



UNIVERSIDAD DE CHILE
FACULTAD DE CIENCIAS FÍSICAS Y MATEMÁTICAS
DEPARTAMENTO DE CIENCIAS DE LA COMPUTACIÓN

LIGHT AND FAST LANGUAGE MODELS FOR SPANISH THROUGH COMPRESSION
TECHNIQUES

TESIS PARA OPTAR AL GRADO DE
MAGÍSTER EN CIENCIAS, MENCIÓN COMPUTACIÓN

MEMORIA PARA OPTAR AL TÍTULO DE
INGENIERO CIVIL EN COMPUTACIÓN

JOSÉ RODOLFO CAÑETE LÓPEZ

PROFESOR GUÍA:
FELIPE BRAVO MÁRQUEZ

MIEMBROS DE LA COMISIÓN:
JOCELYN DUNSTAN ESCUDERO
AIDAN HOGAN
IVÁN SIPIRÁN MENDOZA

Este trabajo ha sido parcialmente financiado por el proyecto FONDECYT de Iniciación
11200290.

SANTIAGO DE CHILE
2023

Modelos de lenguaje ligeros y rápidos para español mediante el uso de técnicas de compresión

Los grandes modelos de lenguaje (LLM, por sus siglas en inglés) se han convertido en un enfoque común y exitoso para abordar tareas de procesamiento de lenguaje natural (NLP, por sus siglas en inglés), incluyendo, pero no limitado a, clasificación de documentos, reconocimiento de entidades nombradas y respuesta a preguntas. A pesar de su notable rendimiento, utilizar estos LLM en entornos con recursos limitados, como aplicaciones web o móviles, es un desafío, especialmente en escenarios en tiempo real que demandan respuestas rápidas. Recientemente han surgido técnicas para comprimir estos LLM en modelos más pequeños y rápidos, particularmente en el caso de modelos en inglés o multilingües, pero aún es un desafío para otros idiomas. De hecho, el español es el segundo idioma con más hablantes nativos pero carece de este tipo de recursos.

En este trabajo, presentamos ALBETO y Speedy Gonzales, dos nuevos recursos para la comunidad de NLP en español que tienen como objetivo cubrir la brecha en términos de modelos más livianos y rápidos para el español. ALBETO es un conjunto de 5 modelos ligeros, con tamaños que van desde 5M a 223M de parámetros, que están pre-entrenados usando exclusivamente datos en español siguiendo la arquitectura de ALBERT. Evaluamos nuestros modelos ALBETO junto con otros modelos disponibles para el español en un conjunto de 6 tareas y luego, mediante el uso de la técnica de *Knowledge Distillation* (KD), presentamos Speedy Gonzales, una colección de modelos de lenguaje, basados en ALBETO, más eficientes en inferencia para tareas en español.

Los resultados de nuestro estudio revelan que nuestros modelos ALBETO tienen un rendimiento en tareas similar a otros modelos con velocidad de inferencia comparable, a pesar de ser más ligeros y tener sustancialmente menos parámetros. Además, nuestro modelo ALBETO *xxlarge* supera a todos los demás modelos pre-entrenados en español que están actualmente disponibles.

En cuanto a nuestros modelos Speedy Gonzales, los resultados indican una mejora en la velocidad de inferencia a expensas de una ligera disminución en el rendimiento en las tareas. Es importante notar que esta disminución es mínima en el caso de nuestros modelos de 8 y 10 capas, mientras que es más pronunciada en los modelos más rápidos con 2-4 capas. Además, nuestro modelo de 10 capas, que llamamos ALBETO *base-10*, proporciona un rendimiento que es generalmente comparable a los modelos de tamaño *base*, al tiempo que demuestra mejor velocidad de inferencia.

Todos nuestros modelos (pre-entrenados, fine-tuneados y destilados) están disponibles públicamente en: <https://huggingface.co/dccuchile>.

Light and Fast Language Models for Spanish Through Compression Techniques

Large language models (LLMs) have become a prevalent and successful approach to address natural language processing (NLP) tasks, including but not limited to document classification, named-entity recognition, and question answering. Despite their remarkable performance, utilizing these LLMs on constrained resources, such as web or mobile applications, is challenging, particularly in real-time scenarios that demand fast responses. Techniques to compress these LLM into smaller and fastest models have emerged for English or Multilingual settings, but it is still a challenge for other languages. In fact, Spanish is the second language with most native speakers but lacks of these kind of resources.

In this work, we present ALBETO and Speedy Gonzales, two new resources for the Spanish NLP community that aim to bridge the gap in terms of lighter and faster models for Spanish. ALBETO is a set of 5 lightweight models, with sizes ranging from 5M to 223M parameters, that are pre-trained exclusively on Spanish corpora following the ALBERT architecture. We evaluate our ALBETO models along with other publicly available models for Spanish on a set of 6 tasks and then, by leveraging on *Knowledge Distillation* (KD), we present Speedy Gonzales, a collection of more inference-efficient task-specific language models based on ALBETO.

The outcomes of our study reveal that our ALBETO models perform at a similar level to other models with comparable inference speed, despite being lighter in weight and having substantially fewer parameters. Moreover, our ALBETO *xxlarge* model outperforms all other pre-trained Spanish models that are currently available.

Regarding our Speedy Gonzales models, the results indicate an enhancement in inference speed at the expense of a slight decline in task performance. Notably, this decay is minimal in the case of our 8 and 10 layer models, while it is more pronounced in the faster models with 2-4 layers. Moreover, our 10-layer model, referred to as ALBETO *base-10*, delivers performance that is generally comparable to *base*-sized models, while also demonstrating improved inference speed.

All of our models (pre-trained, fine-tuned and distilled) are publicly available on: <https://huggingface.co/dccuchile>.

A mi familia, amigos y todos los que han estado en este camino.

Agradecimientos

Me siento muy agradecido de toda la gente bacán que me ha acompañado durante este largo pero también bonito camino que ha sido el paso por la universidad.

Primero a mi familia, que es mi gran apoyo. A mis papás Isabel y José, que siempre han estado para mi. A mis hermanos, Marcela y Daniel que siempre me han acompañado. A mis abuelos, mi tata y mi mama, también a la Myrian, que siempre me han acogido y querido mucho.

Agradezco también los amigos que han estado desde antes: al Nacho y al Joni. También a la Nati, que me acompañó durante gran parte de esta etapa.

También me hace muy feliz pensar en todos los nuevos amigos que hice durante estos años. Agradezco mucho a mis amigos del DCC y programación competitiva, en especial a mi equipo: Thibo y Ale.

A los amigos que hice en el equipo de difusión: Edu, Braian, Gonza, Nico, Gustavo, Karin, Tati, Fran, entre muchos, muchos otros. También a la María José, por siempre recibirnos con su buena onda en la oficina.

Agradezco de todo corazón a todos los que organizaron el ramo de Taller de Programación por encantarme con la computación y la curiosidad y ganas de resolver problemas entretenidos.

Agradezco también a Camilo y Sergio de Adereso por depositar su confianza en mí y enseñarme tanto, especialmente a Sergio por enseñarme de que se trata ser un buen líder.

También agradezco a los amigos que conocí estos últimos años en el magíster y en CENIA, en especial al Seba Donoso, Gabriel, Miguel, Eugenio, Pablo y Lucas, quienes son grandes amigos. También a Felipe Guerrero, del que aprendí mucho.

Finalmente agradezco a mis profesores de la universidad, en especial a Jorge Pérez y Felipe Bravo, de quienes aprendí muchísimo, tanto en lo académico como en la vida. También agradezco a los profesores de mi comisión, Jocelyn Dunstan, Aidan Hogan e Iván Sibirán, por sus valiosos comentarios constructivos sobre este trabajo.

Muchas gracias a todos.

Table of Content

1	Introduction	1
1.1	Research Problem	2
1.2	Hypothesis	4
1.3	Research Questions	4
1.4	Objectives	4
1.4.1	General Objective	4
1.4.2	Specific Objectives	5
1.5	Results	5
1.6	Research Outcome	5
1.7	Thesis Structure	6
2	Background and Related Work	7
2.1	Scientific Disciplines	7
2.1.1	Artificial Intelligence	7
2.1.2	Machine Learning	8
2.1.3	Representation Learning	8
2.1.4	Transfer Learning	9
2.1.5	Natural Language Processing	9
2.1.6	Representations of Text	11
2.2	Transformers	11
2.2.1	Vanilla Transformers	12

2.2.2	BERT	14
2.2.3	RoBERTa	14
2.2.4	ALBERT	16
2.2.5	Multilingual and Monolingual Models	17
2.3	Model Compression and Acceleration	18
2.3.1	Compression Techniques	19
2.3.2	Knowledge Distillation	20
2.3.3	Knowledge Distillation Techniques for BERT models	21
3	Preliminaries	22
3.1	Evaluation Tasks and Metrics	22
3.1.1	Evaluation Metrics	23
3.1.2	Evaluation Tasks	24
3.2	Overview of Pre-trained Models for Spanish	26
3.2.1	BETO	26
3.2.2	DistilBETO	27
3.2.3	RoBERTa-BNE	27
3.2.4	BERTIN	27
3.2.5	Models Not Included	27
4	Proposed Spanish NLP Resources	29
4.1	ALBETO: Lighter Models for Spanish	29
4.1.1	Data Collection and Preprocessing	29
4.1.2	Model Architecture	31
4.1.3	Pre-training Procedure	31
4.1.4	Fine-tuning and Experimental Setup	32
4.2	Speedy Gonzales: Faster Models for Spanish	32
4.2.1	Approach	33

4.2.2	Student Models	34
4.2.3	Applying Knowledge Distillation to Different NLP Tasks	35
4.2.4	Additional Implementation Details and Experimental Setup	37
5	Results and Discussion	40
5.1	Task Performance	40
5.2	Model Efficiency and Inference Speed	43
5.3	Inference Speed on Common Hardware	46
6	Conclusion	48
6.1	Summary of Contributions	48
6.2	Limitations and Future Research Directions	49
	Bibliography	62
	Annexes	63
	Annex A Examples of Considered Tasks	63
	Annex B Training Loss on ALBETO models	66
	Annex C Selected Teacher Models	68
	Annex D Knowledge Distillation Implementation for Models With Shared Vocabulary	69
	Annex E Effect of Caching Teacher Outputs During Training	71
	Annex F Importance of Knowledge Distillation	73

List of Tables

3.1	Details of the datasets used to evaluate our proposed models.	24
3.2	Summary of pre-trained Transformer models for Spanish.	28
4.1	The configurations of each ALBETO model trained in this work.	31
4.2	Training details of all ALBETO models, which were trained using a single TPU v3-8 each one.	32
5.1	Models evaluated on sentence or two sentences classification tasks, results are measured using accuracy on the test set of each dataset.	41
5.2	Models evaluated on sequence tagging tasks, results are measured using the F1 Score on the test set of each dataset.	42
5.3	Models evaluated on question answering datasets, results are noted as F1 Score / Exact Match on the test set of each dataset.	43
5.4	The summary of results of every evaluated model in terms of parameters, inference speedup and overall score across tasks. The speedup is relative to BETO models. The score column shows the average of the metrics on all tasks.	45
5.5	The number of inferences per second of each model on two different hardware settings, CPU and GPU.	47
C.1	The teacher models selected for each task.	68
E.1	Training times when using teacher cache vs not using it. Table report the mean (M) and standard deviation (SD) over three runs.	72
F.1	Comparison of the performance of our proposed models on text classification tasks on two settings: fine-tuning and task-specific knowledge distillation.	74
F.2	Comparison of the performance of our proposed models on sequence tagging tasks on two settings: fine-tuning and task-specific knowledge distillation.	74

F.3	Comparison of the performance of our proposed models on question answering on two settings: fine-tuning and task-specific knowledge distillation.	74
-----	---	----

List of Figures

1.1	The size (number of parameters) and speed (MACs) of every Spanish model evaluated on this work. MACs are measured using a single sequence of length 512, which is the maximum sequence length of all the evaluated models. . . .	2
1.2	A summary of the number of parameters in various modern language models, with increasing model sizes from older to newer models. The models included in this comparison are ELMo [75], GPT [77], BERT [32], GPT-2 [78], MegatronLM [96], T5 [82], Turing-NLG [84], GPT-3 [10], Megatron-Turing NLG [97] and PaLM [23].	3
2.1	The Transformer architecture by Vaswani et al. [105].	12
2.2	The masked language modeling (MLM) task used by BERT as pre-training task.	15
2.3	RoBERTa adopts a dynamic masking strategy, which allows for distinct masking patterns to be applied to a sentence during the pre-training phase.	15
2.4	A design comparison of BERT and ALBERT, focusing on the parameter utilization strategy adopted by each model.	17
2.5	A summary of the number of parameters in various modern multilingual language models, from older to newer. The models included in this comparison are mBERT [32], XLM-R Large [26], mBART [65], mT5 XXL [116], XLM-R XXL [42], XGLM [63] and BLOOM [91].	18
2.6	The figure provides a visual representation of the Knowledge Distillation framework applied in this work. In line with common practices, the framework includes both a distillation loss between the teacher and student models and a cross-entropy loss between the gold labels and the student’s predictions, as indicated by the dashed line.	21
4.1	A broad overview of the process involved in the creation of ALBETO models. Sub-processes relevant to distinct stages are portrayed outside the main frame.	30

4.2	The two stage approach of task-specific knowledge distillation used in this work. (Up) The first stage involves fine-tuning a set of candidate models on a specific dataset, followed by the selection of the best-performing model as the teacher model for that dataset. (Down) The second stage employs the selected teacher model to train a set of student models using knowledge distillation.	33
4.3	Implementation of knowledge distillation for text classifications tasks that use a single sentence as input.	35
4.4	Implementation of knowledge distillation for text classifications tasks that use two sentences as input.	36
4.5	Implementation of knowledge distillation for sequence tagging tasks. The tokens marked with the blue color represents the property of being the first token of a word.	36
4.6	Implementation of knowledge distillation for question answering datasets. The tokens marked with the blue color represents the property of being the first token of a word.	37
5.1	Average of performance in the different tasks. The size of the points represents the size of the model (the number of parameters).	44
A.1	An example of document classification. Taken from the MLDoc [93] dataset.	63
A.2	An example of the paraphrase identification task. Taken from the PAWS-X [117] dataset.	64
A.3	An example of natural language inference. Taken from the XNLI [25] dataset.	64
A.4	An example of the Part-of-Speech tagging task. Taken from the AnCora [101] dataset.	64
A.5	An example of named entity recognition. Taken from the CoNLL2002 NER [103] dataset.	65
A.6	An example of question answering. Taken from the SQAC [44] dataset.	65
B.1	The progression of the training loss on the ALBETO <i>tiny</i> model.	66
B.2	The progression of the training loss on the ALBETO <i>base</i> model.	66
B.3	The progression of the training loss on the ALBETO <i>large</i> model.	67
B.4	The progression of the training loss on the ALBETO <i>xlarge</i> model.	67
B.5	The progression of the training loss on the ALBETO <i>xxlarge</i> model.	67

D.1	Implementation of knowledge distillation for sequence tagging tasks with models that share the same vocabulary.	69
D.2	Implementation of knowledge distillation for question answering datasets with models that share the same vocabulary.	70

Chapter 1

Introduction

The use of dense learned representations of text has become a popular and effective approach in various natural language processing (NLP) tasks, including sentiment analysis [98, 86, 69], sentence [110] and document classification [92], natural language inference [111, 24], among others. These learned representations are usually obtained by training a language model using large collections of texts from the web. This success can be explained by a large number of trainable parameters of these models, a significant amount of training data, and the processing capacity power available nowadays, allowing them to learn complex relationships and cast complex functions.

The size of these models has grown overtime and now very large language models (LLM) are common, with models that range from hundred of millions to billions of parameters. These LLMs are usually trained on English by big technology companies using web-scale datasets and substantial computational resources. Prominent examples include the well-known GPT-3 model [10]. For languages other than English the available models are typically variants of BERT [32] and RoBERTa [64]. In the case of Spanish, which is one of the five most spoken languages in the world and the second with most native speakers, the available models range from 110M to 335M parameters. In Figure 1.1 we showed how different Spanish pre-trained models (including those presented on this work) compare in terms of model size (expressed in the number of parameters) and inference speed (quantified in terms of MACs, that is, Multiply-Accumulate Operations).

Despite the remarkable performance of these LLMs across a range of tasks, it remains a challenge to utilize them effectively in computing environments that are constrained by limited resources, such as web or mobile applications.

New techniques to address this problem have emerged for English [100, 104, 90, 109, 51] or Multilingual [52] models. These typically leverage on more parameter-efficient model architectures [55] and different kinds of Knowledge Distillation (KD) [47] to compress the results of a large and performant model into another one which is typically lighter and more inference efficient. For other languages this is still an open challenge, where we lack from this kind of resources.

In this thesis we approach these challenges for Spanish NLP on two crucial aspects, namely

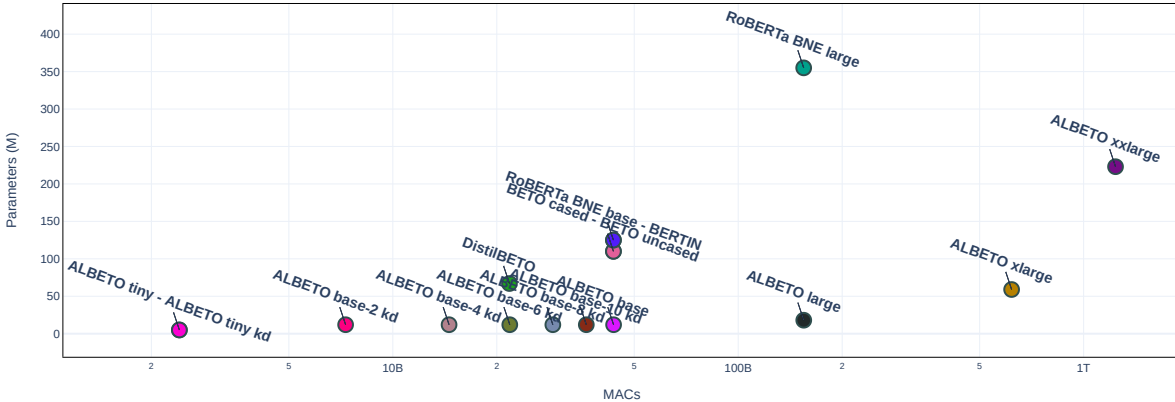


Figure 1.1: The size (number of parameters) and speed (MACs) of every Spanish model evaluated on this work. MACs are measured using a single sequence of length 512, which is the maximum sequence length of all the evaluated models.

model size (number of parameters) and inference speed. To address the issue of model size, we present ALBETO, which is a series of ALBERT [55] models that range from 5M to 223M of parameters and are pre-trained exclusively on Spanish corpora. To tackle the challenge of inference speed, we present Speedy Gonzales, which is a collection of task-specific models based on ALBETO for a diverse set of 6 different tasks, that achieve results comparable to the larger models while exhibiting improved inference speed by utilizing knowledge distillation.

1.1 Research Problem

Language models are rapidly growing in size, this can be visualized in Figure 1.2 where we showed the number of parameters of some models published in the last years. The substantial growth in model size has been facilitated by the availability of significant computing resources, enabling the development of language models with unprecedented capabilities and potential [10].

This remarkable progress in the development of language models has opened up new possibilities for a wide range of natural language processing (NLP) tasks. Language models are becoming increasingly adept at handling complex language tasks, such as text generation [10, 78], machine translation [115, 65, 105], and sentiment analysis [98, 86, 69], among others. In addition, larger models have demonstrated improved performance across a variety of NLP benchmarks [107, 108], suggesting that the expansion of model size has contributed to significant advances in the field.

However, the growing size of language models also poses several challenges. For example, training and fine-tuning these models require large amounts of data and computational resources, which can be prohibitively expensive for many researchers and organizations. The big amount of computational resources that these models require make a special challenge in

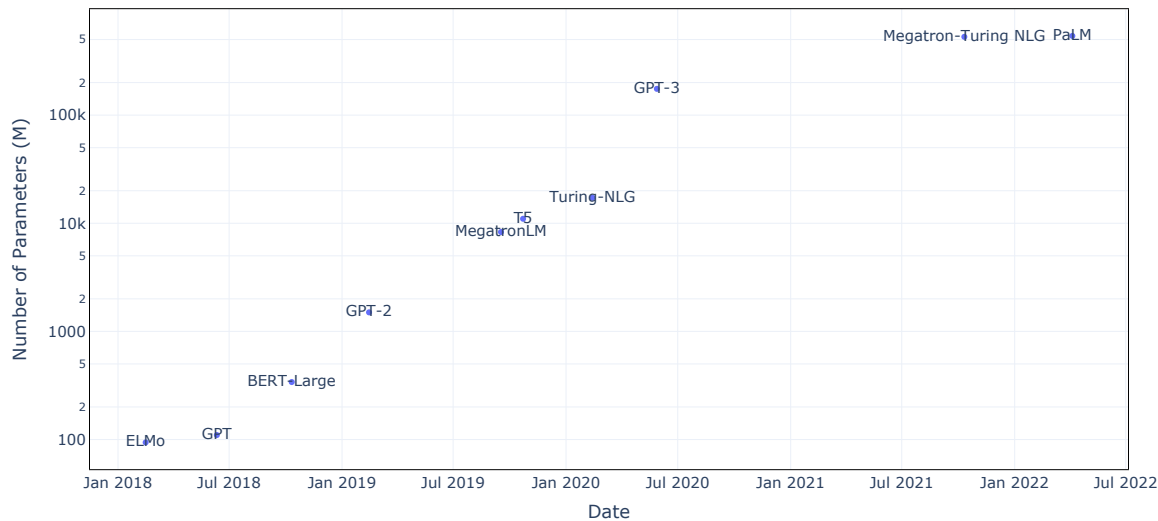


Figure 1.2: A summary of the number of parameters in various modern language models, with increasing model sizes from older to newer models. The models included in this comparison are ELMo [75], GPT [77], BERT [32], GPT-2 [78], MegatronLM [96], T5 [82], Turing-NLG [84], GPT-3 [10], Megatron-Turing NLG [97] and PaLM [23].

deploying these models on constrained resources scenarios like real-time web applications or mobile devices.

As a result, researchers and practitioners are actively exploring new methods to optimize language model performance while reducing their size and resource requirements. These efforts aim to strike a balance between achieving state-of-the-art performance and ensuring that language models remain accessible and efficient for a broader range of users and applications.

One promising approach to address the challenge of deploying language models is the development of more efficient architectures and training algorithms. For example, recent research has focused on developing weight-shared architectures [55] to reduce the overall model size. Other techniques involve parameter pruning [8], quantization [38], and distillation [41], which aim to reduce the number of model parameters and memory requirements as well as the computational complexity of these models, reducing the inference time, while preserving task-performance.

However, deploying these models in languages other than English remains an open problem, as resources for such languages are scarce and the aforementioned techniques have not been fully explored. This issue is especially critical for Spanish, a language spoken widely but lacking adequate resources. Consequently, it is imperative to investigate methods to enhance model efficiency, in terms of size (number of parameters) and inference speed, in order to expand the availability of these models for the Spanish-speaking community.

