



UNIVERSIDAD
DE GRANADA

TRABAJO FIN DE GRADO
GRADO EN INGENIERÍA INFORMÁTICA

Dispositivo para detección de escritura mediante Deep Learning en un sistema empotrado

SmartPen

Autor
Antonio Priego Raya

Directores
Jesús González Peñalver
Juan José Escobar Pérez



GitHub del proyecto



ESCUELA TÉCNICA SUPERIOR DE INGENIERÍAS INFORMÁTICA Y DE
TELECOMUNICACIÓN

—
Granada, Julio de 2022

Antonio Priego Raya



Dispositivo para detección de escritura mediante Deep Learning en un sistema empotrado

SmartPen

Autor
Antonio Priego Raya

Directores
Juan José Escobar Pérez
Jesús González Peñalver

Detección de escritura mediante Deep Learning en un sistema empotrado

SmartPen

Antonio Priego Raya

Palabras clave: TinyML, Machine learning, Deep learning, Sistemas empotrados, Reconocimiento letras, Redes neuronales convolucionales ...

Resumen

Creación de un dispositivo autónomo con forma de lápiz en el que integrar un sistema empotrado, concretamente se hará uso de la *Arduino Nano Sense 33 BLE*. El propósito de este dispositivo será la detección de letras en tiempo real, registrando el movimiento del dispositivo y rasterizando el mismo para procesarlo y clasificarlo.

Para este procesamiento y clasificación de las letras registradas, se empleará *deep learning*, concretamente un modelo de red neuronal convolucional. Con la particularidad de ejecutar el procesamiento del modelo en el propio dispositivo, para dotarlo de autonomía.

Esta autonomía está complementada por una batería que garantiza su independencia. Por tanto se desarrollarán todos los pasos propios del trabajo con redes neuronales: diseño del modelo, recolección de datos, entrenamiento del modelo, su testeo, etc.

Complementario al dispositivo, también se creará un interfaz de usuario donde acceder a las funciones que se desarrollen para este dispositivo.

Deep Learning handwriting detection in an embedded system

SmartPen

Antonio Priego Raya

Keywords: TinyML, Machine learning, Deep learning, Embedded systems, Letter detection, Convolutional neural networks

Abstract

Creation of an autonomous device in the shape of a pencil in which to integrate an embedded system, specifically using the Arduino Nano Sense 33 BLE. The purpose of this device is to detect letters in real time, recording the movement of the device and rasterizing it to be processed and classified.

A Deep Learning model will be used for this letter processing and classification, specifically a convolutional neural network model. With the particularity of executing the processing of the model in the device itself, in order to provide it with autonomy.

This autonomy is complemented using a battery that guarantees its independence. Therefore, all the steps involved in working with neural networks will be developed: design of the model, data collection, model training, testing, etc.

Complementary to the device, a user interface will also be created to access the functions developed for this device.

Yo, **Antonio Priego Raya**, alumno de la titulación *Grado en Ingeniería Informática de la Escuela Técnica Superior de Ingenierías Informática y de Telecomunicación de la Universidad de Granada*, con DNI 31033948W, autorizo la ubicación de la siguiente copia de mi Trabajo Fin de Grado en la biblioteca del centro para que pueda ser consultada por las personas que lo deseen.

Fdo: Antonio Priego Raya

Granada a 8 de Julio de 2022

D. **Jesús González Peñalver**, Catedrático del departamento de Arquitectura y Tecnología de Computadores de la Universidad de Granada.

D. **Juan José Escobar Pérez**, Profesor Sustituto Interino del departamento de Arquitectura y Tecnología de Computadores de la Universidad de Granada.

Informan:

Que el presente trabajo, titulado *Dispositivo para detección de escritura mediante Deep Learning en un sistema empotrado*, ha sido realizado bajo su supervisión por **Antonio Priego Raya**, y autorizamos la defensa de dicho trabajo ante el tribunal que corresponda.

Y para que conste, expiden y firman el presente informe en Granada a X de mes Julio de 2022.

Los directores:

Jesús González Peñalver

Juan José Escobar Pérez

Agradecimientos

A mi familia, sin la que no habría podido llegar a hacer este trabajo, especialmente a mi hermano por colaborar en la recolección de muestras. A mis amigos por animarme cuando la energía flaquea. Y a mi tutor por la ayuda prestada.

Índice general

1. Introducción y Motivación	19
2. Antecedentes y estado actual	21
2.1. Redes neuronales	21
2.2. Microcontroladores	24
2.3. Integración de redes neuronales en microcontroladores	25
3. Especificación del sistema	27
3.1. Requisitos	27
3.1.1. Requisitos funcionales	27
3.1.2. Requisitos no funcionales	27
3.2. Especificación formal	28
3.2.1. Especificación hardware	28
3.2.2. Especificación software	28
3.3. Planificación y presupuestación	28
3.3.1. Planificación	28
3.3.2. Presupuesto	31
4. Diseño del sistema	33
4.1. Estructuración del diseño del sistema	33
5. Microcontrolador	35
5.1. Planificación	35
5.1.1. Elección del microcontrolador	35
5.1.2. Elección del entorno de desarrollo	37
5.1.3. Elección del <i>framework</i> para <i>Deep Learning</i>	37
5.2. Diseño	37
5.2.1. Estructura del firmware del controlador	37
5.2.2. Servicio <i>Bluetooth</i>	38
5.3. Implementación	39
5.3.1. Definición de parámetros de <i>Bluetooth</i>	40

5.3.2. Función de configuración del firmware	40
5.3.2.1. Configuración <i>Bluetooth</i> y de sensores	40
5.3.2.2. Configuración para la red neuronal	41
5.3.3. Función cíclica del firmware	42
6. Red Neuronal	45
6.1. Planificación	45
6.1.1. Elección de framework y librerías	45
6.1.2. Diseño de la red neuronal	46
6.2. Implementación	48
6.2.1. Preparativos	48
6.2.2. Entrenamiento	49
6.2.3. Testeo	51
6.2.4. Transformación a modelo cuantizado	51
6.2.5. Comparación de modelos generados	52
6.2.6. Integración en el microcontrolador	52
6.3. Generación de muestras para el dataset	52
7. Interfaz de usuario	53
7.1. Motivación	53
7.2. Planificación	53
7.2.1. Elección de framework para interfaz gráfica	53
7.2.2. Diseño de la interfaz	54
7.3. Implementación	55
7.4. Traspaso del diseño a <i>QT creator</i>	55
7.5. Configuración de la interfaz	56
7.6. Gestión de lectura del microcontrolador	57
7.6.1. Lectura de la característica <i>Bluetooth Low Energy</i>	57
7.6.2. Lectura del <i>puerto serie</i>	58
8. Encapsulado	59
8.1. Herramientas utilizadas	59
8.2. Implementación	59
9. Validación	61
9.1. Ajuste al presupuesto	61
9.2. Comprobación de objetivos cumplidos	61
10. Trabajos futuros y mejoras	63
11. Conclusiones	65
Apéndices	71

A. Microcontrolador	71
A.1. Resolver problemas de memoria	71
A.2. Instalación de librerías en <i>Arduino IDE</i>	71
A.3. Notificación de conexión <i>Bluetooth</i> del firmware del microcontrolador	71
A.4. Definición de micro-operaciones en el firmware del microcontrolador	72
A.5. Cambiar la orientación de la placa	72
B. Red neuronal	73
B.1. Ajuste para poder utilizar el recolector de muestras de <i>Pete Warden</i>	73
B.2. Asignación incorrecta de índices en el recolector de muestras . . .	73
B.3. El recolector de muestras elimina varias muestras al borrar una .	73
B.4. Descripción de capas <i>Keras</i> empleadas para la implementación de la red neuronal	74
B.5. Experimentación red neuronal	76
B.5.1. Estructura de la red neuronal	76
B.5.2. Entrenamiento de la red neuronal	78
C. Interfaz de usuario	79
C.1. Permisos uso del puerto del microcontrolador (Linux)	79
C.2. Acceso al puerto serie en QT (Linux)	79
C.3. Valores nulos al leer por <i>Bluetooth</i> con <i>QT</i>	79
C.4. Error de reconocimiento de imágenes en QT	80
C.5. Pérdida de valores de <i>características</i>	80
C.6. Error de permisos al lanzar la interfaz de usuario (Linux)	80
C.7. Capturas de la interfaz de usuario	81
D. Encapsulado	83
D.1. Modelos 3D	83
D.2. Resultado de integrar todo en el encapsulado	84

Índice de figuras

2.1. <i>Perceptron</i> frente a <i>Multilayer Perceptron</i>	22
2.2. Estructura simplificada de <i>Red Neuronal Convolucional</i>	23
3.1. Diagrama <i>Gantt</i> para la planificación	29
4.1. Esquema de la estructura del diseño del sistema	34
5.1. Esquema de la estructura planteada para el servicio <i>BLE</i>	39
5.2. Diagrama de flujo de leds	43
5.3. Diagrama de flujo simplificado del firmware del microcontrolador	44
6.1. Esquema de convolución 2D (https://bryanmed.github.io/conv2d/)	47
6.2. Esquema de funcionamiento del la red neuronal diseñada	48
7.1. Diagrama de flujo simplificado del <i>UI</i>	56
B.1. Estructura del modelo generado en <i>TensorFlow</i>	75
B.2. Escalado de <i>accuracy</i> por <i>epochs</i>	78
C.1. Boceto en QT Design Studio	81
C.2. Resultado de la implementación de la interfaz de usuario	82
D.1. Componentes útiles del encapsulado	83
D.2. Decoración del encapsulado	83
D.3. SmartPen	84

Capítulo

1

Introducción y Motivación

Como consecuencia del desarrollo de la informática, y la expansión y filtración del uso de equipos informáticos en la población general, cada vez la escritura tradicional cae más en desuso. Y es que una vez se supera la etapa académica, pocas personas siguen utilizando en su cotidianidad la escritura manual. Incluso para la educación hay un creciente movimiento de adaptación tecnológica que releva cada vez más al lápiz y papel.

No es el objetivo pecar de romanticismo y mirar a través de la lente de la nostalgia, sino avanzar, eso sí, intentando conservar en el proceso de innovación, las técnicas que nos han traído hasta aquí.

Por lo que el motivo de este trabajo siempre ha sido tratar de, creando nueva tecnología, favorecer el uso de la escritura. Ya que el avance y el progreso es no solo imparable sino necesario, la única forma de preservación de la grafía manual es establecer alternativas modernas que complementen a los dispositivos ya constituidos y que empleamos en el día a día.

Gracias al avance tecnológico de décadas, hoy podemos contar con herramientas informáticas de gran capacidad como lo son todos los mecanismos de Inteligencia Artificial. Concretamente el Deep Learning y las redes neuronales son conceptos en gran expansión durante los últimos años. El *Deep Learning* es un campo que está cambiando la informática como la concebíamos, revelándose como una alternativa sobresaliente para problemas que trabajan con grandes volúmenes de información, que presentan una elevada complejidad o simplemente que cumplen mejor de lo que habituaban a hacerlo con técnicas previas. Respondiendo oportunamente al contexto temporal vigente donde el *Big Data*, *Data Science*, automatización de tareas cotidianas, detonación de herramientas y dispositivos inteligentes, humanización de robots y robotización de personas; perfilan y caracterizan la fase en la que nos encontramos. Hecho que nos lleva al interés por un terreno tan intrincado como útil y repleto de potencial. Un potencial evidenciado por las tantas aplicaciones con sobrecededores resultados que hace pocos años casaban más con la ciencia ficción que con algo alcanzable,

y que emplean esta herramienta y que serán citadas a lo largo de este trabajo.

Por sus demostradas altas capacidades para la clasificación en el procesamiento de imagen, por el hito que supuso integrar redes neuronales en sistemas tan reducidos, porque hay algo sugestivo en el hecho de rescatar lo tradicional mediante las técnicas más incipientes, pero por encima de todo, por lo estimulante que resulta trabajar con estos mecanismos y que es algo que siempre ha rondado entre mis pensamientos; este trabajo consistirá en trasladar a la realidad una alternativa moderna a la escritura manual, haciendo uso de Deep Learning en un sistema empotrado para mantener autonomía.

Los usos pueden ser los que se deseen y se alcancen a imaginar; con pocos añadidos podría convertirse en una herramienta para introducir a personas de avanzada edad al manejo de ordenadores, reduciendo la barrera de entrada al tener una forma de interactuar que ya les es familiar; en un instrumento para hacer más ameno y dinámico el proceso de aprender a escribir para niños y niñas; en un cuaderno virtual en el que anotar cuanto queramos sin necesidad de transportar el medio en el que se escribe; incorporando una punta con grafito o tinta, podríamos transcribir digitalmente lo que escribimos en cada momento de manera física, es decir, una copia digital, un registro de lo que hemos escrito; etc.

Capítulo **2**

Antecedentes y estado actual

2.1 Redes neuronales

Pese a que es ahora, en los últimos años cuando, debido a la explosión del fenómeno de la *inteligencia artificial*, comienza a ser más popular todo lo relacionado con la dotación de inteligencia a los dispositivos electrónicos; todo comenzó hace muchas décadas [26]. Ya en 1943, *Warren McCulloch* (neurofisiólogo) y *Walter Pitts* (matemático), escribieron un artículo [24] acerca de las neuronas e incluso en el mismo, fueron capaces de diseñar una red neuronal simple usando exclusivamente circuitos eléctricos y fundamentado en algoritmos de *Lógica de umbral* (*Threshold logic*).

Más tarde, en la década de 1950, en los laboratorios de *IBM* de la mano de *Nathaniel Rochester*, ocurrió el primer intento de simulación de red neuronal; intento que desembocó en fracaso. Sin embargo fue muy estimulante para el campo de la *IA* y motivó el planteamiento de lo que denominaron "máquinas pensantes". También hubo otros acercamientos como la sugerencia del insigne *John Von Neumann* de utilizar relés telegráficos o tubos de vacío para simular el funcionamiento simplificado de las neuronas.

No obstante, no sería hasta 1958 que el neurobiólogo *Frank Rosenblatt* comenzaría a trabajar en el *Perceptron* [35], para muchos el nacimiento de la red neuronal artificial. Como todo precursor, era simple y limitado; hoy se catalogaría de monocapa, algo que evidentemente, ya no se usa en redes neuronales contemporáneas. Como puede observarse en la Figura 2.1a sirviéndose de múltiples entradas binarias, era capaz de producir una única salida, basada ya entonces en la utilización de pesos (número que cuantifica la relevancia de la entrada respecto de la salida), conservada hasta día de hoy, aunque cabe destacar que entonces, los pesos eran directamente atribuidos por el científico al cargo. La salida binaria de esta neurona *Perceptron*, sería como consecuencia de la superioridad o la inferioridad de la suma de la multiplicación de los pesos respecto de un umbral; es por esto que es sabida su influencia del trabajo de *Warren McCulloch* y *Walter*

Pitts anteriormente mencionado. Por tanto se podía destinar a funciones lógicas binarias simples (OR/AND).

El siguiente paso natural era aumentar el número neuronas y capas, llegando en 1965 el *Multilayer Perceptron Perceptron* [48]. Como consecuencia de esta mejora y aumento de la complejidad, nacieron los conceptos de capas de entrada, ocultas y de salida, tal y como se puede observar en la Figura 2.1b. De igual forma y dado que el reparto de pesos todavía no se había automatizado, los valores con los que se trabajaban, seguían siendo binarios.

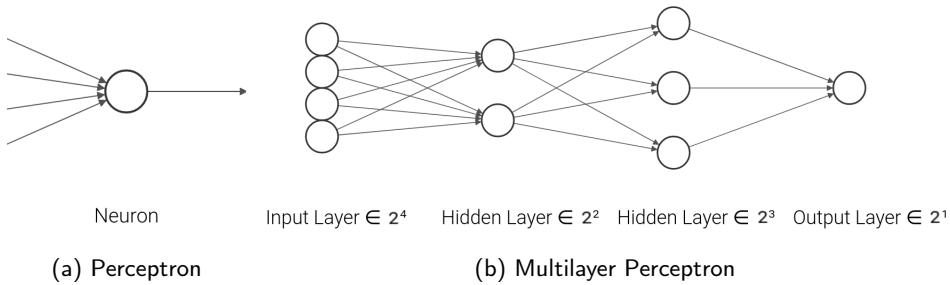


Figura 2.1: *Perceptron* frente a *Multilayer Perceptron*

Y este fue precisamente el siguiente escalón a rebasar, superado ya en la década de 1980, gracias a las *Neuronas Sigmoides Perceptron* [7], afines al *Perceptron* pero con la capacidad de trabajar con números reales. La función de salida ahora sería una sigmoide, a la que deben su nombre y convirtiéndose en la primera función de activación.

A partir de aquí y durante toda la década, comenzaron a aparecer todo tipo de novedades que continúan vigentes: redes *feedforward*, el algoritmo *back-propagation* o la *Red Neuronal Convolucional (Convolutional Neural Network, CNN)* [18]. Las *CNN* son especialmente convenientes para procesamiento de imagen y vídeo, en general información espacial, aunque también se han usado para tareas de procesamiento de lenguaje natural. Esto es debido a que la información se divide en subcampos que sirven como entrada a capas de procesamiento convolucional (ver Figura 2.2), encargadas de apreciar las distintas características que servirán para la clasificación de la información de entrada.

Es llamativo ver cómo las *CNN*, redes que mantienen su vigencia pese a que su origen se remonta a poco antes de los 90. Pero la realidad es que, si bien no lo parece, el campo de las redes neuronales lleva con nosotros mucho tiempo y las *CNN* no son el único ejemplo manifiesto. Las *Recurrent Neural Networks*, originadas en 1989, continúan en uso para procesamiento de datos secuenciales como lo es por ejemplo el texto.

Fue en 2006 cuando *Geoffrey Hinton et al* publicaron un famoso paper [14] presentando una red neuronal profunda, que entrenada, era capaz de reconocer dígitos. Acuñando como *Deep Learning*, a la técnica del *Machine Learning*.

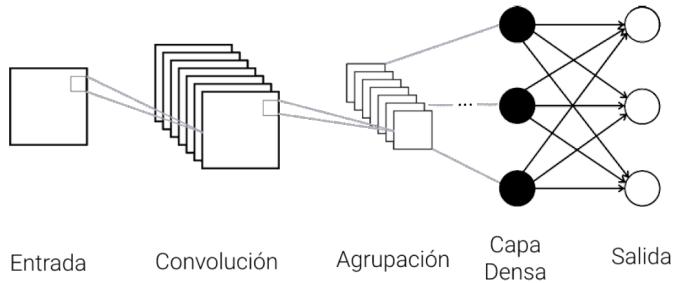


Figura 2.2: Estructura simplificada de *Red Neuronal Convolutinal*

(ya que se basa en el aprendizaje automático), que usa como mecanismo de procesamiento redes neuronales profundas. Se superaba entonces la barrera del entrenamiento de redes neuronales profundas, barrera que había llevado a la comunidad a congelar el avance de esta técnica, y que ahora elevaba al *Deep Learning* un nivel por encima de las técnicas del *Machine Learning*.

