2025年臺灣國際科學展覽會 優勝作品專輯

作品編號 190027

參展科別 電腦科學與資訊工程

作品名稱 自監督學習在臺灣手語辨識上之應用研究

得獎獎項 三等獎

巴西科學博覽會 MOSTRATEC

就讀學校 臺中市立臺中第一高級中學

指導教師 柳佩君

陳志峰

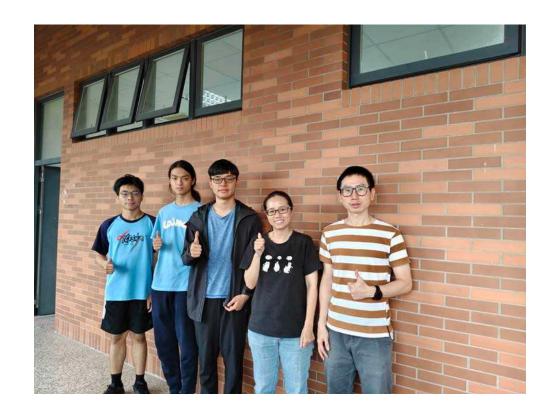
作者姓名 黄政堯

蔡允成

陳宥家

關鍵詞 機器學習、自監督學習、臺灣手語

作者簡介



我是黃政堯,我從以前就特別喜歡人工智慧領域,也正好這次科展,讓我發揮自己的專長做研究。

我是來自台中一中的蔡允成,我們做了一個手語翻譯的研究,手語是一門非常特別的語言,希望我們的研究能在這無聲的世界中做出貢獻。

我是來自台中一中的陳宥家,我的興趣是寫程式、打桌球,很榮幸能來參加 這次的國際科展,希望能藉此學習到更多資訊知識與認識新朋友。

2025 年臺灣國際科學展覽會 研究報告

區別:

科別:電腦科學與資訊工程

作品名稱:自監督學習在臺灣手語辨識上之應用研究

關鍵詞:機器學習、自監督學習、臺灣手語

編號:

摘要

在臺灣手語辨識,先前研究所使用的監督式學習需要大量標記樣本而限制可辨識詞彙量。為此,本研究借鑒自然語言處理領域中 BERT 的遮罩想法,將未標記手語影片隨機遮蓋部分幀數,並讓模型學習預測被遮蓋的幀數以學習臺灣手語的特徵,並透過遷移學習來訓練辨識模型,此作法可克服現有臺灣手語資料缺少的問題。經過實驗,本研究訓練之詞彙辨識模型達成了 242 個詞彙量,94.8%的準確率。

此外,先前研究皆未在手語句子翻譯上有成果。因此本研究基於預訓練模型,整合設計手語翻譯的系統,實驗中,系統在 100 個句子的翻譯表現達到 88%的準,且 BLEU-4 分數取得 20.98,證明自監督學習的方式在手語辨識、翻譯上是有效的。並展現出樣本需求少與辨識詞彙量可輕易擴大的潛力。

Abstract

Previous research of Taiwanese Sign Language (TSL) recognition used supervised learning as their model training method, which required a large number of labeled data, limiting the recognizable vocabulary size. To fix this issue, we took inspiration from the masking concept in a language representation model called BERT. Our idea is to randomly mask certain number of frames in unlabeled TSL videos, allowing the model to learn the features of TSL by predicting the masked frames. Transfer learning is then applied to train the TSL recognition model. The results showed that the TSL recognition model had achieved a recognizable vocabulary size of 242 words with an accuracy of 94.8%.

Moreover, there had been no research about TSL sentence translation. To address that, we designed a TSL translation system based on the TSL recognition model. The system achieved an 88% accuracy in translation for 100 sentences, with a BLEU-4 score of 20.98. This research proved that the self-supervised learning approach is effective in both TSL recognition and translation. With this method, the model requires fewer samples to train, also making the recognizable vocabulary easier to expand.

膏、前言

一、研究動機

在 2020 到 2022 新冠疫情肆虐期間,衛福部每日會定時召開防疫指揮中心記者會,直播說明有關新冠疫情相關之案例報告、政策。研究者時常關注防疫直播,也時常注意到螢幕右下方有一位手語翻譯員,為聾人翻譯與會人士的發言。研究者十分好奇手語翻譯員所比出的手語意義,想要尋求軟體翻譯,但是在搜尋網路之後,發現市面上並沒有臺灣手語的翻譯軟體,因此本研究希望可以自己研究並開發臺灣手語翻譯的系統。

在此以前,不管是國內外,各種的手語都有做過類似的嘗試。國內亦有許多臺灣手語 翻譯的相關研究,研究的主題主要是針對靜態手語詞彙的辨識,以及日常手語詞彙的辨識, 但其所提出的方式均只能辨認少數手語詞彙,且不能翻譯手語句子。

目前研究普遍採用的是監督式學習來進行手語詞彙辨識,這種方式需要準備大量標記 樣本來訓練模型。然而,準備這些樣本的過程非常耗時,且現有資料量十分稀少,這限制模 型能夠識別的手語詞彙數量。而且,現在研究普遍都僅是單詞的辨識,對於在日常情況的手 語句子翻譯未有成果。因此本研究希望能夠解決目前上述臺灣手語辨識研究上所遇到的問題。

二、研究目的

本研究希望透過自監督學習的方式訓練模型,使模型能夠自行學習到手語的特徵。此方案大幅降低所需要的標記樣本的資料,以擴大模型可辨識的手語詞彙。並自行設計手語翻譯系統,達到句子的翻譯。

- (一)探討自監督學習應用在動態手語辨識任務上。
- (二)探討如何用少數的標記資料完成手語辨識任務。
- (三)研究如何將預訓練模型應用於手語翻譯。

貳、文獻回顧

一、遷移學習

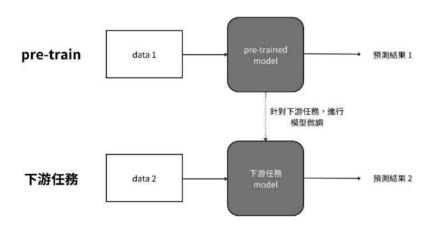


圖 2-1: Fine Tuning 示意圖(來源自行製作)

遷移學習是一種機器學習方法,用於將在一個領域學到的知識應用於相關領域,特別 適用於標註數據稀缺或成本高的情況。其核心思想是利用在大型數據集上預訓練的模型,並 將其應用於特定目標任務,從而節省訓練時間並提升性能。

微調(Fine Tuning)是遷移學習的關鍵技術之一,指的是在預訓練模型的基礎上,對部分層進行進一步訓練,使模型更適應新的任務。典型的預訓練模型如 VGG、ResNet 或 BERT,已在大型數據集上訓練,具備良好的泛化能力。微調可以有效利用這些模型的知識,特別在目標數據集較小的情況下,大幅提升性能,並節省計算資源和時間。

