

# INFORME DETALLADO

SELIN COCARCA

1. INTRODUCCIÓN
2. INSTALACIÓN DE LIBRERÍAS
3. DESCRIPCIÓN DEL DATASET
4. CONFIGURACIÓN DE LA SEMILLA
5. PROCESAMIENTO DE DATOS
6. DATA COLLATOR
7. ENTRENAMIENTO Y EVALUACIÓN DEL MODELO
8. IMPLEMENTACIÓN DE LA SOLUCIÓN

## 1. INTRODUCCIÓN:

En este proyecto que tenemos delante, se detalla el desarrollo de un sistema de traducción automática utilizando técnicas avanzadas de procesamiento de lenguaje natural. El objetivo principal es implementar un modelo preentrenado llamado MarianMT, reconocido por su eficacia en tareas de traducción. A lo largo de este informe, exploraremos todos los aspectos clave del proyecto: desde el conjunto de datos utilizados y la selección del modelo del lenguaje, hasta los métodos de evaluación aplicados y los resultados obtenidos.

## 2. INSTALACIÓN DE LIBRERÍAS:

Para comenzar, este proyecto se han instalado diversas bibliotecas fundamentales que son esenciales para el desarrollo de un sistema de traducción automática avanzado. Estas bibliotecas proporcionan las herramientas necesarias para el procesamiento de lenguaje natural, la implementación de modelos de traducción y la evaluación de resultados. A continuación, se detallan las bibliotecas instaladas:

1. **Langdetect:** Esta biblioteca es útil para detectar automáticamente el idioma de un texto de entrada. En un sistema de traducción automática, permite identificar el idioma de origen de un texto antes de aplicar el modelo de traducción correspondiente. Esto es crucial para dirigir correctamente el texto al modelo adecuado para la traducción.
2. **Transformers y torch:** PyTorch junto con Hugging Face Transformers son esenciales para implementar y desplegar modelos de traducción automática como MarianMT. Estas bibliotecas proporcionan herramientas para cargar modelos preentrenados, realizar inferencias para la traducción y entrenar modelos personalizados si es necesario. Además, Transformers ofrece una interfaz fácil de usar para trabajar con modelos complejos de lenguaje.
3. **tqdm:** Esta biblioteca proporciona barras de progreso durante las iteraciones en el procesamiento de datos. En el contexto de la traducción automática, es útil para monitorear el progreso de la tokenización, el entrenamiento del modelo y la evaluación de los datos. Mejora la visualización y la gestión del tiempo durante estas tareas intensivas en cómputo.
4. **Datasets:** Facilita el acceso a conjuntos de datos estándar utilizados para entrenar y evaluar modelos de traducción automática. Proporciona una interfaz para cargar datos de forma eficiente, realizar divisiones de conjunto de entrenamiento y validación, y gestionar metadatos asociados con los conjuntos de datos.
5. **Sentence-transformers:** Utilizada para obtener representaciones vectoriales de alta calidad de frases y documentos. Estos vectores son útiles en tareas como la recuperación de información y la búsqueda semántica, lo cual puede complementar la traducción automática mejorando la coherencia semántica y la precisión de las traducciones.
6. **Pandas:** Esencial para la manipulación y análisis de datos estructurados. En el contexto de la traducción automática, pandas se utiliza para organizar y visualizar resultados de manera eficiente, lo cual es crucial en la fase de evaluación y análisis posterior al entrenamiento de los modelos.
7. **Umap y hdbscan:** Utilizados en técnicas de reducción dimensional y clustering respectivamente. Son útiles para la visualización y el análisis exploratorio de los datos lingüísticos, permitiendo explorar relaciones entre palabras o documentos que pueden ser útiles para mejorar modelos de traducción automática.
8. **Sacre Moses:** Herramienta de tokenización y normalización de texto utilizada en el preprocesamiento de datos. Prepara el texto de entrada para la tokenización y el procesamiento por parte de modelos de traducción automática, asegurando que los datos estén en el formato adecuado para la tarea.
9. **Accelerate:** Proporciona herramientas para mejorar la velocidad y eficiencia en el entrenamiento de modelos de aprendizaje automático en PyTorch. Esto es crucial en el desarrollo de modelos de traducción automática para reducir el tiempo de entrenamiento y mejorar el rendimiento en grandes volúmenes de datos.