1.2 Hypothesis

In this research, our hypothesis is that adopting more parameter-efficient model architectures and employing knowledge distillation techniques to transfer knowledge from larger models to smaller ones can significantly enhance model compactness and inference speed, while maintaining most of the performance exhibited by larger models on Spanish NLP tasks. We also hypothesize that the size of the smaller models remains a critical factor affecting task performance, where a reduction in model size could lead to a decrease in performance. Hence, our study aims to identify the sweet spot between model size, inference time, and task performance to achieve optimal results.

1.3 Research Questions

The present study aims to address the following research questions:

- Can training models with weight-shared architectures yield benefits in terms of model performance and parameter efficiency?
- Can the knowledge distillation technique be effectively employed to train models with comparable or even better task performance than larger models, while improving their compactness and efficiency?
- Is it feasible to train much smaller and inference-efficient models that exhibit good performance on Spanish language tasks?
- To what extent can the size and computational demands of the model be reduced, and how does this reduction affect the inference speed and task performance?

1.4 Objectives

1.4.1 General Objective

The principal objective of this research is to develop Spanish language models that are more compact and computationally efficient. This will be achieved through the application of parameter-efficient language model architectures and the utilization of task-specific knowledge distillation, which involves transferring the knowledge acquired by larger models that were trained on Spanish NLP tasks to smaller models that are optimized for efficient inference. The goal is to produce models that require fewer computational resources while maintaining high levels of accuracy and performance.

1.4.2 Specific Objectives

To achieve our general objective, we have identified the following specific objectives:

1. Measure the size and inference speed of pre-trained Spanish language models that are currently available.
2. Develop models for the Spanish language that are more parameter-efficient by utilizing weight-shared model architectures.
3. Train models for Spanish that are more inference-efficient by applying task-specific knowledge distillation on Spanish NLP tasks.
4. Evaluate the mentioned techniques on a diverse set of Spanish NLP tasks.
5. Evaluate how the model size impacts the task performance while using these techniques.
6. Release those models publicly as a resource for further research.

1.5 Results

The results of this thesis can be summarized as follows:

- We present ALBETO, a series of 5 models based on the ALBERT architecture pre-trained exclusively on Spanish corpora, with model sizes that range from 5M to 223M parameters.
- We perform a comprehensive evaluation of all publicly available Spanish pre-trained models, which are trained on general-domain corpora, by fine-tuning them across six different tasks and eight datasets.
- We present Speedy Gonzales, a collection of fast models for Spanish trained by selecting the best model on each evaluated dataset and distilling its knowledge into lighter ALBETO models, achieving lighter and inference efficient models, while retaining most of the task performance of the bigger counterparts.
- We make these newly created resources, ALBETO and Speedy Gonzales, consisting of over 140 models (between pre-trained, fine-tuned and distilled models) publicly accessible on the HuggingFace Hub: <https://huggingface.co/dccuchile>

1.6 Research Outcome

In addition to the present thesis, two separate research papers were authored. The first paper, titled “ALBETO and DistilBETO: Lightweight Spanish Language Models” [12] was presented as a conference paper at the Language Resources and Evaluation Conference (LREC) 2022

[11] in Marseille, France. The second paper, entitled “Speedy Gonzales: A Collection of Fast Task-Specific Models for Spanish” [17] has been submitted and is currently under review for another academic conference.

The software code employed in both papers, as well as in this thesis, is publicly available on two GitHub repositories¹². Additionally, the models created in these papers are accessible on the HuggingFace Hub³.

1.7 Thesis Structure

The rest of this thesis is structured as follows:

Chapter 2 provides an overview of the scientific disciplines that are relevant to this research (Section 2.1) and presents the background knowledge necessary to comprehend the thesis. Section 2.2 offers a review of Transformers models, followed by an examination of related research on reducing the size and increasing the speed of these models in Section 2.3.

Chapter 3 provides an overview of the artifacts used to effectively evaluate and compare the lighter and faster models presented in this study. Section 3.1 describes the evaluation tasks and metrics employed, while Section 3.2 gives an overview of the pre-trained models currently available for Spanish.

Chapter 4 offers a comprehensive explanation of the methodology utilized to develop lighter and faster models. Section 4.1 elaborates on the ALBETO approach used to create the lighter models, and Section 4.2 details the methodology used to develop faster models (Speedy Gonzales). Both sections also describe the implementation details and experimental setup used in the study.

Chapter 5 presents the main findings of the research. Section 5.1 presents the outcomes of fine-tuning and task-specific knowledge distillation approaches on various tasks, while Section 5.2 discusses the results related to model size and speed. Additionally, Section 5.3 showcases the inference speed of all evaluated models on common hardware, such as CPU and GPU.

Chapter 6 concludes this research by summarizing the key findings and outlining potential directions for future research.

¹<https://github.com/dccuchile/lightweight-spanish-language-models>

²<https://github.com/dccuchile/speedy-gonzales>

³<https://huggingface.co/dccuchile>

Chapter 2

Background and Related Work

The chapter is divided into three main sections. The first section provides an introduction to the area of knowledge pertinent to the thesis, which works as background to understand the rest. The second section is a review of relevant literature on Transformers, which are neural networks commonly used for natural language processing. The third section focuses on techniques to compress these models, making them lighter and faster. This chapter sets the groundwork for the rest of the thesis, guiding the investigation and analysis of the research questions.

2.1 Scientific Disciplines

2.1.1 Artificial Intelligence

The concept of artificial intelligence (AI) has been subject to various perspectives on its meaning and definition. The notion of “intelligence” is of particular interest, yet intricate to explicate. Russell [89] has gathered several definitions of AI, which generally encompass the terms “reasoning”, “perception”, “learning”, and “acting”, among others. In this thesis, a more practical definition is adopted, as stated by Negnevitsky [70] “*The goal of artificial intelligence (AI) as a science is to make machines do things that would require intelligence if done by humans*”. Therefore, what will matter to us will be computers (machines) solving well defined problems that could be solved by a human.

In recent years, a range of these problems have been explored using AI systems, including visual perception, system movement, language comprehension, and sound analysis, among other domains. These problems are typically addressed using techniques from a successful subfield of AI called Machine Learning.

2.1.2 Machine Learning

Machine Learning is a subfield of Computer Science that studies the question of how to build algorithms that can automatically improve through experience [53]. The goal of machine learning systems is to learn a mathematical model through the use of previous observations (sample data), that can accurately predict outcomes for new, unseen data. This is also known as a model that generalizes.

To create a machine learning model, various components are necessary. Firstly, data samples are required, the type of which may vary depending on the problem at hand. There are at least, two types of learning paradigms, unsupervised and supervised. We call “unlabeled” the type of data needed in the unsupervised paradigm and “labeled” the type of data needed in the supervised paradigm.

According to Jordan and Mitchell [53], unsupervised learning is the “*analysis of unlabeled data under assumptions about structural properties of the data (e.g., algebraic, combinatorial, or probabilistic)*”. The most common example of unsupervised learning is the problem of clustering, where a machine learning algorithm is used to group data samples based on their similarity or dissimilarity, with the aim of grouping together instances of the data that exhibit similar characteristics.

The second learning paradigm is supervised learning. In supervised learning we use a set of data samples with the form of (x, y) , where x is called an example and y is called its label. The goal in supervised learning is to learn a function $f(x)$ that maps from x to y and, as we stated earlier, that generalizes to unseen pairs (x^*, y^*) .

Typically, the function $f(x)$ is a parameterized function that is learned through the optimization of a loss function, that indicates how well is the function doing on learning the mapping from the examples to the labels. There are different types of loss functions, with the most common ones including the mean-squared error loss, the cross-entropy loss and the KL-divergence loss. The optimization algorithm is typically the Stochastic Gradient Descent (SGD) algorithm or a variation of it [40].

There are numerous families of functions $f(x)$ that can be learned, with (artificial) neural networks being the most widely used today. The area that studies large and complex neural networks is called Deep Learning [59] and includes architectures such as Recurrent Neural Networks (RNN) [87, 48], Convolutional Neural Networks (CNN) [58], and Transformers [105], which we will better describe in Section 2.2.

2.1.3 Representation Learning

Traditionally, in classical Machine Learning, the input examples x were represented as feature vectors, which were manually engineered, in a process called “feature engineering”, by domain experts who possessed knowledge on the specific task at hand. These feature vectors were then fed into machine learning algorithms to represent the underlying examples x .

More recently, a different approach has gained significant popularity: the use of learned

representations of the data. This means that not only a function $f(x)$ is learned but also a rich and useful representation x is learned from a simpler representation of the data. Different ways to learn these representations have been proposed and is a very active area of research called Representation Learning [7].

2.1.4 Transfer Learning

The development of Representation Learning has led significant progress in another, very related area, called Transfer Learning. The key idea in Transfer Learning is to reuse the knowledge (or the representation) learned in one very general task for another more specific task. As an example, the representations learned to distinguish a tiger can be then reused to distinguish a different, but related, animal such as a lion or a panther.

Currently, two notable cases of Transfer Learning have gained popularity. The first case pertains to Computer Vision (CV), wherein a model is initially trained on a vast labeled dataset with distinct categories known as ImageNet [31, 88, 85]. The model is then fine-tuned or re-trained to perform other tasks or classify objects in categories not present in ImageNet.

The second example is relevant in Natural Language Processing (NLP), where a model is pre-trained for tasks like Language Modeling [77, 78, 10] or Masked Language Modeling [32, 55, 64]. Subsequently, the pre-trained model is fine-tuned for several other tasks like sentiment analysis, question answering, and document classification. Further elaboration of this approach is presented in the following sections.

2.1.5 Natural Language Processing

Natural Language Processing (NLP) is a subfield of Computer Science that uses techniques from the area of Artificial Intelligence to automate well-defined tasks that require the understanding of written human languages. In his book about the topic, Goldberg [39] defines the area as *“the field of designing methods and algorithms that take as input or produce as output unstructured, natural language data”*.

NLP methods can be classified into two categories based on their expected output. The first category is Natural Language Understanding (NLU) methods, which deal with the meaning or structure of input text. Examples of NLU methods include document classification (e.g., categorizing news articles by topic or tweets by sentiment) and sequence labeling (e.g., named-entity recognition, where the task is to identify entities such as persons, organizations, or places within text).

The second is the category of Natural Language Generation (NLG) methods. These kinds of methods aim to solve tasks that not only require the understanding of human language, but also to generate new, probably different, text from the one that was passed as input. A common example is the task of Machine Translation, where the input is text in one human language and the expected output is the translation of it in another human language.

The recently popular models like GPT-3 [10] and ChatGPT¹ use these kind of approach to simulate a conversation on a chat by generating new text that answers an input text.

There are also tasks that have formulations for both methods, for example question answering and summarization. The NLU formulation of question answering seeks to, given a context text and a question about it, search for the boundaries where the question is completely answered. On the other hand, a NLG formulation seeks to generate the text that answers the question. In the case of summarization, a NLU formulation can be to mark the N most important phrases within the text, while a NLG formulation would be to generate a proper text that summarizes a larger one. Generating text is generally more challenging than understanding it [66]. In both cases, these tasks use the supervised approach from Machine Learning, so the way to train a model to accomplish a task require the use of examples texts accompanied by labels².

Challenges of NLP

Natural Language Processing (NLP) is a complex field that presents several challenges due to the intrinsic properties of language. Some of these properties are:

- Language is symbolic and discrete. This means that language is made by a discrete number of symbols. The basic elements are characters, which forms words. For example, the words “tiger”, “lion” and “panther”. These words, by itself, don’t have a meaning, since this meaning is assigned by us (by the society using the language) and its not embedded to the symbol. We know some relations of these words such as they are both animals and big felines, but that knowledge can not be inferred from the symbols or operations over them. While some properties can be inferred from certain word pairs (e.g., “Chile” and “Chilean”), this is not the case for every pair of words in the language.
- Language is compositional. We form words using characters, sentences using words, paragraphs from sentences and so on. Also, the meaning of these structures is not only the sum of its parts, the order is important to understand the goal of a sentence or a paragraph.
- Language is sparse. From the compositionality of language we can derive that the number of possible well-formed language structures is practically boundless. For instance, the enumeration of natural numbers using natural language exemplifies the notion of an enumerable infinity within language. This leads to a previously mentioned desirable quality of machine learning systems, the ability to generalize. Since it is not possible to learn all possible expressions in the language, a system that understands language has to be capable of learning the rules that govern language.

¹<https://openai.com/blog/chatgpt>

²In some cases, explicit labeling may not be present; instead, an intelligent extraction of labels from the data is performed to generate a synthetic training task. This approach is commonly referred to as self-supervised learning and will be elaborated upon in subsequent subsections.

Moreover, there are other challenges that are inherent to the fact that language is constructed by society. New words are continually introduced into the lexicon, while others fall into disuse. Additionally, words may have multiple meanings depending on the context in which they appear, known as polysemy. For instance, the word “bright” can have different meanings depending on the subject to which it refers; if we talk about a place, it implies light, while if we talk about a person, it implies intelligence.

These challenges, along with others, must be addressed when developing NLP systems. A comprehensive review of the challenges of NLP is provided in the book by Goldberg [39].

2.1.6 Representations of Text

A natural question that follows is how to represent the text examples to feed a machine learning system. While there was an evolution of the approaches through the time, the modern way is to use representations of words, also named as Word Embeddings.

Word Embeddings are a mathematical mapping from a word (a discrete symbol) to a continuous vector of dimensionality d . There are different types of word embeddings; the first ones were sparse vectors in which every word of a vocabulary corresponded to a dimension. Mikolov et al. [68] popularized Word2Vec, a technique to create learned dense vectors that do not depend on the size of the vocabulary and encode some of the semantics of the words. Other approaches were GloVe [73] and FastText [9].

One major limitation of these word representations is that they fail to account for the challenging characteristic of language known as polysemy. Specifically, they are fixed vectors, meaning that a word is represented identically, regardless of its context.

To overcome the issue of polysemy, contextual word representations were introduced. These representations not only used a fixed embedding layer, but also deep neural networks, to account for the complete context of a text in the calculation of a representation of a word. The first contextual representations used RNNs [75] as the neural network architecture and were then replaced in favor of Transformers [32].

2.2 Transformers

Over the past few years, the Transformer architecture proposed by Vaswani et al. [105] has gained considerable importance in multiple areas of machine learning. While originally designed as an architecture for solving sequence-to-sequence problems such as machine translation, it has now applications in other domains, including computer vision (e.g., Vision Transformer (ViT) by Dosovitskiy et al. [34]), speech recognition (e.g., Whisper by Radford et al. [80]), multimodal language-image models (e.g., CLIP by Radford et al. [79]), and reinforcement learning models (e.g., Decision Transformer by Chen et al. [19]), among others. In the following section, we will provide an overview of some of the Transformer-based models that have been used in natural language processing.

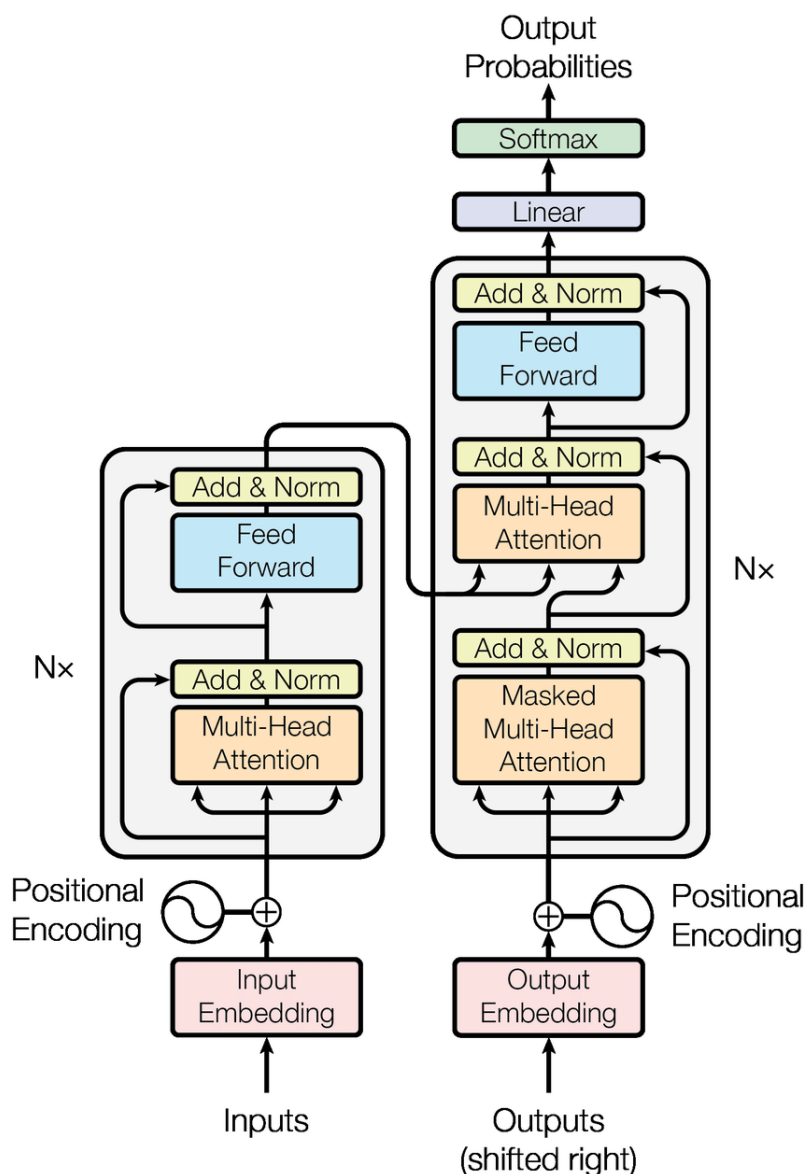


Figure 2.1: The Transformer architecture by Vaswani et al. [105].

2.2.1 Vanilla Transformers

Transformers [105] were first proposed as an Encoder-Decoder architecture that was commonly used in the context of sequence-to-sequence problems [99]. In Transformers, the recurrent architecture is replaced by attention mechanisms, particularly self-attention, which allows the prediction of a token to be dependent on every other token in the same sequence. Moreover, models based on this architecture are able to benefit from the inherent parallelism of the self-attention mechanism, overcoming the main disadvantage of RNNs. This key feature allows a Transformer model to take advantage of training on large collections of data using techniques of self-supervised learning. We will elaborate on these on the following subsections.

Figure 2.1 presents the Vanilla Transformer architecture as originally proposed by Vaswani

et al. [105]. The architecture consists of two primary blocks, which are repeated depending on the number of layers. The Transformer-encoder block is presented on the left side and is composed of two primary components: the multi-head attention, also known as self-attention, and a feedforward layer component. The Transformer-decoder block is presented on the right side and includes a masked multi-head attention, an encoder-decoder multi-head attention, and a feedforward layer component.

The attention mechanism employed in the original work is referred to as ‘‘Scaled Dot-Product Attention’’. Consider an input sequence of n words encoded as one-hot vectors, which is then passed through (multiplied by) an embedding matrix with dimensions (d_{vocab}, d_{model}) . We will denote the resulting input embeddings of size (n, d_{model}) as X . To implement the attention mechanism, three distinct weight matrices are utilized: $W^Q \in \mathbb{R}^{d_{model} \times d_q}$, $W^K \in \mathbb{R}^{d_{model} \times d_k}$, and $W^V \in \mathbb{R}^{d_{model} \times d_v}$. Specifically, the matrix $Q = XW^Q$ is referred to as queries, $K = XW^K$ is referred to as keys, and $V = XW^V$ is referred to as values.

Subsequently, the attention matrix is computed via:

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

In order to compute attention, it is necessary that the dimensionality of d_q be identical to that of d_k to perform the dot product operation between the queries and keys. In practice, Vaswani et al. [105] employ the same dimensions for d_q , d_k , and d_v .

Furthermore, instead of utilizing a single attention mechanism, Vaswani et al. [105] employ multiple attentions, each utilizing different matrices for queries, keys, and values. These attentions are referred to as ‘‘heads’’. The multi-head attention mechanism may be expressed as follows:

$$MultiHeadAttention = Concat(head_1, \dots, head_h)W^O$$

Where, $head_i = Attention(XW_i^Q, XW_i^K, XW_i^V)$ and $W^O \in \mathbb{R}^{hd_v \times d_{model}}$. It is worth mentioning that the employment of the W^O matrix facilitates the output to return to the original dimensions of the input X , thereby enabling the calculation of another multi-head attention if necessary. Such is the case in the Transformer-decoder, where the masked self-attention is followed by an encoder-decoder attention.

We have outlined the self-attention mechanism, which is of significant importance for the Transformer-encoders implemented in this work. The masked self-attention and encoder-decoder attention mechanisms utilized in the Transformer-decoder block are similar, but possess some subtle differences, which are detailed in the original paper by Vaswani et al. [105].

Finally, it is worth noting that the architecture includes an extra component in the input embeddings known as positional embeddings. In contrast to RNNs, where the information for a specific token is reliant on the tokens preceding or following it, self-attention in this case uses all the other tokens in the sequence to compute attention for a given token. As a

result, the order property is lost when using self-attention, which is crucial in language tasks where word order can significantly impact meaning. To overcome this problem, the authors appended positional embeddings to the model to provide information regarding the order of the words.

2.2.2 BERT

BERT [32], which stands for “Bidirectional Encoder Representations from Transformers”, is a Transformer encoder. This means it is constructed using only the encoder part of the Transformer. The self-supervised pre-training techniques used to train BERT were popularized by ULMFit [49], and involve creating a synthetic task using unlabeled data and then training a model on that task in a supervised manner.

There are examples of self-supervised learning in multiple areas. In NLP a way to achieve this is by using a language modeling task, where a model is trained to predict the most possible next word given a previous context. In CV, a way to achieve this is by deleting pixels or patches of pixels in an image, and then train a model to reconstruct them [45]. Similar ideas were applied to audio [80, 4], or even multimodal models [37].

Specifically, ULMFit is first trained for the language modeling task using vast amounts of text. This stage of the training is called pre-training. Then the model can be fine-tuned in a supervised way to solve specific tasks using labeled data. Instead of the language modeling task, BERT uses the masked language modeling (MLM) task, which consists of corrupting an input sequence by arbitrarily deleting some of the tokens and then training the model to reconstruct the original sequence. Figure 2.2 presents an example of the MLM task.

In addition to MLM, a second pre-training task was added to the training of BERT, called next sentence prediction (NSP). In this task, some sentences are changed from the original order and the model is asked to predict whether two sentences were consecutive or not. However, NSP has been found to be an easy task and is prone to overfitting. Consequently, newer variations of BERT, such as RoBERTa and ALBERT, do not use NSP.

Notably, BERT uses byte-pair encoding (BPE) subword tokenization [94], particularly the WordPiece [115] implementation, with a vocabulary of 30K tokens. This is different from previous word representations like Word2Vec and GloVe, which had representations for entire words. The use of BPE subword tokenization allows for the reuse of tokens to form different words, reducing the memory requirements in the embedding layer.