El siguiente hito llegaría a mediados de los 2000 al poder trabajar con redes neuronales profundas (*Deep Learning*), gracias a la introducción de pre-entrenamientos no supervisados para la asignación de pesos previos al usual entrenamiento del modelo. Este avance fue posible debido al desarrollo de la computación con GPUs.

El último acontecimiento o adición reseñable es el de las *Generative Adversarial Networks* (2014) [10], sujeto al empleo de dos redes neuronales complementarias: una se denomina *Generative network* en calidad de modelo generador de muestras y otra *Discriminative network* que evalúa las muestras generadas por la anterior y por el dataset de entrenamiento, es decir, recibe como entrada, la salida de la red anterior y del conjunto de datos de entrenamiento. El propósito de esta simbiosis es que la red generativa consiga reproducir muestras tan válidas como las de entrenamiento, a partir del juicio de la red discriminativa.

De entonces hasta ahora, más que innovación, se han dado muchos avances en términos de implementación, es habitual que cada cierto tiempo salga una nueva aplicación revolucionaria o con mucho potencial que está basada en redes neuronales. Y también es muy frecuente que gigantes tecnológicos como por ejemplo *Google* o *Facebook* compren otros proyectos basados en redes neuronales o directamente las empresas que los llevan a cabo. Siendo una de las más destacables la compra de *DeepMind* por parte de *Google* o la alianza entre *OpenIA* y *Microsoft*.

Algunos ejemplos de estos proyectos son: el tan sonado algoritmo de *AlphaGo* [12] de la empresa *DeepMind*, capaz de vencer al campeón mundial del juego tablero *Go*; el proyecto *DeepFace* de *Facebook* para identificar y automatizar el etiquetado de los usuarios en las imágenes; el *AlphaFold2* [11] de *DeepMind*, capaz de predecir la estructura de las proteínas y que ha sido revolucionario para la resolución del problema del plegamiento de proteínas, lo que antes eran investigaciones del orden de 1 o 2 años, ahora es computable en pocas horas;

GPT-3 [28] de *OpenIA*, un modelo de lenguaje cuya definición podría responder a *chatbot*, es capaz de completar texto, responder preguntas o cualquier tarea que implique interacción con texto; *Copilot* [15] de *OpenIA* y *Github*, *Microsoft*, un sistema construido sobre los cimientos de *GPT-3* capaz de sugerir código autogenerado y comentarios analizando bien las directrices de un comentario o bien directamente interpretando lo que el programador busca; *Nerf* [27] de *Nvidia*, capaz de generar composiciones 3D a partir de imágenes fijas; o por finalizar esta interminable lista de apasionantes ejemplos, el reciente *DALL.E* [28] de *OpenAI*, modelo generador de imágenes a partir de una entrada de texto descriptora.

No se puede quedar sin mencionar las que han sido las dos últimas grandes agitaciones del mundo de las redes neuronales y que están detrás de la mayoría de los ejemplos anteriores: el *Natural Language Processing* y los *Transformers*, aunque en realidad, van de la mano. Van de la mano porque el *Natural Language Processing* [22] ya es en sí mismo una revolución para el mundo de las redes neuronales, ya que el procesamiento de lenguaje, dado que las redes interpretan información numérica, siempre ha sido un desafío para el campo del *Machine Learning*; y no ha sido hasta su llegada, que gracias a lo que propone (vectorización de *tokens*, que son los bloques de datos que se interpretan, ya sean palabras o generalmente en la práctica, subpalabras), que las redes no han empezado a operar de una forma realmente veraz con el lenguaje. Sin embargo no solo ha sido una revolución en sí mismo, sino que ha propiciado el nacimiento de otra como lo son los *Transformers* [23], que parten del progreso conseguido en las redes recurrentes o para ser más precisos, de sus *Mecanismos de Atención* [43], ya que es lo único que mantienen respecto a los modelos recurrentes, es más, se alejan íntegramente pasando a un procesamiento simultáneo y sustituyendo la ordenación recurrente por la vectorización. Aunque pese a renunciar a la recurrencia, y es ahí donde reside su potencial, continúan funcionando con información secuencial.

Esta mejora en la implementación y aparición de tantas aplicaciones puede inferirse que se debe al aumento de la capacidad de cómputo de las GPUs, la llegada de los *Transformers*, la entrada de los mayores gigantes tecnológicos, pero sin duda a la aparición de herramientas de alto nivel e infraestructuras para el trabajo con redes neuronales como lo son *Azure*, *Aporia*, *TensorFlow*, *Keras*, *SciKit-Learn*, etc.

2.2 Microcontroladores

Los microcontroladores son sistemas de dimensiones reducidas y bajo consumo, destinados en su inicio al control de electrodomésticos, pero que a lo largo del tiempo y sobre todo propiciado por la aparición del *Internet of things (IoT)*, y los avances en *IA* han supuesto un cambio en cómo se diseñan e implementan los nuevos dispositivos electrónicos.

El primer microcontrolador fue desarrollado por *Gary Boone* y *Michael Co-*

chran en 1971 y fue bautizado como *TSM 1000* [47], albergando una arquitectura *Harvard* en un mismo circuito contando con el propio microprocesador, memoria ROM, menos de 256 bytes de memoria RAM y el propio reloj del sistema.

Como respuesta, *Intel* comercializó en 1977 su propio sistema para aplicaciones de control, el *Intel 8048* [46], que obtuvo gran popularidad y supuso un pequeño cambio en el paradigma de ventas de *Intel*.

Las memorias que montaban eran *EPROM* en el caso de los microcontroladores reprogramables y *PROM* en el caso de los de bajo presupuesto. Sin embargo esto cambió con la llegada en 1993 de la *EEPROM*, utilizada por primera vez en el *PIC16x64* [47] de *Microchip* y que conllevó un gran avance gracias a la agilización del proceso de creación de prototipos y su programación.

Poco después *Atmel* implementaría por primera vez memoria *flash* en un microcontrolador y se usaría en el *Intel 8051*, que trajo ciertos cambios respecto a su predecesor, como pasar a arquitectura *Von Neumann* o la inclusión de *Universal Asynchronous Receiver-Transmitter (UART)* para el manejo de puertos y dispositivos serie. Además contaba con múltiples compiladores de C para su programación, alternativos al lenguaje *ensamblador*.

Gracias a la inclusión de estas dos últimas memorias en el diseño estandarizado de los microcontroladores, el precio comenzó a ser cada vez más accesible. También se inició la incorporación de periféricos complementarios para dotar a los microcontroladores de más funcionalidad. Periféricos tales como generadores *PWM*, conversores analógicos A/D y D/A, relojes de tiempo real, etc.

Estos complementos lucen ahora arcaicos en comparación con los que se integran en microcontroladores actuales: transceptores 802.15.4, bluetooth, wifi, cámaras, micrófonos, y sensores de todo tipo como de presión, movimiento, orientación, color, brillo, proximidad, humedad, etc. Todos estos acompañados de complejos microprocesadores de 32 bits, cada vez más semejantes a las *CPUs* de equipos de mayores dimensiones. El progreso de los microprocesadores reducidos, ha sido propiciado por el vasto crecimiento del mercado móvil en los últimos años. Y resultando *ARM*, al igual que para los smartphones, una excelente baza para la complejidad que demandan los microcontroladores actuales.

Por un lado, su contenida complejidad frente a equipos de escritorio, hace que económicamente su implementación sea muy viable, aunque esto mismo provoca ciertos inconvenientes a la hora de su empleabilidad para *IA*, mencionados en la siguiente sección.

2.3 Integración de redes neuronales en microcontroladores

Algunos microcontroladores modernos dan soporte a herramientas para la *IA*, como lo es la integración de redes neuronales; sin embargo, el desarrollo de estas sigue siendo dependiente de la asistencia de un PC. De igual forma este soporte a herramientas para la *IA* es aun así sorprendente viendo los resultados que se

pueden obtener de su implementación y siendo este proyecto muestra de ello. Aunque en algunas ocasiones son necesarios ciertos arreglos dadas las carencias de estos dispositivos, como en ciertos casos, la falta de *FPU*s (*Floating-Point Unit*), suponiendo un obstáculo debido a que las redes neuronales realizan su procesamiento en coma flotante.

Se precisa de equipos con mayores prestaciones para la creación y entrenamiento de las redes neuronales que integrarán los microcontroladores, ya que es un proceso complejo y costoso; por suerte no es así para la ejecución, lo que los convierte en grandes candidatos gracias al *CloudML*, *EdgeML* o *TinyML*.

El *CloudML* [8] es la técnica por la que, alojando redes neuronales profundas en la nube, podemos integrar el uso de las mismas en *TPUs* (Unidades de procesamiento tensorial) y *FPGAs*, entre otras. Esta alternativa presenta la capacidad de emancipar el propio procesamiento de los algoritmos de *Machine Learning* fuera del propio dispositivo en el que se integra su implementación. Por otro lado, supone contar con infraestructuras que den soporte a ello y el uso de herramientas y entornos de pago, como entre otros el *Google Cloud ML Engine*.

EdgeML [16] es una librería escrita en *Python* impulsada por *Github*, *Microsoft* que mediante *TensorFlow* (o alternativamente en fase experimental *PYTORCH*), provee de distintas funciones enfocadas al entrenamiento, evaluación y despliegue de algoritmos de *Machine Learning* para sistemas embebidos empleados para labores simples. Por lo que es una elección perfecta para aquellos proyectos en los que se quiera trabajar con *Machine Learning* y utilizar herramientas de apoyo *open source*. Aunque cabe mencionar que, al tratarse de una herramienta para *Machine Learning* en general, las opciones dirigidas a *Deep Learning* no abundan.

Micro-Learn [1] es una librería para *python* que convierte modelos de *machine learning* entrenados con *Scikit-Learn*, a código que virtualiza la ejecución del modelo en cualquier microcontrolador en tiempo real. Es relevante destacar que no existen demasiados proyectos, aunque a cambio ofrece, en principio, soporte para cualquier microcontrolador *arduino*.

También existen infinidad de destacables alternativas para *FPGAs* como *Vitis-AI* o *VTA*, entre otras, aunque no presentan soporte a microcontroladores, por lo que quedan fuera de nuestro espectro de posibilidades.

Y finalmente, *TinyML* [41] se define como el campo que comprende a las tecnologías y aplicaciones relacionadas con el *Machine Learning* y que se fundamenta en la implementación de algoritmos y software capaces de realizar análisis de datos de sensores en un hardware muy limitado y de bajo consumo energético. Encontrando en esta alternativa, numerosos proyectos que consultar, gran actividad de su comunidad y cuantiosa documentación en forma de libros y publicaciones de usuarios.

Capítulo

3

Especificación del sistema

3.1 Requisitos

3.1.1 Requisitos funcionales

- El sistema procesará el movimiento del dispositivo resultando en la identificación de una letra.
- El sistema enviará la letra identificada por el sistema de comunicación pertinente a una interfaz de usuario.
- Al detectar un movimiento despreciable, el dispositivo lo descartará.
- Posterior a la identificación del movimiento, el sistema dejará un periodo suficiente de no registro de movimiento para que el usuario recoloque su postura.
- El sistema debe contar con un programa para ordenador, una interfaz gráfica que medie entre el dispositivo y el usuario.

3.1.2 Requisitos no funcionales

- El tiempo de procesamiento para la identificación de la letra, debe ser inmediato para generar sensación de escritura natural.
- El sistema debe funcionar con conexión por cable e inalámbrica.
- La interfaz de usuario debe ser simple de entender y usar.
- El sistema debe contar con autonomía energética.
- El dispositivo estará integrado en un encapsulado.
- El encapsulado tendrá un tamaño y forma semejante a un lápiz o bolígrafo.
- El microcontrolador debe tener unas dimensiones adecuadas para encajar en el encapsulado.
- El producto debe estar sujeto a un presupuesto no superior a los 65€.

3.2 Especificación formal

3.2.1 Especificación hardware

El dispositivo a crear clasificará, de forma autónoma, los distintos gestos que se realicen como letras, usando como herramienta de procesamiento Deep Learning. Por tanto la placa debe contar con:

- Dada la limitación de presupuesto, se optará por un microcontrolador.
- El microcontrolador debe ser compatible con el procesamiento de tensores para poder trabajar con modelos basados en *Deep Learning*.
- Sensores presentes en el microcontrolador suficientes para hacer reconocible el movimiento con precisión. En su defecto, se integrarán.
- Microcontrolador con tecnología inalámbrica.
- Microcontrolador con dimensiones reducidas.
- Un encapsulado que le dé cabida al propio microcontrolador y a la batería, cableado, etc.

3.2.2 Especificación software

- Interfaz de usuario simple que ofrezca las funcionalidades descritas.
- Framework adecuado para el diseño, entrenamiento, validación y testeo de redes neuronales y su integración en el microcontrolador.
- Firmware para el microcontrolador que integre la red neuronal para la tarea de detección de letras y la recolección del movimiento.
- Un software suficiente para la creación del encapsulado.
- Documentación para que la comunidad pueda hacer uso del producto y aportar nuevas funcionalidades al proyecto.

3.3 Planificación y presupuestación

3.3.1 Planificación

La planificación ha sido esencial para poner en valor los tiempos que manejar y ser consciente de las limitaciones. Es lo primero que se debe plantear unido a una preparación o documentación sobre lo que se va a trabajar para, solo de esta forma, poder estimar de una forma más precisa los plazos de cada elemento ineludible en el desarrollo del producto que se busca.

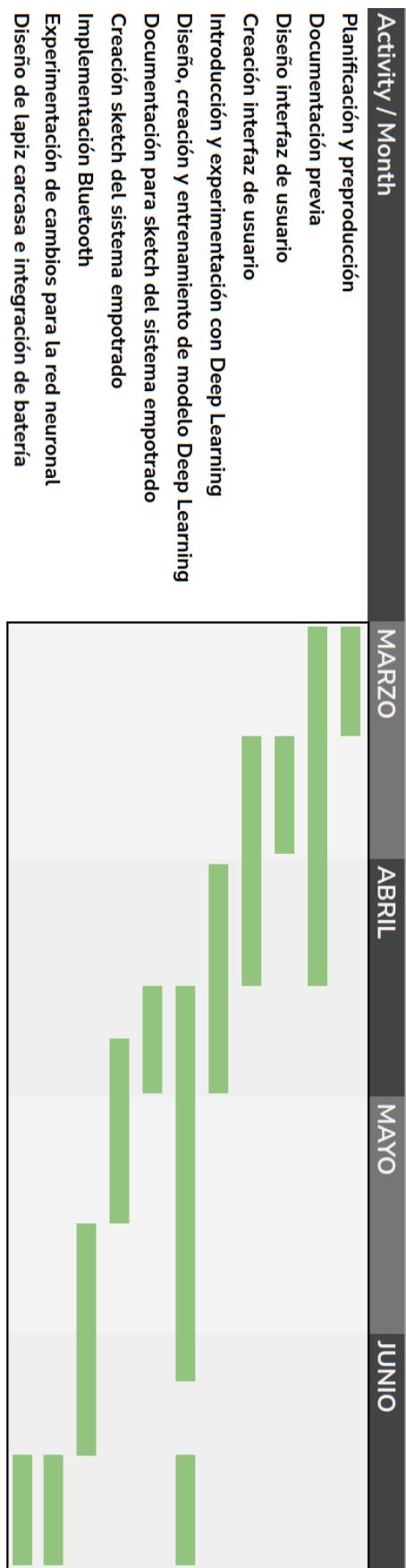


Figura 3.1: Diagrama Gantt para la planificación

Como muestra la Figura 3.1, ha sido una planificación basada en contenidos, sin prácticamente iteraciones de desarrollo, ya que, dados todos los campos que se abarcan y dada la complejidad de alguno de ellos; es improbable poder contar con varias iteraciones. La excepción de este precepto viene con el modelo, con el que se tratará de experimentar, con franqueza por apetencia de estudiar aún más las redes neuronales y su desempeño.

Por lo que podríamos clasificar el desarrollo formalmente como un modelo de prototipos; donde se estudian los requisitos, se documenta al respecto y se crea un prototipo revisable en caso de incumplimiento de requisitos o por errores en la implementación que se reflejen negativamente en el prototipo.

Aunque solo encontramos esa segunda iteración o revisión en la red neuronal, hay un cierto patrón de documentación y ejecución para cada parte del proyecto, por lo que se podría decir que los procesos están fraccionados. Sin embargo no deja de ser, como se ha denominado anteriormente, una planificación basada en contenidos: se ejecuta una parte y se procede con la siguiente.

El orden ha sido relevante y tiene su propósito. En primer lugar está la documentación y planificación, seguida de la creación de la interfaz, no por otra cosa que la carencia del hardware. Estos plazos son los primeros ya que permiten trabajar sin el sistema empotrado; tiempos de selección y obtención del hardware. Complementariamente, el diseño de la interfaz de usuario abre la mente a reflexionar sobre qué podría ofrecer el dispositivo a la persona que hace uso de él, ayuda a pensar en nuevas funcionalidades.

El siguiente bloque de contenido es el del modelo basado en *Deep Learning*, ya que si el modelo no alcanzara en la fase de testeo unos resultados óptimos, todavía podemos contar con algo más de tiempo para solventar su eficiencia.

Y por último la creación del propio firmware del microcontrolador y la integración del bluetooth, dividida en la implementación en la interfaz de usuario y en el firmware; para poder congregar finalmente todos los elementos y probar el resultado del producto.

Una última fase de experimentación en la que también se incluirá la creación y producción del embellecimiento para el dispositivo en forma de lápiz y la integración de su batería en el mismo.

3.3.2 Presupuesto

Para esta sección ha de tenerse en consideración la situación en el momento del desarrollo del proyecto de escasez de silicio, huelgas de transporte, pandemia, etc. Repercutiendo directamente en el precio de la electrónica y en los tiempos de entrega.

Aclarado esto, el presupuesto dada la naturaleza del proyecto, su funcionalidad y que se realiza con fines académicos y experimentales; será uno limitado y acorde a lo que se podría esperar.

En cuanto al coste del trabajo del ingeniero que lo desarrollará, se ha consultado el salario medio de un ingeniero informático en España en 2022 [37]. El pago por hora medio ronda los 13,59€, teniendo en cuenta mi experiencia laboral y mis conocimientos previos acerca de *machine learning*, voy a quedarme con una cifra algo más baja de 12.50€. Por lo que, teniendo en cuenta que este proyecto ocupa 12 créditos y cada crédito equivale a unas 26h de trabajo, este trabajo debería realizarse en unas 312 horas, lo que equivale a 3900€.