二、Transformer 模型

Transformer (Vaswani et al., 2017) 一開始用於自然語言處理任務,如機器翻譯。它的設計摒棄傳統的 RNN 和 LSTM 等序列模型,採用全新的 self-attention 機制,使其在處理長序列時表現更好。

(一) 自注意力機制 (Self-Attention)

Transformer 使用自注意力機制來捕捉輸入序列中不同位置的依賴關係。 每個輸入位置都與其他所有位置建立關聯,這允許模型在處理不同距離的依賴關係時保持高效率。由於 Transformer 沒有明確處理輸入序列的順序訊息,需要添加位置編碼來幫助模型理解單字的 相對位置。

(二)編碼器-解碼器結構 (Encoder-Decoder)

Transformer 由編碼器和解碼器組成,適用於序列到序列的任務,在機器翻譯。編碼器 負責將輸入序列轉換為特徵,解碼器則將這個表示轉換為輸出序列。Transformer 模型的出 現在 NLP 領域引起革命性的變化,它不僅在翻譯任務上取得令人矚目的成果,也成為許多 其他 NLP 工作的基礎模型,如 Bert、GPT 系列等等。

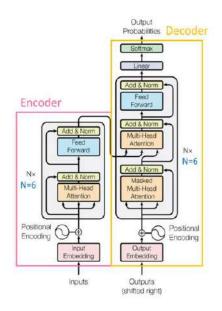


圖 2-2: Transformer 模型架構圖(來源: Vaswani, 2017 [1])

三 · Bidirectional Encoder Representations from Transformers (BERT)

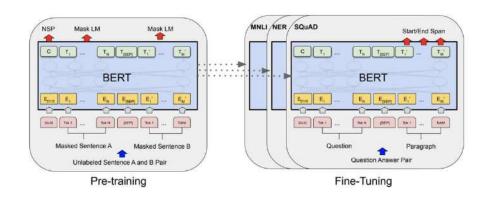


圖 2-3: BERT 示意圖(來源: Jacob Devlin, 2018 [2])

BERT (Bidirectional Encoder Representations from Transformers) 是由 Google 在 2018 年提出的預訓練語言模型,它在 NLP 領域取得巨大成功。 與傳統 NLP 方式不同,BERT 的獨特

之處之一是它採用自監督學習的方法進行預訓練。 自監督學習是一種無監督學習的形式,其中模型從輸入資料中學習,而無需標籤。

(一)預訓練與微調

BERT 首先在大規模文字語料上進行預訓練,學習通用的語言表示。然後,可以透過微調在特定任務上,例如文字分類、命名實體識別等,以適應特定的應用場景。

(二) Transformer Encoder

BERT 模型通常由 Transformer Encoder 組成,每個編碼器層都有多頭自註意力機制和 前饋神經網路。這些層允許模型學習不同層次的語言表示。

(三)掩碼語言模型 (Masked Language Model, MLM) 訓練任務

BERT 在預訓練階段使用一個掩碼語言模型任務,其中一些輸入詞被隨機遮蓋,模型需要預測這些掩蓋詞的標籤。 這鼓勵模型學習更豐富的上下文表示。

自監督學習任務使得 BERT 能夠捕捉大量的語言知識,並且預訓練階段的學到的參數可以在各種 NLP 任務上進行微調,從而獲得更好的性能。自監督學習的想法是透過模型本身產生標籤,因此無需手動標註大量標籤資料。 這種方法使得模型能夠從大規模的未標記資料中學到有用的特徵,然後在特定任務上進行微調。

五、Masked autoencoders (MAE)

Masked Autoencoders(MAE)是一種深度學習模型,主要用於無監督學習,特別是在處理圖像數據方面。這種方法靈感來自於自然語言處理領域的成功技術,例如 BERT。MAE 的核心思想是在輸入數據中隨機遮蓋(mask)一部分內容,然後訓練模型重建被遮蓋的部分。MAE 通常由兩部分組成:一個編碼器(encoder)和一個解碼器(decoder)。

(—) Encoder&Decoder

Encoder 的作用是處理輸入數據,但在此之前,會先隨機選擇並遮蓋數據的一部分。例如,在處理圖像時,會隨機遮蓋圖像的一些像素或區域。編碼器只對未被遮蓋的數據進行處理,從而學習到數據的內在特徵和結構。Decoder 的目標是根據編碼器處理過的數據來重建原始數據的遮蓋部分。這個過程迫使模型學習數據的重要特徵,因為它需要理解和推斷遮蓋部分的內容。

在 MAE 的訓練過程中,首先會隨機選擇並遮蓋輸入數據的一部分,然後編碼器對剩餘的未被遮蓋數據進行處理,提取特徵。接著,解碼器嘗試重建被遮蓋的部分。這一過程涉及到損失函數的計算,用以衡量重建數據與原始數據之間的差異。最後,通過反向傳播和參數更新,模型逐漸學會如何準確重建遮蓋的數據。

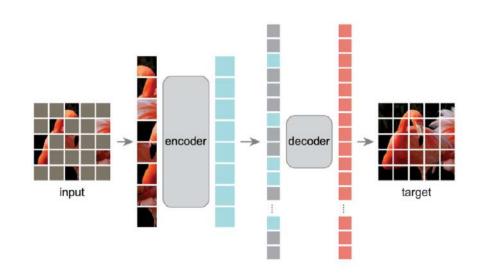


圖 2-5: Masked Autoencoder 架構圖(來源:Kaiming He, 2021[4])

參、研究設備器材

一、設備:

	系統	GPU	СРИ
筆記型電腦	Windows 11	RTX 3060 latop	i7 - 12400H
Google Colab	Ubuntu 22.04.3	Nvidia A100	Intel (R) Xeon

二、軟體:

軟體/套件	python	cuda	pytorch	opency- python	mediapipe	GPT
版本	3.9.4	11.8	2.0	4.7.0	0.10.9	4.0

肆、研究過程與方法

受到前述 Masked Autoencoders 與 Vision Transformer 工作的啟發,本研究認為 Transformer 模型具有處理手語影片序列的潛力。因此,本研究借鑑 Vision Transformer 處理圖片的方法來處理手語序列片段,同時利用 Masked Autoencoders 的自監督訓練方法,使整個模型能夠有效地學習手語的特徵,以應用在詞彙辨識任務中。

為了達到手語翻譯,本研究將預訓練模型進行 Fine tune,得到手語詞彙辨識模型。接下來,將手語影片分段辨識出詞彙,之後應用本研究者設計之滑動窗口演算法得出句子中所含詞彙,再將使用大型語言模型重組文句。