2.2.3 RoBERTa

RoBERTa [64], which is an abbreviation for “Robustly Optimized BERT Pretraing Approach” represents a re-implementation of the BERT model that has advanced the state-of-the-art in natural language processing (NLP) by implementing several modifications to the original model.

Specifically, RoBERTa excludes the Next Sentence Prediction (NSP) task, which has been

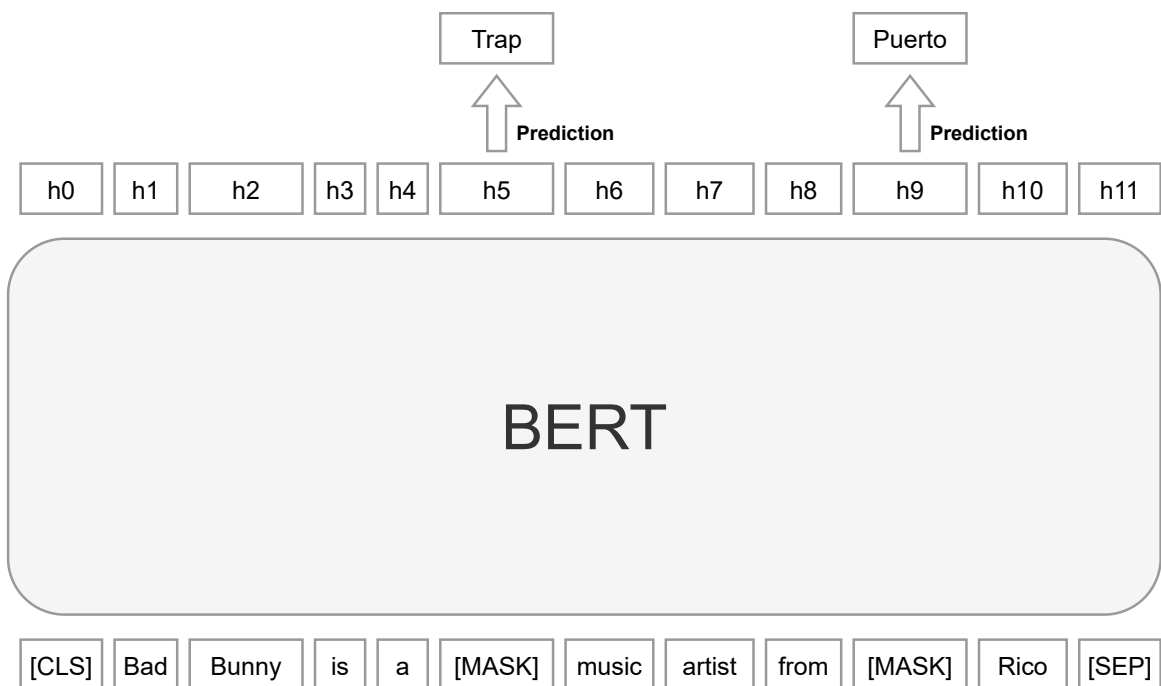


Figure 2.2: The masked language modeling (MLM) task used by BERT as pre-training task.

shown to have a negative effect on the downstream performance of the model, as demonstrated by an ablation experiment conducted by the authors. Additionally, RoBERTa incorporates a dynamic masking technique that provides different masks to the input text instead of using the same mask for every example. An example of the dynamic masking technique used in the MLM task is presented in Figure 2.3.

RoBERTa also uses a larger vocabulary of 50K byte-level BPE tokens, which is a slight modification made by Radford et al. [78] to the original BPE algorithm that, instead of unicode characters, uses bytes as the base subword unit to prevent having unknown tokens. Finally RoBERTa implements modifications to the training process, such as longer training

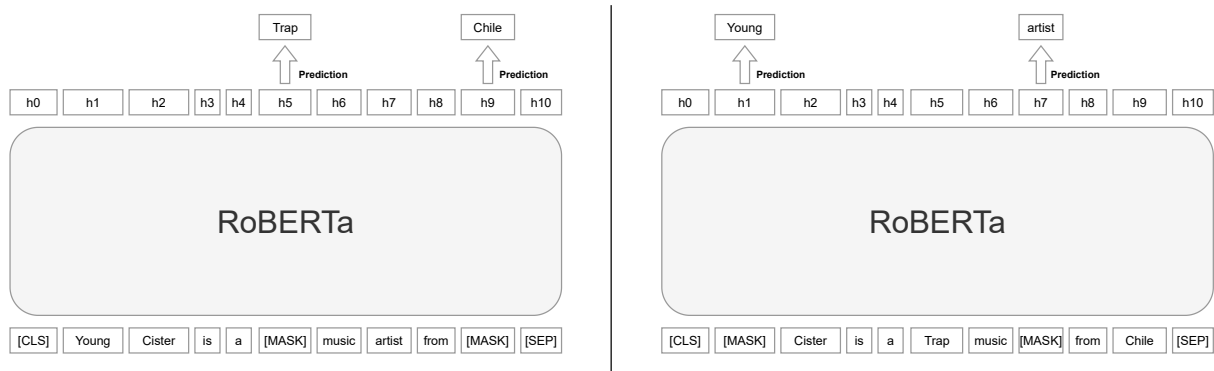


Figure 2.3: RoBERTa adopts a dynamic masking strategy, which allows for distinct masking patterns to be applied to a sentence during the pre-training phase.

times, more data, larger batch sizes, and longer sequences of text. These adjustments have resulted in significant improvements in RoBERTa’s performance compared to the BERT model, making it one of the leading NLP models currently available.

2.2.4 ALBERT

ALBERT [55] is a neural network architecture based on the well-known BERT model. ALBERT aims to optimize the memory usage and training speed of BERT by making two primary modifications to the original architecture.

Firstly, ALBERT decouples the size of the embedding layer from that of the hidden layer. Unlike BERT, which had matching sizes of the embedding and hidden layers, ALBERT’s design is based on the premise that the contextual embeddings (the Transformer layers) are more critical to the model’s efficacy than the non-contextual embeddings (the embedding layer).

To achieve this decoupling, ALBERT factorizes the embedding layer into two smaller matrices, so instead of directly projecting the one-hot vector to the dimension of the hidden layer H , it is projected to a smaller dimension E , followed by a projection from E to H . This reduces the number of parameters from $O(V \times H)$ to $O(V \times E + E \times H)$, with V being the size of the vocabulary. This reduction makes sense when the size of E is much lower than the size of H . As an example, the embedding and hidden layer size of the BERT-*base* model is 768, in the ALBERT-*base* model the embedding size is 128 and the hidden size 768. This modification makes the ALBERT-*base* parameters of the embedding layer only 17% of those in BERT-*base*.

The second modification to the architecture is the use of weight-tied parameters across layers. This means that, for a model of N layers, the model have just one set of parameters, instead of N sets. The architecture of BERT and ALBERT is compared in Figure 2.4, which provides an illustration of the parameter-sharing strategy employed in ALBERT. The fact that all layers have the same weights is a key element we use to propose models that are faster. This is further elaborated in Section 4.2.

In addition to these two architecture modifications, ALBERT also introduces a new pre-training task called sentence order prediction (SOP) that replaces the next sentence prediction (NSP) task used in BERT. In NSP, the task was to predict if two sentences were consecutive or not. In SOP, the task takes two consecutive sentences, and the same sentences with the order swapped and the model is asked to predict which one is the correct order.

Taken together, these modifications enable ALBERT to achieve results in NLP tasks comparable to BERT while using significantly fewer parameters. This improved computational efficiency allows for the training of larger models (with more layers or with the Transformer layer being larger), which can further improve accuracy on NLP tasks.

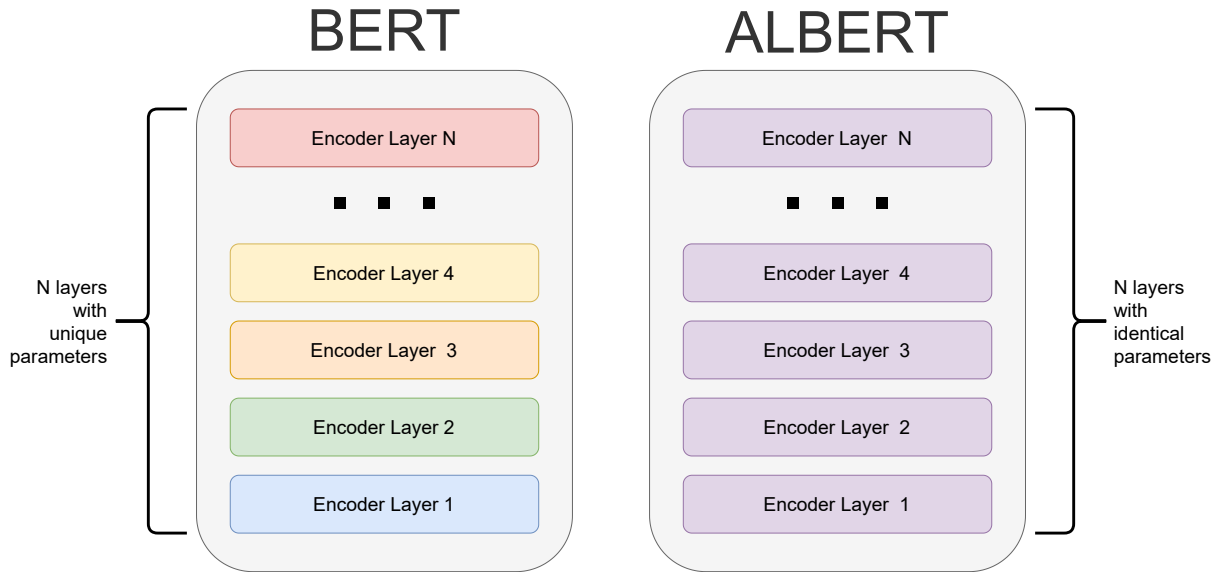


Figure 2.4: A design comparison of BERT and ALBERT, focusing on the parameter utilization strategy adopted by each model.

2.2.5 Multilingual and Monolingual Models

Multilingual models are models that are trained simultaneously using data from several languages. Examples of these models are mBERT [32, 76], a version of the BERT model trained in 104 languages and XLM-R [26] which is based on the RoBERTa architecture and was trained using corpora from 100 languages. Notably, some studies [114, 76] have shown how these multi-language learners set strong baselines for non-English tasks. Furthermore, for tasks that have training data in multiple languages, the performance when training with all the data can be better than the monolingual models trained on a specific language, potentially due to the ability to take advantage of a larger number of training examples [117].

However, these models require larger vocabularies to function effectively since they must represent a much larger range of possible texts, resulting in higher memory requirements that are undesirable for deployment settings. For instance, mBERT uses a vocabulary of 110K tokens, while XLM-R has a vocabulary of 250K tokens, and the recently presented XLM-V [62] employs a vocabulary of 1M tokens.

It is worth noting that multilingual models have also increase in size overtime, both in terms of parameters and required computational resources. This is demonstrated with respect to parameters, as illustrated in Figure 2.5, some examples are: mBART [65] (680M parameters), XGLM [63] (7.5B parameters), XLM-R XXL [42] (10.7B parameters), mT5 XXL [116] (13B parameters) and BLOOM [91] which is the largest, with a size comparable to GPT-3 (176B parameters).

In addition, several single language BERT models have been released. For instance, CamemBERT [67] and FlauBERT [56] for French, BERTje [29] and RobBERT [30] for Dutch,

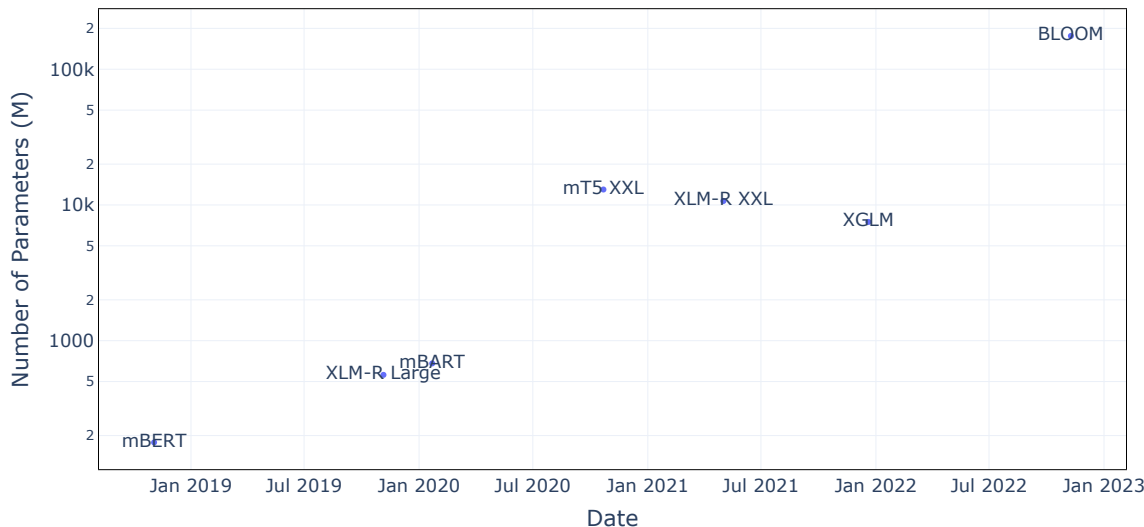


Figure 2.5: A summary of the number of parameters in various modern multilingual language models, from older to newer. The models included in this comparison are mBERT [32], XLM-R Large [26], mBART [65], mT5 XXL [116], XLM-R XXL [42], XGLM [63] and BLOOM [91].

FinBERT [106] for Finish, among others. Those models have generally shown better performance than multilingual ones, highlighting the importance of having language-specific models for language-specific tasks.

For the Spanish language, we can find BETO [18], which is a BERT *base* model, the not publicly available RigoBERTa [95], based on the DeBERTa [46] architecture and a family of RoBERTa models. RoBERTa-BNE [44], which is available in *base* and *large* versions was trained on the corpus crawled by the National Library of Spain. BERTIN [27] is a *base* sized model that was trained on the Spanish portion of mC4 [116]. RoBERTuito [74] is a RoBERTa *base* model trained exclusively on tweets. These models are better described in Section 3.2.

2.3 Model Compression and Acceleration

The notion that neural networks tend to be overparameterized is not novel. In fact, as early as thirty years ago, LeCun et al. [57] demonstrated that it is feasible to eliminate irrelevant weights in some neural networks without any deterioration in the model’s accuracy, resulting in increased speed and efficiency. More recently, Frankle and Carbin [36] established that for the majority of feed-forward neural networks, it is feasible to identify a subnetwork that can achieve comparable or even better accuracy. Follow-up studies by these authors demonstrated that the same is valid for various modern neural network architectures commonly employed in natural language processing [21] and computer vision [22]. These observations, in conjunction with the requirement for lighter and faster models that can be deployed in practical situations, such as real-time web applications and limited devices like mobile phones, have stimulated

numerous research pathways. In this section, we will briefly discuss some of the techniques that are currently employed.

2.3.1 Compression Techniques

Over the years, several techniques have been proposed to compress neural network models, with the aim of reducing their size and computational complexity. The most common ones are pruning [8], quantization [38], and knowledge distillation [47].

It is important to note that these methods are not mutually exclusive and are often combined to achieve better results. For instance, several pruning or quantization techniques are accompanied by a knowledge distillation phase [90, 20, 121].

Pruning

Pruning is a technique that aims to reduce the number of connections (weights) in a neural network by identifying and removing redundant connections, resulting in a lighter model in terms of parameters. Extensive research has been conducted in this area, as documented in the survey by Blalock et al. [8].

There are at least two types of pruning: unstructured and structured. Unstructured pruning involves the removal of individual weights in a neural network. In contrast, structured pruning involves the removal of more complex components, such as entire neurons, attention heads, or even layers [20, 35]. Although both methods can decrease the size of a model in terms of parameters, only structured pruning can reduce computational complexity and runtime of a model since current hardware and libraries are not optimized to speed up the sparse patterns of weights produced by unstructured pruning algorithms [8].

Quantization

Network quantization, on the other hand, compresses the original network by reducing the number of bits required to represent each weight. This results in a lighter model, with a lower memory footprint and faster inference times. Several quantization methods have been proposed for BERT, including TernaryBERT [121] and BinaryBERT [5], which were able to reduce the weight size to 2 and 1 bit (instead of the usual 32-bit floating-point parameter), respectively, while maintaining most of the original BERT performance.

Knowledge Distillation

Knowledge Distillation (KD) [47] is another popular technique for compressing neural networks. In this approach, a big and strong model, or an ensemble of models, called the teacher model, is used to transfer its knowledge to a lighter model, called the student model, by forcing the student to mimic the teacher. One of the advantages of this technique is that

the student model can have a different architecture than the teacher model, allowing for the construction of smaller and more efficient models. A comprehensive review of various knowledge distillation methods can be found in the survey by Gou et al. [41]. In the subsequent subsection, a more comprehensive explanation of KD is presented.

2.3.2 Knowledge Distillation

The technique of Knowledge Distillation (KD) aims to transfer the knowledge learned from a big and capable model, usually called the teacher model, say M_T , to a more restricted model, called the student model, say M_S . To achieve this objective, we train M_S to imitate M_T . There are multiple ways to imitate M_T [41], in this work we use the simple, yet powerful approach, of directly mimic the output of M_T given a input text.

Formally, we define the distillation objective as L_{KD} :

$$L_{KD} = L_O(M_T(x), M_S(x))$$

Where L_O is a loss function that works on the logits of M_T and M_S . The most common choices for this loss are the cross entropy loss, the KL-divergence loss and the mean-squared error loss. In the case of KL-divergence or cross-entropy loss is it a common practice to use soft-targets [47] instead of direct logits, which means to apply a softmax with temperature T (with $T \geq 1$) to $M_T(x)$ and $M_S(x)$ in order to produce a soft probability distribution over the classes.

Also, typically we use not only the output of M_T but also the gold labels from the training dataset. The complete loss, taking into account these labels can be expressed as follows:

$$L = \alpha L_{CE} + (1 - \alpha)L_{KD}$$

Where L_{CE} is the traditional cross-entropy loss against gold labels and $\alpha \in [0, 1]$ defines the weight of each loss.

In practice, KD has demonstrated efficacy as a technique for training smaller yet high-performing models in comparison to training a model from scratch using training data [104, 90, 107, 51]. According to Hinton et al. [47], the success of KD can be attributed to the fact that soft-targets encompass a significantly greater amount of valuable information that cannot be captured by a single hard target, and this information can be acquired by the student model. Notably, KD has been explored not only as a compression technique but also as a form of regularization [120, 119]. Additionally, KD offers the advantage of being applicable in scenarios where labeled data is scarce or nonexistent, allowing direct training based on the teacher model’s output [100].

An overview of the entire framework is shown in Figure 2.6.

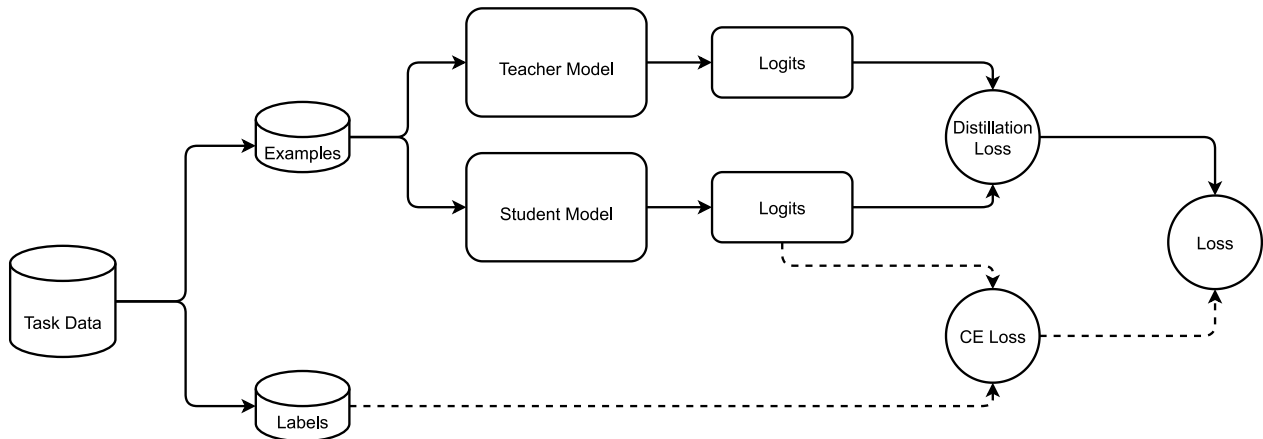


Figure 2.6: The figure provides a visual representation of the Knowledge Distillation framework applied in this work. In line with common practices, the framework includes both a distillation loss between the teacher and student models and a cross-entropy loss between the gold labels and the student’s predictions, as indicated by the dashed line.

2.3.3 Knowledge Distillation Techniques for BERT models

Knowledge Distillation (KD) has emerged as an effective technique for compressing BERT models by transferring knowledge from a larger BERT model to a smaller model. Several studies have explored different approaches to leverage the benefits of KD for BERT models.

For example, Tang et al. [100] use KD to transfer the knowledge from BERT to lighter RNNs. Turc et al. [104] propose pre-training compact BERT models and then using task-specific KD to achieve better results. Sanh et al. [90] introduce a task-agnostic scheme where KD is used on the pre-training task. Wang et al. [109] and Jiao et al. [51] proposed different methods exclusive for Transformers, to directly distill the knowledge from the self-attention layers of the teacher model to the student model.

Our work on faster models (Speedy Gonzales) is similar to Turc et al. [104] by proposing the use of compact Transformers but we use the ALBERT architecture instead of the BERT one. We also use the idea from Sanh et al. [90] of reusing the layers of a pre-trained model, instead of random initializing a new one. Differently from that work, which has to choose which layers to reuse, we only adjust the number of layers (and thus, the inference speed) since all the ALBERT layers are shared. Another difference with those two works is that in our work we skip pre-training (or KD on the pre-training task) and directly apply KD on the task-specific phase. The complete approach is explained in Section 4.2.

Chapter 3

Preliminaries

This chapter presents a comprehensive overview of the artifacts required for the effective evaluation and comparison of the models proposed in our study. The initial focus of the chapter is on the evaluation tasks and metrics employed to assess the models, as described in Section 3.1. Subsequently, a survey of the presently available Spanish models is provided, serving as a reference point for comparison with the proposed ALBETO models and as potential teacher models for the Speedy Gonzales models, which will be presented in Chapter 4.

3.1 Evaluation Tasks and Metrics

Generally, we divide the evaluated tasks in three different types:

1. **Text Classification:** these tasks involve one or two sentences as inputs and require the assignment of these inputs to a specific category, depending on the task at hand. The evaluated tasks associated with this type are Document Classification, Paraphrase Identification, and Natural Language Inference, which are evaluated using Accuracy as a metric.
2. **Sequence Tagging:** these are tasks where the goal is to identify elements within the input sentence that depend on the task. Specifically, the evaluated tasks associated with this type are Part of Speech Tagging and Named-Entity Recognition. These tasks are evaluated using the F1 Score as metric.
3. **Question Answering:** in this case, the task involves finding the boundaries of an answer, specifically the positions where an answer starts and ends, given an input context. For this type of task, the commonly reported evaluation metrics are the Exact Match and the F1 Score.