En resumen, cada una de estas bibliotecas desempeña un papel crucial en diferentes etapas del desarrollo y la evaluación de sistemas de traducción automática, desde la preparación inicial de datos hasta la evaluación final del rendimiento del modelo.

### 3. DESCRIPCIÓN DEL DATASET:

El proyecto se basa en un conjunto de datos multilingüe accesible públicamente, que consta de aproximadamente 997 ejemplos de frases en español, cada una acompañada de sus respectivas traducciones a dos idiomas adicionales: euskera (vasco) e inglés. Este conjunto de datos es fundamental para entrenar y evaluar modelos de traducción automática capaces de manejar varios idiomas de manera efectiva.

Para manipular y procesar estos datos, se importan diversas bibliotecas y módulos clave:

1. **Random** : Se utiliza para generar números aleatorios. Es útil en la inicialización aleatoria de parámetros durante el entrenamiento de modelos de traducción automática o para generar datos de prueba aleatorios que pueden ser útiles en validación o pruebas.

2. **Os** : Proporciona funciones para interactuar con el sistema operativo. Es esencial para acceder a variables de entorno donde se almacenan configuraciones sensibles como claves de API y para gestionar archivos y directorios que contienen datos de entrenamiento y modelos.

3. **Numpy as np** : NumPy se utiliza para operaciones numéricas eficientes en Python. En el contexto de traducción automática, es fundamental para manipulación de matrices durante el procesamiento de datos y cálculos matemáticos avanzados utilizados en algoritmos de optimización y evaluación.

4. **Torch** : PyTorch es un framework de machine learning que proporciona estructuras de datos optimizadas para tensores y soporte para cómputo en GPU. Esencial para ejecutar operaciones eficientes en modelos de deep learning utilizados en traducción automática, especialmente cuando se manejan modelos neuronales grandes y complejos.

5. **Datasets from datasets import load\_dataset, load\_metric** : Estas funciones de la biblioteca Hugging Face datasets permiten cargar conjuntos de datos estructurados y métricas de evaluación asociadas. Simplifican la manipulación y evaluación de datos, proporcionando un acceso rápido y eficiente a los conjuntos de datos multilingües necesarios para entrenar y evaluar modelos de traducción automática.

6. **Transformers from transformers import MarianTokenizer, MarianMTModel, Trainer, TrainingArguments** : Estas clases específicas de Hugging Face Transformers son fundamentales para modelar traducción neuronal y entrenar modelos de NLP utilizando arquitecturas preentrenadas. Permiten la tokenización de texto en varios idiomas y la implementación de modelos de traducción automática avanzados basados en redes neuronales.

8. **Sentence\_transformers from sentence\_transformers import SentenceTransformer** : Esta biblioteca proporciona modelos preentrenados para codificar oraciones y calcular embeddings semánticos. Útil en tareas de traducción automática para capturar la semántica del texto y mejorar la precisión de las traducciones mediante la comparación de similitudes semánticas entre frases.

9. **Umap import umap** : Implementa el algoritmo UMAP para la reducción de dimensionalidad y visualización de datos de alta dimensión. Es útil para explorar y comprender estructuras complejas de datos lingüísticos multilingües, permitiendo una visualización efectiva de la distribución de los embeddings generados por modelos de traducción automática.

10. **Hdbscan import hdbscan** : Implementa el algoritmo HDBSCAN para el agrupamiento de datos basado en densidad en espacios de alta dimensionalidad. Esencial para identificar clusters significativos y potencialmente no lineales en los datos de traducción automática, facilitando el análisis y la interpretación de los resultados obtenidos por los modelos.

#### 4. CONFIGURACIÓN DE LA SEMILLA:

En el proyecto, se implementa la configuración de la semilla (seed) para asegurar la reproducibilidad de los resultados, un aspecto crucial en machine learning y deep learning donde la aleatoriedad puede influir en los resultados experimentales. A continuación, se detalla la función `seed_everything` y su aplicación con una semilla específica:

##### **Función `seed_everything`:**

La función `seed_everything` se utiliza para establecer la semilla en varias bibliotecas y entornos de Python, garantizando consistencia en la generación de números aleatorios:

**1.1 `random.seed(seed)`:** Configura la semilla del generador de números aleatorios de Python (`random.seed`), asegurando la reproducibilidad de secuencias aleatorias.