一、研究及實驗架構流程圖

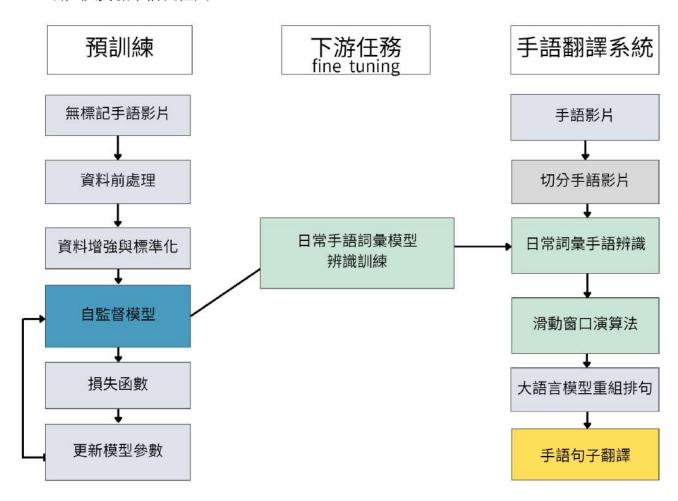


圖 4-1: 研究流程架構圖(來源自行製作)

二、訓練資料的收集與處理

(一)資料來源

本研究發現疫情指揮中心定期召開的防疫記者會 (中央流行疫情指揮中心嚴重特殊傳染性肺炎記者會),旁邊配有手語翻譯人員,提供豐富的資料來源。因此,本研究從防疫指揮中心平台上下載 481 部影片。

(二) MediaPipe 進行影像前處理

本研究將所有影片以 64 幀為一單位進行切分成片段子集 $segment\ subset$,並且將畫面裁切成合適的大小,其中所有片段 v_i 為總長 64 幀,畫面大小為 640x640 的影片。在經過 python 套件 mediapipe 處理,標出翻譯人員手部的關節點。

在上一篇研究當中,研究者發現,如果輸入預訓練模型的資料中僅包含手部的點座標(如圖 4-2),那麼模型將會失去身體骨架的資訊,即無法區分手部在在身體的相對位置,然而,如果訓練資料包含了身體關節座標,則會失去對於手型的精細度,因此在新的研究當中,研究者決定同時訓練兩個模型,一個是訓練資料僅包含手部點座標的手型預訓練模型(Hand-shape Pretrained Encoder),另一個則是骨架預訓練模型(Body Pretrained Encoder),包含手部點座標以及身體骨架。

綜上所述,本研究定義兩種訓練資料,第一種是每筆資料 P_m 是儲存 64 幀中,每幀雙 手 40 個點(一隻手各 20 點)以及關節點 4 點,臉部 1 點的(x,y)座標,共 45 個點。另一種 P_m 則是僅包含每幀雙手 40 個點(一隻手各 20 點)的的座標,共 40 個點的點座標(如圖 4-2)。(在此論文之中的研究過程以及研究結果與討論,皆以優先展示包含 45 個點的 P_m 為 範例)

$$segment\ subset\ =\ \left\{v_1,v_2...v_n\right\}\ , v_i\in R_{64\times 640\times 640\times 3}$$

$$dataset\ =\ \left\{P_{m1},P_{m2}\dots P_{mn}\right\}\ ,\ P_{mi}\in R_{64\times 45\times 2}$$

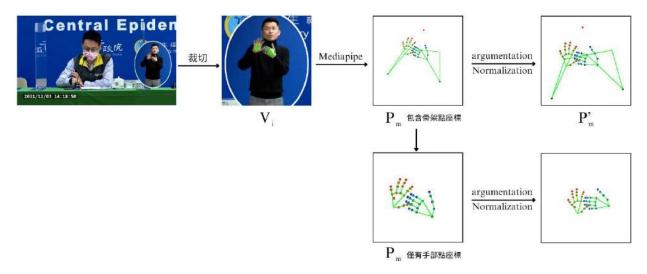


圖 4-2: 資料前處理示意圖(來源取自防疫指揮中心)

三、預訓練過程

(一)資料增強與標準化

每一筆進入模型的資料 P_m ,在進入模型之前須經過縮放(Scaling)、旋轉(Rotation)、平移(Translate)、標準化(Normalization),以進行增強。本研究視 P_m 為 $\left\{P_1,P_2,...,P_{64*45}\right\}$ (其中 P_i 為 P_m 中的所有點)以進行資料增強,具體操作公式如下:

$$\arg u \ mentation(P) = R(\theta) \cdot \left(s \cdot (P-C)\right) + C + v, -20^{\circ} \le \theta \le 20^{\circ}, 0.7 \le s \le 1.5$$

其中,C = (x,y)為錨點, $R(\theta)$ 表示旋轉矩陣,s是隨機的縮放因子,v是隨機的平移向量。這個過程對每一個點 P_i 都進行同樣的變換。

為了消除特徵間的尺度差異以及提高收斂速度, P_m 在經過資料增強後,還需要進行標準化 (Normalization),公式如下:

$$P_{m}^{'} = \frac{P - \mu}{\sigma}, \mu = 335.49, \sigma = 134.28$$

在經過資料增強和標準化後,本研究得到轉換過後的資料 P'_m 。

(二)模型架構

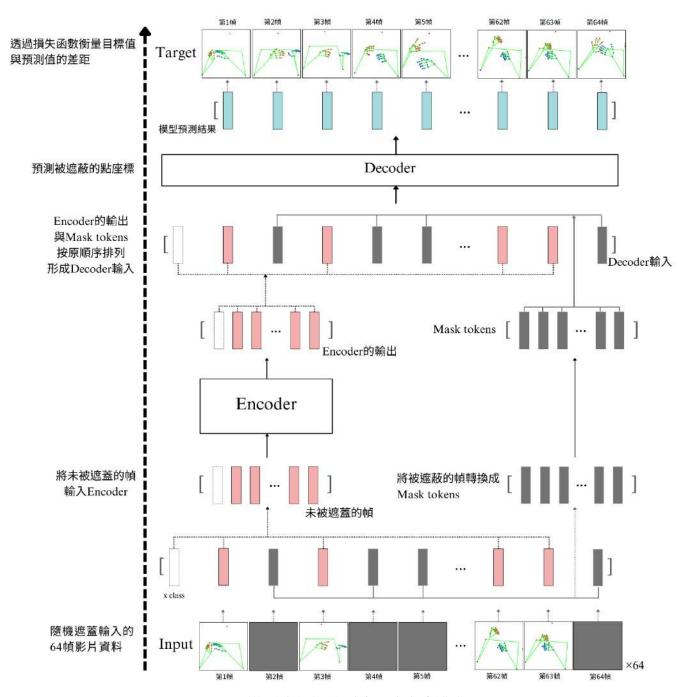


圖 4-3:模型流程概覽(來源自行製作)