3.1.1 Evaluation Metrics

The evaluation metrics mentioned above can be defined as follows:

Accuracy is a metric that calculates the ratio of correct predictions to the total number of predictions made by a model. It can be expressed mathematically as:

$$\text{Accuracy} = \frac{\text{Correct Predictions}}{\text{All Predictions}}$$

To calculate the *F1 Score*, it is first necessary to define two concepts: *Precision* and *Recall*. In the context of binary classification, *Precision* is defined as the proportion of examples classified as positive that are truly positive. This can be expressed as:

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

Recall is defined as the proportion of truly positive examples that are correctly classified. This can be expressed as:

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

The *F1 Score* is then defined as the harmonic mean of *Precision* and *Recall*, given by:

$$\text{F1 Score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

In the case of Question Answering, the Exact Match metric compares the predicted answer string, p_s , with the correct answer string, c_s . The Exact Match for a single example is defined as:

$$\text{Exact Match}_{\text{single}} = \begin{cases} 1, & \text{if } p_s = c_s \\ 0, & \text{otherwise} \end{cases}$$

The Exact Match for a collection of pairs $(p_s, c_s) \in A$ is then defined as the average of the Exact Match for a single example, expressed as:

$$\text{Exact Match} = \sum_{(p_s, c_s) \in A} \frac{\text{Exact Match}_{\text{single}}(p_s, c_s)}{|A|}$$

In the subsequent subsections, a detailed description of the evaluated tasks and the datasets used to perform those tasks will be presented.

Dataset Name	Task Type	Number of Categories	Train Size	Validation Size	Test Size
MLDoc [93]	Text Classification	4	9458	1000	4000
PAWS-X [117]	Text Classification	2	49401	2000	2000
XNLI [25]	Text Classification	3	392702	2490	5010
POS [101]	Sequence Tagging	18	14305	1654	1721
NER [103]	Sequence Tagging	9	8324	1916	1518
MLQA [61]	Question Answering	-	81810	500	5253
SQAC [44]	Question Answering	-	15036	1864	1910
TAR / XQuAD [14, 3]	Question Answering	-	87595	10570	1190

Table 3.1: Details of the datasets used to evaluate our proposed models.

3.1.2 Evaluation Tasks

To evaluate our models, we use datasets that are in Spanish or contain subsets in Spanish. In Table 3.1 we present a summary of the datasets used for evaluation in this work. In addition, Annex A illustrates a representative instance of each task under consideration. Then, in the following subsections, we include a more detailed explanation of each task.

Document Classification

Document Classification is a critical process in natural language processing that involves assigning a document to a specific category based on its underlying semantic meaning. The primary objective of Document Classification is to facilitate efficient information retrieval and management.

To evaluate the effectiveness of the Document Classification task, this study employs the Spanish subset of MLDoc [93], a comprehensive multilingual dataset comprising documents in eight languages. MLDoc is derived from the widely used Reuters Corpus [60], and it comprises documents belonging to four distinct categories: Corporate/Industrial, Economics, Government/Social, and Markets.

Paraphrase Identification

The task of Paraphrase Identification aims to determine whether two given sentences possess the same underlying semantic meaning. This task is a crucial component of many natural language processing applications, including text summarization, information retrieval, and question answering systems.

To evaluate the effectiveness of our models in this task, we have utilized the Spanish subset of PAWS-X [117], which is a multilingual dataset. Specifically, PAWS-X can be regarded as a translation of the PAWS [122] dataset in six different languages. Notably, the training set of PAWS-X has been machine translated, while the validation and test sets were professionally translated by human experts. It is important to note that the quality of professional human translations makes the validation and test sets of the PAWS-X dataset particularly reliable

for assessing the performance of models in Paraphrase Identification.

Natural Language Inference

Natural Language Inference (NLI) is an important task in natural language processing that involves determining the logical relationship between two given sentences, namely a “premise” and an “hypothesis”. Specifically, the task requires inferring whether the premise entails, contradicts, or is neutral to the hypothesis. NLI is a critical component of various natural language applications such as question answering and dialogue systems.

In this study, we have employed the Spanish subset of XNLI [25], a multilingual dataset, to evaluate the effectiveness of our models in this task. The XNLI dataset is similar to PAWS-X in that it offers a machine-translated training set sourced from MultiNLI [112], while the validation and test sets have been professionally translated to 15 different languages, including Spanish. The use of professionally translated data in the validation and test sets ensures reliable performance evaluation of models in NLI.

Part of Speech Tagging

Part of Speech Tagging is a natural language processing task that aims to assign each word in a sentence its corresponding syntactic category. The syntactic categories are based on the grammatical function of the word and include, among others, nouns, verbs, adjectives, adverbs, and pronouns. In this task the dataset used was AnCora [101] which is included in the Spanish part of Universal Dependencies [28] Treebank.

Named Entity Recognition

Named Entity Recognition (NER) is a fundamental sequence labeling task in the field of NLP. It involves identifying and classifying named entities within a text according to their corresponding types. NER is essential in NLP as it enables computers to extract relevant information from unstructured text data, which can be used for a range of downstream applications.

Named entities are typically classified into categories such as people, places, organizations, or miscellaneous entities, and may consist of multiple words. This complexity requires the adoption of the BIO annotation scheme in NER datasets, where each word is labeled as either the beginning (B) of an entity, inside (I) an entity, or outside (O) of any entity.

The CoNLL-2002 shared task dataset [103] is commonly used in the NER research community and serves as a valuable resource for evaluating NER models. In this study, we utilized the Spanish subset of this dataset to evaluate our models.

Question Answering

Question Answering (QA) is a task in natural language processing that involves answering questions posed in natural language. There are various types of QA tasks, as reviewed previously in Section 2.1. In this study, we concentrate on Extractive Question Answering, which aims to extract a span of words from a given context text that fully answers a question posed about that context.

To evaluate our Extractive QA model, we considered four distinct datasets. MLQA [61] is a multilingual dataset that was created by translating English QA instances into six different languages, including Spanish. The dataset provides a validation and a test set for each language, as well as a machine-translated version of the SQuAD v1.1 [83] as a training set. We used the Spanish subsets of the MLQA dataset.

TAR [14] is another machine-translated dataset from SQuAD v1.1 to Spanish. XQuAD [3] provides a test set that was obtained from SQuAD v1.1 and professionally translated into 11 different languages, including Spanish. Following the setup proposed by [18], we combined the train and validation sets from TAR and the Spanish test set from XQuAD as a single evaluation dataset.

The last evaluated dataset is SQAC [44]. It is noteworthy that SQAC stands out among the evaluated datasets as the only one specifically designed for the Spanish language. This unique characteristic suggests that SQAC may offer a more comprehensive and valuable resource for addressing Spanish language-related challenges.

3.2 Overview of Pre-trained Models for Spanish

This section provides an overview of the main and most used pre-trained BERT-like models trained exclusively for the Spanish language.

3.2.1 BETO

BETO by Cañete et al. [18] is the first Transformer encoder that was pre-trained exclusively on Spanish corpora. It is a BERT-base sized model that is offered in two versions, namely *uncased* and *cased*, which depend on whether uppercase characters are included or not. These models have approximately 110 million parameters and each have a vocabulary of 31K BPE [94] subwords which was constructed using SentencePiece [54]. Additionally, 1K place-holder tokens are included for further specialization, resulting in a total vocabulary of 32K subword tokens. Both models were trained for 2M optimization steps on the SUC [16] dataset. Following the work by You et al. [118], the training process is divided into two phases, where the first 900K steps use a batch size of 2048 and maximum sequence length of 128, and the remaining 1.1M steps use a batch size of 256 and maximum sequence length of 512. Unlike BERT, BETO uses dynamic masking from RoBERTa and whole-word masking, which masks several subtokens when a word is split, as introduced in the second version of

BERT by Devlin et al. [32].

3.2.2 DistilBETO

DistilBETO, proposed by Donoso [33] is a compact version of a Transformer encoder, created by utilizing the pre-trained weights of BETO and applying the knowledge distillation technique presented by Sanh et al. [90] on DistilBERT. The training procedure for DistilBERT is a task-agnostic knowledge distillation process, which involves Masked Language Modeling as a distillation task, followed by fine-tuning on downstream tasks. The teacher model, initialization weights, and vocabulary used are from the *uncased* version of BETO. DistilBETO has 6 layers and 67M parameters.

3.2.3 RoBERTa-BNE

Gutiérrez-Fandiño et al. [44] proposed RoBERTa-BNE as two distinct RoBERTa [64] models, which belong to the MarIA collection of models, a collection that also incorporate Transformers models based on the GPT-2 [78] architecture. These models were trained on the Spanish language using the National Library of Spain (BNE) corpus, which was also created by the authors and currently stands as the largest Spanish corpus of its kind, with over 135 billion words and a size of 570GB of clean texts. The models, available in both *base* and *large* sizes, possess 125 million and 355 million parameters, respectively, while sharing a cased subword vocabulary of 50K BPE [94] tokens that were obtained from the BNE corpus. The pre-training procedure was performed with a single-epoch approach over the training dataset, adopting a maximum sequence length of 512 and a batch size of 2048 examples.

3.2.4 BERTIN

BERTIN, introduced by de la Rosa et al. [27] is a RoBERTa-base model trained on the Spanish portion of the mC4 [81] dataset. To reduce the size of the Spanish subset of the corpus from 1TB to 200GB, a novel technique called *perplexity sampling* was developed, aiming to maintain the same distribution of sentence complexities as in the original data. BERTIN is identical to the RoBERTa-BNE *base* model in terms of size, configuration, and vocabulary. The training process of BERTIN consisted of two phases, with a total of 250K steps executed. The first phase involved 230K steps with a batch size of 2048 and a maximum sequence length of 128, followed by a second phase of 20K steps with a batch size of 348 and a maximum sequence length of 512.

3.2.5 Models Not Included

In addition to the models analyzed in the preceding subsections, there exist additional Spanish pre-trained models that were not included in this study. Our aim was to incorporate

Model Name	Architecture	Size	Vocab Size	Vocab Types	Max Seq Length	Parameters	Domain	Availability	Reference
Included									
BETO	BERT	base	32K	uncased, cased	512	110M	General	Public	[18]
DistilBETO	DistilBERT	base	32K	uncased	512	67M	General	Public	[33]
RoBERTa-BNE <i>base</i>	RoBERTa	base	50K	cased	514	125M	General	Public	[44]
RoBERTa-BNE <i>large</i>	RoBERTa	large	50K	cased	514	355M	General	Public	[44]
BERTIN	RoBERTa	base	50K	cased	514	125M	General	Public	[27]
Not Included									
GPT-2-BNE <i>base</i>	GPT-2	base	50K	cased	512	124M	General	Public	[44]
GPT-2-BNE <i>large</i>	GPT-2	large	50K	cased	512	773M	General	Public	[44]
RigoBERTa	DeBERTa	base	50K	-	512	-	General	Private	[95]
RoBERTuito	RoBERTa	base	30K	uncased, cased, deaccented	130	109M	Social Media	Public	[74]
BSC-Bio	RoBERTa	base	50K	cased	514	125M	Biomedical	Public	[15]
RoBERTalex	RoBERTa	base	52K	cased	514	126M	Legal	Public	[43]
Longformer-BNE	Longformer	base	50K	cased	4098	149M	General	Public	¹

Table 3.2: Summary of pre-trained Transformer models for Spanish.

all publicly available Transformer-encoder based models trained on Spanish general domain corpora. However, several models for Spanish do not meet these criteria. Firstly, there are at least two models based on a Transformer-decoder architecture, namely the GPT-2-BNE models trained by the MarIA team [44] that are more suitable for NLG tasks. Secondly, there is RigoBERTa, a DeBERTa-based [46] model developed by Serrano et al. [95], which is not publicly available. Furthermore, there are various models trained on different domains. For example, RoBERTuito trained by Pérez et al. [74] are three RoBERTa models trained on social media data, specifically tweets, while BSC-Bio by Carrino et al. [15] consists of two RoBERTa models trained on biomedical corpora. Additionally, RoBERTalex by Gutiérrez-Fandiño et al. [43] is a RoBERTa model trained on legal domain corpora, among other models. Finally, a Longformer-BNE trained on the BNE corpus by the MarIA team [44] was omitted from this study. This Longformer is derived from the RoBERTa-BNE base model and was recently developed. A thorough comparison of this Longformer model with the other models considered in this work is left for future research.

A summary of all the models mentioned in this section is presented in Table 3.2.

¹<https://huggingface.co/PlanTL-GOB-ES/longformer-base-4096-bne-es>

Chapter 4

Proposed Spanish NLP Resources

This chapter provides a detailed description of the methodology used to develop models that are both lighter and faster for the Spanish language. The first section describe the methodology used to train ALBETO models (which are achieved through the use of the ALBERT architecture) in Section 4.1, and the approach used to develop faster models by employing knowledge distillation with these ALBETO models, which are referred to as Speedy Gonzales, in Section 4.2.

4.1 ALBETO: Lighter Models for Spanish

In this section, we present the elements of ALBETO, which are our ALBERT models designed specifically for the Spanish language. The components discussed include the pre-training data, model architectures, training procedures, as well as the experimental setup utilized for fine-tuning and evaluation of these models on downstream tasks.

The diagram depicted in Figure 4.1 presents an overview of the fundamental stages involved in the development, training, and evaluation of ALBETO models.

4.1.1 Data Collection and Preprocessing

In order to train a language model, it is crucial to utilize a corpus of texts, and the data employed for training all ALBETO models was obtained from the SUC¹ [16] dataset, which is the same dataset used for training the BETO model [18]. SUC is a collection of datasets similar to the corpus proposed by Cardellino [13], which includes additional sources and has a size of around 3 billion words, making it comparable to the corpus utilized for training the original BERT [32]. By utilizing the same dataset as BETO, our setup is similar to that of the original ALBERT, which utilized the same corpus as BERT. The SUC dataset consists of Spanish Wikis (from the April 2019 dump) and almost all of the Spanish portion of the

¹<https://github.com/josecannete/spanish-corpora>

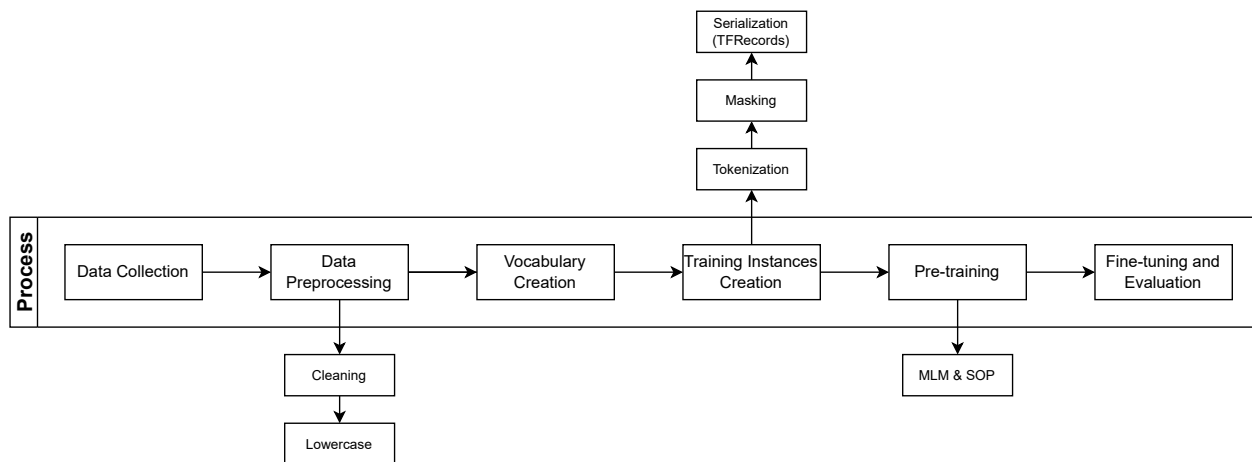


Figure 4.1: A broad overview of the process involved in the creation of ALBETO models. Sub-processes relevant to distinct stages are portrayed outside the main frame.

OPUS Project [102], making it a general domain corpus.

The preprocessing applied to the corpus for our models is identical to that used for BETO, which involves removing URLs and listings. Also, in cases where sentences contained multiple whitespace characters between words, they were replaced with a single whitespace character. Furthermore, all texts were converted to lowercase to accommodate the *uncased* nature of our models.

In addition to the preprocessing stage, a training instance creation step was undertaken to transform the preprocessed data into training examples for our pre-training process. This step involved tokenizing the data based on the model vocabulary and masking 15% of the tokens in each sentence. A dynamic masking approach was employed, utilizing 10 different masks for each sentence, following the whole-word masking (WWM) method introduced in the second version of BERT, which is also employed in BETO. Subsequently, the examples were converted into the efficient serialization TFRecord format developed by the Tensorflow [1] library.

The code² used for this last step and also the code used in pre-training was the same as the original ALBERT, which used Python and Tensorflow [1], with minor adjustments to accommodate Spanish-specific features such as the acceptance of accents and the exclusion of the NFKD normalization, which replaced certain characters utilized in Spanish, including the “ñ” with “n”.

Given the substantial size of the training data, a straightforward parallelization approach was employed. The entire corpus was partitioned into 100,000 smaller sections and processed using diverse Linux processes, which were executed in batches of parallel processes.

²<https://github.com/josecannete/ALBERT>

Model	Parameters	Layers	Hidden	Embedding
ALBETO <i>tiny</i>	5M	4	312	128
ALBETO <i>base</i>	12M	12	768	128
ALBETO <i>large</i>	18M	24	1024	128
ALBETO <i>xlarge</i>	59M	24	2048	128
ALBETO <i>xxlarge</i>	223M	12	4096	128

Table 4.1: The configurations of each ALBETO model trained in this work.

4.1.2 Model Architecture

For our language models, we employed the ALBERT [55] architecture, which is a variant of the BERT model with improved efficiency in terms of parameters. This is achieved through a factorization of the embedding layer and the use of the weight-tied strategy for parameter sharing. Additional details on this architecture are explained in Section 2.2.

We present five ALBETO models, namely, *tiny*³, *base*⁴, *large*⁵, *xlarge*⁶ and *xxlarge*⁷. Except for the *tiny* model, all of them have a similar size to those trained by Lan et al. [55]. The *tiny* model of ALBETO has the same size and configuration as the Chinese ALBERT *tiny*⁸, which has shown relatively good results for Chinese downstream tasks despite being smaller and faster. A comprehensive configuration of each model is provided in Table 4.1.

These five models share a vocabulary of 31K lowercase subword tokens that was constructed using SentencePiece [54], which is a library for fast unsupervised text tokenization that implements a variety of tokenization algorithms. In particular, following the common approach on language models, we used the BPE algorithm over the training dataset (SUC).

4.1.3 Pre-training Procedure

We conducted pre-training of our models using the Masked Language Modeling (MLM) and Sentence Order Prediction (SOP) self-supervised tasks described in Section 2.2. To improve the efficacy of our models, we employed the LAMB optimizer developed by You et al. [118], which has previously shown superior performance in training BERT models compared to other algorithms. We followed the authors’ recommendations for selecting hyperparameters. Each of the ALBETO models was trained using a single TPU v3-8, made available to us through the Google TRC⁹ program, for a limited duration. A maximum sequence length of 512 was used for pre-training, and the largest multiple of 64 that fit in the TPU memory was selected as the batch size. However, during the training process, we experienced divergence

³<https://huggingface.co/dccuchile/albert-tiny-spanish>

⁴<https://huggingface.co/dccuchile/albert-base-spanish>

⁵<https://huggingface.co/dccuchile/albert-large-spanish>

⁶<https://huggingface.co/dccuchile/albert-xlarge-spanish>

⁷<https://huggingface.co/dccuchile/albert-xxlarge-spanish>

⁸<https://github.com/ckiplab/ckip-transformers>

⁹<https://sites.research.google/trc/about/>

Model	Learning Rate	Batch Size	Warmup Ratio	Warmup Steps	Total Steps	Training Time (days)
ALBETO <i>tiny</i>	1.25e-3	2,048	1.25e-2	125,000	8,300,000	58.2
ALBETO <i>base</i>	8.83e-4	960	6.25e-3	53,333	3,650,000	70.4
ALBETO <i>large</i>	6.25e-4	512	3.12e-3	12,500	1,450,000	42.0
ALBETO <i>xlarge</i>	3.12e-4	128	7.81e-4	6,250	2,775,000	64.2
ALBETO <i>xxlarge</i>	3.12e-4	128	7.81e-4	3,125	1,650,000	70.7

Table 4.2: Training details of all ALBETO models, which were trained using a single TPU v3-8 each one.

in the loss on the *large* and *xlarge* models. This issue forced us to stop the training and restart it from an earlier checkpoint with a slightly lower learning rate. Consequently, the training time for each model varied. Detailed information about the training procedure for each model is presented in Table 4.2. We showcase the training loss progress of all models in Annex B.

4.1.4 Fine-tuning and Experimental Setup

To ensure a fair comparison, we fine-tuned our models, as well as the models reviewed in Section 3.2, using identical code. The fine-tuning was conducted using PyTorch [72] and HuggingFace’s Transformers library [113]. We followed the standard procedure proposed by Devlin et al. [32], with the only preprocessing step being tokenization based on the token vocabulary of each model, which converted words into subwords and added the special tokens [CLS], [PAD], and [SEP] (or equivalent tokens depending on the model) to each sentence. The maximum input sentence length was fixed at 512 tokens for all models, with longer sentences being truncated.

We conducted a hyperparameter search on BETO, DistilBETO, RoBERTa-BNE, BERTIN, ALBETO *tiny* and *base*, exploring combinations of batch size {16, 32, 64}, learning rate {1e-5, 2e-5, 3e-5, 5e-5}, and number of epochs 2, 3, 4. However, for the larger ALBETO models (*large*, *xlarge*, and *xxlarge*), we reduced the learning rates to {1e-6, 2e-6, 3e-6, 5e-6} to mitigate numerical instability issues during training.

These fine-tuning procedures were performed on one to two NVIDIA RTX 3090 GPUs, depending on the model and task. To fine-tune the largest models on QA, we utilized two NVIDIA A100 GPUs from the Patagón supercomputer [71]. We used gradient accumulation in situations where the GPU memory was insufficient to reach the target batch size.

Finally, we selected the models that produced the best results on the development set, which we discuss in more detail in Chapter 5.