Descripción	Precio
Microcontrolador	40€
Batería	15€
Adaptación batería	2€
Impresión del encapsulado	1€
Cableado	7€
Tiempo de trabajo	3900€
EQUIPO: 65'00€	
TOTAL: 3965'00€	

Tabla 3.1: Presupuesto estimado para la producción

Capítulo

4

Diseño del sistema

Para poder trabajar con objetivos claros, lo mejor es definir la estructura de nuestro proyecto, constreñir los elementos que constituirán el producto que se trata de alcanzar.

4.1 Estructuración del diseño del sistema

Lo primero es identificar los elementos. Es evidente que nuestro dispositivo estará integrado en un encapsulado, por lo tanto, cuando se haga referencia al *SmartPen*, se estará haciendo alusión al propio encapsulado que contiene toda la electrónica agregada. Dicha electrónica consta de dos partes físicamente separadas: el microcontrolador y la batería que lo alimentará cuando se haga uso de su característica inalámbrica (bluetooth), dotándolo de la autonomía necesaria. Nuestro *SmartPen* necesitará de un equipo en el que mostrar las funcionalidades que ofrece el producto, podría tratarse de, por ejemplo un smartphone, pero en este caso se ha optado por hacerlo en un ordenador para agilizar el desarrollo. El ordenador es una herramienta que será utilizada no solo cuando se concluya el desarrollo, como interfaz para el *SmartPen*, sino como herramienta para la propia producción de todo el software requerido durante el desarrollo del proyecto: interfaz de usuario, firmware del microcontrolador y creación de la red neuronal. A su vez, el firmware hará uso de varios elementos clave para su funcionamiento; ya que la red neuronal tomará como entrada imágenes, se necesitará de una parte del firmware destinada a rasterizar el movimiento, movimiento que por otro lado deberá registrarse por medio de los sensores presentes en el microcontrolador. Finalmente, el dispositivo requerirá según lo especificado (sección 3.1.2), de mínimo dos canales de comunicación uno inalámbrico y otro por cable. Canales de comunicación que no solo se emplearán con fines de conexión con la interfaz de usuario, sino que también servirán para generar las muestras con las que se desarrollará (validación, entrenamiento y testeo) la red neuronal.

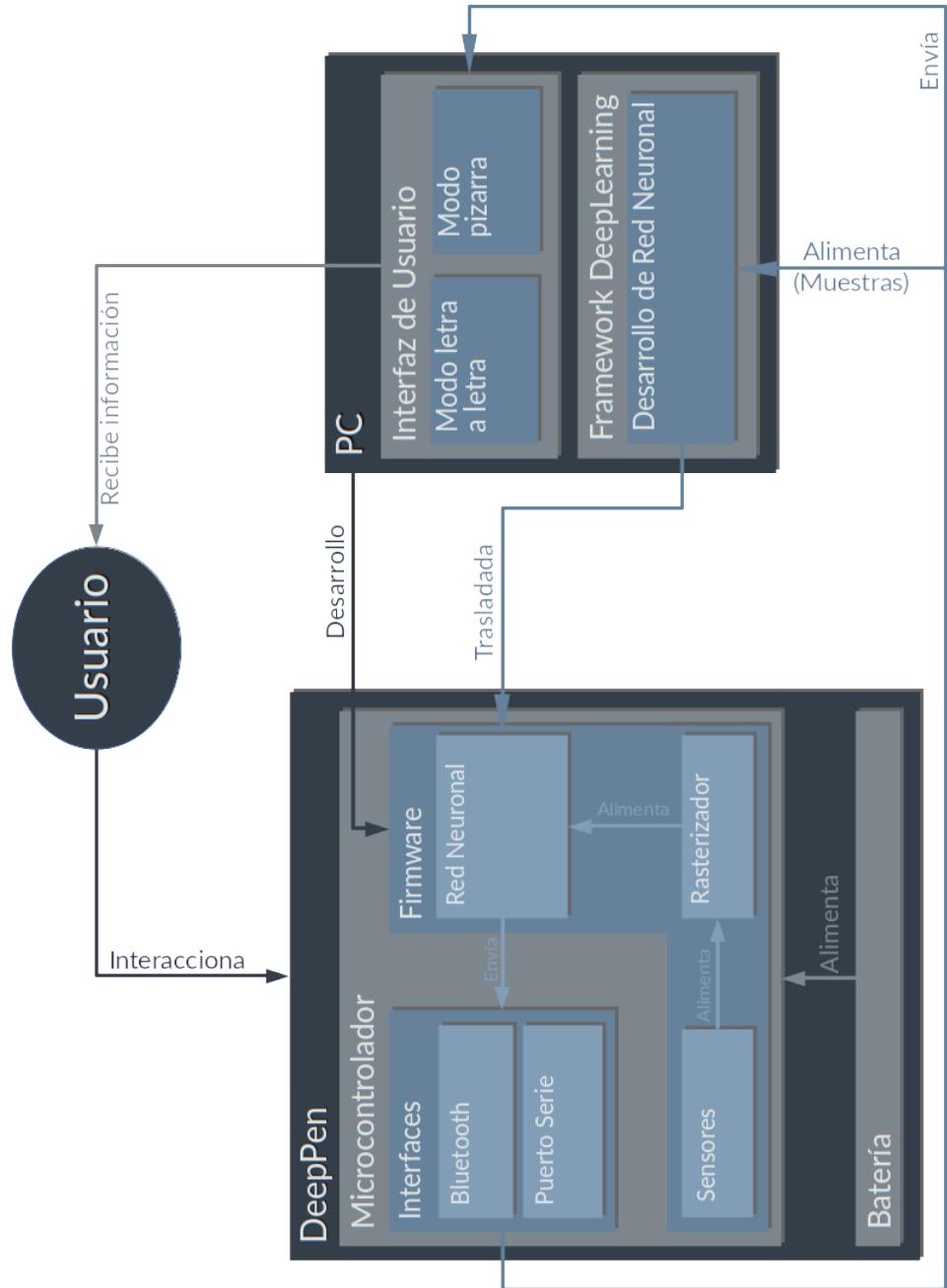


Figura 4.1: Esquema de la estructura del diseño del sistema

Capítulo

5

Microcontrolador

5.1 Planificación

En la planificación se recogerá la elección de equipo y herramientas de desarrollo previas al diseño.

5.1.1 Elección del microcontrolador

Como ya se especificó en la Sección 3.2.1 *Requisitos hardware*, son necesarias ciertas características inmutables. Para poder elegir un microcontrolador, antes se debe hacer el ejercicio de búsqueda de algunos candidatos y analizar sus características. La búsqueda se ha sustentado en encontrar dispositivos compatibles con *TensorFlow Lite*.

Microcontrolador	Precio	Documentación	Dimensiones	Sensores	Disponibilidad	Soporte para DL
SparkFun Edge	✓	✓	-	✓	-	✓
Arduino Nano Sense 33 BLE	✓	✓	✓	✓	-	✓
STM32F746G	✗	-	✗	✓	✓	✓
Adafruit EdgeBadge	✓	✓	✗	✓	-	✓
STM32	✓	-	✓	✗	✓	✓
ESP32	✓	-	✓	✗	✓	✓

Tabla 5.1: Tabla de requisitos para elección de microcontrolador

La *ESP32* al igual que la *STM32* pueden descartarse debido a que carecen de

sensores relacionados con el movimiento, se les podrían integrar periféricamente, pero sería algo más molesta su inserción en una carcasa. Y dado que tenemos alternativas que cuentan con estos sensores, podemos permitirnos descartarlas.

Por otro lado la *Adafruit EdgeBadge* es una placa muy llamativa, potente y llena de posibilidades, pero cuenta con unas dimensiones superiores a lo que se busca.

Respecto a la *STM32F746G* el inconveniente es la bajísima disponibilidad y los plazos de envío desorbitados. Por tanto tampoco podemos contar con ella.

Por lo que la disyuntiva se plantea entre la *Arduino Nano Sense 33 BLE* y la *SparkFun Edge*. En este caso la decisión no está motivada por características técnicas, que además son muy parecidas en ambos dispositivos, sino que una de ellas ofrece algo fundamental cuando se trabaja por primera vez en un campo y más aún cuando se cuenta con limitación de tiempo. El motivo taxativo es la documentación que provee Arduino, así como el apoyo de su activa comunidad y el gran volumen de proyectos que podemos consultar y que hacen uso de esta misma placa. Sumado a que este microcontrolador es prácticamente el más extendido para *Machine Learning*, por tanto no solo encontraremos multitud de proyectos en los que apoyarnos, sino que muchos de ellos o la práctica mayoría estarán enfocados al *Machine Learning*. Incluso se ha convertido en el abanderado de los proyectos basados en *TinyML*, el mayor valor añadido con el que puede contar un dispositivo que se empleará partiendo de pocos conocimientos y con restricción temporal. Cabe destacar también, que dada la complicada situación respecto al mercado de la electrónica al momento del desarrollo de este proyecto, la escasez de silicio, retrasos en la producción por la pandemia, huelgas de transporte, etc; la prontitud de entrega del propio microcontrolador, ya era en sí mismo el factor más determinante. Por lo cual me vi obligado a provisionarme de ambas y comenzar a trabajar con el microcontrolador que antes llegara. Por suerte pude hacerme primero con la *Arduino Nano Sense 33 BLE*, que era la prevista para el proyecto por lo expuesto.

De su propia nomenclatura podemos extraer todos los elementos precisados para este proyecto:

- **Arduino:** Garantiza que encontraremos documentación, y asistencia y proyectos de otros usuarios.
- **Nano:** Posee unas dimensiones convenientes para poder incorporarlo en un encapsulado adecuado para la escritura.
- **Sense:** Cuenta con diversos sensores, concretamente con una *IMU (Inertial Measurement Unit)* que provee de *giroscopio* y *acelerómetro*.
- **BLE:** *Bluetooth Low Energy*, un *Bluetooth* de bajo consumo que proporcionará autonomía y libertad de movimiento.

5.1.2 Elección del entorno de desarrollo

El entorno de desarrollo escogido será el propio *Arduino IDE*, debido a que nos facilita mucho el trabajo en ciertas tareas como el acceso al puerto serie para labores de depuración, la instalación de librerías para Arduino y sus dispositivos, o la carga del firmware en la placa con un solo click. Adicionalmente, se hará uso de *Visual Studio Code* en los períodos de programación sin interacción con el microcontrolador, dado que es un entorno más cómodo para gestionar varios archivos simultáneamente y programar durante sesiones algo más largas.

5.1.3 Elección del *framework* para *Deep Learning*

A la hora de trabajar con *Deep Learning* y *redes neuronales*, es importante apoyarse en herramientas de alto nivel, ya que si desarrollaramos la red neuronal a bajo nivel, necesitaríamos de un nivel de documentación que llevaría mucho más tiempo del que tenemos. Y no solo eso, sino que no podríamos integrar el modelo en el microcontrolador. Por tanto, necesitamos de un marco de trabajo, un *framework*, que nos facilite el trabajo, y nos brinde la infraestructura y herramientas esenciales para poder desarrollar nuestro modelo basado en *Deep Learning*.

Ya vimos en la sección 2.3 (*Integración de redes neuronales en microcontroladores*) las alternativas que se nos presentan a la hora de integrar redes neuronales en microcontroladores; de todas ellas y debido a la actividad de su comunidad y la tendencia a exponer sus proyectos, se optará por *TinyML*, el cual está complementado a la excelencia con *TensorFlow Lite*. Por ello y por otras múltiples razones como que es de código abierto, gratuito, forma una gran sinergia al agregar algunas otras herramientas de alto nivel (como por ejemplo, *Keras* o *Scikit Learn*), etc; se ha optado por *TensorFlow Lite*. Pero al igual que en las elecciones anteriores, lo que más decanta la balanza es siempre la expansión desde su inicio y la cuota de utilización frente a sus alternativas. Ya que esto se traduce en, generalmente, mayor documentación, mayor interacción de la comunidad, más proyectos que poder consultar y más experiencias de otros usuarios que pueden ser de interés.

5.2 Diseño

En esta sección se planteará estructuralmente y a nivel de funcionamiento, el firmware del controlador.

5.2.1 Estructura del firmware del controlador

El firmware se distribuirá en diferentes secciones. *smart_pen.ino* (extensión propia de *Arduino*, empleada en el archivo principal de sus proyectos) incluirá todo lo relativo a la configuración previa de las características del microcontrolador

y las habituales funciones `setup()` y `loop()`.

En `smart_pen_model_data` encontraremos exclusivamente el modelo de la red neuronal entrenado y listo para funcionar en formato binario, dada la carencia de sistema de archivos.

`labels` será la sección que ocupe la gestión de las etiquetas del modelo; definición y traducción etiqueta a letra.

El apartado `rasterize_stroke` rasterizará el movimiento recogido, transformándolo en las imágenes que sirven como entrada para el modelo.

Por último, `stroke_collector` estará reservado a la recolección del movimiento y la configuración del servicio *Bluetooth (BLE)* vinculado a la recolección de muestras para el modelo.

5.2.2 Servicio *Bluetooth*

Pese a que se implementarán dos servicios, uno de ellos es parte del *Data collector*, sección extraída del proyecto de *Pete Warden* [44] y por tanto no la desarrollaré más allá de una breve explicación en su apéndice. Se describirá, por tanto, solamente el servicio implementado de cero: *letterSenderService*.

Previo a la descripción del diseño, debemos entender cómo funciona esta versión bluetooth de bajo consumo.

Teoría 5.2.1: Estructura del *BLE(Bluetooth Low Energy)* [3] [6] [42]

El funcionamiento de este bluetooth de bajo consumo es notoriamente disidente de la versión general, tanto que tenemos que hablar de una estructura propia y que será clave para poder hacer uso de esta herramienta. Esta estructura jerárquica está definida por *Atributos* (Attributes). Cada uno de los elementos a continuación enumerados son *Atributos*, todos ellos identificados por un *UUID*(Universally Unique Identifier):

1. *Servicios* (Services)

Agrupaciones de características. Un servicio suele componerse de características vinculadas al ámbito del servicio. Generalmente cada servicio corresponde a una prestación del dispositivo

2. *Características* (Characteristics)

Cada característica contiene un tipo(*UUID*) de característica, sus propias propiedades y sus propios permisos. Y continuando con la disposición jerárquica, cada característica está formada por ninguno, uno o múltiples descriptores.

Representan estados del dispositivo, datos de la configuración del mismo o simplemente un dato correspondiente a alguna función del servicio.

3. *Descriptores* (Descriptors)

La unidad mínima de la estructura. Es la que contiene la información transmitida por cada comportamiento de una característica y sus metadatos asociados.

El servicio (*letterSenderService*) está compuesto por dos características: rx (*rxChar*) y tx (*txChar*). Tomaremos estas características como canales de comunicación unidireccionales. Han sido denominados teniendo en cuenta la placa como sistema de referencia; *rx* será la característica receptora de datos y *tx* la característica transmisora.

Utilizaremos el canal *tx* para transmitir la letra y el canal *rx* a modo de gestor de flujo; para la comunicación con el programa de usuario. Cuando el interfaz de usuario reciba la letra y la almacene, escribirá en el canal *rx* la correspondiente señal para que el canal *tx* se borre y pueda dar paso a una nueva letra.

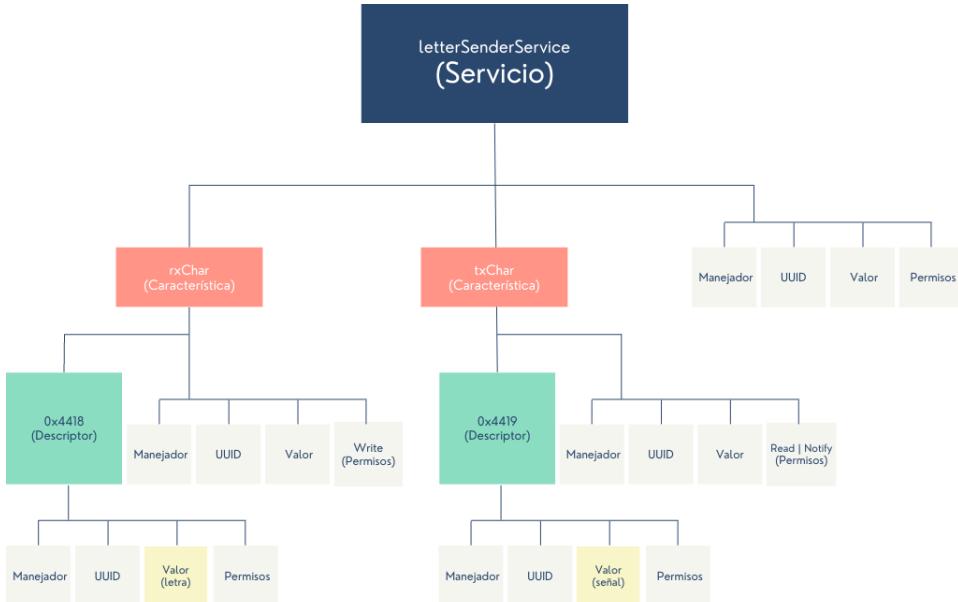


Figura 5.1: Esquema de la estructura planteada para el servicio *BLE*

5.3 Implementación

Previo a la implementación del código, se deben instalar ciertas librerías (proceso sencillo gracias a la decisión de utilizar el *Arduino IDE* pero que puede consultarse en el apéndice A.2) para el trabajo con algunas características de la *Arduino Nano Sense 33 BLE* como lo son la librería *Arduino_LSM9DS1*, para el uso de la IMU y su acelerómetro, magnetómetro y giroscopio; y la librería *ArduinoBLE*, para el uso del *Bluetooth* de bajo consumo. Por otro lado, también es necesario instalar la librería *Arduino_TensorFlowLite*, para trabajar con el modelo de la red neuronal en el firmware del microcontrolador.

Además, en ocasiones, es necesario dar permisos al puerto del microcontrolador como reza el apéndice C.1

5.3.1 Definición de parámetros de *Bluetooth*

En el primer bloque del firmware, encontramos definición de los parámetros que más tarde se utilizarán para configurar el *BLE* (Bluetooth Low Energy) [20].

La configuración de permisos de las características son consecuentes con su cometido. La característica de lectura, tendrá permisos de escritura para que la *central* pueda escribir los valores que la placa leerá. Análogamente la característica de escritura tendrá permisos de lectura y notificación.

Teoría 5.3.1: Roles de las partes en conexiones *bluetooth* [42]

Al darse una conexión bluetooth, existen ciertas implicaciones inherentes a la propia conexión. Y es que siempre habrá una de las partes que se lucra de los servicios de la(s) otra(s).