本研究參考 MAE 研究的架構,在預訓練期間,將輸入的資料 P_m 中的幀依遮蓋率(mask ratio)隨機遮蓋,其中未被遮蓋的幀輸入進 Encoder,被遮蓋的幀轉換成 Mask tokens,隨後將 Encoder 的輸出與 Mask tokens 按原順序排列形成 Decoder 的輸入資料,將之輸入Decoder 以預測被遮蓋的點座標,以此重建原始手語幀的座標形成預測值。最後透過損失函

數衡量目標值與預測值的差距,更新模型參數以最小化損失,使模型能夠學習到手語的特徵。在預訓練之後,Decoder 被移除,而未被隨機遮蓋的手語序列幀則輸入進 Encoder 以進行識別任務。

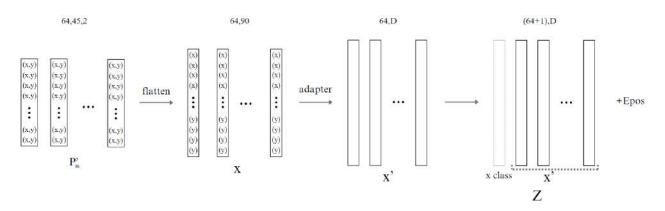


圖 4-4:步驟一的流程圖(來源自行製作)

如圖 4-4,本研究將資料 $P_m^{'}\in R_{64\times45\times2}$ 轉換成序列 $x\in R_{64\times90}$ 以符合 Encoder 所要求的 shape,而為了使輸入的x投射到設定的 Hidden size D,本研究通過設計一個轉接層 Linear Projection Adaptor 使 $x\in R_{64\times90}$ 變成 $x'\in R_{64\times D}$,公式如下:

$$x^{'}=adaptor(x)=xW+b$$
 , $x\in R_{64 imes 90}$, $W\in R_{90 imes D}$, $b\in R_{64 imes D}$

在 position embedding 的部分上,本研究選擇與 vision transformer 一樣的方式,在序列x' 前連接一個 lernable embedding x_{class} ,再加上 position embedding E_{pos} ,使整個序列保留位置的資訊,公式如下:

$$z = \left[x_{class}; x^{'}\right] + E_{pos}, E_{pos} \in R_{(64+1)\times D}$$

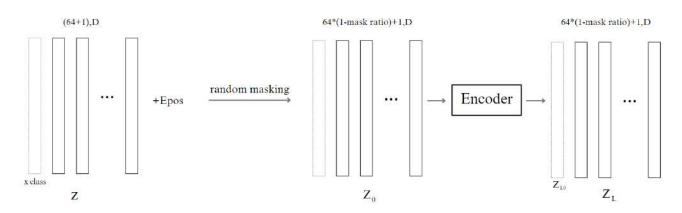


圖 4-5: 步驟二的流程圖(來源自行製作)

如圖 4-5,在經過 position embedding 後,本研究遵循 Masked auotencoder 論文的方式,透過生成一個抽樣表,裡面包含隨機的編號,按照編號依遮蓋率(masked ratio)隨機將序列z中的部分 token 抽離,剩餘的形成新的序列 $z_0 \in R_{1+64*(1-mr)\times D}$ (其中mr為設定的 masked ratio)作為 Encoder 的輸入。

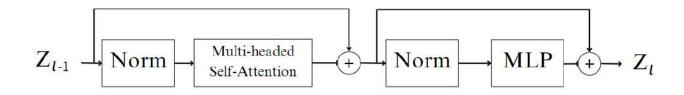


圖 4-6: Pre-Norm Transformer 的架構(來源自行製作)

如圖 4-6 所示,本研究選擇採用改良版的 Transformer encoder,稱為 Pre-Norm Transformer,與原始論文中使用的 Post-Norm Transformer 有所不同,其特性是在訓練的過程中對 learning rate 不那麼敏感,較為穩定。這個 Encoder 結構包括交替排列的 Multiheaded Self-Attention 和 MLP 層(參見公式 $1 \cdot 2$)[3]。在每一層的前面,本研究進行 Layer Normalization (LN),並在每一層的後面加入 Residual Connection。接下來,本研究將被隨機遮蓋的序列 z_0 輸入進 Encoder,並且得到 Encoder 的輸出 $z_L \in R_{1+64\cdot(1-mr)\times D}$ 。 而 x_{class} 在 Encoder 輸出端的狀態 z_{L0} (z_L 的第 1 項)在經過 Layer Norm(LN)得到的預測輸出 $y_{feature} \in R_{1\times D}$ 將作為手語的特徵向量(參見公式 3)。

公式 1
$$z'_l = MSA(LN(z_{l-1})) + z_{l-1}$$
, $l = 1,2...,L$

公式 2
$$z_l = MLP(LN(z'_l)) + z'_l$$
, $l = 1,2...,L$

公式 3
$$y_{feature} = LN(z_{L0})$$

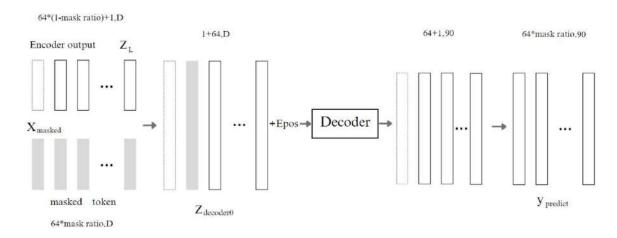


圖 4-7:步驟三的流程圖(來源自行製作)

以 $x_{masked} \in R_{1 \times D}$ 與 Encoder 的輸出 z_L ,按原先的順序排列後組成 $Z_{decoder0} \in R_{(1+64) \times D}$,加上 E_{pos} 以保留位置資訊,再輸入進 Decoder。 $Z_{decoder0}$ 包含完整長度的序列(如圖 4-7 所示),其中 x_{masked} 為共享的向量,表示需要預測的缺失訊息的存在。

本研究的 Decoder 也採用 Pre-Norm Transformer encoder, 其中的層數比起 Encoder 較少。
Decoder 的輸出在經過 prediction layer 後,按照抽樣表編號抽樣,將得到遮蔽部分的預測結果 $y_{predict}$,公式如下。

$$y_{predict} = prediction \ layer(x) = xW \ \ x \in R_{64 \times D} \ , \ W \in R_{D \times 90}$$

其中x為 Decoder 的輸出

(三)損失函數 (Loss Function)

為了衡量模型預測與實際目標之間差異,本研究選擇MPJPE()[6]作為指標,公式如下,其中 \hat{P}_i 為模型預測值的點, P_i 為目標值的點。

MPJPE loss
$$(\widehat{P}, P) = \frac{1}{N} \sum_{i=1}^{N} ||P_i - \widehat{P}_i||_2$$