4.2 Speedy Gonzales: Faster Models for Spanish

In pursuit of our goal to have efficient models for Spanish in various tasks, we employ the method of Knowledge Distillation. In this section our approach to use KD will be further

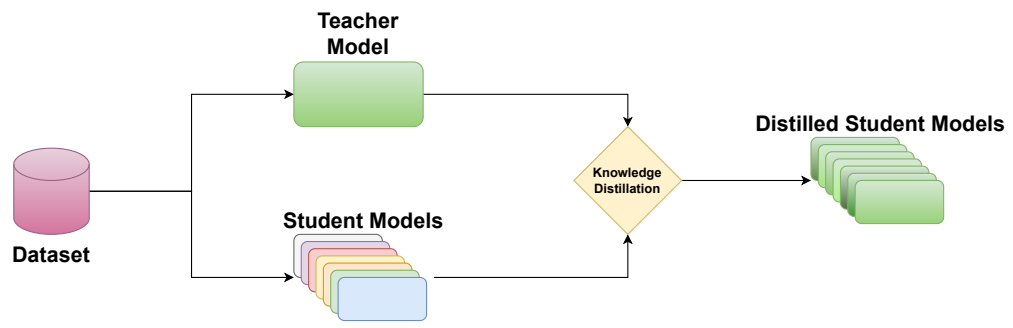
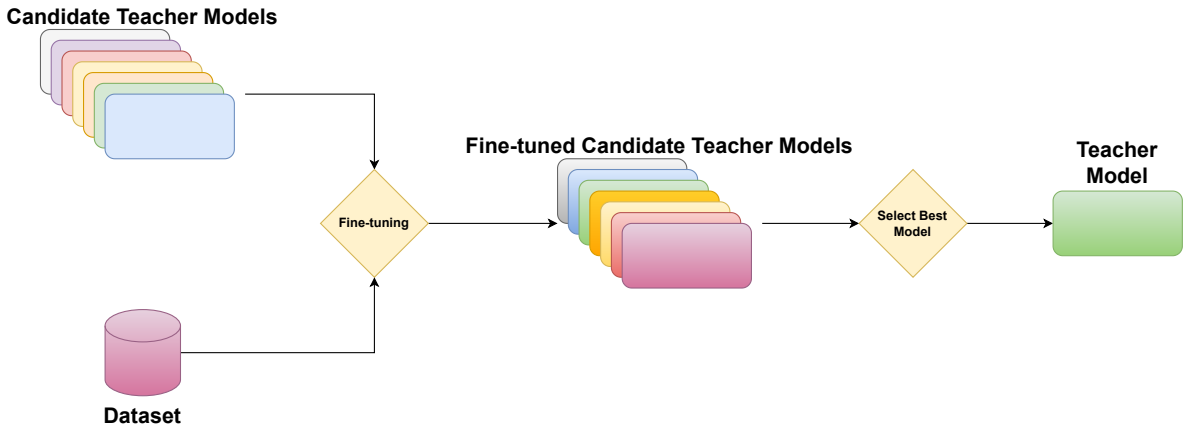


Figure 4.2: The two stage approach of task-specific knowledge distillation used in this work. (Up) The first stage involves fine-tuning a set of candidate models on a specific dataset, followed by the selection of the best-performing model as the teacher model for that dataset. (Down) The second stage employs the selected teacher model to train a set of student models using knowledge distillation.

elaborated.

4.2.1 Approach

Our approach is straightforward. In contrast to the methodology used by Sanh et al. [90] in DistilBERT and Donoso [33] in DistilBETO, which employs a task-agnostic approach involving knowledge distillation (KD) on a self-supervised task such as masked language modeling followed by standard fine-tuning, our methodology utilizes task-specific KD. Specifically, we employ KD to distill the knowledge of teacher models that have already been fine-tuned in our target tasks. Unlike the previous method that utilizes KD on the masked language modeling task with a substantial amount of data, our approach solely relies on a fine-tuning dataset, typically of smaller scale, thereby reducing the overall training time.

To implement our approach, we utilize a two-stage process, which is illustrated in Figure 4.2. In the first stage, we fine-tune a set of candidate teacher models for a given task. We then select the best candidate teacher model as the model with the minimum validation loss

among all candidates, to serve as the teacher model for that task. In the second stage, we employ KD using the selected teacher model and a set of student models that vary in size or computational complexity. This process is repeated for each of the tasks of interest, which are those specified in Section 3.1.

It is noteworthy that our approach possesses the advantage of reutilizing models that have already been fine-tuned for tasks relevant to our interests. The widespread availability of such models that have been fine-tuned for diverse tasks owes much to the open-source community and platforms like the HuggingFace Model Hub¹⁰, which presently hosts over 150,000 models that have been trained across a broad spectrum of domains and tasks by individuals and organizations.

In our case, the first stage of our approach corresponds to the evaluation stage of the ALBETO models presented in the preceding section. The ability to utilize such fine-tuned models in our approach offers the significant advantage of reducing the overall computational cost and time required for model training, while enhancing the overall efficiency and effectiveness of our methodology.

Therefore, the set of candidate teacher models utilized in our work consists of Transformer-encoder models that were pre-trained exclusively for Spanish in general domain corpora, and which are publicly available, as described in Section 3.2, and our recently introduced ALBETO models, which are described in Section 4.1. In the subsequent subsection, we describe the set of student models employed in our approach.

4.2.2 Student Models

For our student models we employed the ALBERT [55] architecture. Specifically, we utilized some of the previously introduced ALBETO models, which are characterized by their smaller model size. In particular, we considered the ALBETO *tiny* model, which is the lightest and fastest among all ALBETO models. Additionally, inspired by Sanh et al. [90], we developed models with fewer layers to achieve faster inference times while maintaining a configuration similar to the ALBETO *base* model, except for the number of layers. These lighter ALBERT models were initialized with the weights of the ALBETO *base* model. From another perspective, this approach can be seen as a form of layerwise pruning of the ALBETO *base* model, similarly to the technique introduced by Fan et al. [35], where entire layers of the model are pruned after training to facilitate faster inference.

In order to obtain a broad range of models with varying inference speeds, we created five models based on the ALBETO *base* model, ranging from 2 to 10 layers, with a difference of 2 layers between each of them. These models are identified in the tables as ALBETO *base-n*, where n represents the number of layers in the model.

¹⁰<https://huggingface.co/models>

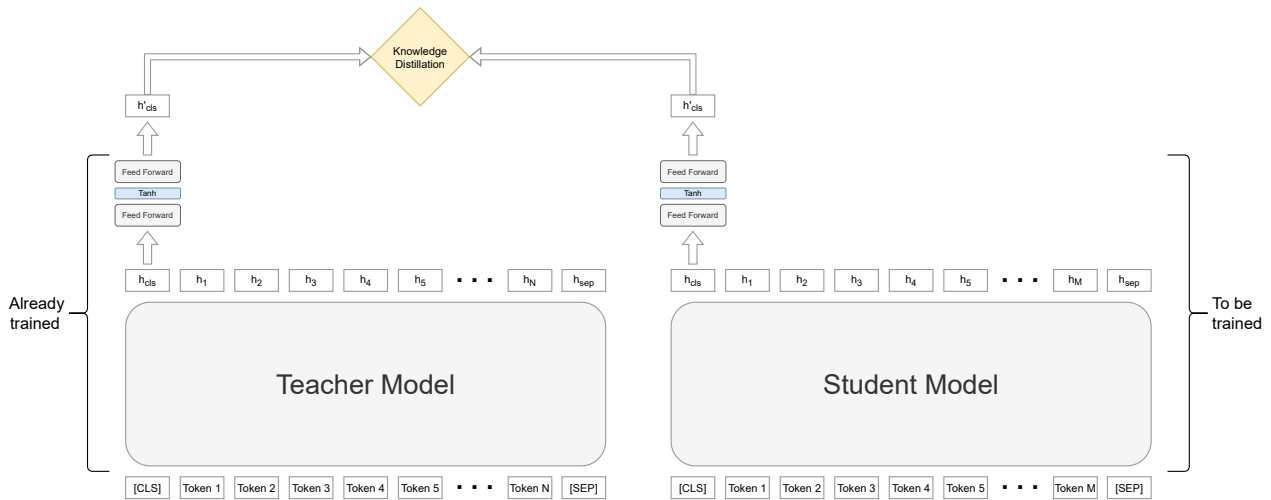


Figure 4.3: Implementation of knowledge distillation for text classifications tasks that use a single sentence as input.

4.2.3 Applying Knowledge Distillation to Different NLP Tasks

The procedure for fine-tuning BERT and Transformer encoders for specific tasks varies depending on the task’s characteristics. Similarly, the implementation of knowledge distillation varies across different tasks.

This subsection details the application of knowledge distillation to the three distinct task categories explored in our research: text classification, sequence tagging, and question answering.

Figures 4.3 and 4.4 serve to illustrate the implementation of knowledge distillation for text classification tasks involving a singular sentence (e.g. MLDoc) and two sentences (e.g. PAWS-X and XNLI), respectively. This represents a relatively straightforward scenario, as knowledge distillation is applied only to the last representation of the [CLS] token (or an equivalent, as determined by the model). The only discrepancy between the one and two sentence cases lies in the initial encoding, where the special [SEP] token is inserted between the tokenized sentences in the latter instance.

Figure 4.5 depicts the implementation of knowledge distillation for sequence tagging tasks such as Named Entity Recognition (NER) and Part-of-Speech (POS) tagging. We consider a general scenario where the vocabularies of the teacher and student models are distinct. Consequently, tokenization of the same sentence could vary, and hence, the representations of the two models may not align for straightforward application of knowledge distillation. To address this issue, an alignment strategy is employed, where KD is applied between the representations of the first token of each word. This approach is analogous to the fine-tuning technique utilized by Devlin et al. [32] in sequence tagging, where the first token of each word is employed for word classification. To achieve this alignment, the first token of each word is marked in the tokenization step, represented by the blue color in the figure. After obtaining the final representations of all tokens, a mask is employed to filter out tokens that are not the first token of a word. This process is applied to both teacher and student models, and

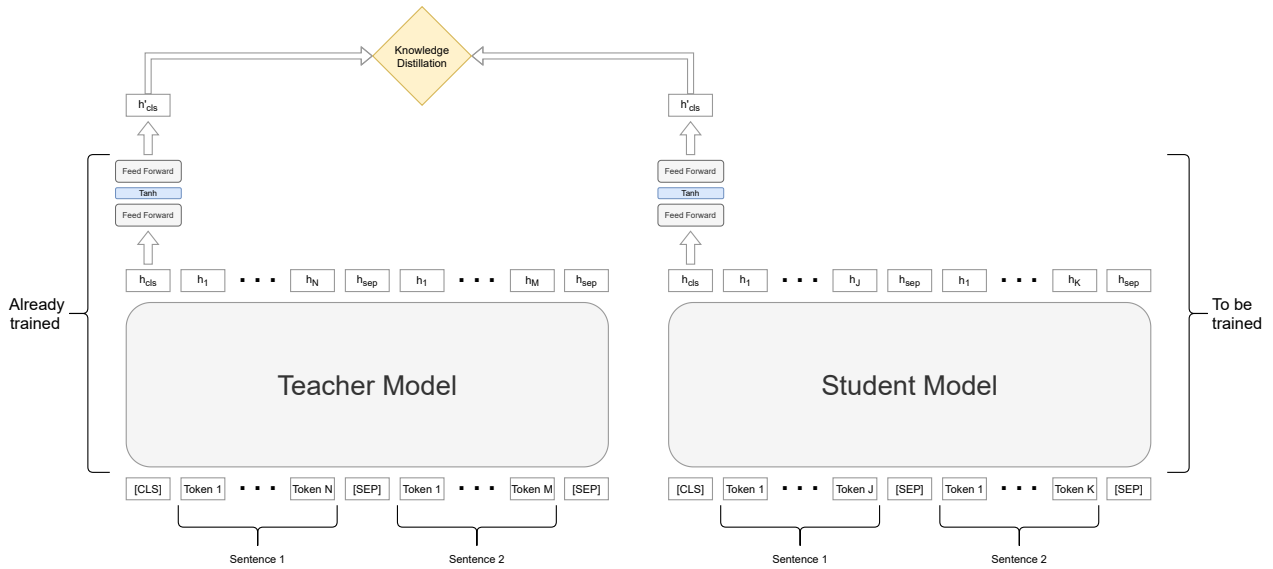


Figure 4.4: Implementation of knowledge distillation for text classifications tasks that use two sentences as input.

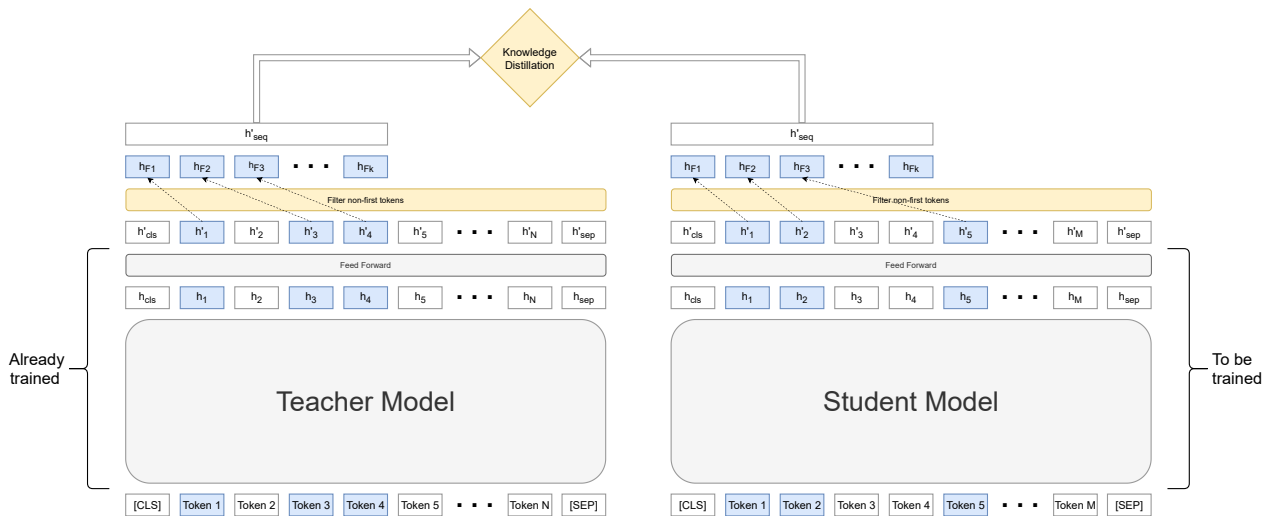


Figure 4.5: Implementation of knowledge distillation for sequence tagging tasks. The tokens marked with the blue color represents the property of being the first token of a word.

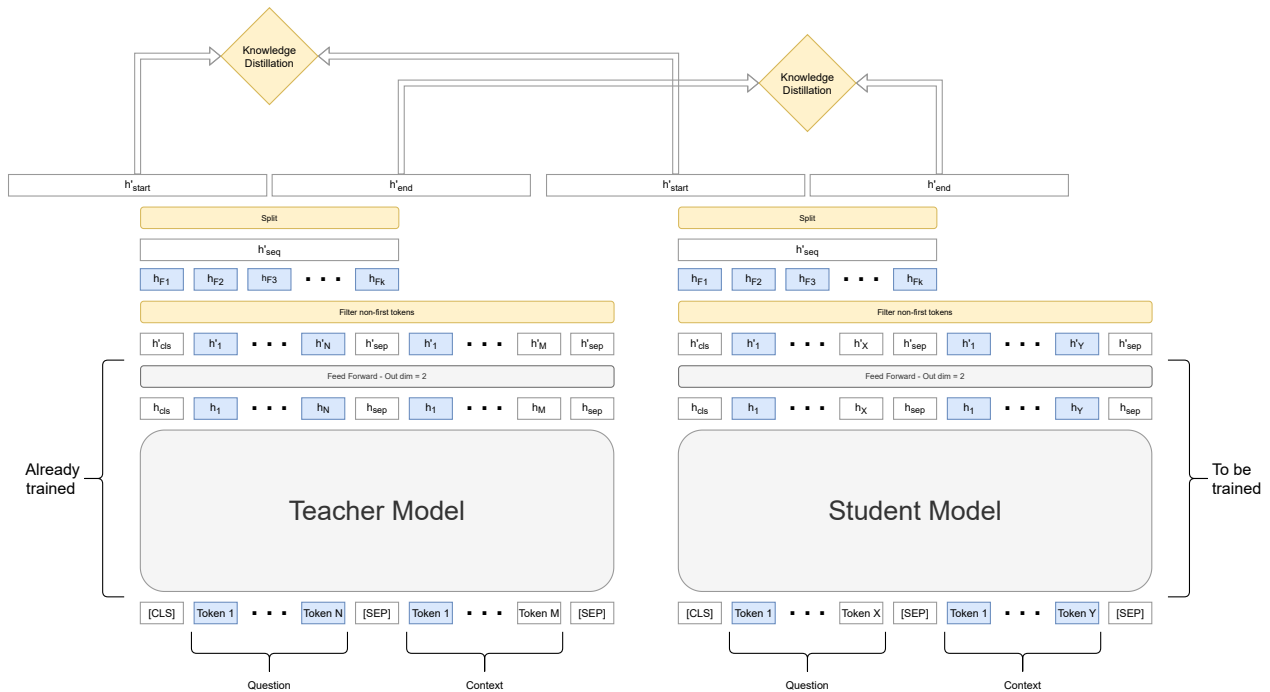


Figure 4.6: Implementation of knowledge distillation for question answering datasets. The tokens marked with the blue color represents the property of being the first token of a word.

the resulting aligned representations are utilized for KD.

Figure 4.6 depicts the knowledge distillation approach for question answering datasets such as MLQA, SQAC, and TAR-XQuAD. This scenario is akin to the previous sequence tagging case, wherein it is necessary to align the representations of the teacher and student models that may differ due to the use of different vocabularies. The alignment strategy is the same as before. Subsequently, the representations of each first token, which are of dimension 2, are split to obtain the probability of a token being the start or end of an answer. Consequently, we use KD twice, once for the representations of the start of an answer and once for the representations of the end of an answer. In this case the loss optimized during training is just the average of the two KD losses.

It is worth noting that these last two cases are the general scenarios where the vocabularies of the teacher and student models may differ. The cases where the vocabularies are the same are relatively straightforward and are illustrated in Annex D.

4.2.4 Additional Implementation Details and Experimental Setup

Code and Libraries

All our code uses Python and PyTorch [72] as the machine learning framework and is publicly available on GitHub¹¹.

¹¹<https://github.com/dccuchile/speedy-gonzales>

Measuring Inference Metrics

The evaluation of the inference speed of the proposed models, which is presented in Section 5.2, is performed through the utilization of the Multiply-Accumulate (MACs) metric, which provides a hardware-agnostic evaluation and is thus considered to be a more robust evaluation criterion. This measurement is conducted using the THOP¹² library, which operates on PyTorch models, to accurately measure MACs. In addition, to provide a more intuitive understanding of the models’ performance, actual inference speeds on commonly used hardware configurations are also reported in Section 5.3.

Improving Training Speed on Knowledge Distillation Experiments by Caching Teacher Outputs

Typically, the expense in a knowledge distillation experiment lies in employing a larger teacher model to guide the training of the smaller student model. Despite the teacher model being frozen (i.e., the model weights are not updated), it still accounts for the majority of the compute used during training.

To enhance this process, a teacher cache is introduced. Since the teacher model is static, the model’s predictions remain constant when an input is repeatedly passed through it. Therefore, the teacher model’s predictions can be computed once during the initial epoch and reused in subsequent epochs. In our experiments, where the number of training epochs is fixed at 50, this approach significantly enhances the experimentation speed and reduces training costs.

In practice, the implementation is accomplished by fixing the batch order passed to the models (i.e., the data is not shuffled between epochs) and utilizing Joblib¹³ as a disk-caching tool.

Annex E further emphasizes on the theoretical and practical impact of the implemented cache.

On the Importance of Knowledge Distillation

In Section 2.3, we elucidated the empirical benefits of employing Knowledge Distillation (KD) in prior studies. However, it remains uncertain whether these advantages persist when applied to our lighter and faster models. To address this inquiry, we conduct an ablation test, comparing the two scenarios, as outlined in Annex F.

¹²<https://github.com/Lyken17/pytorch-OpCounter>

¹³<https://joblib.readthedocs.io/en/latest/>

Selecting Knowledge Distillation Hyperparameters

Regarding KD, as formally introduced in Section 2.3, we conducted initial experiments utilizing three distinct loss functions: mean-squared error loss, cross-entropy loss, and KL-divergence loss. We varied the parameters α and T across these losses using Optuna [2]. The outcomes of these experiments revealed that the optimal settings were $\alpha = 0$ and $T = 1$. Although all three losses yielded satisfactory outcomes with this configuration, KL-divergence produced marginally superior results. Consequently, we proceeded to perform subsequent experiments using KL-divergence as loss function.

Preprocessing

For both stages of our approach, the only preprocessing applied was tokenization of the input texts according to the subword vocabulary of every model. In this tokenization step was also included the process of identifying the first tokens of each word for the case in which it was necessary (sequence tagging).

Experimental Setting

For the first stage, which is fine-tuning of the candidate teacher models, the experimental setup is detailed in Section 4.1. Also, it is important to note that the selected teacher model are not always the best models reported in the fine-tuning section of the tables in Chapter 5. These teacher models were selected based on the lowest loss on the *validation* set while those tables report results on the *test* set. While the expectation is that results obtained from *validation* correlates to *test*, it is not always the case. The list of selected models can be found in Annex C.

For the second stage, the implementations of KD for text classification and sequence tagging are as described in the previous section. However, for question answering, we adopted the implementation outlined in Annex D since the chosen teacher models had the same vocabulary as the student models.

For the experiments on this second stage we conducted a grid search using the hyperparameters: learning rate = $\{5e-5, 1e-4\}$, batch sizes = $\{16, 32, 64\}$ and epochs = 50. We also implemented early stopping with a tolerance of 10 epochs without improving.

The results of our experiments with task-specific KD are presented in Chapter 5. It is important to note that, in the same way we selected models to be teacher models, the best distilled models were selected based on the best results on the *validation* set among the grid search experiments, while the results reported in Chapter 5 refer to the *test* set of each dataset. These models are also the ones publicly available on the HuggingFace Hub.

Chapter 5

Results and Discussion

This chapter presents and discusses the main findings of the study. Initially, the results are presented with respect to the overall performance of the models in the evaluated tasks under two settings, namely fine-tuning and task-specific knowledge distillation. Subsequently, the efficiency of the models is discussed in terms of their size (number of parameters) and inference speed (MACs per inference).

5.1 Task Performance

Tables 5.1, 5.2 and 5.3 present the results of each model across all evaluated tasks. A general observation is that there are two distinct behaviors among the tasks. Firstly, there is minimal variation in performance between smaller and larger models in certain tasks, as evidenced by the comparable high scores achieved by all models in the MLDoc and POS tasks. It is hypothesized that these tasks are relatively simple, and as a result, the utilization of larger models results in overparameterization.

Secondly, there are tasks where there is a notable difference in performance between smaller and bigger models. This is evident in tasks such as Paraphrase Identification (PAWS-X), Natural Language Inference (XNLI), Named Entity Recognition (NER) and Question Answering (MLQA, SQAC, TAR/XQuAD), where the larger models tend to outperform the smaller models. This suggests that these tasks are more complex and require a greater model capacity. Overall, the results of this evaluation demonstrate the importance of considering the appropriate model size for a given task, as overparameterization can lead to suboptimal inference performance.