Pues bien, los dispositivos que se gestionan o se benefician de los servicios de otros, se conocen como *central* y los que proveen, son los *periféricos*.

El resultado de la configuración del servicio *BLE* obedece completamente al diseño planteado.

Finalmente se deben definir una serie de manejadores para gestionar los distintos eventos relacionados con el servicio *BLE*. Estos manejadores se activarán como consecuencia de suceso y se encargarán de gestionar la respuesta a dicho suceso; concretamente respuestas a conexiones, desconexiones o recibo de señales por el canal de lectura *rx*.

5.3.2 Función de configuración del firmware

En esta sección (denominada función '*setup()*' en el diseño de *Arduino*), como de constumbre, se deberá llevar a cabo el ajuste propio para la ejecución de nuestro código [2] [44].

Es reseñable mencionar que el microcontrolador con el que estamos trabajando solo posee la *memoria flash* y *SRAM* características de los dispositivos de estas prestaciones. Sin embargo tenemos que almacenar alguna información imprescindible como la red neuronal o la reserva de algunas secciones de memoria. Para ver la solución implementada con un mayor detalle técnico, consultar el apéndice A.1. Más adelante se hará uso de la misma técnica para resolver el problema de la integración de la red neuronal en el microcontrolador, entre otros.

5.3.2.1 Configuración *Bluetooth* y de sensores

En esta primera etapa de configuración, se definirá el comportamiento de algunas características de la placa; en primer lugar, el de la *IMU* (*Inertial Measurement Unit*), la unidad con la que trabajamos para obtener los datos del movimiento sirviéndose de un *giroscopio* y un *acelerómetro*, y por otro lado, la

configuración del ya descrito, definido y diseñado *BLE* (*Bluetooth Low Energy*) con los parámetros que se definieron en la sección 5.3.1.

Teoría 5.3.2: Por qué es necesario configurar la *IMU* [5]

La *IMU* es una parte fundamental de este proyecto, ya que es el dispositivo contenido en la placa que gestiona las mediciones de aceleración y velocidad del movimiento y que consta para ello de acelerómetro y giroscopio. Con la captura de estos datos, podrá rasterizarse una imagen que contenga el trazo descrito y con la que podremos alimentar la red neuronal.

Problemas 5.3.1: Librería Arduino _ LSM9DS1

Uno de los problemas con esta parte del código, fue que la librería *Arduino _ LSM9DS1* necesaria para poder trabajar con la *IMU*, tiene varias versiones. Al haber leído para algunos de los proyectos *TFLite de arduino*, que era recomendable utilizar su primera versión, fue la elegida. Sin embargo esta primera versión, no posee una de las funciones necesarias para trabajar con la *IMU* en nuestro caso, que es el llenado continuo de la FIFO de lectura de medidas recogidas, que por defecto funciona en *oneShotMode*, es decir, llenado a ráfagas. Es necesario poder disponer de la función para trabajar en tiempo real con la predicción de letras y también es imprescindible para la recolección de muestras, como se verá más adelante.

Por lo que la solución es o bien añadir manualmente la función en la librería, o bien actualizarla a la versión 1.1.0. Yo me he decantado por actualizar la librería, ya que es algo más limpio que editar una librería que no has programado tú mismo.

Por lo que solo resta fijar los parámetros creados para *BLE*, vincular sus señales con manejadores de los eventos y hacer lo propio con la *IMU*.

Cabe destacar que adicional a la configuración *BLE* descrita, también se incorpora otro servicio con la característica *strokeCharacteristic*, menos interesante de explicar, ya que utilizaré el método de recolección de muestras de *Pete Warden* para uno de los ejemplos *TFLite de Arduino*: *magic_wand* [44].

5.3.2.2 Configuración para la red neuronal

Parte de máxima trascendencia, ya que, es imprescindible configurar manera adecuada los parámetros del modelo para que el reconocimiento se de manera óptima.

Una vez obtenemos el modelo definido en *deep_pen_model_data.cpp*, proceso que se explicará en la sección 6.2.6.

Se deben establecer las micro-operaciones que se darán en el modelo para tener definido el repertorio en tiempo de ejecución que utilizará nuestro interprete [40]. Existe una alternativa que añade todas las micro-operaciones de forma genérica, a costa de un mayor uso de memoria. Para mayor profundización en la implementación, consultar el apéndice A.4.

Problemas 5.3.2: Al cargar el firmware la placa deja de ser detectada

En las primeras cargas del firmware experimentando con el modelo, la placa dejaba de ser detectada. Lo primero que pensé es que el bootloader se había bloqueado. Sin embargo al restaurar la placa manualmente (Pulsación del botón reset justo al conectar la placa), el 'L' led de la placa, comenzó a parpadear; indicativo de que la placa se había restaurado. Por tanto solo cabía que el programa cargado era erróneo. Dado que tanto la compilación como la ejecución no informaban de errores, fue complicado dar con que este error se debía a una mala configuración de las microoperaciones definidas.

Ya que al hacer cambios en el diseño del modelo, es imperativo añadir las microoperaciones ampliadas. Un error de principiante que llevó mucho tiempo arreglar.

Para finalizar la configuración de *TensorFlow Lite* definimos el interprete del modelo, que hará uso del repertorio de micro-operaciones especificadas.

Y por último, inicializamos los led pins como salida, para poder hacer uso del led como indicativo de estado.

5.3.3 Función cíclica del firmware

En esta función (denominada 'loop()' en el diseño *Arduino*), será donde se establezca el código que se ejecutará persistentemente mientras la placa esté alimentada [2] [44].

En primer lugar, se tratará la lógica para los leds es simple: En estado idle el led encendido será el verde y en estado de detección de letra, será azul (quedará en este estado durante 800ms, tiempo durante el que no se detectará nada, para que el usuario pueda repositionar su postura para volver a escribir otra letra); en caso de que la *IMU* no esté disponible, se encenderá el led rojo.

Para notificar una vez la conexión a un dispositivo, pese a que la ejecución es cíclica, se implementa una sencilla solución descrita en el apéndice A.3.

En cuanto al registro de movimiento, se han utilizado algunas de las funciones del trabajo de *Pete Warden* y su *magic_wand* [44] a modo de librería para mi propio proyecto y así aligerar el desarrollo de la recolección de trazos y su rasterización.

Durante este proceso, se hará uso del giroscopio para determinar los cambios

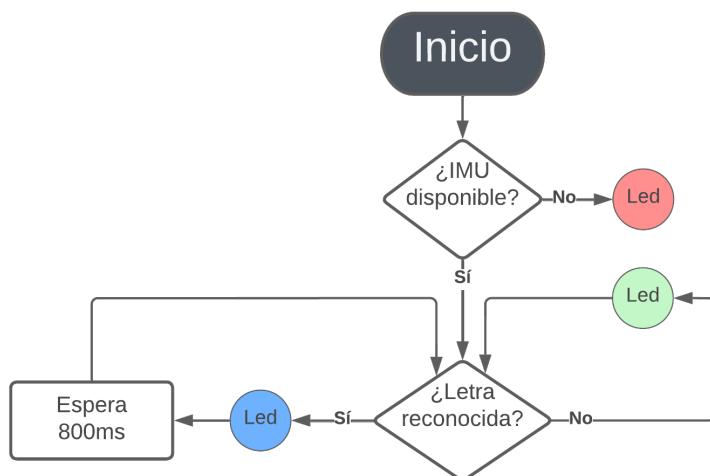


Figura 5.2: Diagrama de flujo de leds

del trazo y el acelerómetro para pequeños ajustes sobre los resultados obtenidos con el giroscopio, por ejemplo cálculos de gravedad y parámetros de velocidad del cambio de trayectoria del trazo. Gracias a contar con las funciones antes citadas, tomar los valores de los sensores es tan fácil como llamar a la función *ReadAccelerometerAndGyroscope*. Con los valores leídos, se hace un pequeño arreglo de estimación de desvío del giroscopio y se envían los valores leídos como característica del servicio para *Data Collector*, donde se recogerán las muestras para entrenar la red neuronal. Y con el trazo completamente construido y corregido, se rasteriza el movimiento para obtener la imagen que pasará por el modelo.

Con todo lo anterior definido, todo está listo para dar comienzo con el procesamiento en la red neuronal. Para llamar a su ejecución, se invoca el interprete, que arrojará los resultados en un puntero anteriormente definido. Este puntero de salida, contiene los datos asociados al tensor, es decir, el producto de que el trazo rasterizado haya pasado por la red neuronal. En nuestro caso, lo que se obtiene de la red neuronal, es una valoración de ajuste de afinidad del trazo rasterizado a lo que el modelo ha sido entrenado para reconocer como letras (nuestros *labels*). Sintetizando, lo que se obtiene como salida de la red neuronal, es una valoración de semejanza a cada letra, codificada como un índice.

Por tanto, lo que resta es trivial, solo tenemos que obtener el *label* con mayor valoración. Y como consecuencia, la letra que el modelo ha estimado más posible respecto a su entrenamiento.

Obtenida la letra, es enviada al puerto serie, para cuando se trabaje con conexión física; y se introduce en la característica de transferencia(*tx*) para cuando se trabaje con el servicio *Bluetooth (BLE)*.

Como comprobación concluyente, se verifica si la característica de lectura

(rx) ha sido escrita por el programa de usuario, suponiendo esto una señal de que ya ha sido leída la letra actual en tx y como consecuencia, haciendo que se restaure su valor a uno por defecto; ya que de no hacerlo, la letra permanecería inmutable en la característica y el programa de usuario leería las mismas letras reiteradamente hasta escribir otras. Esto se debe a que las *características* en BLE funcionan con valores constantes hasta que se modifique el estado actual.

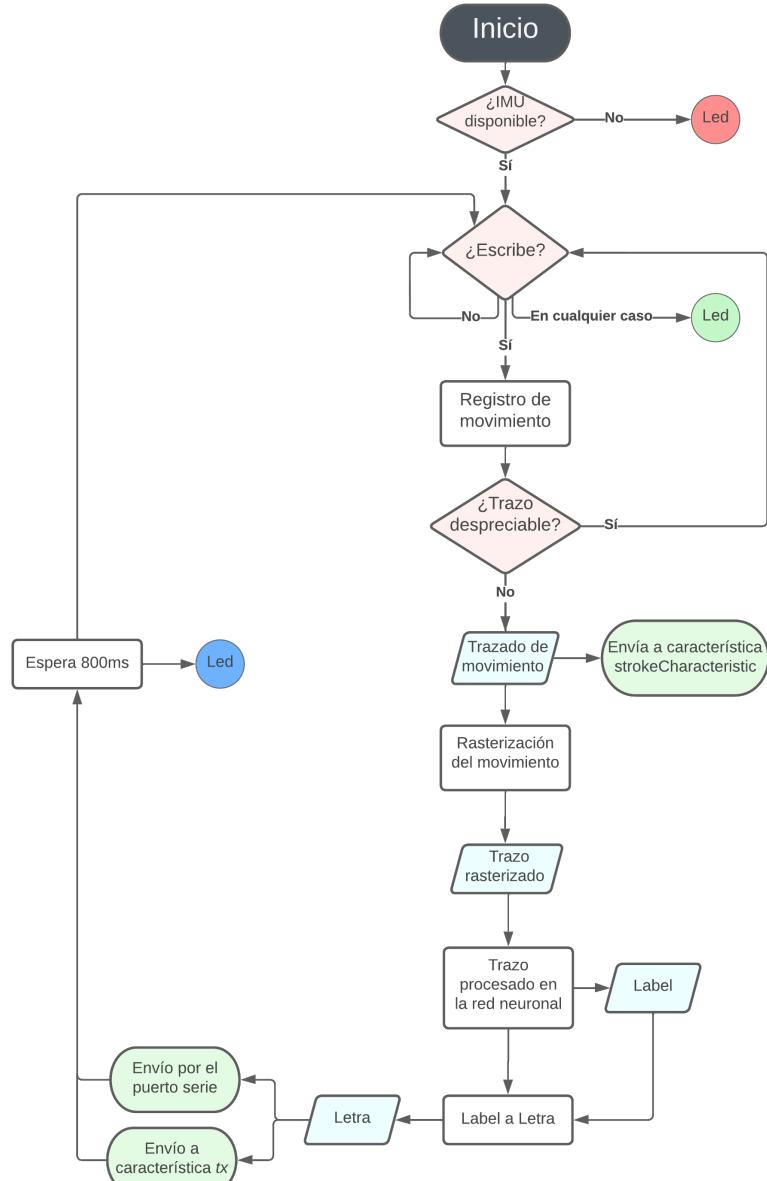


Figura 5.3: Diagrama de flujo simplificado del firmware del microcontrolador

Red Neuronal

6.1 Planificación

6.1.1 Elección de framework y librerías

En esta sección se defenderá la elección de las herramientas utilizadas para el desarrollo de la red neuronal, resultando una ampliación de la [sección 5.1.3](#), donde se ha tratado más la materia en términos del modelo apartándose de su integración en el microcontrolador.

TensorFlow Lite se utilizará para la creación del modelo optimizado para microcontroladores, pero esto no perturba en absoluto el desarrollo ordinario en *TensorFlow*. Es decir, se trabajará de igual forma que si la red neuronal se utilizara en un equipo convencional, con la única salvedad de la conversión del modelo ya generado a una versión optimizada para microcontroladores; así es la dinámica de trabajo con *TinyML*.

Como complemento a *TensorFlow* y para facilitar el desarrollo de la red neuronal, se empleará *Keras*, una *API* de alto nivel que, desde hace unos años, pertenece a *TensorFlow*; sirviendo para actuar como interfaz a nivel de capa para *TensorFlow*, o expuesto de otra forma, facilita el desarrollo al simplificar el trabajo con las capas del modelo. Además también se utilizará para el estudio de los resultados del entrenamiento.

Como entorno de desarrollo, se utilizará *Google Colab*, se podrían emplear alternativas que se ejecutan localmente como *PyCharm*, *Jupyter* o *Eclipse*; sin embargo con *Google Colab* podemos aprovechar la capacidad de cómputo de sus *GPUs*, lo cual agiliza mucho los tiempos de *entrenamiento y validación*. Si bien, el tiempo de uso del que se dispone de estas *GPUs*, es limitado, aunque siempre podemos continuar haciendo uso de las *CPUs*. Además, al tratarse de una alternativa en la nube, la etapa de configuración es mínima: apenas instalando las librerías de las que vamos a hacer uso, ya tenemos desplegado el entorno de trabajo.

6.1.2 Diseño de la red neuronal

El diseño de un modelo es determinante para que se desenvuelva apropiadamente a la hora de realizar su función. Por más que se dote al modelo durante el entrenamiento de un gran volumen de datos, será vano si el diseño no está enfocado a la labor que tiene que desempeñar. Por lo tanto, para el diseño de este modelo, se han estudiado otros muchos de clasificación de formas y procesamiento de imagen; la mayoría hacían uso del célebre *MNIST*: una base de datos que cuenta con 60.000 imágenes de entrenamiento y 10.000 de testeo, es un gran dataset de imágenes de dígitos. Aunque en general se han consultado todo tipo de modelos para el procesamiento de imágenes y clasificación, destacando el referente de este proyecto *Magic Wand* [44] de *Pete Warden*. Para examinar con mayor profundidad la experimentación respecto al diseño del modelo, observe el *apéndice B.5.1*.

El modelo que mejor funciona y más se utiliza en trabajos simples de procesamiento de imagen para clasificación, es el de *redes neuronales convolucionales* (*CNN*); definidas en una estructura muy marcada como ya se describió en la Figura 2.2.

Teoría 6.1.1: Redes neuronales convolucionales

Estas *redes neuronales convolucionales* no son más que un tipo de red neuronal que, reduciendo a un nivel muy básico, utiliza un tipo de capa que realiza una operación matemática llamada convolución.

Lo cual nos lleva al siguiente paso, la definición de convolución: operador matemático que haciendo uso de dos funciones f y g , genera una nueva a partir de estas, que representa la magnitud de su superposición. Concretamente en 2 dimensiones, podemos entenderlo, tal y como ilustra la Figura 6.1 como el producto del *kernel* o *filtro* y una subventana de la matriz, generalmente una imagen; al repetir este producto para todas las subventanas, obtenemos una nueva matriz resultado de la convolución.

Encontrar los valores del *kernel* y por tanto la magnitud resultante del análisis de la imagen de entrada, será la principal labor de la capa de convolución.

Al resultado de la convolución, se le denomina *mapa de características*, y su función es evidenciar dónde se encuentra la característica buscada por el *kernel*. Estos *mapas de características* pueden reflejar cambios de contraste, texturas, superficies planas, etc.

Pues bien, la red neuronal realizará esta operación secuencialmente, es decir, el resultado de una capa, alimentará a la siguiente. Lo cual provoca que, unido a un proceso de *agrupación*, cada vez la información relevante se condense y refine más y más, permitiendo extraer patrones cada vez más complejos.

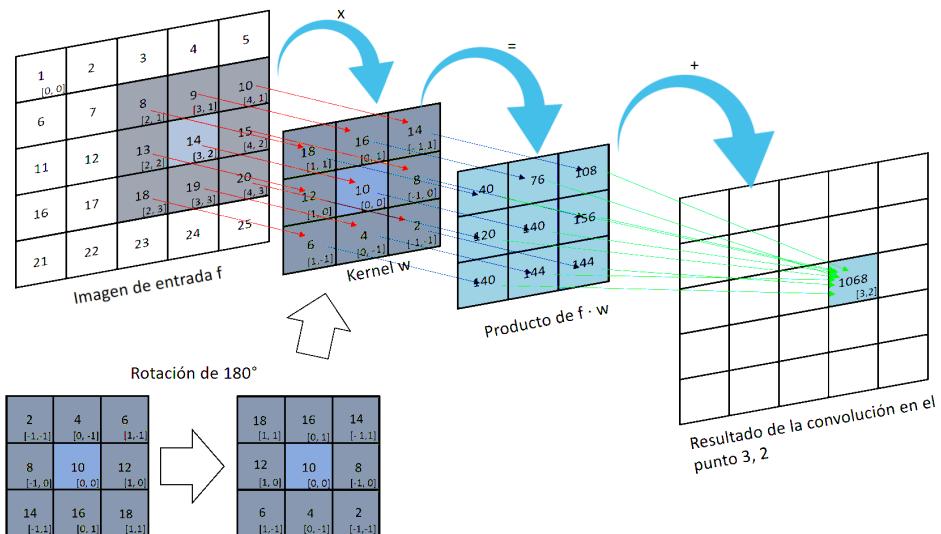


Figura 6.1: Esquema de convolución 2D (<https://bryanmed.github.io/conv2d/>)

Conociendo la teoría de la *red neuronal convolucional*, ya es posible comentar el diseño concreto implementado. El modelo cuenta con una estructura de *Sequential* de Keras; gracias a esta utilidad, es posible añadir capas al modelo de forma cómoda, además de permitir el acceso a algunas funciones para el entrenamiento. Se describirá la estructura en términos generales, para encontrar una descripción a nivel de implementación real con Keras, véase el *apéndice B.4*

El modelo utilizado va a contar con tres bloques de procesamiento tras la entrada, formados por una capa de *convolución*, otra de *normalización*, *activación* y *fropout*. Estos bloques deben su composición a que el aporte de cada una de las capas, complementa al procesamiento convolucional, en el caso de nuestro modelo, que trabaja con imágenes tan pequeñas, son prácticamente irremplazables. Comenzando por la inherente capa de *convolución*, donde se dará el procesamiento anteriormente ilustrado. Tras esta, se normaliza la salida mediante una capa para este propósito, lo cual garantiza una mayor estabilidad y eficiencia en el proceso de aprendizaje. Se define como función de activación, *relu*; no por otra cosa que porque es la más utilizada, la que mejor funciona para prácticamente todas las labores sin tener que profundizar en el estudio del mismo y por su bajo coste computacional; destacable debido a la naturaleza del dispositivo en el que se ejecutará la red neuronal. Y por último la capa de *Dropout* para mitigar el overfitting, un problema que se ha podido experimentar debido a que es una red neuronal de reducida complejidad.