模型在預訓練中的目標是最小化這個差異。MPJPE()計算每個關節的預測位置與真實位置之間的歐幾里德距離,然後對所有關節的這些距離求平均,得到的結果是模型在二維空間預測關節位置的平均誤差。為了計算模型預測的 Loss,本研究將預測值 $y_{predict} \in R_{64\cdot mr \times 90}$ Reshape 成符合計算 $y_{predict} \in R_{64\cdot mr \times 45 \times 2}$,並將目標值 P'_m 並按抽樣表的編號抽

樣組成 $t \arg e \ t \in R_{64*mr \times 45 \times 2}$,接下來計算 $MPJPE\left(y_{predict}^{'}, t \arg e \ t\right)$,之後 Optimizer 會進行梯度下降,更新模型參數以最小化損失,使模型能夠學習到手語的特徵。

(四)預訓練實驗設置

在預訓練階段,本研究將*dataset*分成訓練集與測試集兩部分,其中訓練集有 78.8 萬筆,測試集有 4.1 萬筆。此外,在遮蔽率 mask ratio 設置上,研究者在上一篇研究中已經證明 50%是最佳的數值,因此本研究在這選用了 50%進行預訓練。

	Encoder	Decoder
Layers	12	4
Hidden size	768	512
MLP size	3027	2048
Heads	12	16
Params	86M	40M

learning rate	1e-6
mask ratio	50%
batch size	128
optimizer	AdamW

表 4-1:模型訓練各項超參數設置

四、下游任務

預訓練過後,本研究將模型的中的 Decoder 移除,模型僅保留 Pretrained Encoder。一樣將單詞手語片段經過 Mediapipe 處理得到 P_m ,並進行標準化後輸入進 Pretrained Encoder ,此時,Encoder 的輸入是完整的 64 幀,不會被隨機遮蓋。

接著,在下游任務進行手語辭彙 fine tuning 階段,本研究實驗了四種不同的模型,以分析比較實驗數據。

(一)手型預訓練模型與骨架預訓練模型

本研究將手型、骨架預訓練模型分別接上 MLP 層與 Softmax 層進行 fine tuning (如圖 4-8),比較辨識準確率。

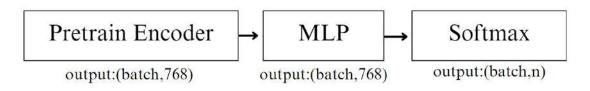


圖 4-8:模型架構圖(來源自行製作)

(二)融合模型

為進一步探索預訓練模型的潛力,本研究認為,骨架預訓練模型(Body Pretrained Encoder)較有能力關注到整體,而手型預訓練模型(Hand-shape Pretrained Encoder)較有能力關注到手型的細節,如果能夠整合此二模型,將會達到更好的辨識準確率。

因此,本研究將 Body Pretrained Encoder 以及 Hand-shape Pretrained Encoder 所輸出的 Feature 取出並且連接在一起組成 connective feature ,然後再接上 GELU Activation function 跟 MLP 層,以及最後的 Softmax 層以用於分類(如圖 4-9)。

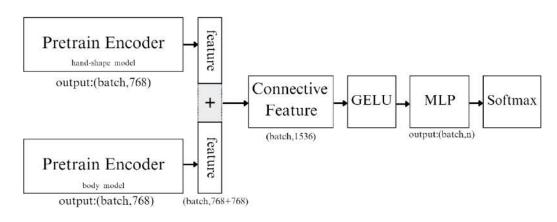


圖 4-9:模型架構圖(來源自行製作)

(三)照組組:在此模型,本研究將不會使用預訓練模型,僅使用隨機初始化的 Encoder, 以此來對比是否使用預訓練模型的準確度差距。

五、日常詞彙手語辨識實驗

在辨識實驗中,本研究評估模型在4個日常詞彙手語辨識的能力,從中正大學手語辭典 所提供的500多個例句中,挑選了100個句子,其中包含了242個日常使用字彙(表4-2), 因此本研究選擇了這242個詞彙進行實驗,實驗流程如下(圖4-10):

1. 請實驗者為 242 個日常詞彙手語,每一個類別錄製 5 次作為訓練集,並將四個模型設定 batch size 為 64,進行 fine tuning。

2. 請受試者為 242 個日常詞彙手語,每一個類別錄製 1 次作為測試集,輸入進四個模型,以 分析比較模型性能。

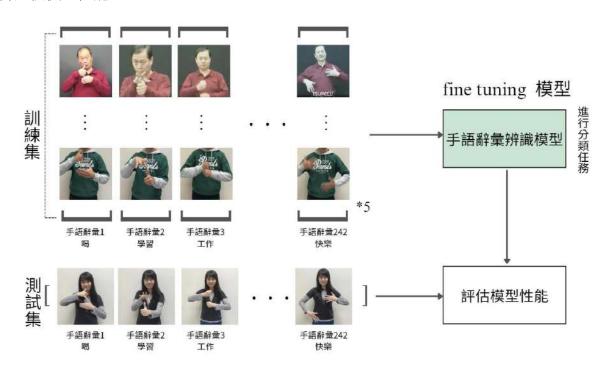


圖 4-10:日常詞彙實驗流程圖(來源自行製作)

2	夏天	記得	朋友	幫_N	錢	每天	如果	直	我們兩個	清楚	硬_A	英文	晴天_A
3	互相	完了	卡片	床	胖_A	家_B	昨天	不要_s	欣賞	刷牙	金	介紹	也許
4	可以	跑_B	旅行	他們	老師_N	整理	結束	停	作弊	答案	臉	大樓_A	
5	森林	安靜	出爾反於	照相_A	走	標準	奇怪_S	工作	開花	做	時間	老師_s	
6	油_s	剛剛	環境	他們兩個	好_A	寫	安排	抱	今天	學校	禁止	裙子	
7	綠	每	問題	夫妻_B	很	長大	還沒	明天	有	一定	叫	像	
8	回家_B	讀書	下雨	兒子	名字_s	加	快樂	平靜	真	鈴_A	客廳_A	眼鏡	
9	小孩	原來	幾時_B	忙	人_A	幸福	殘忍	倒	游泳_A	電影	結婚	過去的量	最近
10	危險	拉	自己N	腳踏車	爸爸	不要 N	立刻	上癮B	手語	我們	她	漂亮 S	
11	生氣_A	旅館_A	負責B	少	未來_B	回答	希望	地方	雞	碰見	知道_s	變	
12	菜_N	好_B	愛	肯	想	去_B	難	不好_B	生病	怕_N	外面	妹妹	
13	棒_A	其中	∦I_S	告訴	參加	兩個	更	見面_B	衣服	奶_s	動物	吃	
14	什麼_A	努力	結果	放	香煙	那	皮	困難	行動	比賽	身體B	抽煙	
15	馬上	種類	見	會_N	早上 B	媽媽	獲得	經過	提供 B	繼續	學生	會議	
16	正確	日本	要_s	決定	現在	哭_A	是	作業	保護	報紙	玩	考試_S	
17	睡_A	他_A	完全	責任	出	有沒有	你	再_A	台北	快_s	健康	開車	
18	同學_A	舒服_A	去_A	嬰兒	代替	我_B	邀請	依然_A	失約_A	燈	這	棒球_B	
19	鐓	認真	張	西瓜	亮 A	收集	以後 B	貴A	目的	吃飯 A	桌子	輸	
20	牙齒	第一名_	郵票	不能	講	炒	不知道_	眼睛	加入_B	送	近	來	
21	世界	學習	哪裡	整天 B	關心	還不錯	答應	看N	有錢	一起A	事情	豐富	