The results of our experiments on text classification tasks are detailed in Table 5.1. Firstly, on the fine-tuning setting we observed that the best models are RoBERTa-BNE *large* or ALBETO *xxlarge*. In PAWS-X, RoBERTa-BNE *large* outperforms every other model, followed by a small difference by the *base* version of RoBERTa-BNE, ALBETO *xxlarge* and the *cased* version of BETO. In XNLI, the best model in this setting is ALBETO *xxlarge*, followed by BETO *cased* and the *xlarge* size of ALBETO, with a difference of 0.4 and 0.7

Model	MLDoc	PAWS-X	XNLI
Fine-tuning			
BETO <i>uncased</i>	96.38	84.25	77.76
BETO <i>cased</i>	96.65	89.80	81.98
DistilBETO	96.35	75.80	76.59
ALBETO <i>tiny</i>	95.82	80.20	73.43
ALBETO <i>base</i>	96.07	87.95	79.88
ALBETO <i>large</i>	92.22	86.05	78.94
ALBETO <i>xlarge</i>	95.70	89.05	81.68
ALBETO <i>xxlarge</i>	96.85	89.85	82.42
BERTIN	96.47	88.65	80.50
RoBERTa BNE <i>base</i>	96.82	89.90	81.12
RoBERTa BNE <i>large</i>	97.00	90.00	51.62
Task-specific Knowledge Distillation			
ALBETO <i>tiny</i>	96.40	85.05	75.99
ALBETO <i>base-2</i>	96.20	76.75	73.65
ALBETO <i>base-4</i>	96.35	86.40	78.68
ALBETO <i>base-6</i>	96.40	88.45	81.66
ALBETO <i>base-8</i>	96.70	89.75	82.55
ALBETO <i>base-10</i>	96.88	89.95	82.26

Table 5.1: Models evaluated on sentence or two sentences classification tasks, results are measured using accuracy on the test set of each dataset.

Model	POS	NER
Fine-tuning		
BETO <i>uncased</i>	97.81	80.85
BETO <i>cased</i>	98.95	87.14
DistilBETO	97.67	78.13
ALBETO <i>tiny</i>	97.34	75.42
ALBETO <i>base</i>	98.21	82.89
ALBETO <i>large</i>	97.98	82.36
ALBETO <i>xlarge</i>	98.43	83.06
ALBETO <i>xxlarge</i>	98.43	83.06
BERTIN	99.02	85.66
RoBERTa BNE <i>base</i>	99.00	86.80
RoBERTa BNE <i>large</i>	61.83	21.47
Task-specific Knowledge Distillation		
ALBETO <i>tiny</i>	97.36	72.51
ALBETO <i>base-2</i>	97.17	69.69
ALBETO <i>base-4</i>	97.60	74.58
ALBETO <i>base-6</i>	97.82	78.41
ALBETO <i>base-8</i>	97.96	80.23
ALBETO <i>base-10</i>	98.00	81.10

Table 5.2: Models evaluated on sequence tagging tasks, results are measured using the F1 Score on the test set of each dataset.

percentual difference (pd) respectively.

Secondly, in the task-specific KD setting, we observed that models with a depth of 8 or more layers exhibit performance comparable to the best larger models, while also demonstrating significant improvements in inference time. Specifically, for the XNLI dataset, we found that the ALBETO *base-8* model outperforms all other models evaluated in our study.

Table 5.2 presents the results obtained on sequence tagging tasks. The Part-of-Speech (POS) tagging task exhibits a fairly consistent performance across most models, except for RoBERTa BNE *large*. The top-performing models in this task include BERTIN, RoBERTa-BNE *base*, and BETO *cased*, with only slight variations in their performances.

On NER we observe a significant difference between our lighter models (ALBETO), our faster models (those trained in the task-specific KD setting) and the *cased* models (BETO, BERTIN, RoBERTa-BNE), especially with BETO *cased*, which was the best model on the task. Furthermore, we observe a difference of almost 4.1 pd between ALBETO *xxlarge*, and BETO *cased*, even though ALBETO *xxlarge* is one of the largest models in the fine-tuning setting. Additionally, we find a difference of almost 6.3 pd between the *cased* and *uncased* versions of BETO. Based on these observations, we posit that the difference in performance between *cased* and *uncased* models can be attributed to the additional hints provided by capitalization for solving the NER task. Specifically, the names of persons, organizations, and places typically begin with a capital letter. Furthermore, our results from models trained

Model	MLQA	SQAC	TAR, XQuAD
Fine-tuning			
BETO <i>uncased</i>	64.12 / 40.83	72.22 / 53.45	74.81 / 54.62
BETO <i>cased</i>	67.65 / 43.38	78.65 / 60.94	77.81 / 56.97
DistilBETO	57.97 / 35.50	64.41 / 45.34	66.97 / 46.55
ALBETO <i>tiny</i>	51.84 / 28.28	59.28 / 39.16	66.43 / 45.71
ALBETO <i>base</i>	66.12 / 41.10	77.71 / 59.84	77.18 / 57.05
ALBETO <i>large</i>	65.56 / 40.98	76.36 / 56.54	76.72 / 56.21
ALBETO <i>xlarge</i>	68.26 / 43.76	78.64 / 59.26	80.15 / 59.66
ALBETO <i>xxlarge</i>	70.17 / 45.99	81.49 / 62.67	79.13 / 58.40
BERTIN	66.06 / 42.16	78.42 / 60.05	77.05 / 57.14
RoBERTa BNE <i>base</i>	67.31 / 44.50	80.53 / 62.72	77.16 / 55.46
RoBERTa BNE <i>large</i>	67.69 / 44.88	80.41 / 62.14	77.34 / 56.97
Task-specific Knowledge Distillation			
ALBETO <i>tiny</i>	54.17 / 32.22	63.03 / 43.35	67.47 / 46.13
ALBETO <i>base-2</i>	48.62 / 26.17	58.40 / 39.00	63.41 / 42.35
ALBETO <i>base-4</i>	62.19 / 38.28	71.41 / 52.87	73.31 / 52.43
ALBETO <i>base-6</i>	66.35 / 42.01	76.99 / 59.00	75.59 / 54.95
ALBETO <i>base-8</i>	67.39 / 42.94	77.79 / 59.63	77.89 / 56.72
ALBETO <i>base-10</i>	68.29 / 44.29	79.89 / 62.04	78.21 / 56.21

Table 5.3: Models evaluated on question answering datasets, results are noted as F1 Score / Exact Match on the test set of each dataset.

using knowledge distillation (KD) suggest that this hint is not easily replicable in an uncased model.

The performance on Question Answering datasets, as indicated in Table 5.3, follows a pattern similar to that observed in text classification tasks. The larger models, specially ALBETO *xxlarge* and *xlarge*, exhibit higher performance, while our proposed models on the task-specific KD setting featuring 8 or more layers present results similar to those of the base-sized models.

It should be noted that some models performed significantly worse than the others. Specifically, the utilization of RoBERTa-BNE *large* on XNLI, POS, and NER tasks produced sub-par results. This deviation from the performance of the same model on other tasks, as well as the results reported by Gutiérrez-Fandiño et al. [44], suggests that RoBERTa-BNE *large* may be particularly sensitive to hyperparameter selection and may benefit from additional hyperparameter tuning.

5.2 Model Efficiency and Inference Speed

In Figure 5.1 we aggregated the results by type of task and visualized them taking into account the size of the model (the number of parameters) and the inference speed (the MACs of a



Figure 5.1: Average of performance in the different tasks. The size of the points represents the size of the model (the number of parameters).

single inference). When compared at equal inference speed, our ALBETO fine-tuned models present similar performance to other models, while exhibiting improved parameter-efficiency. This is the case when we compare our ALBETO *base* model with BETO models, BERTIN or RoBERTa-BNE *base*. The results of our model are relatively close to those models, despite having significantly fewer parameters. Another interesting observation is that our ALBETO *base*, which is an uncased model, outperforms BETO *uncased* in almost every evaluated task, but is outperformed by the *cased* version of BETO, which may imply that the use of a vocabulary with upper and lowercases is generally a better design choice.

It is also important to notice that, while our larger ALBETO models (*large*, *xlarge* and *xxlarge*) are light in terms of parameters because they share parameters across layers, their inference speed is not necessarily faster. This can be attributed to the larger size of the Transformer layers, with larger hidden dimensions, and the increased number of layers, with the *large* and *xlarge* models consisting of 24 layers, and the *xxlarge* model containing 12 layers. As a result, greater computational resources are required for inference.

In the task-specific KD setting, a general progression in performance of our proposed models is demonstrated as the number of layers increases. For these models, a clear trade-off between task performance and inference speed is observed, with a more pronounced effect in text classification and question answering tasks, and a weaker effect in sequence tagging. Additionally, at equal inference speed, our models trained with task-specific distillation exhibit improved performance compared to DistilBETO, which was trained with task-agnostic distillation, despite having significantly fewer parameters.

A similar effect can be observed when comparing ALBETO *base*-{8-10} to the original 12-layer ALBETO *base* fine-tuned using standard techniques. The former exhibits improved

Model	Parameters	Speedup	Score
Fine-tuning			
BETO <i>uncased</i>	110M	1.00×	81.02
BETO <i>cased</i>	110M	1.00×	84.82
DistilBETO	67M	2.00×	76.73
ALBETO <i>tiny</i>	5M	18.05×	74.97
ALBETO <i>base</i>	12M	0.99×	83.25
ALBETO <i>large</i>	18M	0.28×	82.02
ALBETO <i>xlarge</i>	59M	0.07×	84.13
ALBETO <i>xxlarge</i>	223M	0.03×	85.17
BERTIN	125M	1.00×	83.97
RoBERTa BNE <i>base</i>	125M	1.00×	84.83
RoBERTa BNE <i>large</i>	355M	0.28×	68.42
Task-specific Knowledge Distillation			
ALBETO <i>tiny</i>	5M	18.05×	76.49
ALBETO <i>base-2</i>	12M	5.96×	72.98
ALBETO <i>base-4</i>	12M	2.99×	80.06
ALBETO <i>base-6</i>	12M	1.99×	82.70
ALBETO <i>base-8</i>	12M	1.49×	83.78
ALBETO <i>base-10</i>	12M	1.19×	84.32

Table 5.4: The summary of results of every evaluated model in terms of parameters, inference speedup and overall score across tasks. The speedup is relative to BETO models. The score column shows the average of the metrics on all tasks.

performance. This underscores the vital role of task-specific knowledge distillation in obtaining improved performance for these faster models.

Table 5.4 summarizes our findings. Following the methodology of GLUE [107], we compute a global score that encompasses all tasks, which is displayed in the third column. The score is the simple mean of the individual task results. In the instance of Question Answering, which provides two metrics, we opted for the F1 Score as the representative score for the task. The ALBETO *xxlarge* model achieved the best overall performance, although it was also the slowest and had the second largest number of parameters. With a mere 0.35 performance drop from the top model, the RoBERTa BNE *base* and BETO *cased* models exhibited comparable results. The ALBETO *base-10*, exhibiting a 19% improvement in speed compared to BETO models, is our strongest proposed model in the task-specific KD setting, with a performance drop of approximately 0.5 compared to the aforementioned models. A negative result can be observed in the case of our faster model derived from ALBETO *base*, which is the *base-2* model, with only two layers. We observe a significant difference between the *base-2* model and the rest, which may show a limit of our approach in terms of how much layers can be removed. Our remaining models on the task-specific KD setting display varying degrees of improved inference speed, at the expense of slight reductions in task performance.

5.3 Inference Speed on Common Hardware

In order to assess the performance of deep learning models, it is important to evaluate the inference speed, which is measured in terms of Multiply-Accumulate (MAC) operations in our work. This metric is advantageous as it is agnostic to hardware variations. However, it can be useful to also report the actual inference speed of models on common hardware, as this can provide a more intuitive understanding of their performance.

Table 5.5 presents the average number of inferences per second that can be achieved on two different hardware platforms, a CPU with an Intel Core i7-11700K and a GPU with a NVIDIA GeForce RTX 3090. To account for variance in the measurements, we first conducted 10 warm-up inferences followed by 100 real measures for each model. We then applied an aggressive outlier filtering method based on the modified Z-Score [50] with a threshold of 0.75, which resulted in the removal of approximately 40-45% of the measures. The remaining 55-60% of the measures were used to calculate, with very low variance, the average inference speed (in milliseconds) and the number of inferences that could be performed in one second, which serves as a clearer illustration of the model’s inference speed.

It is worth noting that the difference in speed between the larger models and the proposed models trained using task-specific KD is substantial. Specifically, on the CPU setting, which is representative of popular serverless platforms used in industry, the best model found in this study in terms of task performance, ALBETO *xxlarge*, would take several seconds for a single inference, making it unsuitable for real-time user-facing applications. On the other hand, if we consider our proposed faster models, we can observe that ALBETO *base-6* is capable of executing more than 10 inferences per second, which is a much more acceptable latency for a real-time application.

Model	Inferences per second	
	CPU	GPU
Fine-tuning		
BETO <i>uncased</i>	3.96	107.19
BETO <i>cased</i>	4.26	109.02
DistilBETO	9.12	217.40
ALBETO <i>tiny</i>	32.53	539.61
ALBETO <i>base</i>	4.50	108.62
ALBETO <i>large</i>	1.29	33.62
ALBETO <i>xlarge</i>	0.35	11.72
ALBETO <i>xxlarge</i>	0.14	6.60
BERTIN	3.99	109.39
RoBERTa BNE <i>base</i>	3.82	107.77
RoBERTa BNE <i>large</i>	1.18	33.65
Task-specific Knowledge Distillation		
ALBETO <i>tiny</i>	32.53	539.61
ALBETO <i>base-2</i>	31.08	625.30
ALBETO <i>base-4</i>	15.16	319.32
ALBETO <i>base-6</i>	10.45	213.53
ALBETO <i>base-8</i>	6.82	160.66
ALBETO <i>base-10</i>	6.01	128.38

Table 5.5: The number of inferences per second of each model on two different hardware settings, CPU and GPU.

Chapter 6

Conclusion

6.1 Summary of Contributions

In this work, we introduced ALBETO and Speedy Gonzales, which are two novel resources for the Spanish NLP community that were created to improve two key aspects of machine learning models, namely model size and inference speed.

We have developed a series of language models, known as ALBETO, which have been pre-trained exclusively for the Spanish language. These models are built upon the ALBERT architecture, which utilizes the weight-tying strategy to achieve greater efficiency in terms of model parameters. ALBETO models were trained in five different sizes: *tiny*, *base*, *large*, *xlarge* and *xxlarge*. The first four mentioned models are lighter when compared with base-sized models for Spanish such as BETO, BERTIN, and RoBERTa-BNE *base*. The larger ALBETO *xxlarge* is lighter in terms of parameters to the *large* version of RoBERTa-BNE.

The successful performance of our ALBETO models in the fine-tuning setting allows us to address our first research question, which investigates the feasibility of developing lightweight and parameter-efficient models through the utilization of the weight-tying technique while still achieving comparable performance to larger models. This is substantiated by the fact that our ALBETO *base* model outperforms or is only marginally inferior to other *base*-sized models for Spanish such as BETO, BERTIN, and RoBERTa-BNE *base*. Notably, our *base* model, which is an uncased model, outperforms the *uncased* version of BETO on almost all the evaluated tasks.

Additionally, we conducted an investigation into the feasibility of utilizing Task-Specific Knowledge Distillation to efficiently transfer knowledge from a larger, high-performing model to a lighter and faster model across a range of tasks in the Spanish language. Our findings demonstrate that this approach is effective. Specifically, we developed a suite of fast models for Spanish, called Speedy Gonzales, which are built upon our previously introduced ALBETO models. These models achieved task performance levels that are comparable to those of most base-sized models in the fine-tuning context while also exhibiting enhanced inference speed due to the application of the task-specific KD technique.

Our faster models were derived from two of the ALBETO models, the *tiny* and *base*. In the case of the *tiny*, we utilized the pre-trained model without any modifications, providing evidence of the effectiveness of the task-specific KD technique over simple fine-tuning. Our results showed that the model trained using KD performed better than the one trained using fine-tuning in most tasks. In the case of the *base* model, we introduced 5 models derived from it, with the unique difference of having fewer layers of the model. Our evaluation of these models in the KD setting led to the observation of a hypothesized behavior, namely the existence of a trade-off between efficiency of the model and task performance. The models that are faster are typically worse in terms of task performance than those that have more layers and that are slower, but this difference is generally small. Of these models, one model was considerably worse, the model with two layers, which can give us a hint on the limits of our method.

In addition, another interesting observation was that there exist tasks that are presumably easier ones, where smaller and larger models perform almost equally good. This finding suggests that larger models may be excessively overparameterized for these tasks, and that lighter and faster models can be effectively constructed for these tasks. On the other hand, the opposite behaviour can be observed with tasks that are probably more difficult such as Question Answering. These tasks benefit from the use of models that are larger in parameters and are more computationally complex. Nevertheless, the benefit of using larger models is not substantial and must be viewed as a trade-off based on the use case.

As a concluding contribution, we have made available to the research community all of our models, comprising more than 140 models in total, including pre-trained, fine-tuned, and distilled models, to facilitate further investigations in the field.

6.2 Limitations and Future Research Directions

Despite our contributions, our work has some limitations that can be addressed in future research. First, we only evaluated our models on a limited set of tasks. It would be beneficial to evaluate our models on other Spanish-language datasets, including intrinsic and extrinsic tasks. We also believe that the availability of these models and the expansion of the Knowledge Distillation method to additional tasks will drive the widespread utilization of language models in the Spanish speaking community, particularly for individuals and organizations seeking to tackle crucial NLP challenges, such as question answering, text similarity and semantic search, in both academic and industrial settings.

Secondly, our knowledge distillation method can be further improved to produce more efficient task-specific language models. Specifically, it would be interesting to explore alternative KD approaches, such as distilling intermediate layers of the teacher model, in addition to its output. Additionally, a multi-teacher approach could be studied, in which the models learn from a collection of teacher models rather than just one.

Finally, there exists a trade-off between model size, inference speed, and task performance, making it challenging to choose an appropriate model without context. Therefore, it is important to develop metrics to formally assess this balance.

Bibliography

- [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: a system for large-scale machine learning. In *Osdi*, volume 16, pages 265–283. Savannah, GA, USA, 2016.
- [2] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A next-generation hyperparameter optimization framework. In Ankur Tere-desai, Vipin Kumar, Ying Li, Rómer Rosales, Evimaria Terzi, and George Karypis, editors, *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2019, Anchorage, AK, USA, August 4-8, 2019*, pages 2623–2631. ACM, 2019. doi: 10.1145/3292500.3330701. URL <https://doi.org/10.1145/3292500.3330701>.
- [3] Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. On the cross-lingual transferability of monolingual representations. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4623–4637, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.421. URL <https://aclanthology.org/2020.acl-main.421>.
- [4] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460, 2020.
- [5] Haoli Bai, Wei Zhang, Lu Hou, Lifeng Shang, Jin Jin, Xin Jiang, Qun Liu, Michael Lyu, and Irwin King. BinaryBERT: Pushing the limit of BERT quantization. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4334–4348, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.334. URL <https://aclanthology.org/2021.acl-long.334>.
- [6] Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- [7] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828, 2013. doi: 10.1109/TPAMI.2013.50.

- [8] Davis Blalock, Jose Javier Gonzalez Ortiz, Jonathan Frankle, and John Guttag. What is the state of neural network pruning? In I. Dhillon, D. Papailiopoulos, and V. Sze, editors, *Proceedings of Machine Learning and Systems*, volume 2, pages 129–146, 2020. URL <https://proceedings.mlsys.org/paper/2020/file/d2ddea18f00665ce8623e36bd4e3c7c5-Paper.pdf>.
- [9] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. Enriching word vectors with subword information. *Transactions of the association for computational linguistics*, 5:135–146, 2017.
- [10] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf>.
- [11] Nicoletta Calzolari, Frédéric Béchet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, H el ene Mazo, Jan Odijk, and Stelios Piperidis, editors. *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, Marseille, France, June 2022. European Language Resources Association. URL <https://aclanthology.org/2022.lrec-1.0>.
- [12] Jos e Ca nete, Sebastian Donoso, Felipe Bravo-Marquez, Andr es Carvallo, and Vladimir Araujo. ALBETO and DistilBETO: Lightweight Spanish language models. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 4291–4298, Marseille, France, June 2022. European Language Resources Association. URL <https://aclanthology.org/2022.lrec-1.457>.
- [13] Cristian Cardellino. Spanish Billion Words Corpus and Embeddings, March 2016. URL <https://crscardellino.github.io/SBWCE/>.
- [14] Casimiro Pio Carrino, Marta R. Costa-juss a, and Jos e A. R. Fonollosa. Automatic Spanish translation of SQuAD dataset for multi-lingual question answering. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 5515–5523, Marseille, France, May 2020. European Language Resources Association. ISBN 979-10-95546-34-4. URL <https://aclanthology.org/2020.lrec-1.677>.
- [15] Casimiro Pio Carrino, Joan Llop, Marc P amies, Asier Guti errez-Fandi no, Jordi Armengol-Estap e, Joaqu ın Silveira-Ocampo, Alfonso Valencia, Aitor Gonzalez-Agirre, and Marta Villegas. Pretrained biomedical language models for clinical NLP in Spanish. In *Proceedings of the 21st Workshop on Biomedical Language Processing*, pages 193–199, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.bionlp-1.19. URL <https://aclanthology.org/2022.bionlp-1.19>.

- [16] José Cañete. Compilation of Large Spanish Unannotated Corpora, May 2019. URL <https://doi.org/10.5281/zenodo.3247731>.
- [17] José Cañete and Felipe Bravo-Marquez. Speedy gonzales: A collection of fast task-specific models for spanish, 2023.
- [18] José Cañete, Gabriel Chaperon, Rodrigo Fuentes, Jou-Hui Ho, Hojin Kang, and Jorge Pérez. Spanish pre-trained bert model and evaluation data. In *PML4DC at ICLR 2020*, 2020.
- [19] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems*, 34:15084–15097, 2021.
- [20] Shi Chen and Qi Zhao. Shallowing deep networks: Layer-wise pruning based on feature representations. *IEEE transactions on pattern analysis and machine intelligence*, 41(12):3048–3056, 2018.
- [21] Tianlong Chen, Jonathan Frankle, Shiyu Chang, Sijia Liu, Yang Zhang, Zhangyang Wang, and Michael Carbin. The lottery ticket hypothesis for pre-trained bert networks. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 15834–15846. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/b6af2c9703f203a2794be03d443af2e3-Paper.pdf>.
- [22] Tianlong Chen, Jonathan Frankle, Shiyu Chang, Sijia Liu, Yang Zhang, Michael Carbin, and Zhangyang Wang. The lottery tickets hypothesis for supervised and self-supervised pre-training in computer vision models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16306–16316, 2021.
- [23] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. Palm: Scaling language modeling with pathways. *CoRR*, abs/2204.02311, 2022. doi: 10.48550/arXiv.2204.02311. URL <https://doi.org/10.48550/arXiv.2204.02311>.
- [24] Alexis Conneau, Guillaume Lample, Ruty Rinott, Adina Williams, Samuel R Bowman, Holger Schwenk, and Veselin Stoyanov. Xnli: Evaluating cross-lingual sentence representations. *arXiv preprint arXiv:1809.05053*, 2018.

