancia de las frases incluidas en el resumen.

En su página *web* no se explicita el algoritmo concreto que se emplea, pero prestando atención a la descripción del proceso proporcionada [48], parecen emplear igualmente PageRank.

Se puede acceder a Smmry en: <https://smmry.com/>.

Tabla comparativa

Caraterísticas	JIZT	Bert Extractive Summarizer	ExplainToMe	SMMRY
Tipo de resumen ¹	Abstractivo	Extractivo	Extractivo	Extractivo
Tiempo resumen corto ²	TODO	6 seg.	9 seg.	TODO
Tiempo resumen largo ³	4 min.	No disponible ⁴	Error	∞
Ajustes básicos	✓	✓	✓	✓
Ajustes avanzados	✓	×	×	✓
Entrada: texto plano	✓	✓	×	✓
Entrada: URL	Próximamente	×	✓	✓
Soporte multi-modelo ⁵	Próximamente	×	×	×
Soporte multi-tarea ⁶	Próximamente	×	×	×
API REST	✓	✓	✓	✓
Arquitectura	Microservicios	Monolítica	Monolítica	?
Plataforma	Multiplataforma ⁷	Web	Web	Web
Open-source	✓	✓	✓	×
Gratis	✓	Limitado	✓	Limitado
Proyecto activo	✓	✓	×	✓

1. En los resúmenes *abstractivos*, se toman las frases literales del texto original. En los *extractivos*, se añaden palabras o expresiones nuevas.
2. Texto de entrada con ~6.500 caracteres.
3. Texto de entrada con ~90.000 caracteres.
4. La versión gratuita está limitada. No hemos tenido acceso a la versión completa.
5. Capacidad de generar resúmenes utilizando diferentes modelos.
6. Capacidad de realizar otras tareas de NLP diferentes a la generación de resúmenes.
7. Soporte nativo para Android, iOS y *web*. Pronto, soporte para Linux, macOS y Windows.

Tabla 6.1: Comparativa de las características ofrecidas por las diferentes alternativas para la generación de resúmenes.

Conclusiones y Líneas de trabajo futuras

Todo proyecto debe incluir las conclusiones que se derivan de su desarrollo. Éstas pueden ser de diferente índole, dependiendo de la tipología del proyecto, pero normalmente van a estar presentes un conjunto de conclusiones relacionadas con los resultados del proyecto y un conjunto de conclusiones técnicas. Además, resulta muy útil realizar un informe crítico indicando cómo se puede mejorar el proyecto, o cómo se puede continuar trabajando en la línea del proyecto realizado.

Bibliografía

- [1] “Daniel Crevier. AI: The tumultuous history of the search for artificial intelligence. NY: Basic Books, 1993. 432 pp. (Reviewed by Charles Fair)”. En: *Journal of the History of the Behavioral Sciences* 31.3 (1995), págs. 273-278. DOI: [https://doi.org/10.1002/1520-6696\(199507\)31:3<273::AID-JHBS2300310314>3.0.CO;2-1](https://doi.org/10.1002/1520-6696(199507)31:3<273::AID-JHBS2300310314>3.0.CO;2-1). eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/1520-6696%28199507%2931%3A3%3C273%3A%3AAID-JHBS2300310314%3E3.0.CO%3B2-1>. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/1520-6696%28199507%2931%3A3%3C273%3A%3AAID-JHBS2300310314%3E3.0.CO%3B2-1>.
- [2] A. M. Turing. “Computing Machinery and Intelligence”. En: *Mind* LIX.236 (oct. de 1950), págs. 433-460. ISSN: 0026-4423. DOI: [10.1093/mind/LIX.236.433](https://academic.oup.com/mind/article-pdf/LIX/236/433/30123314/lix-236-433.pdf). eprint: <https://academic.oup.com/mind/article-pdf/LIX/236/433/30123314/lix-236-433.pdf>. URL: <https://doi.org/10.1093/mind/LIX.236.433>.
- [3] Pamela McCorduck. *Machines Who Think*. USA: W. H. Freeman y Co., 1979. ISBN: 0716710722.
- [4] Joachim Rahmfeld. *Recent Advances in Natural Language Processing*. Sep. de 2019. URL: <https://venturebeat.com/2021/01/06/ai-models-from-microsoft-and-google-already-surpass-human-performance-on-the-superglue-language-benchmark/>. Último acceso: 26/01/2021.
- [5] Thang Luong, Hieu Pham y Christopher D. Manning. “Effective Approaches to Attention-based Neural Machine Translation”. En: *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. Lisbon, Portugal: Association for Computational