表 4-2: 日常詞彙手語表(來源自行製作)

六、手語翻譯實驗

為了在真實生活進行手語翻譯,本研究整合了手語詞彙辨識模型,自行設計了手語句子翻譯系統,並測試系統的整體準確度,實驗流程如下(圖 4-11)。

- 1.請受試者錄製句子影片,並將影片以每 20 幀切分輸入進手語詞彙辨識模型後,得到每個片段所代表的手語詞彙。
- 2.將步驟 1 輸出的所有手語詞彙,輸入進本研究設計的滑動窗口演算法,得出影片中所包含的所有中文詞彙。
- 3.將每個中文詞彙輸入進大型語言模型,重新排列成中文句子,與原句子比較分析,並計算 準確度與 BLEU 分數以量化模型翻譯的表現。

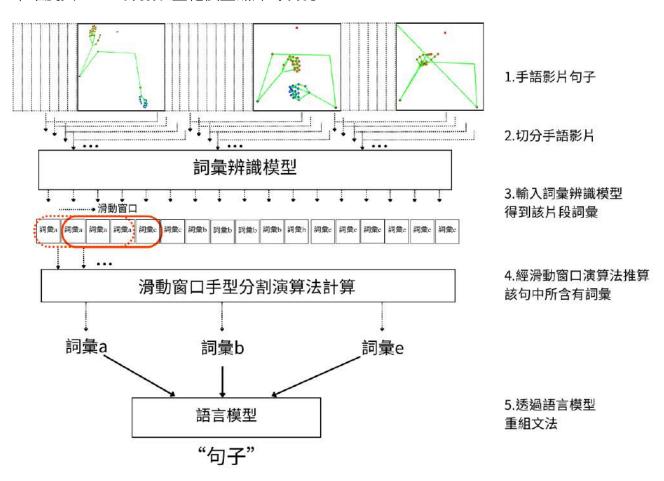


圖 4-11:手語翻譯流程圖(來源自行製作)

(一)滑動窗口

由於手語句子是由許多手語詞彙所組成的,本研究透過切分影片後,輸入進辨識模型 得出可能詞彙,並且設計了滑動窗口演算法來判斷哪些詞彙存在於手語句子中,實現方法 如下:

- 1.將每20幀的手語影片輸入詞彙辨識模型並記錄辨識結果,作為詞彙序列(如圖4-12)。
- 2.設定 18 大小的滑動窗口,搜尋詞彙序列,若窗口內某詞彙的總數大於 75%就記錄該詞彙, 若沒有詞彙的總數低大於 75%就記為-1(如圖 4-13)。
- 3.將每個窗口的標記整理成序列後,找尋連續相同的標記,合併成結果。(如圖 4-14)。

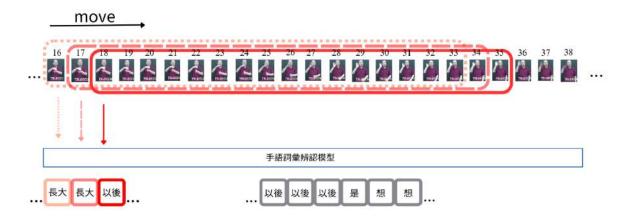


圖 4-12:影片切分辨識示意圖(來源自行製作)

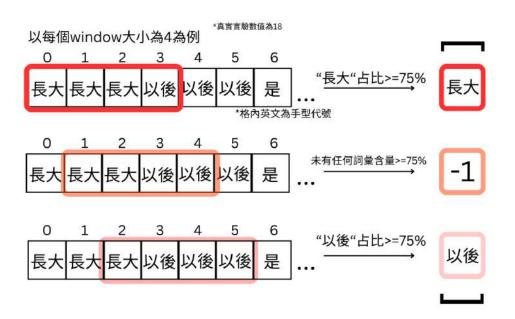


圖 4-13: 滑動視窗示意圖(以每個 window 取 4 幀為範例)(來源自行製作)



圖 4-14:合併相同標記示意圖(來源自行製作)

(二)大型語言模型翻譯

透過上述方式即可得出每個手語詞彙,最後將詞彙辨識結果輸入進大語言模型重組成句子,達成手語翻譯。

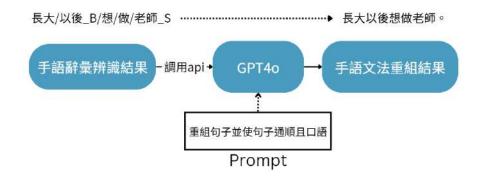


圖 4-15: 大型語言模型調用流程圖(來源自行製作)

(三)評估模型翻譯性能

除了人工判讀翻譯準確率以外,在本研究中,我們使用了 BLEU 分數作為翻譯系統性能的評估指標。BLEU 常用於量化機器翻譯結果的性能。它通過計算機器生成翻譯與參考翻譯之間的 n-gram 匹配度來評估翻譯的準確性。其中, P_n 是 n-gram 精確度,翻譯間匹配的n-gram 比例, W_n 是權重、N是最大 n-gram 長度。

BLEU = BP · exp
$$\left(\sum_{n=1}^{N} w_n log p_n\right)$$

在我們的實驗中,我們採用了 BLEU-1 至 BLEU-4 的評估,以全面衡量模型在不同層次語言結構上的表現。BLEU-1 和 BLEU-2 更偏向於詞和短語的精確性評估,而 BLEU-3 和 BLEU-4 則側重於句子的語言連貫性。

伍、研究結果與討論

一、預訓練結果

在預訓練中,本研究隨機遮蓋訓練資料中的座標輸入進模型,模型將預測遮蓋部分的座標。本研究使用 MPJPE loss function 衡量預測值與目標值的差異以更新模型參數,實驗訓練兩種不同輸入資料的模型,一個是包含骨架座標的模型 (Body Pretrained Encoder);另一個是包含手部點座標的模型(Hand-shape Pretrained Encoder)。