**1.2 `os.environ['PYTHONHASHSEED'] = str(seed)`:** Establece la semilla para la generación de hashes en Python a través de la variable de entorno `PYTHONHASHSEED`, manteniendo la consistencia en la creación de hashes entre ejecuciones.

**1.3 `np.random.seed(seed)`:** En NumPy, `np.random.seed(seed)` fija la semilla para la generación de números aleatorios, asegurando resultados consistentes en operaciones numéricas complejas.

**1.4 `torch.manual_seed(seed)`:** En PyTorch, `torch.manual_seed(seed)` establece la semilla para la generación de números aleatorios en la CPU de manera determinística.

**1.5 `torch.cuda.manual_seed(seed)` y `torch.backends.cudnn.deterministic = True`:** Estas líneas aseguran la reproducibilidad en operaciones GPU. `torch.cuda.manual_seed(seed)` fija la semilla para la GPU, mientras `torch.backends.cudnn.deterministic = True` asegura operaciones determinísticas en la librería cuDNN de PyTorch.

##### **Llamada a `seed_everything(357)`:**

Al llamar `seed_everything(357)`, se especifica la semilla (357 en este caso). Esta acción garantiza que todas las operaciones que dependen de números aleatorios utilicen esta semilla específica. Como resultado:

**Reproducibilidad:** Todas las ejecuciones del código con esta semilla producirán los mismos resultados, esencial para validar y comparar experimentos de machine learning de manera coherente.

**Consistencia:** Se evitan variaciones no deseadas debido a la aleatoriedad, proporcionando estabilidad y fiabilidad en los resultados del modelo.

En resumen, la configuración adecuada de la semilla mediante `seed_everything` asegura la coherencia y la reproducibilidad en todas las operaciones que involucran aleatoriedad, facilitando un desarrollo robusto y confiable en proyectos de machine learning.

## 5. PROCESAMIENTO DE DATOS:

El procesamiento de datos es una etapa crítica en cualquier proyecto de procesamiento de lenguaje natural (NLP). En este proyecto, hemos utilizado diversas bibliotecas y funciones para cargar, limpiar, explorar y preparar los datos en inglés, español y euskera en formato Latin. A continuación, se explica detalladamente el propósito y el funcionamiento de cada paso del procesamiento de datos y las razones detrás de las elecciones hechas.

Para comenzar, hemos utilizado la función `'load_dataset'` del módulo `'datasets'` para cargar conjuntos de datos en los idiomas mencionados. Esta función es fundamental porque facilita la carga de conjuntos de datos predefinidos o personalizados, permitiendo un acceso rápido y eficiente a los datos necesarios para entrenar y evaluar modelos de NLP. Específicamente, hemos cargado los siguientes conjuntos de datos:

- `'load_dataset('facebook/flores', 'eus_Latn')`: Carga el conjunto de datos en euskera (vasco) escrito en caracteres latinos.
- `'load_dataset('facebook/flores', 'spa_Latn')`: Carga el conjunto de datos en español escrito en caracteres latinos.
- `'load_dataset('facebook/flores', 'eng_Latn')`: Carga el conjunto de datos en inglés escrito en caracteres latinos.

La elección de estos conjuntos de datos específicos del proyecto 'facebook/flores' se debe a su calidad y a la riqueza de sus datos etiquetados, lo que es crucial para tareas de NLP como la traducción automática y la generación de texto.

Una vez cargados los conjuntos de datos, el siguiente paso es explorar y limpiar estos datos. Esto se realiza mediante una función personalizada que itera a través de las particiones 'dev' de cada conjunto de datos, verifica la existencia de traducciones y recopila las frases disponibles. Este proceso es esencial para asegurar que los datos sean de alta calidad y estén listos para ser utilizados en el entrenamiento de los modelos. A continuación se describe el funcionamiento detallado:

- 1. Iteración a través de los conjuntos de datos:** Se recorre cada entrada en la partición 'dev' de los conjuntos de datos en inglés, español y euskera.
- 2. Verificación de la existencia de traducción:** Se comprueba si cada ejemplo contiene información de traducción válida.
- 3. Impresión de las frases:** Se imprime el texto correspondiente si se encuentra una traducción válida; de lo contrario, se indica que el ejemplo no contiene información de traducción.