En cada uno de los tres bloques, lo único que varía es el número de filtros de la capa de convolución, duplicándose respecto al anterior.

Comentada la estructuración de estos bloques y su fundamento, ya es posible desarrollar el modelo completo. Al inicio, preliminar a los tres bloques citados, hay una capa de *reescalamiento*, que va a servir meramente para normalizar los valores los píxeles de las imágenes a una escala [0,1]. Tras esta, se encuentran

los tres bloques descritos y posterior a estos, una capa de *agrupación* por simple coherencia estructural debido a que reduce la dimensionalidad pero conserva la mayor parte de información relevante (utilizada en la práctica totalidad de CNNs). Otra capa *dropout* por el mismo motivo que en los bloques y finalmente una capa densa *softmax* con la que obtendremos una distribución de probabilidad y así adquirir un output clasificatorio del mismo tamaño que número de letras reconoce la red neuronal. Toda esta arquitectura se ilustra, simplificada, en la Figura 6.2.

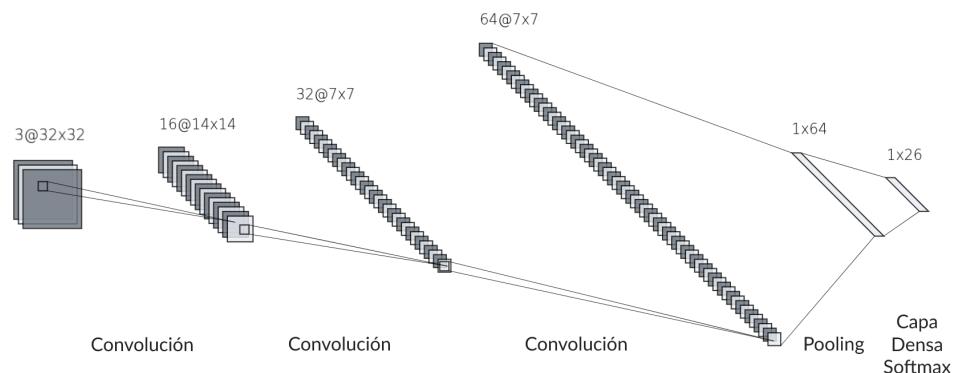


Figura 6.2: Esquema de funcionamiento del la red neuronal diseñada

6.2 Implementación

6.2.1 Preparativos

Previo al diseño, se han de instalar ciertas dependencias y definir algunos parámetros asistentes para el código. La mayoría de dependencias son para trabajar con estructuras de datos en python, interfaces con el OS, o librerías para *TensorFlow* y *Keras*. También se ha hecho uso de algunas funciones del trabajo de *Pete Warden* [44].

Otra dependencia imprescindible es `xxd`.

Teoría 6.2.1: Uso de `xxd` [38]

Es esencial contar con `xxd`, esto se debe a que la mayoría de microcontroladores no tienen soporte para sistema de archivos nativo. Con `xxd` obtenemos el modelo en formato matriz de chars directamente compatible con C/C++ e integrable en cualquier microcontrolador.

Es recomendable establecer esta matriz como constante por cuestiones de eficiencia en el acceso a memoria.

También necesitaremos el dataset para el entrenamiento, validación y testeo del modelo. Se ha decidido que en pos de la experimentación y de un mejor desempeño de la red neuronal al entrenarla con muestras generadas de la misma forma que se generarán las que se procesen cuando se haga uso de esta en el microcontrolador; que se generarán las muestras que utilizaremos como dataset, es decir, se va a generar un dataset propio. Este dataset se irá actualizando a medida que se toman más muestras y que se podrá descargar desde el *Google Drive* institucional. Se acompaña la descarga de una comprobación de existencia del fichero descargado para concluir esta sección.

Con el dataset descargado, se almacenan los trazos en un array, en el que cada elemento del array será un trazo; de momento de todos los labels.

Los trazos en este momento están cargados como conjunto de coordenadas que conforman, unas seguidas de otras, un recorrido; presumiblemente una letra. Pero como ya es sabido, nuestro modelo necesita como entrada imágenes; por lo que el siguiente paso es rasterizar los trazos para producir una imagen.

Al igual que en el firmware del microcontrolador, se hará uso de la función de rasterización creada por *Pete Warden* [44]. Definidas las funciones de rasterización, es posible rasterizar todos los trazos del dataset y destinarlos a las distintas fases del desarrollo del modelo como se verá en las siguientes secciones.

6.2.2 Entrenamiento

El entrenamiento del modelo será el segundo puntal que, junto con el diseño del propio modelo, dotará de estabilidad a nuestra red neuronal; siendo ambas dos determinantes para que esta responda de forma óptima a la tarea para la que a la que ha sido dispuesta. Por lo que la configuración del entrenamiento debe ser lo mejor posible, ajustando *epochs* (iteraciones en el proceso de entrenamiento), *learning_rate* (escalabilidad del aprendizaje), *optimizadores*, etc. Para ver el ajuste experimental de estos valores, véase el *apéndice B.5.1*.

Las imágenes de los datasets ocuparán 32x32 píxeles. Este es otro parámetro que puede ser estudiado, no obstante, estas dimensiones han sido las que han arrojado mejores resultados de forma homogénea. Como de costumbre, para más información a este respecto, consultar el *apéndice B.5.2*.

Se utilizará el mismo conjunto de datos aleatoriamente distribuido para cada uno de los tres dataset. Cada dataset contará con un porcentaje del conjunto de datos total. Estos porcentajes han sido estudiados y a priori no suponen extrema relevancia más allá de que el de entrenamiento debe ser ampliamente mayor. Ha sido fijado un 10% para test, otro 10% para validación y el restante 80% para entrenamiento.

Teoría 6.2.2: Uso del dataset en Deep Learning(1) [25]

Los tres datasets que se usan para Deep Learning son:

■ Validation

En Deep Learning se usan datos de validación para corroborar durante el entrenamiento, que el ajuste se está dando de forma óptima.

Dilatando un poco más esta sencilla explicación, el *validation dataset* es un conjunto de datos imperativamente distinto del *training dataset*, que sirve para estimar la eficacia de la red en tiempo de entrenamiento.

En general se suele usar la validación para hacer estudios del ajuste del modelo, para evitar sobreajustes (*overfitting*) y subajustes (*underfitting*).

Teoría 6.2.3: Overfitting y Underfitting

- ***Overfitting***

Es un fenómeno que se da cuando el modelo reconoce peculiaridades demasiado específicas como distintivo para la evaluación. Estas peculiaridades no serían los rasgos o características que constituyen a los elementos que estudiados y por lo tanto se produce un sobreajuste; un ajuste por encima de lo óptimo.

- ***Underfitting***

Término análogo y opuesto al anterior, el ajuste se presentaría laxo y falto de rigurosidad; ajuste por debajo de lo óptimo.

■ Training

Este dataset es el más simple de entender por mera inmediación semántica. Es el conjunto de datos que se utiliza en tiempo de entrenamiento para balancear los pesos de las capas. En cada iteración de entrenamiento, se calcula la pérdida con los datos de entrenamiento introducidos y se da el ajuste de pesos en base a la pérdida. Esto supone que, cada vez la pérdida sea menor y generalmente la eficacia, o en términos más comunes a este ámbito, la *precisión* (*accuracy*), sea mayor.

■ Test

El conjunto de datos que se utilizará posterior al entrenamiento, para validar la efectividad del entrenamiento. Es el dataset con el que se pone a prueba el modelo entrenado.

6.2.3 Testeo

El testeo del modelo es útil para garantizar que funciona correctamente, sin embargo podemos obtener valores muy buenos sin resultar en un funcionamiento adecuado, debido al ya mencionado *overfitting* (*Teoría 6.2.2*). Por tanto solo debemos tomarlo como una herramienta más, siempre evaluando adecuadamente los resultados.

La fase de testeo consiste en simular lo que será una ejecución habitual, pero conociendo los inputs de la red neuronal, que serán del dataset de testeo. Obteniendo imágenes de cada letra (*label*), se evalúa la clasificación de la red neuronal de la misma, junto con la precisión de estimación; si la clasificación coincide y además lo hace con valores de precisión próximos a 1, significará, en principio, un buen desempeño de la red neuronal para esta letra. El objetivo es conseguir esto para todas las letras.

Es importante separar los dataset de testeo de los de entrenamiento y validación, ya que si se emplearan los mismos, el modelo ha sido entrenado con este mismo input, por tanto siempre se obtendrían buenos resultados y el testeo perdería validez.

Para un mayor acercamiento a la experimentación con el testeo de los modelos probados, visite el *apéndice B.5.2*

6.2.4 Transformación a modelo cuantizado

El hecho de cuantizar el modelo basado en *Deep Learning* cuando se plantea su integración en microcontroladores y equipos de bajo rendimiento, suele ser imperativo.

Teoría 6.2.4: Cuantización

Cuantizar es desvirtuar la naturaleza continua de un conjunto de valores continuos, restringiéndolos a un conjunto de valores discretos.

En general esta práctica viene propiciada debido a que la mayoría de microcontroladores no cuentan con *FPU* (*Floating Point Unit*). Sin embargo no es este el caso, ya que en este proyecto se escogió un microcontrolador que sí dispone de esta unidad. Aun así se ha tenido que recurrir a la cuantización del modelo, como consecuencia de su limitación de memoria, ya que el hecho de cuantizar el modelo, supone reducir su precisión y por tanto reducir la memoria que este requiere.

Sin entrar en muchos detalles, ya que es parte de la siguiente sección, el modelo generado estándar queda muy lejos de poder integrarse en el microcontrolador, gracias a la conversión a un modelo *TensorFlow Lite*, se reduce el tamaño, pero sigue sin ser suficiente, por tanto es conveniente cuantizar el modelo *TensorFlow Lite* para obtener su versión más compacta posible.

6.2.5 Comparación de modelos generados

Los modelos generados para la red neuronal, se muestran en la Tabla 6.1.

Modelo	Tamaño	Reducción TF	Reducción TFLite
TensorFlow	756009 Bytes	0%	-
TensorFlow Lite	140300 Bytes	81'4%	0%
TF Lite cuantizado	41136 Bytes	94'5%	70'6%

Tabla 6.1: Tabla de comparación de modelos generados

Los resultados son impresionantes en sí mismos, pero lo son aún más cuando ponemos en contexto los valores de precisión que se alcanzan para cada uno de ellos en la clasificación; variando en el peor de los casos del orden de 0,1% de *TensorFlow* a *TensorFlow Lite* y del orden de 0,2% de *TensorFlow* a *TensorFlow Lite* cuantizado. Por lo que la compactación del modelo no tiene un impacto notorio durante su ejecución.

6.2.6 Integración en el microcontrolador

Esta sección sirve ahora a modo de recopilación de las anteriores, ya que el procedimiento prácticamente ha sido expuesto.

La cuantización del modelo para optimizar el uso de memoria es el primer paso; una vez el modelo cuenta con un tamaño apto para el microcontrolador, lo que resta es solucionar el problema de la carencia de sistema de archivos necesario para el manejo de la red neuronal construida con *TensorFlow*. Para lo cual se hace uso de *xxd*, herramienta que automatiza el proceso de conversión del modelo a una estructura de datos (*vector de chars*), agregable al firmware del microcontrolador.

Cuando pasamos el modelo por *xxd* y obtenemos el código en *C*, es suficiente con añadirlo al código del *firmware* y este ya es funcional en el microcontrolador. Con el que se trabajará gracias a las librerías de *TensorFlow Lite*.

6.3 Generación de muestras para el dataset

Por suerte se podremos contar con una herramienta creada por *Pete Ward* [44] para tomar muestras de trazados con el microcontrolador que estamos empleando. Se ha modificado superficialmente y se han hecho algunos cambios estéticos ([DataCollector.html](#)). Este proceso supone, para el sujeto que genera las muestras, un tiempo de adaptación al movimiento que se debe ejecutar para que se recoja un buen trazado, añadido al hecho de que las muestras que no sean óptimas, deben eliminarse; lo convierten en un proceso cargante y lento para quien está creando las muestras y para quien tiene que supervisarlas.

A lo largo de la toma de muestras se han detectado ciertos fallos que se han resuelto y son consultables en los Apéndices B.1, B.3 y B.2.

Capítulo

7

Interfaz de usuario

7.1 Motivación

Antes de comenzar con la planificación, debe ponerse en valor la conveniencia de esta interfaz de usuario. El principal motivo es que lo que se busca con este proyecto es crear un producto real, y como tal, no podemos valernos de un simple lector de puerto serie para el dispositivo, como podría ser el integrado en el *Arduino IDE*. Además, esta interfaz no solo está creada con la funcionalidad de este proyecto en mente de servir de mediador con el usuario para la lectura de las letras escritas con el *SmartPen*, sino que busca ser el nexo de unión de todas las funcionalidades que se implementen para este. Durante el propio desarrollo de esta interfaz, se ha reflexionado sobre ciertas funcionalidades, como la de un modo pizarra donde reflejar el trazado íntegro que se realice, y que se ha propuesto finalmente como una meta secundaria.

7.2 Planificación

En esta sección se recogerán los motivos por los que se necesita de esta interfaz de usuario, la elección de herramientas para su desarrollo y el bocetado de la distribución de sus elementos.

7.2.1 Elección de framework para interfaz gráfica

Las alternativas con las que vamos a partir, tras una fase de documentación sobre frameworks para desarrollo de interfaces gráficas, preferiblemente multiplataforma, y sumarlas a las que ya conocía; se recogen en la Tabla 7.1, junto a la valoración de algunas características importantes.

Electron y *GTK* quedan descartados porque, si bien existen procedimientos para poder gestionar *Bluetooth Low Energy* con estos, no tienen soporte nativo

Alternativas	Documentación	Usada anteriormente	Comunidad	Soporte para BLE	Lectura P.Serie nativa	Multiplataforma
Flutter	✓	✓	✓	✓	✓	✓
GTK+	✓	-	✓	✗	✗	✓
QT	✓	✓	✓	✓	✓	✓
PyQT	✓	✗	✓	✓	✓	✓
Electron	✓	✗	✓	✗	✓	✓

Tabla 7.1: Tabla de elección de *frameworks* para desarrollo de *UIs*

ni librerías para ello, por tanto, teniendo la posibilidad de hacer uso de otros frameworks que faciliten esta tarea, es factible descartar estas opciones.

Ante *QT (C++)* y *PyQT*, debido a que se ha trabajado anteriormente con *QT (C++)* y el tiempo de desarrollo disponible es un factor limitante, es razonable excluir *PyQT*.

Solo quedarían *QT* y *Flutter*, se ha escogido ante estas dos posibilidades *QT* por dos razones. La primera razón es que *Flutter* está más dirigido a interfaces móviles. Y la segunda razón es que pese a haber trabajado con ambas, a título personal, me encuentro más habituado a *QT*.

Por lo que el framework seleccionado para realizar la interfaz de usuario, será *QT* ya que ofrece soporte para las tecnologías que en principio van a utilizarse y porque es una herramienta con la que ya se ha trabajado y esto resulta en un menor tiempo de documentación, causa de peso por el tiempo tan ajustado.

7.2.2 Diseño de la interfaz

Para el bocetado de la interfaz previo a la implementación, se hará uso de *QT Design Studio*, herramienta de *QT* para diseñar interfaces. En el bocetado de la interfaz se fijarán los elementos que la formarán y su disposición procediendo razonadamente. Para entender visualmente lo que se va a describir, puede examinar el *apéndice C.7*.

Como se ha razonado en la Sección 7.1 anterior, queremos que este programa sea el *hub* de funcionalidades para el *SmartPen*, por lo que se hará uso de una barra lateral para acceder a estas. Queremos que el protagonismo se encuentre en la funcionalidad, por lo que esta barra lateral será desplegable, para reducir el tamaño que esta ocupe en la interfaz. Alojará las funcionalidades y, dada la naturaleza académica del proyecto, una sección sobre el desarrollador, para más

información sobre el trabajo.

Cada sección de la interfaz, se mostrará en la pantalla central, quedando las barras superior y lateral, constantemente a la vista ya que la barra lateral sirve para acceder a otras secciones de la interfaz y la barra superior muestra información vital.

Es necesario un indicador de conexión y un botón de conexión inalámbrica, que se ubicarán a la derecha en la barra superior para evitar sobrecargar la zona izquierda de la interfaz.

En los primeros bocetos se evidenciaba un vacío en la parte central de la barra superior y también se echaba en falta saber qué pantalla se estaba mostrando en cada momento, por lo que uniendo ambas carencias, se solventó disponiendo un indicador de pantalla en dicha región.

7.3 Implementación

En cuanto a la implementación software, al margen de aplicar el diseño descrito en la Sección 7.2.2 anterior, los retos planteados son principalmente dos. El primero de ellos: la gestión simultánea de la lógica de la interfaz, la lectura del puerto serie y la lectura de *características* del servicio *Bluetooth Low Energy* (definido en la Sección 5.2.2). El segundo: la implementación de lectura, especialmente la lectura haciendo uso del servicio *Bluetooth*.

Para afrontarlos, se hará uso de tecnologías que nos proporciona el propio *framework*. Para gestionar simultáneamente varios procesos, es posible hacer uso de *hilos* (o *threads*) gracias a los *QThreads* de *QT*, que proveen de una clase para gestionar *hilos* independientemente del sistema. Por lo que crearemos un hilo concurrente a la ejecución de la lógica de la interfaz, que se hará cargo de los procedimientos de lectura, tanto *Bluetooth* como del puerto serie; la clase *ReaderThread*.