- [25] Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel Bowman, Holger Schwenk, and Veselin Stoyanov. XNLI: Evaluating cross-lingual sentence representations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2475–2485, Brussels, Belgium, October–November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1269. URL <https://aclanthology.org/D18-1269>.
- [26] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.747. URL <https://aclanthology.org/2020.acl-main.747>.
- [27] Javier de la Rosa, Eduardo G Ponferrada, Paulo Villegas, Pablo González de Prado Salas, Manu Romero, and Maria Grandury. Bertin: Efficient pre-training of a spanish language model using perplexity sampling. *Procesamiento del Lenguaje Natural*, 68(0):13–23, 2022. ISSN 1989-7553. URL <http://journal.sepln.org/sepln/ojs/ojs/index.php/pln/article/view/6403>.
- [28] Marie-Catherine de Marneffe, Christopher D. Manning, Joakim Nivre, and Daniel Zeman. Universal Dependencies. *Computational Linguistics*, 47(2):255–308, June 2021. doi: 10.1162/coli.a.00402. URL <https://aclanthology.org/2021.cl-2.11>.
- [29] Wietse de Vries, Andreas van Cranenburgh, Arianna Bisazza, Tommaso Caselli, Gertjan van Noord, and Malvina Nissim. BERTje: A Dutch BERT Model. arXiv:1912.09582, December 2019. URL <http://arxiv.org/abs/1912.09582>.
- [30] Pieter Delobelle, Thomas Winters, and Bettina Berendt. RobBERT: a Dutch RoBERTa-based Language Model. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3255–3265, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.292. URL <https://aclanthology.org/2020.findings-emnlp.292>.
- [31] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- [32] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
- [33] Sebastián Alejandro Donoso. Entrenamiento y evaluación de modelos pequeños de lenguaje natural basado en métodos de autoatención. 2021. URL <https://repositorio.uchile.cl/handle/2250/183444>.

- [34] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [35] Angela Fan, Edouard Grave, and Armand Joulin. Reducing transformer depth on demand with structured dropout. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=Sy102yStDr>.
- [36] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=rJl-b3RcF7>.
- [37] Xinyang Geng, Hao Liu, Lisa Lee, Dale Schuurams, Sergey Levine, and Pieter Abbeel. Multimodal masked autoencoders learn transferable representations. *arXiv preprint arXiv:2205.14204*, 2022.
- [38] Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W. Mahoney, and Kurt Keutzer. A survey of quantization methods for efficient neural network inference. *CoRR*, abs/2103.13630, 2021. URL <https://arxiv.org/abs/2103.13630>.
- [39] Yoav Goldberg. Neural network methods for natural language processing. *Synthesis lectures on human language technologies*, 10(1):1–309, 2017.
- [40] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [41] Jianping Gou, Baosheng Yu, Stephen J. Maybank, and Dacheng Tao. Knowledge distillation: A survey. *Int. J. Comput. Vision*, 129(6):1789–1819, jun 2021. ISSN 0920-5691. doi: 10.1007/s11263-021-01453-z. URL <https://doi.org/10.1007/s11263-021-01453-z>.
- [42] Naman Goyal, Jingfei Du, Myle Ott, Giri Anantharaman, and Alexis Conneau. Larger-scale transformers for multilingual masked language modeling. In *Proceedings of the 6th Workshop on Representation Learning for NLP (Repl4NLP-2021)*, pages 29–33, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.repl4nlp-1.4. URL <https://aclanthology.org/2021.repl4nlp-1.4>.
- [43] Asier Gutiérrez-Fandiño, Jordi Armengol-Estapé, Aitor Gonzalez-Agirre, and Marta Villegas. Spanish legalese language model and corpora, 2021.
- [44] Asier Gutiérrez-Fandiño, Jordi Armengol Estapé, Marc Pàmies, Joan Llop Palao, Joaquin Silveira Ocampo, Casimiro Pio Carrino, Carme Armentano Oller, Carlos Rodriguez Penagos, Aitor Gonzalez Agirre, and Marta Villegas. Maria: Spanish language models. *Procesamiento del Lenguaje Natural*, 68, 2022. ISSN 1135-5948. doi: 10.26342/2022-68-3. URL <https://upcommons.upc.edu/handle/2117/367156#.YyMTB4X9A-0.mendeley>.
- [45] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022.

- [46] Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. DeBERTa: Decoding-enhanced bert with disentangled attention. *arXiv preprint arXiv:2006.03654*, 2020.
- [47] Geoffrey Hinton, Oriol Vinyals, Jeff Dean, et al. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2(7), 2015.
- [48] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [49] Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. *arXiv preprint arXiv:1801.06146*, 2018.
- [50] Boris Iglewicz and David C Hoaglin. *Volume 16: how to detect and handle outliers*. Quality Press, 1993.
- [51] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. TinyBERT: Distilling BERT for natural language understanding. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4163–4174, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.372. URL <https://aclanthology.org/2020.findings-emnlp.372>.
- [52] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. Lightmbert: A simple yet effective method for multilingual bert distillation. *arXiv preprint arXiv:2103.06418*, 2021.
- [53] Michael I Jordan and Tom M Mitchell. Machine learning: Trends, perspectives, and prospects. *Science*, 349(6245):255–260, 2015.
- [54] Taku Kudo and John Richardson. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium, November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-2012. URL <https://aclanthology.org/D18-2012>.
- [55] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. ALBERT: A lite BERT for self-supervised learning of language representations. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL <https://openreview.net/forum?id=H1eA7AetvS>.
- [56] Hang Le, Loïc Vial, Jibril Frej, Vincent Segonne, Maximin Coavoux, Benjamin Lecouteux, Alexandre Allauzen, Benoit Crabbé, Laurent Besacier, and Didier Schwab. FlauBERT: Unsupervised language model pre-training for French. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 2479–2490, Marseille, France, May 2020. European Language Resources Association. ISBN 979-10-95546-34-4. URL <https://aclanthology.org/2020.lrec-1.302>.
- [57] Yann LeCun, John Denker, and Sara Solla. Optimal brain damage. *Advances in neural information processing systems*, 2, 1989.

- [58] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [59] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553): 436–444, 2015.
- [60] David D. Lewis, Yiming Yang, Tony G. Rose, and Fan Li. RCV1: A new benchmark collection for text categorization research. *J. Mach. Learn. Res.*, 5:361–397, 2004. URL <http://jmlr.org/papers/volume5/lewis04a/lewis04a.pdf>.
- [61] Patrick Lewis, Barlas Oguz, Ruty Rinott, Sebastian Riedel, and Holger Schwenk. MLQA: Evaluating cross-lingual extractive question answering. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7315–7330, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.653. URL <https://aclanthology.org/2020.acl-main.653>.
- [62] Davis Liang, Hila Gonen, Yuning Mao, Rui Hou, Naman Goyal, Marjan Ghazvininejad, Luke Zettlemoyer, and Madian Khabsa. Xlm-v: Overcoming the vocabulary bottleneck in multilingual masked language models. *arXiv preprint arXiv:2301.10472*, 2023.
- [63] Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, Ramakanth Pasunuru, Sam Shleifer, Punit Singh Koura, Vishrav Chaudhary, Brian O’Horo, Jeff Wang, Luke Zettlemoyer, Zornitsa Kozareva, Mona Diab, Veselin Stoyanov, and Xian Li. Few-shot learning with multilingual generative language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9019–9052, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. URL <https://aclanthology.org/2022.emnlp-main.616>.
- [64] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized BERT pretraining approach. *CoRR*, abs/1907.11692, 2019. URL <http://arxiv.org/abs/1907.11692>.
- [65] Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. Multilingual denoising pre-training for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8: 726–742, 2020. doi: 10.1162/tacl_a_00343. URL <https://aclanthology.org/2020.tacl-1.47>.
- [66] Abhishek Mahajani, Vinay Pandya, Isaac Maria, and Deepak Sharma. A comprehensive survey on extractive and abstractive techniques for text summarization. In Yu-Chen Hu, Shailesh Tiwari, Krishn K. Mishra, and Munesh C. Trivedi, editors, *Ambient Communications and Computer Systems*, pages 339–351, Singapore, 2019. Springer Singapore. ISBN 978-981-13-5934-7.
- [67] Louis Martin, Benjamin Muller, Pedro Javier Ortiz Suárez, Yoann Dupont, Laurent Romary, Éric de la Clergerie, Djamé Seddah, and Benoît Sagot. CamemBERT: a

- tasty French language model. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7203–7219, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.645. URL <https://aclanthology.org/2020.acl-main.645>.
- [68] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26, 2013.
- [69] Saif Mohammad, Felipe Bravo-Marquez, Mohammad Salameh, and Svetlana Kiritchenko. Semeval-2018 task 1: Affect in tweets. In *Proceedings of the 12th international workshop on semantic evaluation*, pages 1–17, 2018.
- [70] Michael Negnevitsky. *Artificial intelligence: a guide to intelligent systems*. Pearson education, 2005.
- [71] Austral University of Chile. Patagón supercomputer, 2021. URL <https://patagon.uach.cl>.
- [72] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdf>.
- [73] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, 2014. URL <http://www.aclweb.org/anthology/D14-1162>.
- [74] Juan Manuel Pérez, Damián Ariel Furman, Laura Alonso Alemany, and Franco M. Luque. RoBERTuito: a pre-trained language model for social media text in Spanish. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 7235–7243, Marseille, France, June 2022. European Language Resources Association. URL <https://aclanthology.org/2022.lrec-1.785>.
- [75] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1202. URL <https://aclanthology.org/N18-1202>.

- [76] Telmo Pires, Eva Schlinger, and Dan Garrette. How multilingual is multilingual BERT? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4996–5001, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1493. URL <https://aclanthology.org/P19-1493>.
- [77] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.
- [78] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019.
- [79] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [80] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. *arXiv preprint arXiv:2212.04356*, 2022.
- [81] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140): 1–67, 2020. URL <http://jmlr.org/papers/v21/20-074.html>.
- [82] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(1), jun 2022. ISSN 1532-4435.
- [83] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas, November 2016. Association for Computational Linguistics. doi: 10.18653/v1/D16-1264. URL <https://aclanthology.org/D16-1264>.
- [84] Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 3505–3506, 2020.
- [85] Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik-Manor. Imagenet-21k pretraining for the masses, 2021.
- [86] Sara Rosenthal, Noura Farra, and Preslav Nakov. Semeval-2017 task 4: Sentiment analysis in twitter. In *Proceedings of the 11th international workshop on semantic evaluation (SemEval-2017)*, pages 502–518, 2017.
- [87] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *nature*, 323(6088):533–536, 1986.

- [88] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115: 211–252, 2015.
- [89] Stuart J Russell. *Artificial intelligence a modern approach*. Pearson Education, Inc., 2010.
- [90] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of BERT: smaller, faster, cheaper and lighter. *CoRR*, abs/1910.01108, 2019. URL <http://arxiv.org/abs/1910.01108>.
- [91] Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*, 2022.
- [92] Holger Schwenk and Xian Li. A corpus for multilingual document classification in eight languages. In Nicoletta Calzolari (Conference chair), Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Koiti Hasida, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, H el ene Mazo, Asuncion Moreno, Jan Odijk, Stelios Piperidis, and Takenobu Tokunaga, editors, *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Paris, France, may 2018. European Language Resources Association (ELRA). ISBN 979-10-95546-00-9.
- [93] Holger Schwenk and Xian Li. A corpus for multilingual document classification in eight languages. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan, May 2018. European Language Resources Association (ELRA). URL <https://aclanthology.org/L18-1560>.
- [94] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/v1/P16-1162. URL <https://aclanthology.org/P16-1162>.
- [95] Alejandro Vaca Serrano, Guillem Garcia Subies, Helena Montoro Zamorano, Nuria Aldama Garcia, Doaa Samy, David Betancur Sanchez, Antonio Moreno Sandoval, Marta Guerrero Nieto, and Alvaro Barbero Jimenez. Rigoberta: A state-of-the-art language model for spanish. *arXiv preprint arXiv:2205.10233*, 2022.
- [96] Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-lm: Training multi-billion parameter language models using model parallelism. *CoRR*, abs/1909.08053, 2019. URL <http://arxiv.org/abs/1909.08053>.
- [97] Shaden Smith, Mostofa Patwary, Brandon Norick, Patrick LeGresley, Samyam Rajbhandari, Jared Casper, Zhun Liu, Shrimai Prabhumoye, George Zerveas, Vijay Korthikanti, Elton Zheng, Rewon Child, Reza Yazdani Aminabadi, Julie Bernauer, Xia Song, Mohammad Shoeybi, Yuxiong He, Michael Houston, Saurabh Tiwary, and

- Bryan Catanzaro. Using deepspeed and megatron to train megatron-turing NLG 530b, A large-scale generative language model. *CoRR*, abs/2201.11990, 2022. URL <https://arxiv.org/abs/2201.11990>.
- [98] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642, 2013.
- [99] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112, 2014.
- [100] Raphael Tang, Yao Lu, Linqing Liu, Lili Mou, Olga Vechtomova, and Jimmy Lin. Distilling task-specific knowledge from BERT into simple neural networks. *CoRR*, abs/1903.12136, 2019. URL <http://arxiv.org/abs/1903.12136>.
- [101] Mariona Taulé, M. Antònia Martí, and Marta Recasens. AnCora: Multilevel annotated corpora for Catalan and Spanish. In *Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC’08)*, Marrakech, Morocco, May 2008. European Language Resources Association (ELRA). URL http://www.lrec-conf.org/proceedings/lrec2008/pdf/35_paper.pdf.
- [102] Jörg Tiedemann. Parallel data, tools and interfaces in opus. In *Lrec*, volume 2012, pages 2214–2218. Citeseer, 2012.
- [103] Erik F. Tjong Kim Sang. Introduction to the CoNLL-2002 shared task: Language-independent named entity recognition. In *COLING-02: The 6th Conference on Natural Language Learning 2002 (CoNLL-2002)*, 2002. URL <https://aclanthology.org/W02-2024>.
- [104] Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Well-read students learn better: On the importance of pre-training compact models. *arXiv preprint arXiv:1908.08962v2*, 2019.
- [105] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [106] Antti Virtanen, Jenna Kanerva, Rami Ilo, Jouni Luoma, Juhani Luotolahti, Tapio Salakoski, Filip Ginter, and Sampo Pyysalo. Multilingual is not enough: Bert for finnish. *arXiv preprint arXiv:1912.07076*, 2019.
- [107] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium, November 2018. Association for Computational Linguistics. doi: 10.18653/v1/W18-5446. URL <https://aclanthology.org/W18-5446>.

- [108] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. *SuperGLUE: A Stickier Benchmark for General-Purpose Language Understanding Systems*. Curran Associates Inc., Red Hook, NY, USA, 2019.
- [109] Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS’20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- [110] Alex Warstadt, Amanpreet Singh, and Samuel R Bowman. Neural network acceptability judgments. *Transactions of the Association for Computational Linguistics*, 7: 625–641, 2019.
- [111] Adina Williams, Nikita Nangia, and Samuel R Bowman. A broad-coverage challenge corpus for sentence understanding through inference. *arXiv preprint arXiv:1704.05426*, 2017.
- [112] Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1101. URL <https://aclanthology.org/N18-1101>.
- [113] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online, October 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-demos.6. URL <https://aclanthology.org/2020.emnlp-demos.6>.
- [114] Shijie Wu and Mark Dredze. Beto, bentz, becas: The surprising cross-lingual effectiveness of BERT. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 833–844, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1077. URL <https://aclanthology.org/D19-1077>.
- [115] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*, 2016.
- [116] Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. mT5: A massively multilingual pre-trained

- text-to-text transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.41. URL <https://aclanthology.org/2021.naacl-main.41>.
- [117] Yinfei Yang, Yuan Zhang, Chris Tar, and Jason Baldridge. PAWS-X: A cross-lingual adversarial dataset for paraphrase identification. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3687–3692, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1382. URL <https://aclanthology.org/D19-1382>.
- [118] Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep learning: Training bert in 76 minutes. *arXiv preprint arXiv:1904.00962*, 2019.
- [119] Li Yuan, Francis EH Tay, Guilin Li, Tao Wang, and Jiashi Feng. Revisiting knowledge distillation via label smoothing regularization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [120] Sukmin Yun, Jongjin Park, Kimin Lee, and Jinwoo Shin. Regularizing class-wise predictions via self-knowledge distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [121] Wei Zhang, Lu Hou, Yichun Yin, Lifeng Shang, Xiao Chen, Xin Jiang, and Qun Liu. TernaryBERT: Distillation-aware ultra-low bit BERT. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 509–521, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.37. URL <https://aclanthology.org/2020.emnlp-main.37>.
- [122] Yuan Zhang, Jason Baldridge, and Luheng He. PAWS: Paraphrase adversaries from word scrambling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1298–1308, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1131. URL <https://aclanthology.org/N19-1131>.

Annex A

Examples of Considered Tasks

This section provides illustrations of all tasks that have been taken into consideration. Specifically, Figure A.1, Figure A.2, and Figure A.3 exemplify Document Classification, Paraphrase Identification, and Natural Language Inference, respectively. Moreover, Figure A.5 and Figure A.4 depict typical examples of Sequence Tagging tasks, such as Named-Entity Recognition and Part-of-Speech. The final task, Question Answering, is demonstrated in Figure A.6.

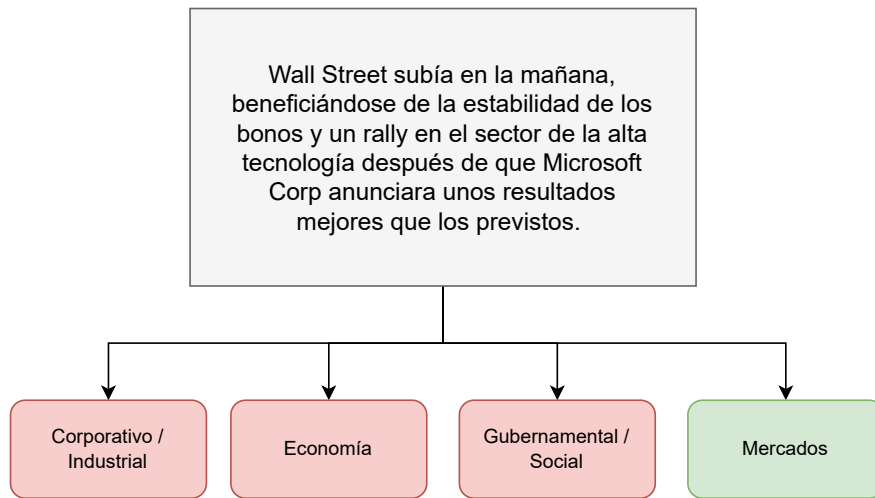


Figure A.1: An example of document classification. Taken from the MLDoc [93] dataset.

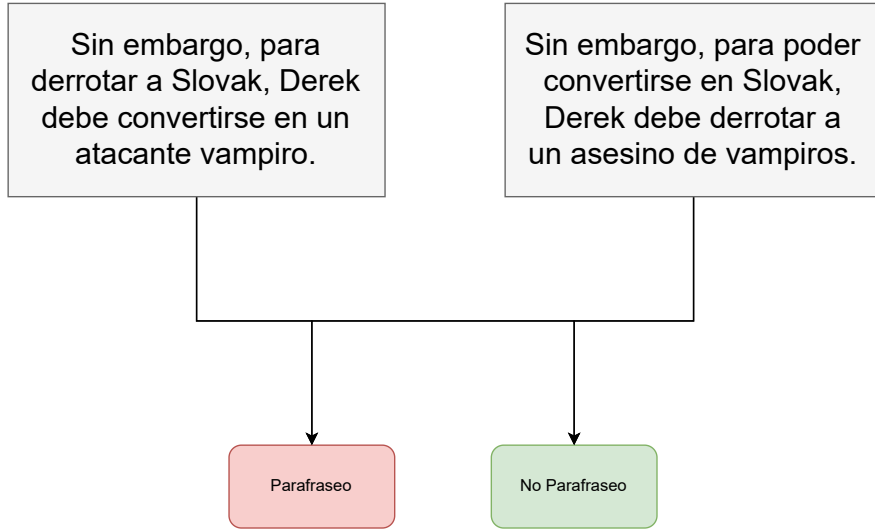


Figure A.2: An example of the paraphrase identification task. Taken from the PAWS-X [117] dataset.

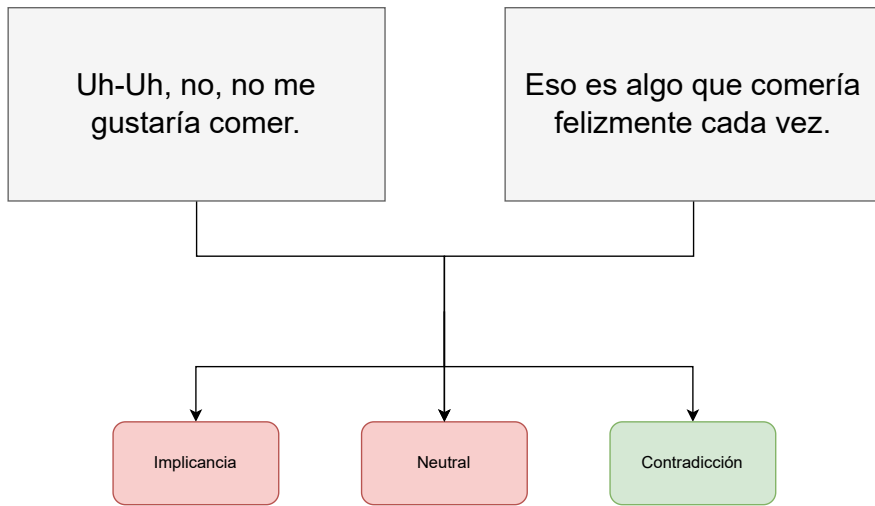


Figure A.3: An example of natural language inference. Taken from the XNLI [25] dataset.

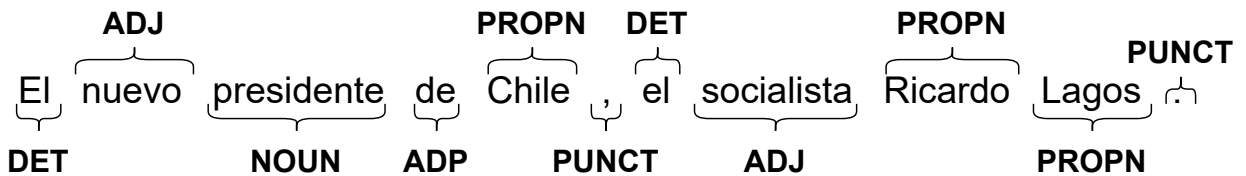


Figure A.4: An example of the Part-of-Speech tagging task. Taken from the AnCora [101] dataset.

En su lugar entró el chileno Iván Luis Zamorano

PER

Figure A.5: An example of named entity recognition. Taken from the CoNLL2002 NER [103] dataset.