- Linguistics, sep. de 2015, págs. 1412-1421. DOI: [10.18653/v1/D15-1166](https://doi.org/10.18653/v1/D15-1166). URL: <https://www.aclweb.org/anthology/D15-1166>.
- [6] Dzmitry Bahdanau, Kyunghyun Cho y Yoshua Bengio. *Neural Machine Translation by Jointly Learning to Align and Translate*. 2016. arXiv: [1409.0473](https://arxiv.org/abs/1409.0473) [cs.CL].
- [7] Ashish Vaswani y col. “Attention Is All You Need”. En: *CoRR* abs/1706.03762 (2017). arXiv: [1706.03762](https://arxiv.org/abs/1706.03762). URL: <http://arxiv.org/abs/1706.03762>.
- [8] Thomas Macaulay. *Someone let a GPT-3 bot loose on Reddit — it didn't end well*. Oct. de 2020. URL: <https://thenextweb.com/neural/2020/10/07/someone-let-a-gpt-3-bot-loose-on-reddit-it-didnt-end-well>. Último acceso: 26/01/2021.
- [9] Kyle Wiggers. *AI models from Microsoft and Google already surpass human performance on the SuperGLUE language benchmark*. Ene. de 2021. URL: <https://venturebeat.com/2021/01/06/ai-models-from-microsoft-and-google-already-surpass-human-performance-on-the-superglue-language-benchmark/>. Último acceso: 26/01/2021.
- [10] Google. *Cloud Natural Language*. URL: <https://cloud.google.com/natural-language>. Último acceso: 26/01/2021.
- [11] IBM. *Watson*. URL: <https://www.ibm.com/watson/about>. Último acceso: 26/01/2021.
- [12] Amazon. *Comprehend*. URL: <https://aws.amazon.com/es/comprehend>. Último acceso: 26/01/2021.
- [13] Microsoft. *Text Analytics*. URL: <https://azure.microsoft.com/es-es/services/cognitive-services/text-analytics>. Último acceso: 26/01/2021.
- [14] Raconteur. *A Day in Data*. 2019. URL: <https://www.raconteur.net/infographics/a-day-in-data/>. Último acceso: 26/01/2021.
- [15] Abigail See, Peter J. Liu y Christopher D. Manning. “Get To The Point: Summarization with Pointer-Generator Networks”. En: *CoRR* abs/1704.04368 (2017), pág. 1. arXiv: [1704.04368](https://arxiv.org/abs/1704.04368). URL: <http://arxiv.org/abs/1704.04368>.
- [16] Colin Raffel y col. “Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer”. En: *CoRR* abs/1910.10683 (2019), pág. 11. arXiv: [1910.10683](https://arxiv.org/abs/1910.10683). URL: <http://arxiv.org/abs/1910.10683>.

- [17] Microsoft. *Defining Cloud Native*. Mayo de 2020. URL: <https://docs.microsoft.com/en-us/dotnet/architecture/cloud-native/definition>. Último acceso: 27/01/2021.
- [18] John Arundel y Justin Domingus. *Cloud Native DevOps with Kubernetes*. O'Reilly Media, Inc., mar. de 2019. ISBN: 9781492040767.
- [19] Adam Bellemare. *Building Event-Driven Microservices: Leveraging Organizational Data at Scale*. O'Reilly Media, Inc., 2020. ISBN: 9781492057895.
- [20] Crunchy Data. *Crunchy PostgreSQL Operator*. Mayo de 2021. URL: <https://access.crunchydata.com/documentation/postgres-operator/latest>. Último acceso: 27/01/2021.
- [21] Robert Martin. *Clean Architecture: A Craftsman's Guide to Software Structure and Design*. Pearson Education, 2015. ISBN: 9780134494166.
- [22] Daniel Sauble. *Offline First Web Development*. Packt, 2015. ISBN: 9781785884573.
- [23] Mike Lewis y col. "BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension". En: *CoRR* abs/1910.13461 (2019). arXiv: [1910.13461](https://arxiv.org/abs/1910.13461). URL: <http://arxiv.org/abs/1910.13461>.
- [24] Wikipedia. *Reconocimiento de entidades nombradas - Wikipedia, La enciclopedia libre*. 2020. URL: https://es.wikipedia.org/wiki/Reconocimiento_de_entidades_nombradas. Último acceso: 27/01/2021.
- [25] Christopher Manning - Stanford University. *Stanford CS224N: NLP with Deep Learning. Winter 2019. Lecture 13. Contextual Word Embeddings*. 2019. URL: <https://www.youtube.com/watch?v=S-CspeZ8FHc>. Último acceso: 28/01/2021.
- [26] Linlin Hou y col. *Method and Dataset Entity Mining in Scientific Literature: A CNN + Bi-LSTM Model with Self-attention*. 2020. arXiv: [2010.13583](https://arxiv.org/abs/2010.13583) [cs.AI].
- [27] Tomas Mikolov y col. *Efficient Estimation of Word Representations in Vector Space*. 2013. arXiv: [1301.3781](https://arxiv.org/abs/1301.3781) [cs.CL].
- [28] Tomás Mikolov y col. "Distributed Representations of Words and Phrases and their Compositionality". En: *CoRR* abs/1310.4546 (2013). arXiv: [1310.4546](https://arxiv.org/abs/1310.4546). URL: <http://arxiv.org/abs/1310.4546>.