El propósito de este paso es visualizar rápidamente las frases disponibles y evaluar la calidad y la cobertura de los datos etiquetados antes de proceder con tareas más avanzadas de NLP.

Luego, se realiza una limpieza y exploración más profunda de los datos mediante una función dedicada. Esta función itera a través de las entradas en la partición 'dev', verifica la existencia de la clave 'sentence' en cada ejemplo y recopila las frases en una lista. Este paso es crucial para asegurar que solo se utilicen datos válidos y bien estructurados en las siguientes fases del proyecto.

Para una mejor visualización y análisis, las frases limpias se almacenan en un DataFrame utilizando la biblioteca Pandas. Esto proporciona una estructura clara y legible para los datos, facilitando su revisión y análisis inicial. El DataFrame contiene columnas para cada idioma y muestra las primeras filas de los datos para una inspección rápida.

En la etapa de tokenización, se utiliza la función `'tokenize_text'` que emplea modelos preentrenados específicos para cada idioma, cargados con la ayuda de `'MarianTokenizer'` y `'AutoTokenizer'` de la biblioteca `'transformers'`. La tokenización es un paso crucial porque

convierte el texto en una secuencia de tokens que el modelo puede procesar. Se seleccionan modelos adecuados para cada idioma y se tokenizan las frases con opciones de padding y truncamiento, retornando los resultados como tensores de PyTorch.

Finalmente, para asegurar que la tokenización se ha realizado correctamente, se utiliza una función de verificación. Esta función comprueba la forma y el contenido de los tokens generados, asegurando que estén listos para su uso en el entrenamiento y evaluación de los modelos de NLP. Este paso garantiza que los datos tokenizados sean adecuados y estén en el formato esperado para las tareas posteriores de procesamiento de lenguaje natural.

En resumen, el procesamiento de datos en este proyecto incluye la carga, limpieza, exploración, tokenización y verificación de datos en inglés, español y euskera. Cada paso es esencial para asegurar la calidad y preparación de los datos antes de su uso en modelos avanzados de NLP, lo que facilita el desarrollo de aplicaciones precisas y eficientes en el ámbito del procesamiento de lenguaje natural.

## 6. DATA COLLATOR:

El código proporcionado define `MiDataCollator`, una clase que extiende `DataCollatorForSeq2Seq` de la biblioteca Transformers de Hugging Face. A continuación, se presenta una explicación detallada sobre su propósito y funcionamiento, así como las razones detrás de la elección de cada componente utilizado en el proyecto.

Primero, se importan las bibliotecas necesarias: `torch` de PyTorch, `DataCollatorForSeq2Seq` para el manejo de lotes en tareas de secuencia a secuencia (seq2seq), y herramientas de tokenización como `AutoTokenizer` y `BatchEncoding`. La elección de estas bibliotecas y herramientas se debe a su eficacia y amplia adopción en la comunidad de NLP, lo que garantiza un rendimiento sólido y una integración fluida con otros componentes del proyecto.

La clase `MiDataCollator` se define para mejorar la preparación de los datos de entrada para modelos de seq2seq. Esto es esencial porque los modelos de secuencia a secuencia, como los utilizados en traducción automática y otras tareas de NLP, requieren que las secuencias de entrada sean de longitud uniforme para un procesamiento eficiente en GPU. A continuación se explica cada componente de `MiDataCollator` y su importancia en el proyecto.

El método `__init__` actúa como constructor de la clase. Inicializa `MiDataCollator` con un tokenizer, una herramienta crucial que convierte texto en tokens numéricos que los modelos pueden procesar. Se pasa `None` para el modelo en este constructor para permitir una mayor flexibilidad en la adaptación del collator a diferentes modelos de seq2seq. Esta flexibilidad es vital para experimentar con diferentes configuraciones y optimizar el rendimiento del modelo.