Por otro lado, para resolver el problema de la lectura de la placa, tanto la lectura por cable como la inalámbrica, se emplearán librerías que asisten para estas tareas. Ya se planteó como requisito que el *framework* de creación de la interfaz, contara con soporte para las comunicaciones *Bluetooth* y *puerto serie*. Estas librerías son *QLowEnergy** y *QBluetooth**, y *QSerialPort** y son parte de las dependencias que hay que añadir al proyecto.

El resultado de todo lo implementado será lo que refleja la Figura 7.1.

7.4 Traspaso del diseño a *QT creator*

Siguiendo el esquema de diseño de la Sección 7.2.2 (*Diseño de la interfaz*) y gracias a la herramienta *Design* de *QT creator*, el traspaso es trivial; crear la estructura, añadir botones, widgets y demás elementos de la interfaz. Con estos dispuestos, ya es posible comenzar a trabajar con la lógica detrás estos y su configuración.

Para ver el resultado obtenido, consulte el apéndice C.7.

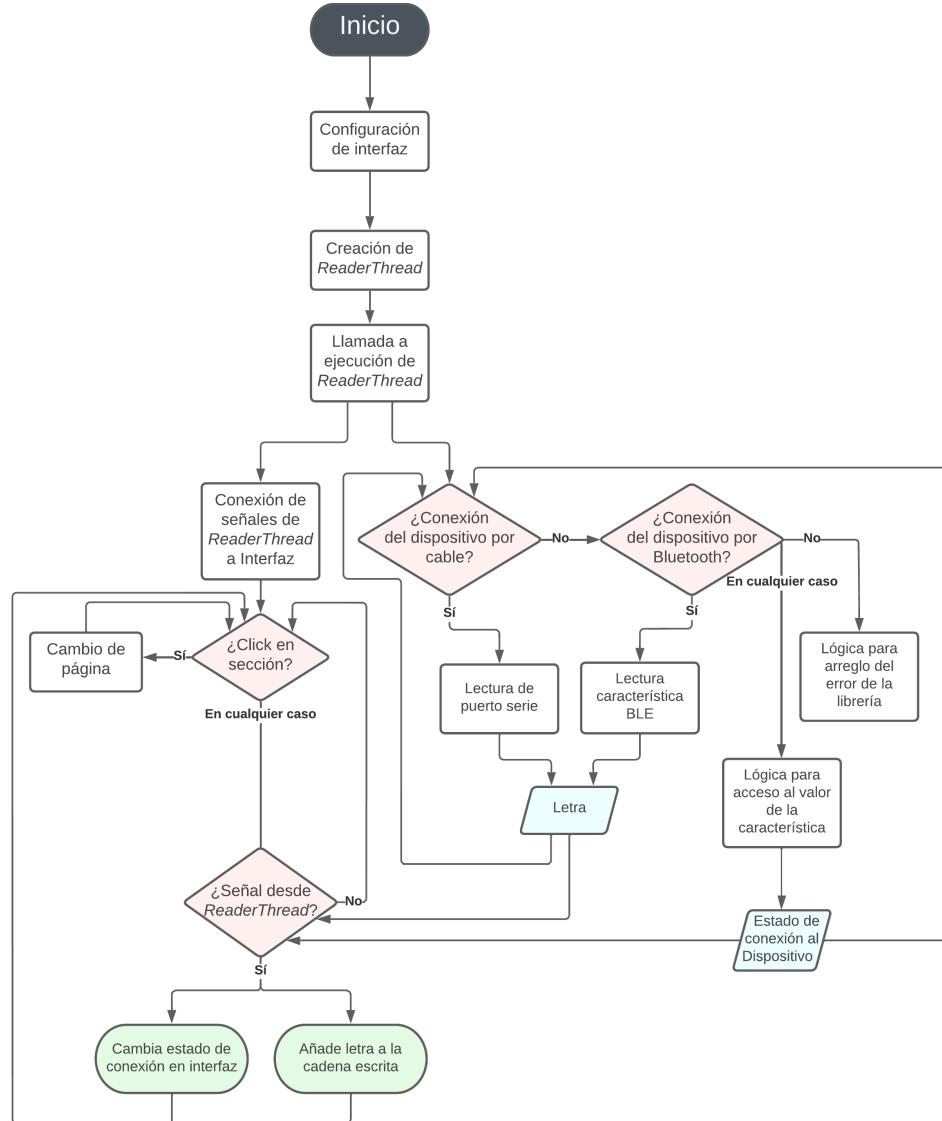


Figura 7.1: Diagrama de flujo simplificado del UI

7.5 Configuración de la interfaz

La configuración de la interfaz consistirá en inicializar los parámetros con los que parten los elementos creados en la Sección 7.4 anterior y definir su comportamiento durante la ejecución.

Para el acceso a cada sección, se plantea un widget que varía su contenido en función del índice seleccionado. De esta forma se guarda el resto de la interfaz

inmutable. La configuración quedará de tal forma que al pulsar cada botón de sección, cambiará título, sección seleccionada, y contenido de este widget.

Esta sección se basa de forma general, en fijación la trivial de parámetros de estilo y establecer el comportamiento de los elementos dispuestos en pantalla.

7.6 Gestión de lectura del microcontrolador

El thread anteriormente citado en la Sección 7.3, *ReaderThread* será el encargado de esta labor, aunando en una única señal, las lecturas de *Bluetooth* y *puerto serie*, y definiéndose de esta forma como una interfaz de lectura para el programa. Esta interfaz de comunicación entre *ReaderThread* e interfaz se llevará a cabo mediante *señales* y *slots* [31] de QT. Y será así porque es una forma de comunicación que permite la interacción entre entidades concurrentes de una manera cómoda e intuitiva.

7.6.1 Lectura de la característica *Bluetooth Low Energy*

Para esta sección será determinante una buena documentación, ya que la implementación es un poco complicada cuando nunca se ha trabajado antes con esta tecnología. Gracias a contar con ello como requisito para la elección del *framework* de creación de la *UI*, podemos contar con documentación al respecto [33]. Aunque es reseñable que bien por complejidad de la arquitectura *BLE* o bien porque la documentación es sucinta en exceso, el proceso de aprendizaje para estas librerías, puede llegar a ser más lento de lo que cabría esperar.

Se creará una clase auxiliar (*Device*) que lidie con la gestión *bluetooth* del programa: búsqueda de dispositivos, conexión al *SmartPen*, notificación de estados de conexión, obtención de servicios, características, descriptores y valores, y evidentemente lectura de las letras. La comunicación entre *Device* y el *ReaderThread* se llevará a cabo, al igual que entre *ReaderThread* e interfaz, mediante *señales* y *slots* [31].

Problemas 7.6.1: Error de permisos para *qt.bluetooth*

En ocasiones se experimentan crashes durante la ejecución, unido a un error derivado del uso de *bluetooth* en los *logs*, lo que me llevó a investigar este problema para resolverlo. Llegando a que se debía a un error de permisos con el protocolo que gestiona el *bluetooth* en *Linux* (*BlueZ*), y que más adelante seguirá causando problemas, pero en este caso, pudo resolverse como indica el Apéndice C.6.

Problemas 7.6.2: Errores de desconexiones *Bluetooth Low Energy*

Possiblemente el error que más tiempo ha ocupado y que lamentablemente como desarrollador, solo se puede mitigar. Y es que es un problema aparentemente popular entre los desarrolladores que implementan mecanismo *bluetooth LE* para *linux*, ya que se experimentan desconexiones del dispositivo sin aparente causa, ya que no son debidas a la implementación del desarrollador, sino de la librería que gestiona las comunicaciones *bluetooth* en *QT* haciendo uso de los protocolos provistos por *linux*.

Resulta especialmente problemático, no por el hecho de la desconexión, que también, sino porque se pierden valores durante esta; provocando sensación de que el dispositivo no funciona como debe. La pérdida de valores, tras mucha investigación y el aporte de la comunidad, tiene solución editando la propia librería que causa el problema; desarrollado en el apéndice C.5.

Desafortunadamente, el problema de las desconexiones persiste y aparece deberse a *BlueZ*, el *stack* de *bluetooth* de *linux*, que provee de los protocolos y capas necesarias para trabajar con este. Y por tanto, su solución queda fuera de mi margen de actuación, o al menos en un tiempo tan limitado. Sin embargo sí que he implementado un ajuste para que de cara al usuario, solo suponga un pequeño retardo, en ocasiones, de la letra escrita y que se basa en el número de iteraciones en el proceso de conexión sin que se notifique señal por parte del dispositivo periférico; de forma que sí se notifique una desconexión persistente real, pero no las desconexiones provocadas por la falla descrita.

Pese a los problemas que se han encontrado, y que no podían sospecharse hasta adentrarse en la implementación, el resultado es satisfactorio y la experiencia de usuario es prácticamente la misma que si no existieran estos problemas.

7.6.2 Lectura del *puerto serie*

La implementación para la lectura del *puerto serie* es mucho más sencilla que la de *bluetooth*. Como la gestión es mucho más simple y apenas es necesario trabajar con una sola clase que media con el *puerto serie*, se implementará en una simple función de lectura en el propio *ReaderThread*, la cual se ejecutará en caso de que se detecte conexión en el puerto serie y que será, el caso predeterminado cuando exista posibilidad simultánea de conexión por cable e inalámbrica.

Véase el Apéndice C.2 para acceder a detalles de la implementación para la lectura del puerto serie.

Capítulo

8

Encapsulado

Para esta parte, dada la menor relevancia al poner en contexto todo el proyecto, se contará con poco tiempo. Esto supondrá decisiones tomadas con esta limitación,

8.1 Herramientas utilizadas

Para la creación de los modelos 3D del encapsulado, se han utilizado principalmente dos herramientas distintas para modelar y una para imprimir.

Para el modelado se utilizará *Blender*, un software de modelado muy potente y lleno de posibilidades, aunque considerablemente complejo al empezar y que se utilizará para no limitar el modelo a algo simple. Fuera de lo planificado, también se ha utilizado *SketchUp*, por los motivos que se desarrollan en Problemas 8.2.2.

Para la impresión, de la que yo no me encargaré, se hará uso del *slicer* (software para impresión 3D) *Ultimaker Cura*.

8.2 Implementación

Por la limitación de tiempo descrita al comienzo de este Capítulo 8, no se dispone de tiempo para una fase de diseño previa a implementar el modelo. El modelo se ha construirá con las especificaciones más básicas en mente.

El encapsulado debe contener de forma firme el microcontrolador para que no se mueva, ya que se espera que el usuario lo mueva con determinación para trazar una letra. Por otro lado, también será necesario incorporar una batería para alimentar el dispositivo sin conexión directa al ordenador. Con estas dos únicas limitaciones y con una forma semejante a la de un bolígrafo, se construirá el modelo.

El modelo estará dividido en componentes por motivos de la logística de la impresión, tanto por comodidad para tratar con la persona que lo imprimirá, como por cuestiones de que los modelos deben guardar ciertas limitaciones

estructurales para poder imprimirse. Por ejemplo, en lugar de imprimir el [cilindroBajo](#) ya con la [punta](#), para aunar toda la parte inferior del encapsulado, se imprimen por separado, ya que se necesita de una base estable para la impresión que no se lograría con los dos elementos unidos.

Problemas 8.2.1: Modelo incompatible con la impresión

Durante la impresión surgieron varios problemas, la mayoría menores. Sin embargo hubo uno que se resistió, aunque no entrañaba realmente ninguna complejidad. Este era que una de las partes del encapsulado, el [cilindroBajo](#), la parte que contendrá al microcontrolador, mostraba en el [slicer](#) (software para la impresión), un error de formato en la base. Este error finalmente se debía a un grosor por debajo del mínimo. Problema ocasionado por haber construido este componente a partir del escalado vertical de otro ([cilindroAlto](#)); al escalar verticalmente, el grosor de la pared del modelo, se conserva, sin embargo no ocurre lo mismo para la base, disminuyendo su grosor y provocando este problema. La solución fue simplemente redefinir el grosor de la base.

Con todos los componentes impresos, surge un problema con la batería con la que contaba para incorporar en el encapsulado

Problemas 8.2.2: Alimentación intermitente de la batería

La alimentación resultaba deficiente, por lo tanto solo pude hacerme con otra batería y adaptar el encapsulado a este nuevo componente. Para poder crear el modelo complementario para adaptar la batería, hice uso de *SketchUp*, un software muchísimo más simple que *Blender*; esta sencillez ha sido determinante porque el tiempo era muy restrictivo en ese momento. Resultando en el componente [slotBatería](#), para poder incluir una pila de 9V recargable. La alimentación con esta pila, al ser de más de 5V, deberá ser vía pines de la placa, que admiten hasta 21V.

Con esta nueva batería, se logra una autonomía de más de dos días en el primer ciclo de carga. Según el fabricante, los primeros ciclos de carga suelen ser los que menor energía suministran, por lo que podría incluso mejorar la duración.

La batería se conecta directamente a los pines de alimentación (*V/N*) y masa (*GND*) del microcontrolador.

En los Apéndices D.1 y D.2 pueden encontrarse ilustrados, respectivamente, los modelos creados y el resultado final del encapsulado con todo el hardware integrado.

Capítulo **9**

Validación

9.1 Ajuste al presupuesto

Respecto a como se estimó en la sección 3.3.2 (*Presupuesto*), se ha ajustado satisfactoriamente el coste del dispositivo, resultando todos los costes los reflejados en la Tabla 9.1.

Descripción	Precio
Arduino Nano Sense 33 BLE	35'80\$~33'82€
Pila 9V Recargable	10'99€
Adaptador pila 9V	3€
Impresión	1€
Interruptor	0'05€
Adaptador microUSB a USB	1€
Cable MicroUSB Datos	6€
Tiempo de trabajo	~3900€
EQUIPO: 55'86€	
TOTAL: 3955'86€	

Tabla 9.1: Costes de producción del *SmartPen*

Respecto a las horas de trabajo, se ha respetado la estimación de la planificación, ya que he tratado de ceñirme a estas durante todo el desarrollo del proyecto.

9.2 Comprobación de objetivos cumplidos

Los requisitos para el proyecto y sus derivadas especificaciones, han sido satisfactoriamente cumplidas.

El presupuesto se ha ajustado a lo estimado tanto en su aspecto de trabajo como en el de ajuste para los fondos empleados en el *hardware*.

El microcontrolador efectivamente, como se propuso, hace uso de un modelo basado en *Deep Learning* y además lo hace de manera eficaz. Para alimentar al modelo de procesamiento se hace uso de los sensores que ofrece el microcontrolador, el cual cuenta con dimensiones muy reducidas para poder acoplarse en un encapsulado embellecedor que dará cabida también a una batería que lo alimentará para, como se especifica, dotarlo de autonomía. Autonomía necesaria ya que el dispositivo funciona tanto por cable como de forma inalámbrica, tal y como se demandaba.

Todo el desarrollo es transparente y de código abierto, por lo que está acondicionado para que otros desarrolladores y usuarios hagan uso del progreso logrado en este despliegue del producto y aporten funcionalidades al producto si así lo consideran.

Todo el *software* utilizado cumple con las especificaciones, porque de otra forma, no podría haberse llevado a cabo el proyecto. La interfaz de usuario cumple con la funcionalidad mínima propuesta; el *framework* para el diseño, entrenamiento, validación y testeo de la *red neuronal* es el mejor que podría haberse empleado; el *firmware* del microcontrolador ejecuta exactamente lo que se ideó al comienzo del proyecto; y respecto al *software* para el modelado del encapsulado, se han trabajado con varias alternativas en función de lo que se requería en cada momento.

En general el producto no solo realiza su cometido, sino que lo hace de forma satisfactoria.

Capítulo 10

Trabajos futuros y mejoras

Si bien el dispositivo realiza lo que se postulaba al comienzo como requerimientos. Sin embargo quedan muchas mejoras y ampliaciones posibles, en su mayoría por carencia de tiempo.

Una de las que más me incordian, es el no haber podido añadir la funcionalidad de 'modo pizarra', ya que si bien no estaba planteado como requisito, era algo que pretendía hacer desde el planteamiento del trabajo, ya que me resulta especialmente útil como complemento a una herramienta que pretende ser un híbrido entre lápiz y teclado. De hecho la solución era integrar el módulo de exposición del trazado de la herramienta de recolección de datos (*SmartPen_DataCollector*) en el propio interfaz de usuario, haciendo uso de las herramientas de *QT* para visualización web. Sin embargo al emplear la web *bluetooth* para la recolección del trazo generado por el microcontrolador, dificultaba mucho la implementación y viendo que se alargaría más de lo esperado, se tomó la decisión de posergarlo.

Otro añadido que habría dado más autonomía al dispositivo, es el uso de un buffer de memoria en el microcontrolador para almacenar las letras que se escribieran previas a la conexión con la interfaz de usuario. Y que es en realidad una funcionalidad muy fácil de implementar ya que solo habría que crear el propio buffer en el microcontrolador y cambiar la lógica de la comunicación interfaz-microcontrolador de una letra a varias letras consecutivamente enviadas.

Esta sin embargo sería una mejora mucho más compleja y que llevaría un proceso de documentación a bajo nivel, un tanto distendido. El problema de desconexiones descrito en Problemas 7.6.1. La solución sería encontrar el problema que desemboca en la desconexión y que parte de la librería mencionada y el stack de protocolos *bluetooth* de *linux*.

Ocasionalmente la interfaz de usuario sufre problemas de corrupción de memoria. Son tan raros que trazarlos es muy complicado y no he sido capaz de entender la causa. Este es un problema que queda pendiente a resolver al no disponer de más tiempo. Para resolverlo, el primer paso podría ser utilizar *Valgrind*

para dar con el origen del problema, que con total seguridad, será de nuevo, algo respectivo a la librería mencionada o a la gestión de *QThread*.

Lo cual nos lleva a otra posible mejora y es la implementación de la interfaz para *Windows* o *macOS*. Me gustaría poder haberla realizado, pero el tiempo no alcanza para más y dado que trabajo habitualmente en *linux*, era más inteligente comenzar implementando la interfaz solo para este sistema. En principio la solución pasaría por, simplemente, añadir la configuración de puertos correspondientes para cada sistema.

Una mejora que significaría un rango mucho mayor de detección para el *SmartPen* y algo más de precisión, sería aumentar el número de muestras con el que entrenar el modelo. Pero como se expuso en la Sección 6.3, es un proceso que lleva muchísimo tiempo. Aunque sencilla, esta mejora supone emplear tiempo del que no se dispone.

En realidad hay un sinfín de posibles mejoras y nuevas funcionalidades, pero por esto mismo y dado que es un proyecto interesante y lleno de posibilidades, se ha planteado su desarrollo abierto y libre a aportaciones de la comunidad.