- [29] Jeffrey Pennington, Richard Socher y Christopher Manning. “GloVe: Global Vectors for Word Representation”. En: *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Doha, Qatar: Association for Computational Linguistics, abr. de 2014, págs. 1532-1543. DOI: [10.3115/v1/D14-1162](https://doi.org/10.3115/v1/D14-1162). URL: <https://www.aclweb.org/anthology/D14-1162>.
- [30] Matthew E. Peters y col. “Deep contextualized word representations”. En: *CoRR* abs/1802.05365 (2018). arXiv: [1802.05365](https://arxiv.org/abs/1802.05365). URL: <http://arxiv.org/abs/1802.05365>.
- [31] Jacob Devlin y col. “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”. En: *CoRR* abs/1810.04805 (2018). arXiv: [1810.04805](https://arxiv.org/abs/1810.04805). URL: <http://arxiv.org/abs/1810.04805>.
- [32] Bożena Cetnarowska. “Ingo Plag, Word-formation in English (Cambridge Textbooks in Linguistics). Cambridge: Cambridge University Press, 2003. Pp. xiv 240.” En: *Journal of Linguistics* 41.1 (2005). DOI: [10.1017/S0022226704303233](https://doi.org/10.1017/S0022226704303233).
- [33] Dipanjan Sarkar, Raghav Bali y Tamoghna Ghosh. *Hands-On Transfer Learning with Python*. Packt Publishing, 2018. ISBN: 9781788831307.
- [34] Manoj Kumar y col. *ProtoDA: Efficient Transfer Learning for Few-Shot Intent Classification*. 2021. arXiv: [2101.11753](https://arxiv.org/abs/2101.11753) [cs.CL].
- [35] Nuredin Ali. *Exploring Transfer Learning on Face Recognition of Dark Skinned, Low Quality and Low Resource Face Data*. 2021. arXiv: [2101.10809](https://arxiv.org/abs/2101.10809) [cs.CV].
- [36] Yi Liu y Shuiwang Ji. *A Multi-Stage Attentive Transfer Learning Framework for Improving COVID-19 Diagnosis*. 2021. arXiv: [2101.05410](https://arxiv.org/abs/2101.05410) [eess.IV].
- [37] Patrick von Platen. *How to generate text: using different decoding methods for language generation with Transformers*. Mar. de 2020. URL: <https://huggingface.co/blog/how-to-generate>. Último acceso: 31/01/2021.
- [38] Ashwin K. Vijayakumar y col. “Diverse Beam Search: Decoding Diverse Solutions from Neural Sequence Models”. En: *CoRR* abs/1610.02424 (2016). arXiv: [1610.02424](https://arxiv.org/abs/1610.02424). URL: <http://arxiv.org/abs/1610.02424>.
- [39] Louis Shao y col. “Generating Long and Diverse Responses with Neural Conversation Models”. En: *CoRR* abs/1701.03185 (2017). arXiv: [1701.03185](https://arxiv.org/abs/1701.03185). URL: <http://arxiv.org/abs/1701.03185>.

- [40] Lucian Vlad Lita y col. “TRuEcasIng”. En: *Proceedings of the 41st Annual Meeting on Association for Computational Linguistics - Volume 1*. ACL '03. Sapporo, Japan: Association for Computational Linguistics, 2003, págs. 152-159. DOI: [10.3115/1075096.1075116](https://doi.org/10.3115/1075096.1075116). URL: <https://doi.org/10.3115/1075096.1075116>.
- [41] The Pallets Projects. *Flask*. 2021. URL: <https://palletsprojects.com/p/flask>. Último acceso: 29/01/2021.
- [42] Flask-RESTful Community. *Flask-RESTful*. 2021. URL: <https://flask-restful.readthedocs.io/en/latest>. Último acceso: 29/01/2021.
- [43] Docker. *Why Docker?* 2021. URL: <https://www.docker.com/why-docker>. Último acceso: 29/01/2021.
- [44] Jacob Devlin y col. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. 2019. arXiv: [1810.04805](https://arxiv.org/abs/1810.04805) [cs.CL].
- [45] Derek Miller. “Leveraging BERT for Extractive Text Summarization on Lectures”. En: *CoRR* abs/1906.04165 (2019). arXiv: [1906.04165](https://arxiv.org/abs/1906.04165). URL: <http://arxiv.org/abs/1906.04165>.
- [46] Rada Mihalcea y Paul Tarau. “TextRank: Bringing Order into Text.” En: jul. de 2004.
- [47] Lawrence Page y col. *The PageRank Citation Ranking: Bringing Order to the Web*. Technical Report 1999-66. Previous number = SIDL-WP-1999-0120. Stanford InfoLab, nov. de 1999. URL: <http://ilpubs.stanford.edu:8090/422/>.
- [48] Smmry Team. *Smmry*. 2021. URL: <https://smmry.com/about>. Último acceso: 31/01/2021.