El método `__call__` permite que `MiDataCollator` sea llamado como una función. Recibe `features`, una lista de diccionarios que contienen `input_ids` y `attention_mask`, que son representaciones tokenizadas de las secuencias de entrada y sus respectivas máscaras de atención. Esto permite que `MiDataCollator` maneje dinámicamente diferentes lotes de datos de entrada, lo cual es crucial para adaptarse a la variabilidad inherente en las secuencias de texto natural.

Dentro del método `__call__`, se utiliza `BatchEncoding` para crear una codificación de lotes con las características proporcionadas. Esto estructura los datos de manera que puedan ser procesados uniformemente. Luego, se aplica padding dinámico utilizando el tokenizer para asegurar que todas las secuencias en el lote tengan la misma longitud. Este padding es dinámico porque se adapta a la longitud de la secuencia más larga en cada lote, minimizando así el desperdicio de memoria y mejorando la eficiencia computacional durante el entrenamiento o la evaluación.

El método finalmente devuelve un diccionario con `'input_ids'` y `'attention_mask'` después de aplicar el padding. Esta salida está lista para ser utilizada directamente en el entrenamiento o evaluación de modelos de secuencia a secuencia, asegurando que los datos estén en el formato y longitud adecuados para el procesamiento eficiente en GPU. Este enfoque no solo mejora la uniformidad de las secuencias, sino que también optimiza el uso de recursos computacionales, lo cual es esencial para trabajar con grandes volúmenes de datos y modelos complejos en NLP.

Para ilustrar el uso de `'MiDataCollator'`, se proporciona un ejemplo donde se carga un tokenizer preentrenado y se instancia `'MiDataCollator'` con dicho tokenizer. Luego, se crea una lista de `'features'` con datos de entrada y máscaras de atención de secuencias. Al llamar a `'MiDataCollator'` con estas características, se obtiene un lote con padding aplicado, demostrando la utilidad práctica de esta clase en la preparación de datos para modelos de seq2seq.

En resumen, `'MiDataCollator'` facilita la preparación de lotes para modelos de secuencia a secuencia al aplicar padding dinámico, asegurando uniformidad en la longitud de las secuencias y eficiencia en el procesamiento en GPU. Esto es esencial para tareas de NLP como la traducción automática, donde la coherencia en la longitud de las secuencias mejora el rendimiento del modelo y optimiza el uso de recursos computacionales. La elección de `'DataCollatorForSeq2Seq'`, junto con herramientas de tokenización avanzadas, asegura que el proyecto se beneficie de las mejores prácticas y tecnologías disponibles en el campo del procesamiento de lenguaje natural.

## 7. ENTRENAMIENTO Y EVALUACIÓN DEL MODELO:

En este fragmento de código, creamos y dividimos un conjunto de datos combinado que contiene tokens de varios idiomas: español, inglés y euskera. Primero, importamos la clase `'Dataset'` de la biblioteca `'datasets'`, lo cual nos permite manipular conjuntos de datos de manera eficiente en Python. La elección de esta biblioteca se debe a su robustez y facilidad de uso para manejar datos en tareas de procesamiento de lenguaje natural (NLP).

Creamos un `'combined_dataset'` utilizando la función `'from_dict'`, donde cada idioma tiene asociado un conjunto de tokens. Este paso nos permite consolidar datos multilingües en un único dataset, lo cual es crucial para trabajar con datos estructurados en diferentes idiomas. Esto es especialmente importante en proyectos de NLP que requieren manejar múltiples idiomas simultáneamente, como la traducción automática. Al combinar los datos, facilitamos la integración y el procesamiento uniforme de los mismos en etapas posteriores.

Dividimos este `'combined_dataset'` en conjuntos de entrenamiento y validación utilizando el método `'train_test_split'`. Específicamente, asignamos el 80% de los datos al conjunto de entrenamiento (`'train_dataset'`) y el 20% restante al conjunto de validación (`'val_dataset'`). Esta división es fundamental para evaluar la capacidad de generalización de nuestro modelo después del entrenamiento. Al reservar un conjunto de validación, podemos medir el rendimiento del modelo en datos no vistos, lo que es esencial para evitar el sobreajuste y asegurar que el modelo funcione bien en datos del mundo real.