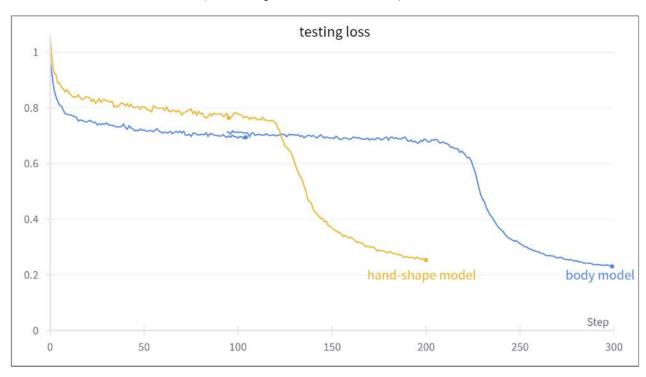


圖 5-1:不同 Batch Size,模型訓練的平均 testing loss(來源自行製作)

從圖 5-1 可以觀察在預訓練過程中,兩個模型分別在 step120 以及 step220 時, loss 值出 現一次巨大的轉折,急遽下降,模型是在這期間學習到了手語的特徵。

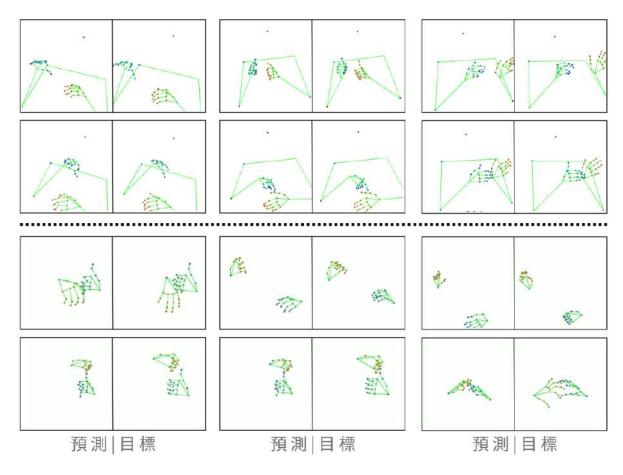


圖 5-2:模型預測的結果(左),目標值(右),上半部分為骨架預訓練模型,下半部分為 手型預訓練模型。(來源自行製作)

圖 5-2 很好的展示出兩個模型對於遮蓋部分的預測與目標值相當接近,可以發現模型在 雙手位置與關節點的部份預測得相當精準,而儘管在每個點預測上與目標值仍有些差異,不 過仔細觀察,這是合理的誤差,模型仍然有預測出與目標值相同的手勢與姿態,足以說明模 型學習到手語的特徵。

二、手語詞彙辨識實驗

在詞彙辨識實驗中,本研究 fine tune 了四種模型,在 242 個日常詞彙手語資料中進行訓練,並且透過測試集來評估模型的準確率。

模型/準確率	手型模型	骨架模型	融合模型	對照組
Testing accuracy	74.3%	82.6%	94.8%	12.3%

表 5-1:四種詞彙辨識模型最終在測試集辨識準確率(來源自行製作)

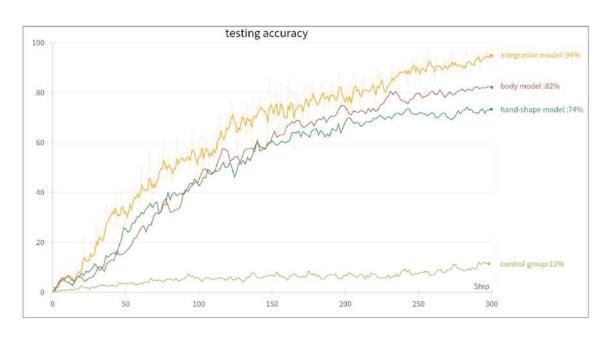


圖 5-3:四種詞彙辨識模型在訓練階段,測試集辨識準確率(來源自行製作)

從表 5-1 可以看到,四種模型裡以融合模型的辨識準確率是最高的,相較於對照組提升了 82.5%的辨識準確率,也比手型、骨架模型的辨識準確率來的更高,足以說明融合模型可以將手型、骨架模型各自的優勢很好的加在一起,完成 94.85%的辨識準確率。因此,後續的手語句子翻譯系統,本研究都選用表現最好的融合模型。

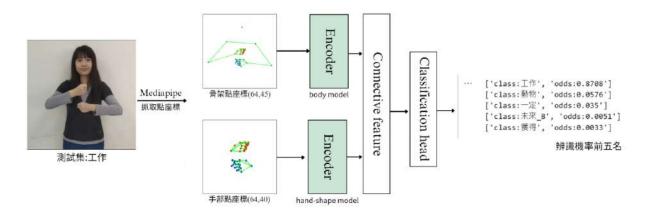


圖 5-4: 詞彙辨識實驗圖(來源自行製作)

圖 5-4 受試者比出手語,經過 mediapipe 處理得到手部點座標,輸入進融合模型辨識,以"工作"為例。

三、滑動窗口演算法

本研究為了將每個手語詞彙從句子中辨識出來與去除雜訊,開發了滑動窗口句子分割 演算法。實現方法為:每一幀往後取 20 幀輸入詞彙辨識模型,紀錄辨識結果,並將辨識結 果由滑動窗口進行篩選,最後將連續詞彙進行合併,得到句子中所含詞彙。(如圖 5-5)。

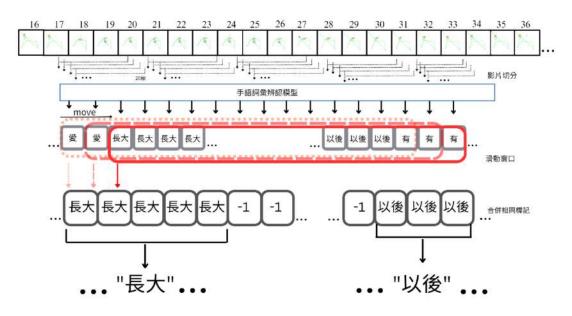


圖 5-5: 滑動視窗手型分割演算法實驗圖(來源自行製作)

四、手語文法及大型語言模型實驗

臺灣自然手語的文法結構與中文差異顯著,且包含大量省略詞彙和需意會的動作。由於現有手語與中文的對照資料不足,本研究另闢蹊徑。測試結果顯示,LLM 能有效解決此問題,特別是在處理省略和需意會部分時,優於傳統翻譯模型。其中,GPT-4 表現最佳,因此本研究使用了GPT-4 的 API 進行手語文法翻譯。

圖 5-6: GPT4-API 調用 (來源自行製作)