Capítulo 11

Conclusiones

Las conclusiones tras el desarrollo del proyecto y ver los resultados, son que es algo de lo que estaría orgulloso si lo hubiese visto al comienzo del mismo. Sin embargo no puedo evitar sentir que, por todo lo que tenía en mente haber implementado, podría haber sido incluso mejor. Aunque ya pude suponer desde el principio que el tiempo del que disponía no era suficiente para llenar de funcionalidades el proyecto, por eso planteé lo que queda descrito en esta memoria de la forma que lo hice y he cumplido con ello de forma satisfactoria.

Por otro lado, el haber trabajado con *Deep Learning, redes neuronales, capas, entrenamientos* y otros tantos conceptos que me resultaban tan llamativos e interesantes, y haberlo podido unir al campo del *hardware* que tanto me apasiona, mediante la integración del procesamiento en un microcontrolador, es otra de las razones por las que me complace haberme decidido por este proyecto.

Ha sido un trabajo complicado desde su propio planteamiento, los complejos mecanismos que se emplean, el no haberlos utilizado nunca antes y por los problemas que han ido surgiendo durante el desarrollo, fruto de la implementación. Pero el resultado y el aprendizaje surgido del esfuerzo para poder completarlo, han hecho que valga la pena.

Finalmente, es necesario concluir que este no es el final de este proyecto, planeo continuar con su desarrollo como pasatiempo y me llenaría de orgullo que otras personas interesadas en estos campos, participaran también agregando sus aportaciones como se plantea al hacer todo el proyecto público y abierto.

Bibliografía

- [1] Adarsh1001. Repositorio micro-learn.
<https://github.com/adarsh1001/micro-learn>.
- [2] Andriyadi. andriyadi magic wand repository.
<https://github.com/andriyadi/MagicWand-TFLite-Arduino>.
- [3] Developers Android. Descripción general del bluetooth de bajo consumo para android. <https://developer.android.com/guide/topics/connectivity/bluetooth-le?hl=es-419>.
- [4] Arduino. Project hub. https://create.arduino.cc/projecthub?by=part&part_id=108462&sort=trending.
- [5] Aprendiendo Arduino. Aprendiendo a manejar arduino en profundidad.
<https://aprendiendoarduino.wordpress.com/tag imu/>.
- [6] Bluetooth. A developer's guide to bluetooth technology.
<https://www.bluetooth.com/blog/a-developers-guide-to-bluetooth/>.
- [7] Fernando Sancho Caparrini. Redes neuronales: una visión superficial.
<http://www.cs.us.es/~fsancho/?e=72>.
- [8] CloudML. Cloudml webpage. <https://cloudml.io/>.
- [9] QT Community. Codereview de qt-project. <https://codereview.qt-project.org/c/qt/qtconnectivity/+/233087>.
- [10] Antonia Creswell, Tom White, Vincent Dumoulin, Kai Arulkumaran, Biswa Sengupta, and Anil A. Bharath. Generative adversarial networks: An overview. 2018.
- [11] DeepMind. Alphafold2 webpage. <https://www.deepmind.com/research/highlighted-research/alphafold>.

- [12] DeepMind. Alphago webpage. <https://www.deepmind.com/research/highlighted-research/alphago>.
- [13] Free Software Foundation. GNU General Public License. <http://www.gnu.org/licenses/gpl.html>.
- [14] Hinton G.E., Osindero S., and Teh Y. Movies of the neural network generating and recognizing digits. 2006.
- [15] GitHub. Github copilot webpage.
<https://github.com/features/copilot/>.
- [16] GitHub/Microsoft. EdgeML webpage.
<https://microsoft.github.io/EdgeML/>.
- [17] Aurélien Géron. *Hands-on Machine Learning with Scikit-Learn, Keras, and TensorFlow*. O'Reilly, 2019.
- [18] lars.geo. Breve historia de las redes neuronales artificiales.
<https://steemit.com/spanish/@iars.geo/breve-historias-de-las-redes-neuronales-artificiales-articulo-1>.
- [19] Keras. Keras webpage. <https://keras.io/>.
- [20] Ladvien. Repository named 'arduino_ble_sense'. https://github.com/Ladvien/arduino_ble_sense.
- [21] Alex Lenail. Creación de esquemas de redes neuronales. <http://alexlenail.me/NN-SVG/index.html>.
- [22] Ben Lutkevich. Natural language processing (NLP).
<https://www.techtarget.com/searchenterpriseai/definition/natural-language-processing-NLP>.
- [23] Maxime. What is a transformer?
<https://medium.com/inside-machine-learning/what-is-a-transformer-d07dd1fbec04>.
- [24] Warren S. McCulloch and Walter Pitts. A logical calculus of the ideas immanent in nervous activity. 1943.
- [25] Andreas C. Müller and Sarah Guido. *Introduction to Machine Learning with Python*. O'Reilly, 2017.
- [26] Na8. Breve historia de las redes neuronales artificiales.
<https://www.aprendemachinelearning.com/breve-historia-de-las-redes-neuronales-artificiales/>.

- [27] Jonathan Stephens (Nvidia). Getting started with nvidia instant nerfs. <https://developer.nvidia.com/blog/getting-started-with-nvidia-instant-nerfs/>.
- [28] OpenAI. Openai webpage. <https://openai.com/>.
- [29] Peltarion. Global average pooling 2d. <https://peltarion.com/knowledge-center/documentation/modeling-view/build-an-ai-model/blocks/global-average-pooling-2d>.
- [30] QT. Documentación Bluetooth LE. <https://doc.qt.io/qt-6/qtbluetooth-le-overview.html>.
- [31] QT. Documentación de señales y slots en qt. <https://doc.qt.io/qt-6/signalsandslots.html>.
- [32] QT. Documentación QThread. <https://doc.qt.io/qt-6/threadstechologies.html>.
- [33] QT. Implementaciones qt para demostrar el uso de ble. <https://doc.qt.io/qt-6/qtbluetooth-lowenergyscanner-example.html>.
- [34] Partha Pratim Ray. A review on tinyml: State-of-the-art and prospects. 2022.
- [35] F. Rosenblatt. The perceptron, a perceiving and recognizing automaton project para. 1957.
- [36] Sierdzio. Foro qt. https://forum.qt.io/topic/114503/missing-cap_net_admin_permission/2.
- [37] Talent.com. Salario medio de un ingeniero informático en españa en 2022. <https://es.talent.com/salary?job=ingeniero+inform%C3%A1tico>.
- [38] TensorFlow. Build convert. https://www.tensorflow.org/lite/microcontrollers/build_convert.
- [39] TensorFlow. Github ejemplo magic_wand de tensorflow. https://github.com/tensorflow/tflite-micro/tree/main/tensorflow/lite/micro/examples/magic_wand.
- [40] TensorFlow. Introducción a los microcontroladores por tensorflow. https://www.tensorflow.org/lite/microcontrollers/get_started_low_level.
- [41] TinyML. Tinyml webpage. <https://www.tinyml.org/>.
- [42] Kevin Townsend, Carles Cufí, Akiba, and Robert Davidson. Getting started with bluetooth low energy. <https://www.oreilly.com/library/view/getting-started-with/9781491900550/ch04.html>.

- [43] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. 2017.
- [44] Pete Warden. Pete warden magic wand repository. https://github.com/petewarden/magic_wand.
- [45] Pete Warden and Daniel Situnayake. *TinyML: Machine Learning with TensorFlow Lite on Arduino and Ultra-Low-Power Microcontrollers*. O'Reilly, 2019.
- [46] Wikipedia. Intel mcs-48. https://en.wikipedia.org/wiki/Intel_MCS-48.
- [47] Wikipedia. Microcontrolador. <https://es.wikipedia.org/wiki/Microcontrolador>.
- [48] Wikipedia. Multilayer perceptron. https://en.wikipedia.org/wiki/Multilayer_perceptron.
- [49] Orhan G. Yalçın. Image classification in 10 minutes with MNIST dataset. <https://towardsdatascience.com/image-classification-in-10-minutes-with-mnist-dataset-54c35b77a38d>.

Apéndice **A**

Microcontrolador

A.1 Resolver problemas de memoria

A lo largo del desarrollo del firmware, será necesaria la gestión de áreas de memoria reservadas para ciertas funcionalidades, por ejemplo en el bloque de configuración del firmware, *setup()*; será necesario reservar un área en memoria para las labores de E/S y memoria intermedia de las acciones con tensor.

Para ello se contará con la propia *flash* del dispositivo a modo de región de almacenamiento. Creando arrays en la propia memoria *flash* que harán las veces de mecanismo de almacenamiento como se expone en el siguiente fragmento de código:

Fragmento de *deep_pen.ino*

```

62 // **** SETUP FUNCTION ****/
63 // Setting the area of memory reserved to tensor input-output actions.
64 constexpr int TensorAreaSize = 30 * 1024;
65 uint8_t tensor_arena[TensorAreaSize];

```

A.2 Instalación de librerías en *Arduino IDE*

La instalación de librerías es muy sencilla haciendo uso del *Arduino IDE*. En el menú: *Tools->Manage Libraries...* y se abrirá una ventana donde podemos gestionar las librerías instaladas, instalar otras versiones e instalar nuevas librerías.

A.3 Notificación de conexión *Bluetooth* del firmware del microcontrolador

Para solucionar el problema de la notificación de conexión *Bluetooth* se crea una pequeña estructura para consultar el estado de conexión de la iteración

anterior (*last_connection*). De esta forma es posible mantener consistencia en la conexión pese a estar dentro de un bucle.

A.4 Definición de micro-operaciones en el firmware del microcontrolador

Las micro-operaciones que se definirán en el firmware que se emplearán en la red neuronal son:

- Conv2: Para el procesamiento de la capa homónima.
- Mean: Para promedios como el que se necesita en la capa *GlobalAveragePooling2D*
- FullyConnected: Para las capas densas, entre otros.
- SoftMax: Como función de activación.

Como alternativa más cómoda, pero a costa de un mayor uso de memoria, del cual no podemos abusar dada la naturaleza de nuestro dispositivo; es plausible usar *tflite::AllOpsResolver*, que cargará todas las operaciones disponibles para *TFLite*.

A.5 Cambiar la orientación de la placa

Para poder hacer uso del *SmartPen* en una posición de escritura natural, vertical, hay que hacer una serie de cambios. Los mejores resultados se han conseguido cambiando en la biblioteca de los sensores, en *LSM9DS1.cpp* en todas los sensores:

Fragmento de *LSM9DS1.cpp* de la librería homónima

```

121 // Original          // Orientacion cambiada
122 x = data[0] * 4.0 / 32768.0; // z = -data[0] * 4.0 / 32768.0;
123 y = data[1] * 4.0 / 32768.0; // y = data[1] * 4.0 / 32768.0;
124 z = data[2] * 4.0 / 32768.0; // x = data[2] * 4.0 / 32768.0;
```

Aunque continúan sin obtenerse trazados correctos del movimiento.

Red neuronal

B.1 Ajuste para poder utilizar el recolector de muestras de *Pete Warden*

Es necesario hacer uso del recolector de muestras, utilizar *Google Chrome* y activar la flag para desarrolladores '*Experimental Web Platform features*' activa en '*chrome://flags/*' [44].

B.2 Asignación incorrecta de índices en el recolector de muestras

Como se ha citado en el *apéndice B.3* anterior, al borrar muestras se borran más trazos de los debidos y eso genera una redistribución de los índices de cada muestra, resultando incorrectos.

Para resolverlo, se ha creado un script que reasigna los índices con la distribución adecuada: [sort_dataset.py](#)

B.3 El recolector de muestras elimina varias muestras al borrar una

Este fallo aparentemente se produce al borrar instancias por debajo de la última. En ocasiones, se produce un pequeño error en las comprobaciones de índices de muestras, por lo que se borran varias muestras y estas terminan con índices desordenados.

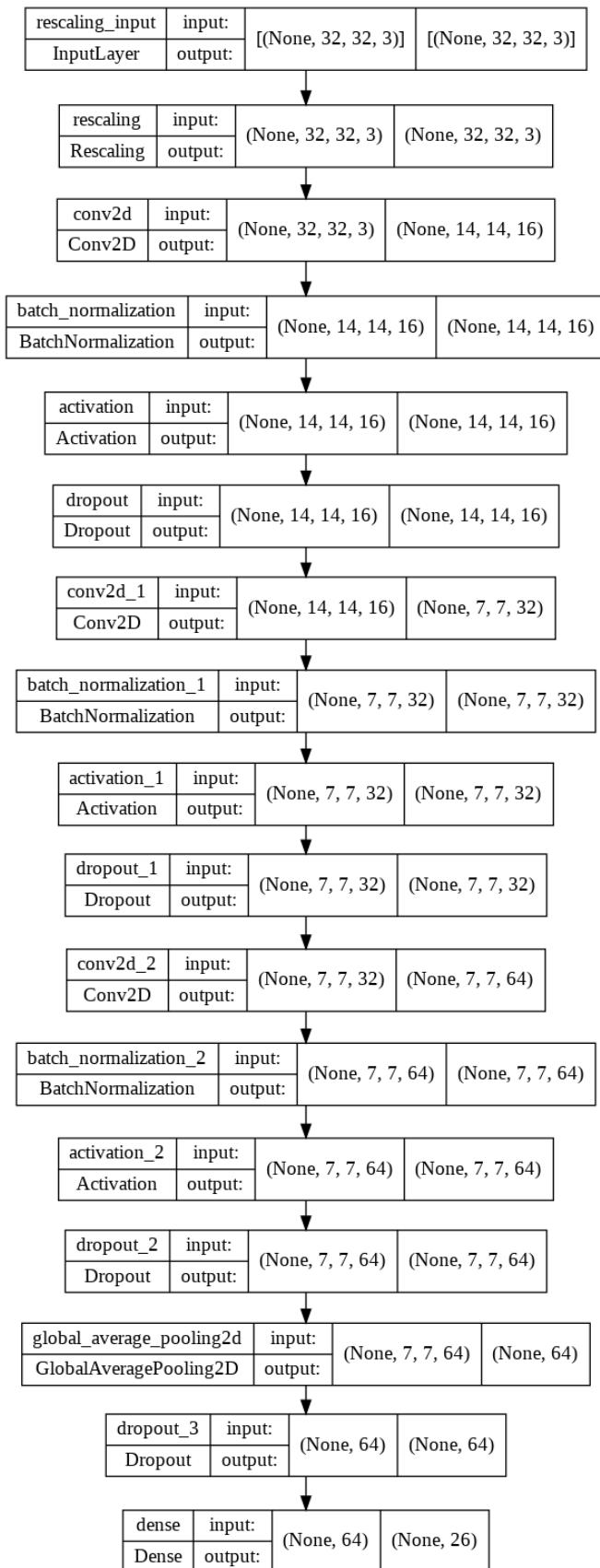
Para paliar este problema y comprobar si había el número de trazos esperado, se hacía uso de un comando en bash cada vez que se generara un *JSON*:

```
1 ~$ grep -o index <nombre_archivo>.json | wc -l
```

B.4 Descripción de capas Keras empleadas para la implementación de la red neuronal

Las capas utilizadas durante la implementación del modelo, son las siguientes [19]:

- Rescaling: Como resultado de esta capa, se consigue reescalar o normalizar los valores de los inputs, en nuestro caso las imágenes. Con esto definimos una escala uniforme y es un procedimiento típico en procesamiento de imagen: *image normalization*, donde se normalizará el valor de los píxeles de las imágenes a una escala [0,1].
- Conv2D: Capa para procesamiento convolucional, es la capa que realmente procesará dentro del modelo. Se crea un kernel que convoluciona con la entrada de la capa. Se define la capa en base ciertos parámetros que determinarán su funcionamiento y complejidad, siendo en nuestro caso los parámetros:
 - filters: el número de filtros de salida de la convolución.
 - kernel_size: el tamaño de ventana de convolución, cuando se especifica un solo número, se interpreta una ventana cuadrada de ese tamaño.
 - strides: define el tamaño de los tramos de salto de la ventana de convolución.
 - input_shape: establecer el tamaño de la entrada. En nuestro caso imágenes de 32x32.
- BatchNormalization: Normaliza sus entradas por lotes. Aplica transformaciones que conservan la media de la salida cercana a 0 y la desviación estándar a 1.
- Activation: aplica una función de activación a una salida. Las funciones de activación son las que arbitran la activación de las neuronas de la red neuronal y por ello, repercute en su salida. Existen diversas funciones de activación, como lo son *rectified linear unit(relu)*, *sigmoid*, *softmax*(función de distribución de probabilidad), etc.
- Dropout: Capa que introduce cierta entropía, de forma que se descarta la contribución de ciertas neuronas de forma estadística. Se suelen implementar para soslayar el *overfitting*.
- GlobalAveragePooling2D: Esta capa tomará un tensor de dimensión $x*y*z$ y calculando el valor medio de los valores x e y , producirá una salida basada en z elementos.
- Dense: Es una capa común de red neuronal, solo que está *densamente* conectada, es decir, cada neurona de esta capa está conectada a todas las neuronas de la capa anterior. Se usa, como es nuestro caso, en redes clasificadoras.

Figura B.1: Estructura del modelo generado en *TensorFlow*

La anterior *figura B.1* es consecuencia del siguiente código:

Fragmento de *Train.ipynb*

```

1 ##### MAKING THE MODEL #####
2
3 def make_model(input_shape, num_classes):
4     model = models.Sequential()
5
6     # Rescaling
7     model.add( layers.Rescaling(1.0 / 255) )
8     # Block 1
9     model.add( layers.Conv2D(16, 5, strides=2, input_shape=input_shape) )
10    model.add( layers.BatchNormalization() )
11    model.add( layers.Activation("relu") )
12    model.add( layers.Dropout(0.45) )
13    # Block 2
14    model.add( layers.Conv2D(32, 5, strides=2, padding="same") )
15    model.add( layers.BatchNormalization() )
16    model.add( layers.Activation("relu") )
17    model.add( layers.Dropout(0.45) )
18    # Block 3
19    model.add( layers.Conv2D(64, 3, strides=1, padding="same") )
20    model.add( layers.BatchNormalization() )
21    model.add( layers.Activation("relu") )
22    model.add( layers.Dropout(0.45) )
23    # Pooling + another Dropout
24    model.add( layers.GlobalAveragePooling2D() )
25    model.add( layers.Dropout(0.45) )
26    # Softmax
27    model.add( layers.Dense(num_classes, activation="softmax") )
28
29
30    return model
31
32 model = make_model(input_shape=(IMAGE_WIDTH, IMAGE_HEIGHT, 3),
33                     num_classes=NUM_CLASSES)
34 model.build(input_shape=(None, IMAGE_WIDTH, IMAGE_HEIGHT, 3))
35 model.summary()
36 keras.utils.plot_model(model, show_shapes=True)

```

B.5 Experimentación red neuronal

B.5.1 Estructura de la red neuronal

Cada uno de los modelos o cambios con los que se ha experimentado, son pequeñas iteraciones o pequeñas variaciones respecto del modelo base. Creado a partir del estudio de otros modelos para tareas parecidas.