Para entrenar el modelo de traducción automática, utilizamos el transformer MarianMT de Hugging Face. Este modelo preentrenado está específicamente diseñado para tareas de traducción entre idiomas romances e inglés. La elección de MarianMT se basa en su alta precisión y capacidad de manejar múltiples pares de idiomas, lo que lo hace adecuado para nuestro proyecto multilingüe. Utilizamos `'MarianTokenizer'` para tokenizar el texto en un formato compatible con el modelo, asegurando que las secuencias de entrada sean adecuadamente procesadas y que el modelo pueda interpretar correctamente los datos.

Para la carga y preparación de datos, utilizamos conjuntos de datos de traducción en varios idiomas del conjunto de datos Flores de Facebook. Los conjuntos de datos de desarrollo (dev) y prueba de desarrollo (devtest) se concatenan para cada idioma. Esto permite maximizar la



cantidad de datos disponibles para el entrenamiento y evaluación, mejorando así la calidad del modelo. Al extraer frases únicas del dataset combinado, eliminamos duplicados y aseguramos que los datos sean variados y representativos de los idiomas de interés.

Definimos una clase de `'Dataset'` personalizada llamada `'CustomDataset'` para manejar pares de texto fuente y destino. Esta clase tokeniza los textos utilizando `'MarianTokenizer'`, asegurando que ambos textos estén alineados en longitud y formato antes de ser pasados al modelo. Este paso es crucial para mantener la integridad de los datos y garantizar que el modelo reciba entradas coherentes y adecuadamente formateadas para el entrenamiento.

Preparamos un `'DataLoader'` para manejar el conjunto de datos completo durante el entrenamiento. Este `'DataLoader'` agrupa las instancias en batches, maneja el padding para asegurar que todos los textos tengan la misma longitud y aleatoriza los batches para mejorar la generalización del modelo. El uso de `'DataLoader'` es esencial para manejar eficientemente grandes volúmenes de datos y asegurar que el modelo se entrene de manera robusta y eficiente.

Para el entrenamiento del modelo, utilizamos el optimizador Adam (`'optim.Adam'`), conocido por su eficacia en ajustar los parámetros del modelo durante el entrenamiento. Utilizamos una tasa de aprendizaje de 0.001 para balancear entre la velocidad de convergencia y la estabilidad del entrenamiento. Entrenamos el modelo durante 10 épocas, iterando a través de los batches de datos de entrenamiento. En cada época, calculamos la pérdida (loss) de predicción y ajustamos los gradientes utilizando backpropagation para optimizar el modelo. Este proceso iterativo es fundamental para mejorar la precisión del modelo y reducir los errores de predicción.

Al finalizar el entrenamiento, guardamos el estado del modelo utilizando `'torch.save'`. Esto nos permite almacenar el modelo entrenado en un archivo (`'modelo_entrenado.pth'`), facilitando su carga y uso posterior para traducción o generación de nuevas predicciones. Guardar el modelo es una práctica estándar que asegura la reproducibilidad y permite reutilizar el modelo entrenado sin necesidad de reentrenarlo desde cero.

Finalmente, el código evalúa el modelo utilizando un `'DataLoader'` para la validación y calcula métricas como precisión, recall y F1-score. Durante la evaluación, generamos predicciones y comparamos los resultados con las etiquetas reales para medir el rendimiento del modelo. Estas métricas son cruciales para evaluar la efectividad del modelo y ajustar su rendimiento. Posteriormente, simulamos varias épocas de entrenamiento y almacenamos las métricas calculadas en DataFrames de pandas para su visualización. Esto proporciona una manera completa de entrenar y evaluar modelos de aprendizaje automático, mostrando cómo calcular y mostrar métricas importantes como parte del proceso de evaluación del modelo.

En resumen, cada componente y paso del código está diseñado para manejar eficientemente datos multilingües, entrenar un modelo de traducción automática de alta calidad y evaluar su rendimiento de manera rigurosa. La integración de estas técnicas y herramientas en el proyecto asegura que podamos desarrollar y refinar un modelo capaz de manejar tareas complejas de NLP en múltiples idiomas.