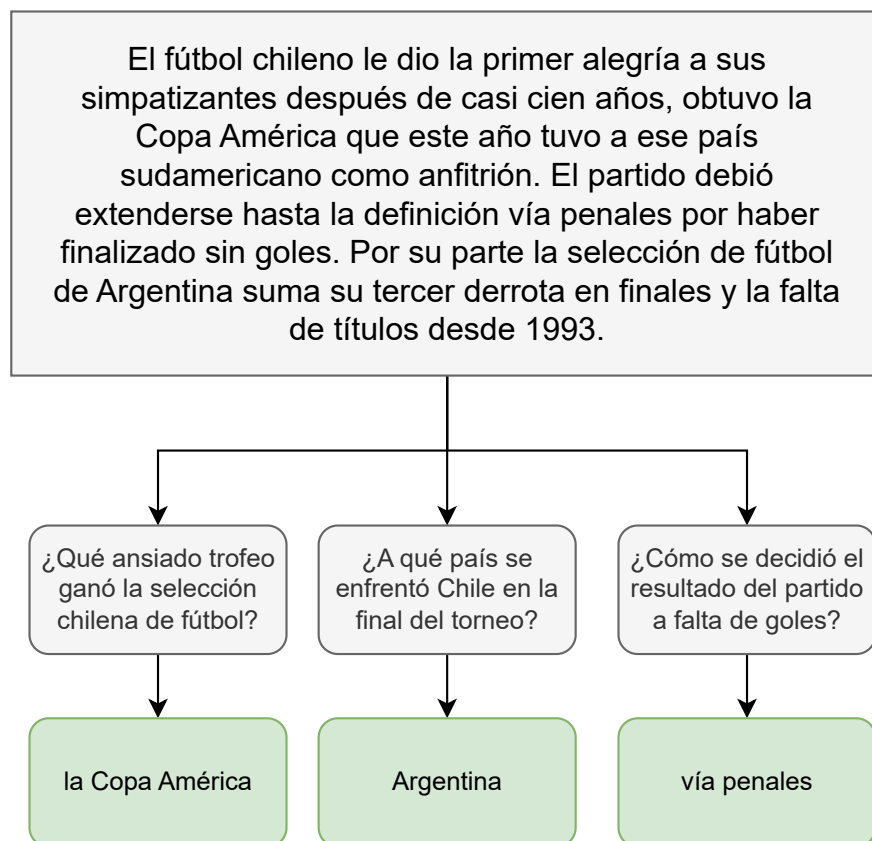


Figure A.6: An example of question answering. Taken from the SQAC [44] dataset.

Annex B

Training Loss on ALBETO models

In this section, the loss curves of all the ALBETO models trained in this study are presented. These models are arranged in ascending order, starting from the smaller ALBETO *tiny* model and finishing with the larger ALBETO *xxlarge* model.

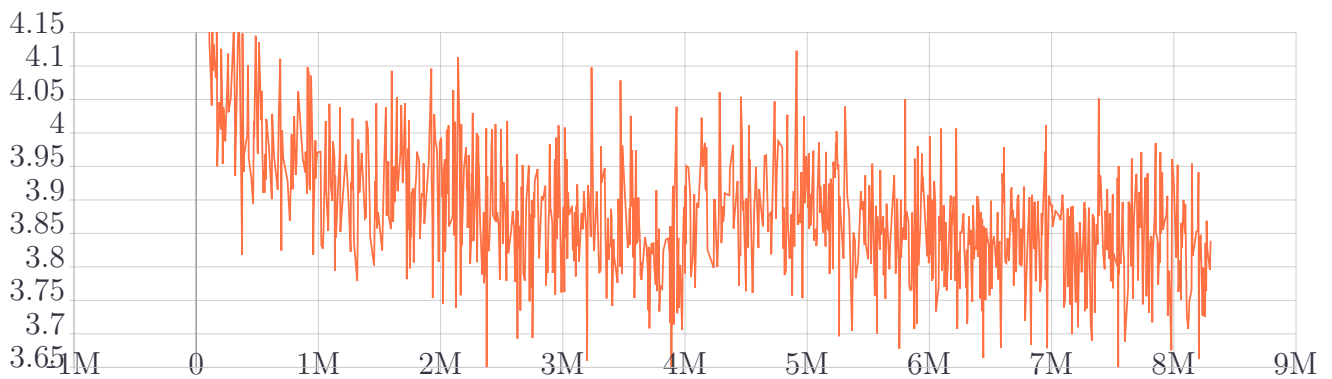


Figure B.1: The progression of the training loss on the ALBETO *tiny* model.

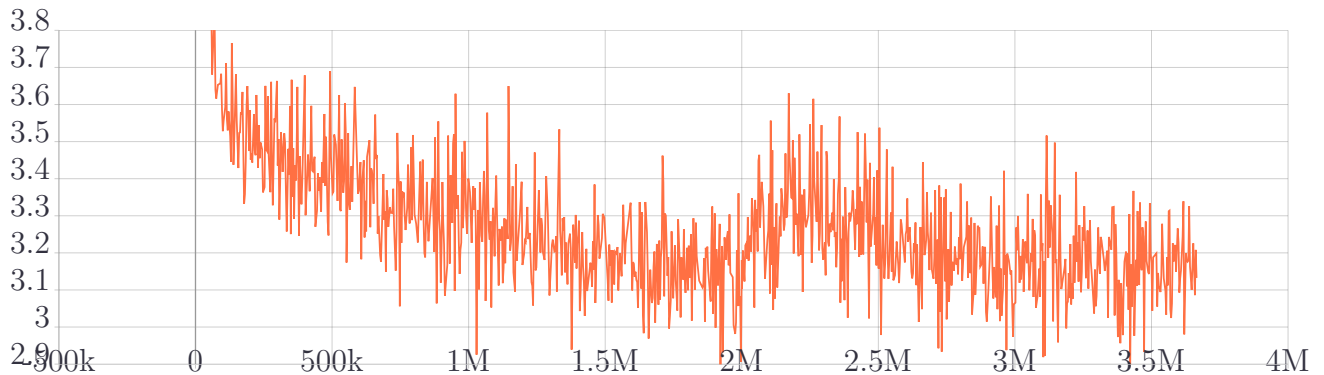


Figure B.2: The progression of the training loss on the ALBETO *base* model.

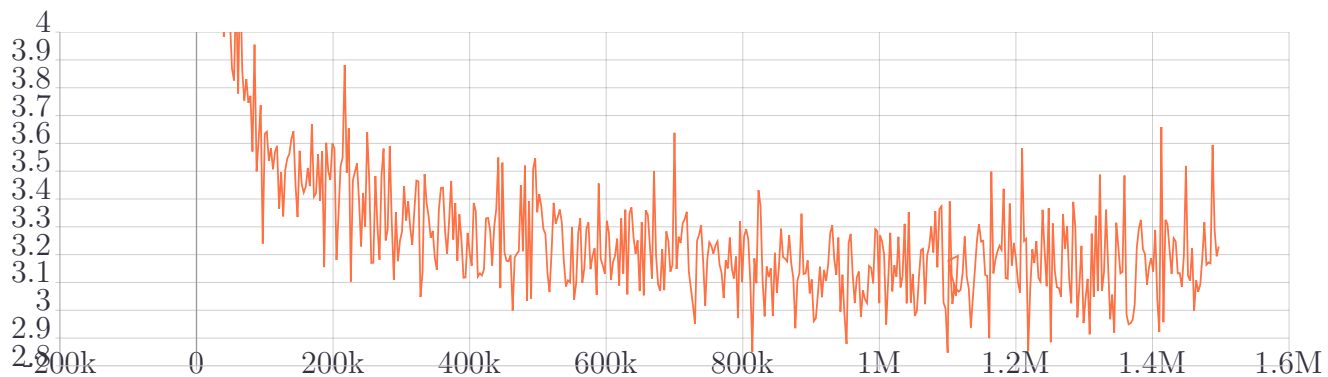


Figure B.3: The progression of the training loss on the ALBETO *large* model.

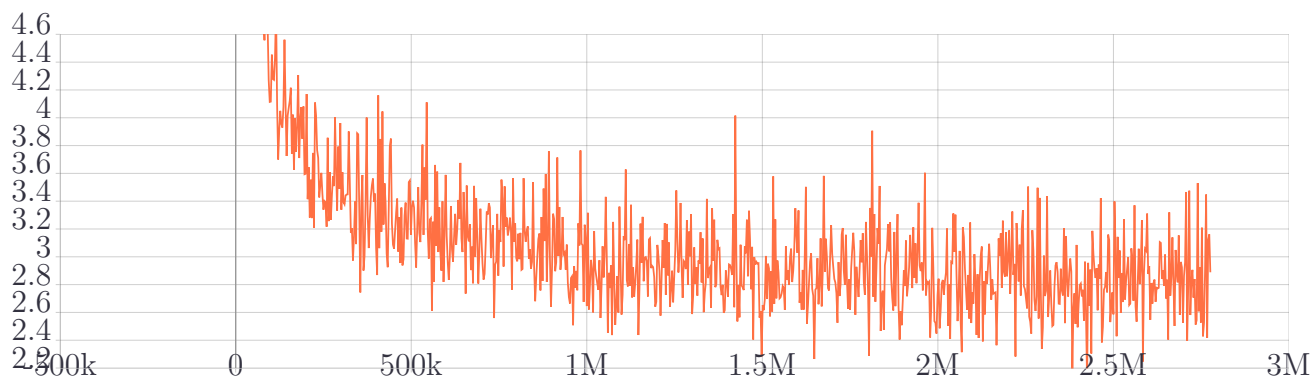


Figure B.4: The progression of the training loss on the ALBETO *xlarge* model.

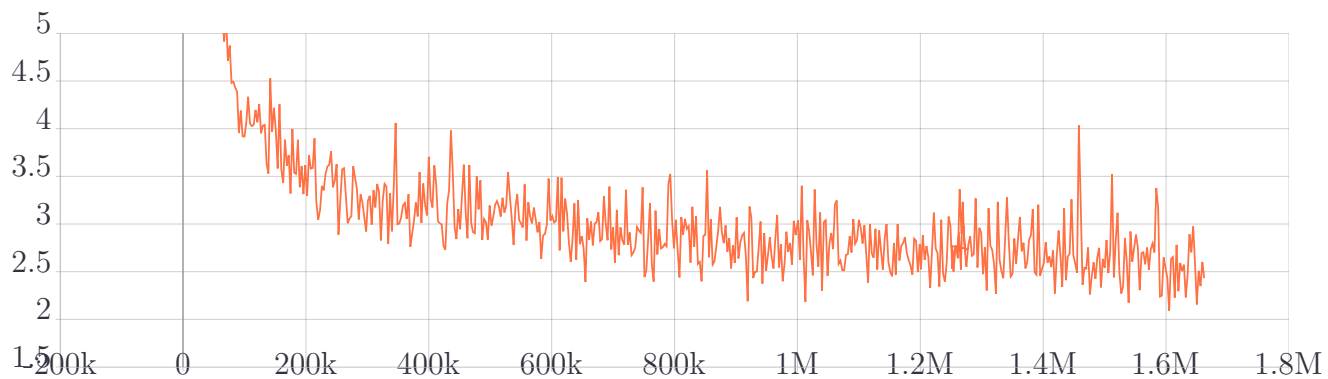


Figure B.5: The progression of the training loss on the ALBETO *xxlarge* model.

Annex C

Selected Teacher Models

Table C.1 presents the teacher models selected for each task. The selection process is based on the lowest validation loss achieved among the candidate teacher models that were fine-tuned for each task.

Dataset	Teacher Model
MLDoc	RoBERTa BNE <i>large</i>
PAWS-X	ALBETO <i>xxlarge</i>
XNLI	ALBETO <i>xxlarge</i>
POS	RoBERTa BNE <i>base</i>
NER	RoBERTa BNE <i>base</i>
MLQA	ALBETO <i>xxlarge</i>
SQAC	ALBETO <i>xxlarge</i>
TAR / XQuAD	ALBETO <i>xxlarge</i>

Table C.1: The teacher models selected for each task.

Annex D

Knowledge Distillation Implementation for Models With Shared Vocabulary

This section provides visual representations of the simpler scenarios of knowledge distillation, where the teacher and student models employ an identical vocabulary, and thus, alignment between tokenized sentences is not required.

Figure D.1 presents an exemplar of such implementation in the context of sequence tagging tasks. Moreover, Figure D.2 illustrates a similar approach for the question answering task.

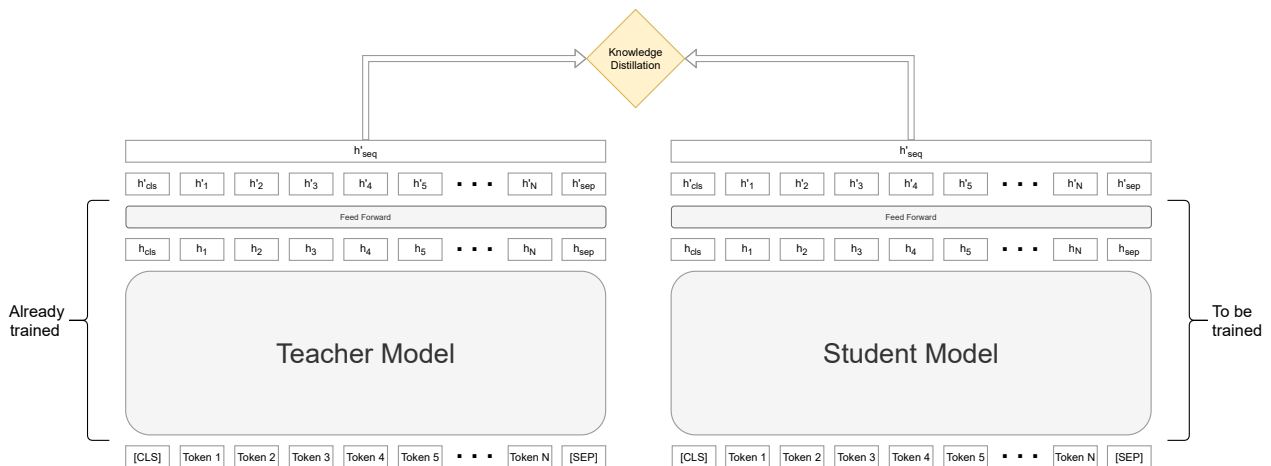


Figure D.1: Implementation of knowledge distillation for sequence tagging tasks with models that share the same vocabulary.

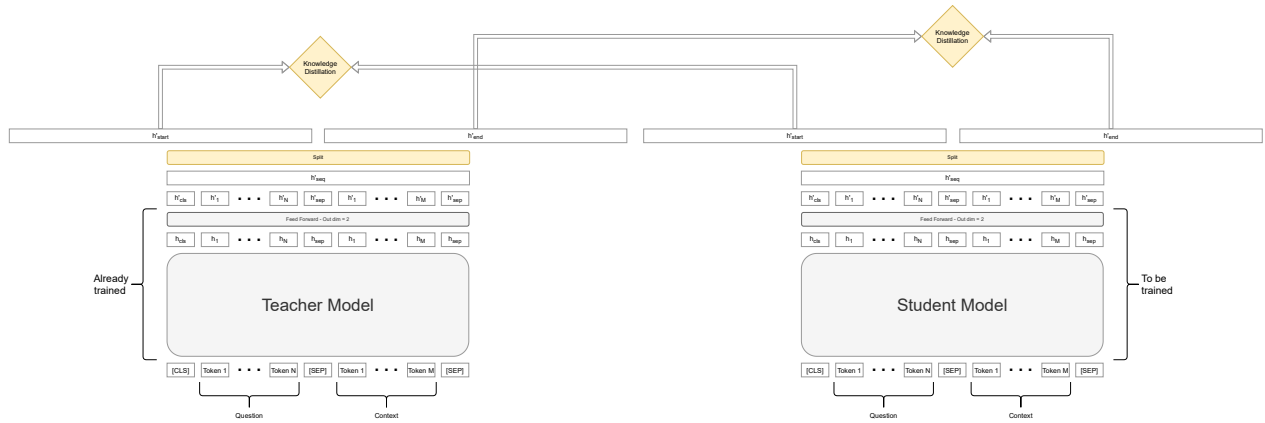


Figure D.2: Implementation of knowledge distillation for question answering datasets with models that share the same vocabulary.

Annex E

Effect of Caching Teacher Outputs During Training

A significant challenge in our experimental study is the use of large and costly language models as teacher models for our faster and lighter models. Despite this, as discussed in Appendix F, the importance of these teacher models is essential for achieving better results with our proposed models.

Thus, the use of these teacher models poses challenges in terms of experimentation, particularly when working with restricted budgets, as is often the case in research outside big tech companies. To mitigate this issue, we implement a cache for the outputs of the teacher model, which allows us to train and experiment more efficiently.

The idea behind this approach is straightforward: since the teacher model is fixed during training, its outputs on an input x remain unchanged during different epochs, allowing us to compute them once and reuse them in subsequent epochs.

Formally, suppose F_t and F_s represent the computational cost of the forward pass for the teacher and student models, respectively, on an entire dataset, and E is the number of epochs used to train our proposed models. By caching the teacher’s output, the total cost of computing the forward pass reduces from $O(E \cdot (F_t + F_s))$ to $O(F_t + E \cdot F_s)$.

It is worth noting that typically $F_t \gg F_s$, and the number of epochs used in knowledge distillation is often higher than that used in simple fine-tuning. To illustrate, our fine-tuning experiments employ between 2 and 4 epochs, while our knowledge distillation experiments use a maximum of 50 epochs.

To evaluate the impact of our cache implementation, we compare the training times of our proposed models on the XNLI dataset, which is the largest dataset considered in this study, for only 5 epochs (1/10 of the epochs used in our primary experiments) when using the cache and when not using it. Table E.1 reports the results of this experiment, presenting the mean (noted as M) and standard deviation (noted as SD) over three runs. As expected, the use of the cache reduces the training time significantly, with results indicating that training time is approximately 1/4 of the time required to train without a cache. This reduction in training

Model	Training Time (hours)			
	Cache		No Cache	
	M	SD	M	SD
ALBETO <i>tiny</i>	3.8	3.1×10^{-2}	16.2	3.1×10^{-3}
ALBETO <i>base-2</i>	3.8	1.6×10^{-3}	16.3	3.6×10^{-3}
ALBETO <i>base-4</i>	4.2	3.3×10^{-4}	16.6	2.6×10^{-3}
ALBETO <i>base-6</i>	4.5	1.5×10^{-3}	17.0	1.5×10^{-3}
ALBETO <i>base-8</i>	4.8	1.9×10^{-4}	17.3	5.8×10^{-3}
ALBETO <i>base-10</i>	5.3	9.6×10^{-3}	17.6	5.6×10^{-3}

Table E.1: Training times when using teacher cache vs not using it. Table report the mean (M) and standard deviation (SD) over three runs.

time is expected since the forward pass of the teacher model is the most costly operation and is computed only in the first epoch and then retrieved in the next 4 epochs. Furthermore, this difference will increase as the number of epochs increases.

In conclusion, while our cache implementation is a simple engineering trick, it has a significant impact on our experimentation phase in terms of training time and required compute.

Annex F

Importance of Knowledge Distillation

In addition to other experiments, we conducted ablation experiments to evaluate the contribution of Task-Specific Knowledge Distillation to the results of our faster models based on ALBETO.

Tables F.1, F.2, and F.3 compare the performance of each of our proposed models under two training settings: regular fine-tuning (FT) and task-specific knowledge distillation (KD). For fine-tuning, we followed the experimental setup described in Section 4.1, while we followed the setup described in Section 4.2 for KD.

Overall, our results indicate that training using KD generally yields better results than simple fine-tuning, except for sequence tagging tasks (POS, NER), where the results are mixed.

Table F.1 presents the results of text classification tasks, where we observe that KD outperforms fine-tuning. In MLDoc, which is hypothesized as an easier task, the performance is similar for both training schemes and different models. However, in PAWS-X and XNLI, we observe a significant difference between the fine-tuning and KD training schemes.

Table F.2 presents the results for sequence tagging tasks, where the performance of models under the KD and fine-tuning settings are mixed. Unlike other types of tasks, where the KD training method is the clear winner, the results here vary. In the case of NER, faster models perform better under the fine-tuning setting, while those with larger compute requirements perform better under the KD setting.

Finally, Table F.3 presents the results for question answering, where we observe that models trained using KD generally exhibit better performance than those trained using simple fine-tuning, with a significant difference of around 3-4 percentage points, depending on the model and dataset.

In summary, our results underscore the significance of KD, particularly for harder tasks where the effect is more pronounced, allowing for lighter and faster models to achieve better task performance.

Model	MLDoc		PAWS-X		XNLI	
	FT	KD	FT	KD	FT	KD
ALBETO <i>tiny</i>	95.82	96.40	80.20	85.05	73.43	75.99
ALBETO <i>base-2</i>	94.67	96.20	73.45	76.75	72.08	73.65
ALBETO <i>base-4</i>	95.88	96.35	82.90	86.40	75.83	78.68
ALBETO <i>base-6</i>	95.88	96.40	85.20	88.45	78.42	81.66
ALBETO <i>base-8</i>	95.82	96.70	87.30	89.75	79.44	82.55
ALBETO <i>base-10</i>	95.65	96.88	88.80	89.95	79.62	82.26

Table F.1: Comparison of the performance of our proposed models on text classification tasks on two settings: fine-tuning and task-specific knowledge distillation.

Model	POS		NER	
	FT	KD	FT	KD
ALBETO <i>tiny</i>	97.34	97.36	75.42	72.51
ALBETO <i>base-2</i>	97.46	97.17	71.70	69.69
ALBETO <i>base-4</i>	97.87	97.60	76.18	74.58
ALBETO <i>base-6</i>	98.03	97.82	78.10	78.41
ALBETO <i>base-8</i>	98.18	97.96	79.46	80.23
ALBETO <i>base-10</i>	98.17	98.00	80.46	81.10

Table F.2: Comparison of the performance of our proposed models on sequence tagging tasks on two settings: fine-tuning and task-specific knowledge distillation.

Model	MLQA		SQAC		TAR, XQuAD	
	FT	KD	FT	KD	FT	KD
ALBETO <i>tiny</i>	51.84 / 28.28	54.17 / 32.22	59.28 / 39.16	63.03 / 43.35	66.43 / 45.71	67.47 / 46.13
ALBETO <i>base-2</i>	45.97 / 23.60	48.62 / 26.17	53.32 / 34.34	58.40 / 39.00	61.82 / 40.67	63.41 / 42.35
ALBETO <i>base-4</i>	59.99 / 35.69	62.19 / 38.28	65.66 / 45.54	71.41 / 52.87	68.91 / 49.07	73.31 / 52.43
ALBETO <i>base-6</i>	63.75 / 38.58	66.35 / 42.01	72.22 / 53.61	76.99 / 59.00	74.33 / 52.68	75.59 / 54.95
ALBETO <i>base-8</i>	64.99 / 40.58	67.39 / 42.94	75.22 / 56.43	77.79 / 59.63	75.47 / 54.11	77.89 / 56.72
ALBETO <i>base-10</i>	66.29 / 41.69	68.29 / 44.29	77.14 / 59.21	79.89 / 62.04	77.06 / 56.47	78.21 / 56.21

Table F.3: Comparison of the performance of our proposed models on question answering on two settings: fine-tuning and task-specific knowledge distillation.