Una de las abundantes variables con las que experimentar, es el tamaño de las imágenes de entrada para la red neuronal. Se ha percibido una notoria mejora en el reconocimiento con letras complejas, como por ejemplo la 'k', a medida

que las dimensiones aumentan. De forma análoga, con letras simples, como por ejemplo la 'c', a razón de una menor resolución, mejores resultados se obtenían. Esto cuadra con lo esperable, ya que las letras complejas necesitan de un análisis más preciso para obtener mejores predicciones que el resto y equivalentemente, las letras más sencillas obtienen mejores predicciones sobre el resto, cuando se analizan muestras sencillas.

También se ha estudiado el comportamiento del tamaño del *kernel* en las capas *Conv2D*. Se ha experimentado para tamaños (cuadrados) de entre 3 a 5, ya que menos resultaría poco conveniente y más sería desproporcionado teniendo en cuenta que trabajamos con información espacial igual o inferior a 32x32. En general se ha observado que a mayor tamaño de kernel, se obtiene una mayor precisión en el testeo y si siente al probar el modelo, que funciona ligeramente mejor. A costa evidentemente de una mayor complejidad del modelo y por tanto de un mayor tamaño en memoria. Por lo que se ha optado por mantener en el modelo final tamaño de *kernel* 5x5 en los dos primeros bloques de *convolución* y de 3x3 en el último, ya que a esta capa llega información de dimensiones mucho menores (7x7). Sin embargo, la decisión de aumentar la complejidad del modelo, también implica aumentar el número de *epochs* en la etapa de entrenamiento, para desentrañar mejor el comportamiento más adecuado de una estructura más compleja, aunque esto se tratará en la siguiente sección.

Otro parámetro que se ha analizado es el número de bloques de convolución (explicados en la Sección 6.1.2), con solo dos, los resultados no son nada buenos, y añadiendo un cuarto, se obtienen resultados muy buenos pero concretamente en algunas letras hay clasificaciones erróneas, posiblemente debido a *overfitting* (concepto introducido en Teoría 6.2.2) y se obtienen buenos resultados en general, pero no suficientes para una buena experiencia.

Se experimentó con la capa *BatchNormalization* de los bloques de convolución. Al eliminar esta capa de los bloques, se obtienen en fase de testeo, clasificaciones correctas, pero con una precisión mucho menor, llegando a alcanzar en algunos casos, la mitad de precisión que con el modelo base. Como cabía esperar, la normalización de las entradas resulta muy útil para mejorar la precisión, y aunque no se aprecie en el testeo, también proporciona un entrenamiento más rápido.

Con el resto de capas del bloque de convolución, no se ha experimentado, porque son capas manifiestamente indispensables.

Aunque si bien, sí que se ha experimentado con los valores de *Dropout*, estudiándolos para valores de 0,4 a 0,6, ya que para el resto de valores, resulta un *dropout* demasiado restrictivo o demasiado laxo. Obteniendo los mejores resultados para 0.45, aunque si bien, las diferencias eran prácticamente indistinguibles durante su uso.

B.5.2 Entrenamiento de la red neuronal

En esta fase, los cambios siguen siendo muy relevantes para el desempeño de la red neuronal, sin embargo hay menos parámetros que estudiar.

Yo me he centrado principalmente en las *epochs* y el *learning rate*, aunque también he probado someramente los *optimizadores*. Las *epochs* son las iteraciones que se darán en el algoritmo de entrenamiento. A más *epochs*, mayor será la profundidad en el entrenamiento. Esto no quiere decir que si empleamos el doble de *epochs*, consigamos el doble de precisión, no es un fenómeno lineal, sino que el modelo presenta un techo alcanzable y por más que se entrene, no van a obtenerse mejores resultados. Por lo que la clave es encontrar un número de *epochs* que resulte en un buen entrenamiento, pero sin emplear tiempo de más. En nuestro caso, para los modelos experimentales más complejos, se ha empleado un número mayor de *epochs*, pero en general, ningún modelo ha seguido presentando mejoras notorias más allá de las 100 *epochs*. Por tanto, se usará este valor.

El segundo parámetro de estudio, ha sido el *learning rate*, que simplificándolo, sería el valor que determina el progreso de entrenamiento que se hace en cada *epoch*. Un *learning rate* más alto, se traduce en un entrenamiento que escala mejor con pocas *epochs*, pero que está más limitado a alcanzar el resultado óptimo del modelo; análogamente, un *learning rate* más bajo, resulta en un entrenamiento que escalará peor, pero que puede alcanzar un resultado potencial mejor, como ilustra la Figura B.2. En nuestro caso se ha experimentado con valores de *learning rate* entre 0,00095 y 0,0015 ya que en muchos modelos estudiados y parecidos a este, se usa un valor cercano a 0,001. y como nuestro modelo no es muy complejo, podemos permitirnos un escalado menor, ya que no llevará mucho tiempo entrenar al modelo y hay margen para aumentarlo sin problema. Por lo que nos quedamos con un *learning rate* de 0,00095.

Respecto a los *optimizadores*, se ha experimentado muy laxamente, así que para asegurar, nos quedamos con el *optimizador* de referencia del modelo base.

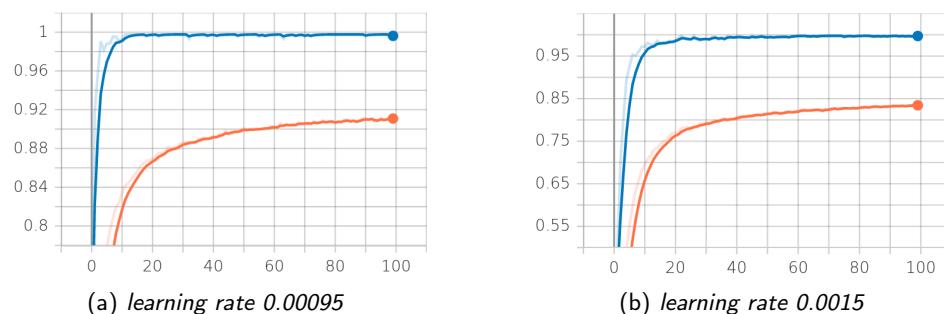


Figura B.2: Escalado de *accuracy* por *epochs*

Azul: Validación | **Naranja:** Entrenamiento. Mismo modelo, distinto *learning rate*.

Interfaz de usuario

C.1 Permisos uso del puerto del microcontrolador (Linux)

En algunos casos, como ocurrió en el trámite de este trabajo, si queremos acceder al puerto serie de la placa o usar el IDE de Arduino, debemos conceder permisos al microcontrolador:

```
1 ~$ ls -l /dev/ttyACM* # En mi caso
2 ~$ sudo usermod -a -G dialout <usuario>
```

Como respuesta al error “*can't open device /dev/ttyACM0*”.

C.2 Acceso al puerto serie en QT (Linux)

Hay que añadir al *.pro del proyecto QT:

```
1 QT += serialport
```

Y tras esto, añadir la librería con normalidad y acceder al puerto con el nombre, en nuestro caso, “*ttyACM0*”.

C.3 Valores nulos al leer por *Bluetooth* con *QT*

Este fue un error complejo de depurar como suele ser habitual con los problemas derivados de versiones de librerías. No se leía ningún valor de las *características* del dispositivo pese a estar detectándolas y estar correctamente conectado, se debía a dos factores.

El primero estar haciendo uso de una librería anterior a la documentación con la que estaba trabajando (librería de QLowEnergyService). Hay grandes diferencias en el comportamiento de algunos métodos de esta librería de las versiones

5.x a la 6.x, aunque estas no provocan errores, sí que provocan que el código no funcione como se esperaría (*Enums* con valores diferentes o inexistentes, etc).

En segundo lugar se estaba llamando a un método cuando todavía no se había recibido la *característica*. Por tanto esta aparentaba estar bien registrada, ya que podía obtenerse su *Uuid*, pero no contenía ningún valor. Se estaba leyendo el valor en *connectToService()*, cuando debería hacerse en *serviceDetailsDiscovered()*.

C.4 Error de reconocimiento de imágenes en QT

Al importar imágenes en QT tras haber exportado el programa a Win o Linux para corroborar que todo continuara funcionando correctamente, las imágenes dejaron de mostrarse. Esto se debe a que QT cambia la ruta del proyecto al directorio en el que se exporta el programa. Por tanto las rutas especificadas para las imágenes, dejan de ser válidas. Para corregir este problema, basta con cambiar el '*Build directory*' del proyecto (Desde '*Projects*' en el panel de la derecha de QT creator).

C.5 Pérdida de valores de *características*

Este es teóricamente el fix para las desconexiones aleatorias, sin embargo en mi caso y en el de otros desarrolladores, no termina de funcionar, si bien sí que permite no perder las características durante las desconexiones. [9]

Este parche hace que el código del callback de la librería, no asuma que hay un periférico conectado, que es lo que presuntamente, provoca la desconexión.

Debe editarse la librería *QLowEnergyControllerPrivateBluezDBus* con los cambios resaltados en el siguiente enlace: https://codereview.qt-project.org/c/qt/qtconnectivity/+/233087/3/src/bluetooth/qlowenergycontroller_bluezdbus.cpp#563.

Pese a que no haya funcionado en mi caso, sí que ha sido muy importante ya que ha corregido en parte un mal comportamiento de la librería. Estas son las causas por las que es mucho mejor elegir herramientas con una comunidad amplia y activa detrás.

C.6 Error de permisos al lanzar la interfaz de usuario (Linux)

Al ejecutar el *UI*, aparece este error en los logs: '*qt.bluetooth.bluez: Missing CAP_NET_ADMIN permission. Cannot determine whether a found address is of random or public type.*' Para ello hay que modificar la configuración del *dbus* y dotar de permisos a nuestro usuario. [36]

```
1 <policy user="yourUserName">
2   <allow own="org.bluez"/>
3   <allow send_destination="org.bluez"/>
4   <allow send_interface="org.bluez.Agent1"/>
5   <allow send_interface="org.bluez.MediaEndpoint1"/>
6   <allow send_interface="org.bluez.MediaPlayer1"/>
7   <allow send_interface="org.bluez.Profile1"/>
8   <allow send_interface="org.bluez.GattCharacteristic1"/>
9   <allow send_interface="org.bluez.GattDescriptor1"/>
10  <allow send_interface="org.bluez.LEAdvertisement1"/>
11  <allow send_interface="org.freedesktop.DBus.ObjectManager"/>
12  <allow send_interface="org.freedesktop.DBus.Properties"/>
13 </policy>
```

Si persiste el problema, otra forma de solucionarlo es mediante el comando:

```
1 sudo setcap CAP_NET_ADMIN=eip <path hasta el ejecutable>
```

C.7 Capturas de la interfaz de usuario

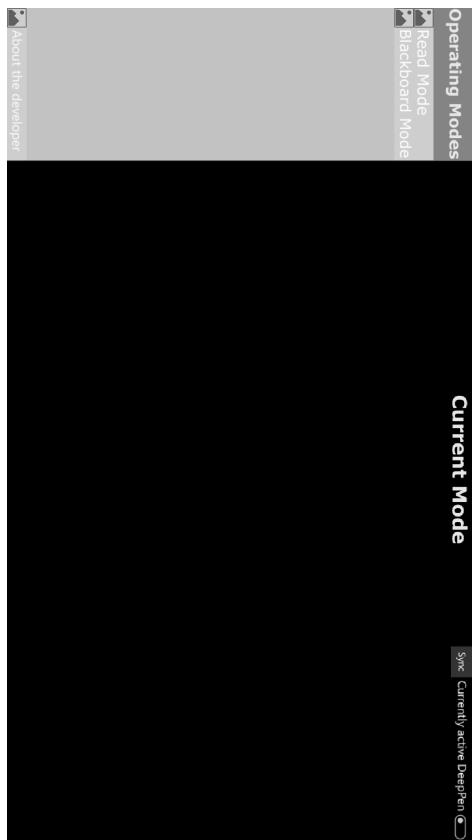


Figura C.1: Boceto en QT Design Studio

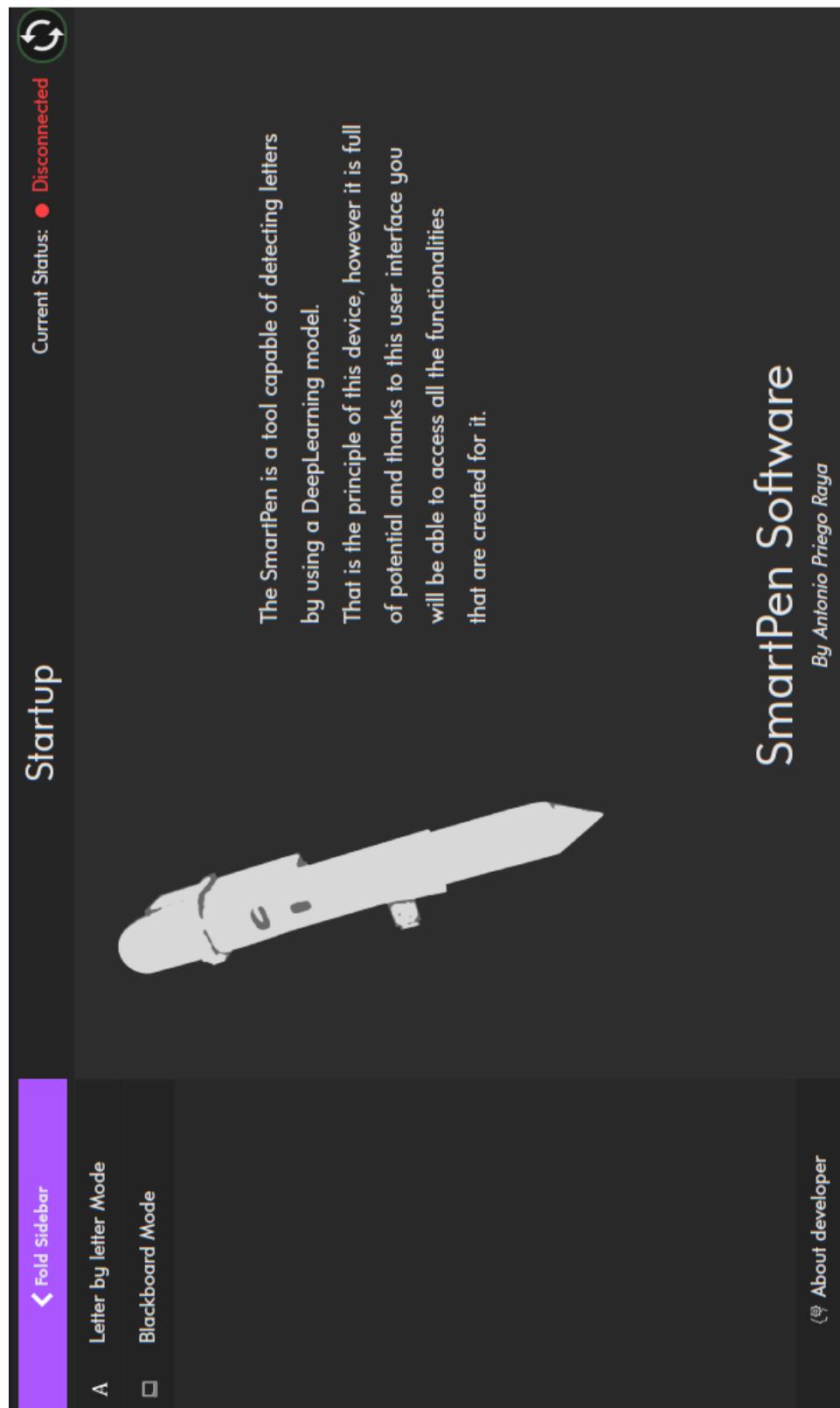


Figura C.2: Resultado de la implementación de la interfaz de usuario

Encapsulado

D.1 Modelos 3D

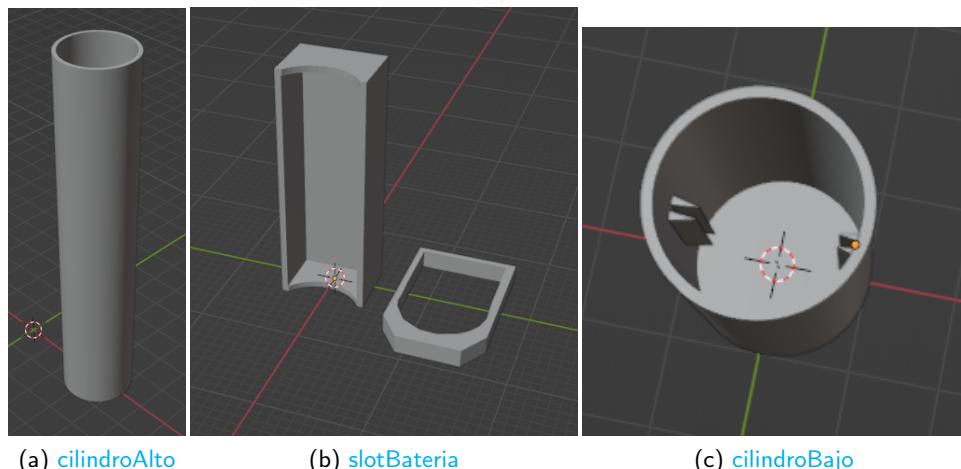


Figura D.1: Componentes útiles del encapsulado

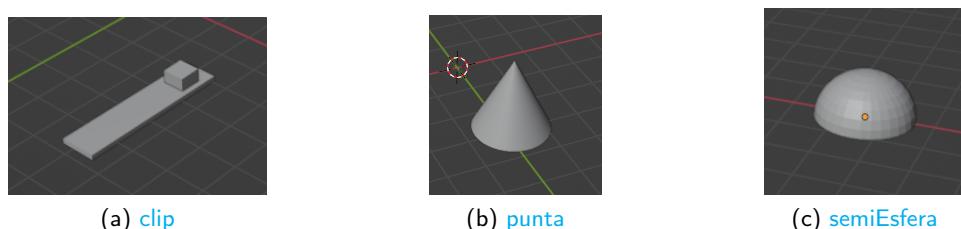


Figura D.2: Decoración del encapsulado

D.2 Resultado de integrar todo en el encapsulado



(a) SmartPen vista frontal

(b) SmartPen vista trasera

Figura D.3: SmartPen



Github del proyecto
<https://github.com/AntonioPriego/SmartPen>