## 8. IMPLEMENTACIÓN DE LA SOLUCIÓN:

Este código Python utiliza modelos preentrenados de Hugging Face para realizar traducciones automáticas entre varios idiomas. Define funciones para cargar modelos y tokenizadores, traducir texto utilizando estos modelos, y presenta una interfaz gráfica interactiva en Jupyter/Colab para ingresar texto y seleccionar los idiomas de origen y destino. Al hacer clic en "Traducir", muestra la traducción resultante en pantalla, permitiendo ajustes dinámicos del tamaño del área de texto según el contenido ingresado.

Primero, importamos las bibliotecas necesarias como ``sys``, ``os`` y ``contextlib``, y de la biblioteca ``transformers`` importamos ``MarianMTModel`` y ``MarianTokenizer`` para cargar los modelos y tokenizadores preentrenados. También se importan widgets de ``ipywidgets`` para crear una interfaz gráfica interactiva y ``warnings`` para suprimir advertencias.

Definimos un context manager llamado ``suppress_output`` para redirigir ``stdout`` y ``stderr`` a ``os.devnull``, lo cual nos permite suprimir las salidas innecesarias durante la carga de los modelos y evitar que se muestren mensajes de advertencia o errores en la interfaz.

Luego, creamos un diccionario ``model_names`` que mapea pares de idiomas (origen y destino) a nombres de modelos específicos de Hugging Face. Esto nos permite seleccionar el modelo correcto para la traducción en función de los idiomas seleccionados por el usuario. Este enfoque asegura que se utilice el modelo adecuado para cada par de idiomas, optimizando la precisión de las traducciones.

Definimos la función ``get_model_name`` que toma los idiomas de origen y destino como argumentos y devuelve el nombre del modelo correspondiente del diccionario ``model_names``. Esto simplifica la lógica de selección del modelo en las funciones posteriores.

La función ``load_model_and_tokenizer`` carga el modelo y el tokenizador adecuados basándose en los idiomas de origen y destino proporcionados. Utiliza el context manager ``suppress_output`` para suprimir salidas durante la carga. Si no se puede encontrar o cargar el modelo, se lanza una excepción con un mensaje de error descriptivo. Esta función es esencial para preparar el modelo y el tokenizador necesarios para la traducción.

Para traducir texto, definimos la función ``translate_text``. Esta función toma el texto de entrada y los idiomas de origen y destino, carga el modelo y el tokenizador adecuados, y realiza la traducción. Primero, se verifica que el texto de entrada no esté vacío. Luego, se tokeniza el texto de entrada, se generan las traducciones utilizando el modelo y finalmente se decodifica la traducción generada para convertirla de nuevo a texto legible. Esta función maneja errores y excepciones, devolviendo mensajes de error en caso de fallos.

La función ``translate_and_display`` llama a ``translate_text`` y muestra la traducción resultante en la interfaz. Esto asegura que el proceso de traducción se realice de manera fluida y que los resultados se presenten claramente al usuario.

Para mejorar la experiencia del usuario, definimos la función ``update_textarea_size`` que ajusta dinámicamente el tamaño del área de texto en función de la cantidad de líneas de texto ingresadas. Esto permite que el área de texto se expanda o contraiga automáticamente, proporcionando una interfaz más amigable y adaptable.

Creamos varios widgets de ``ipywidgets`` para la interfaz gráfica. ``text_input`` es un área de texto para ingresar el texto a traducir, ``src_lang_dropdown`` y ``tgt_lang_dropdown`` son menús desplegables para seleccionar los idiomas de origen y destino, y ``translate_button`` es un botón para iniciar la traducción. ``output_area`` se utiliza para mostrar la traducción resultante. Estos widgets permiten al usuario interactuar con la aplicación de manera intuitiva.

Definimos la función ``on_translate_button_clicked`` que se ejecuta cuando se hace clic en el botón de traducir. Esta función limpia el área de salida, llama a ``translate_and_display`` y maneja posibles errores, mostrando mensajes descriptivos si algo sale mal.

Finalmente, observamos cambios en el ``text_input`` para ajustar dinámicamente su tamaño y organizamos todos los widgets en una disposición vertical utilizando ``VBox`` de ``ipywidgets``. Esta disposición se muestra al usuario, proporcionando una interfaz interactiva y fácil de usar para realizar traducciones automáticas entre varios idiomas.