五、手語翻譯實驗結果

為了在真實生活進行手語翻譯,本研究整合了手語詞彙辨識模型,自行設計了手語句 子翻譯系統,並邀請受試者錄製手語句子影片來測試系統的整體準確度。

評估指標	BLEU-1	BLEU-2	BLEU-3	BLEU-4
scores	49.21	35.53	27.53	20.98

表 5-2:翻譯系統在不同評估指標中取得的分數

在表 5-2,我們採用了 BLEU-1~BLEU-4 的評估,衡量模型從詞彙辨識的準確率到句子翻譯的連貫性。結果顯示,模型在 BLEU-1 上達到了 49.21,這表示系統在詞彙辨識的層面上準確率十分優秀,然而,隨著 n-gram 的數量增加,分數逐漸下降到 BLEU-4 的 20.98,顯示在更長語境的翻譯中,模型的語言連貫性尚有提升空間。我們認為,這樣的結果主要是因為本研究中的翻譯實驗集中在生活中的短語,這使得在 BLEU-4 評估中未能充分展現模型的優勢。短語的簡單結構限制了 BLEU-4 分數在更高 n-gram 評估中的發揮。

圖 5-7 展示了手語句子"我長大以後想當老師"的翻譯實驗圖,將影片切分為 300 個 20 幀的片段輸入進詞彙辨識模型得到對應的詞彙,並經滑動窗口演算法處理,隨後得到"我_B、長大、以後_B、想、工作、老師_S"。將此輸入語言模型翻譯,最後模型翻譯結果為"我長大以後想做老師"。其結果顯現出,手語詞彙辨識模型辨識的相當準確,且大型語言模型是能夠理解手語中意會的部分,儘管與原句不是一字不漏地翻譯,其翻譯結果仍與原句意思相同。

表 5-3 展示了手語句子翻譯的實驗最後的統計表,全部實驗了 100 個手語句子,並取得 **88%**的翻譯準確率,其中詞彙辨識模型失誤率 9%,大型語言模型失誤率 6%。

"我長大以後想當老師"手語影片

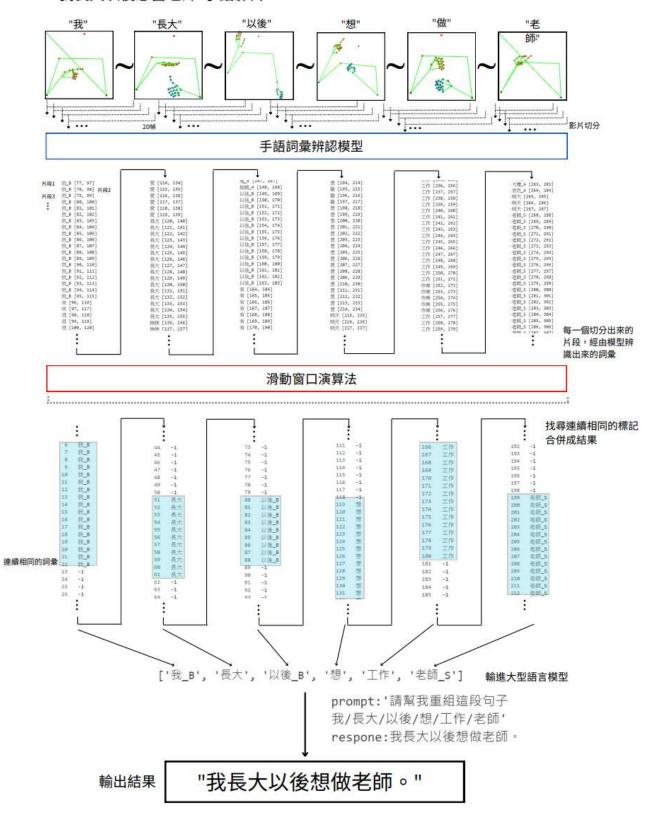
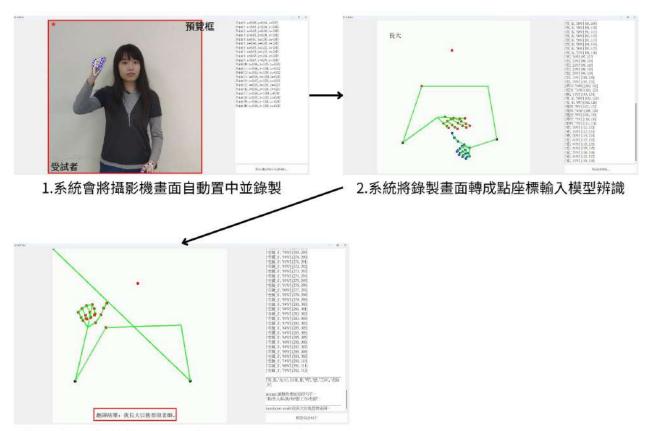


圖 5-7:手語翻譯實驗圖(來源自行製作)

原句	原詞	模型煙行輸出	大型語言模型重組輸出
有一件相子	我_8/裙子/有/一	我_0/裙子/青/一	我有一根裙子。
1星標準苦底 作要試真	道/答案/標準 工作/認真/要 S	這/苦虧/標準 工作/認真/要 S	這是標準苦務。 工作要認真。
- TF 表现具 F南不能出門	下而/出/不能	下開/出/不能	下南時不能出去。
也跑得很快	他_A/跑_B/更/换_S	他_A/她_B/夏/探_S	他跑得很快。
5們領工作忙,很少見賣	他俩/工作/忙/晃面_8/少	他俩/工作/忙/晃面_8/少	他們工作忙,所以見面少。
金能要我馬上回家	爸爸/告訴我/馬上/回家/馬上	爸爸/告訴我/馬上/回家/馬上	爸爸告訴我馬上回家
青自我介紹一下 50%の特別性達70M	(叫人)/我_0/自己/介紹 语/陈/陈系。(陈/陈天林	我_8/自己/介绍 性/第/体统 / / / / / / / / / / / / / / / / / /	我可以自己來介紹自己
直離肉飯勉強還可以 等優地方都有好人和壞人	這/維/吃板_A/吃/達不錯 各些/人_A/好/人_A/不好_B/有	這/雜/吃飯_A/吃/還不維 人_A/好/人_A/有	這雜吃起來還不繼 人很好,
ラ天皇皇朝後昭	今天/星期機	今天/星期機	今天是星期幾?
見伝是賴點呢	現在/競時_B	現在/幾時_8	現在檢點了?
工理停,排燈行	红_5/億/停//排/億/走	红_5/电/停//教/电/走	紅燈停・我走・
牙齒是很緊硬的	牙磨/硬_A	牙齒/睫_A	牙齒很硬。
陸 契 前要別等 也顧慶和你結婚嗎	醛_A/前/耐牙/要_5 他 A/你/结婚/他 A/詞	醛_A/束/则牙/要_5 他 A/你/結婚/他 A/期	要經覺了,來則牙· 你覺得他發意和你結婚嗎?
(日間を付ける場合) (日間を付ける場合)	作業/真/医力/加/會 N	作集/寫/图力/加/會_N	京作業會增加壓力
能能很生意	能能/生氣 A/很	爸爸/主题 A/很	能能很生順 。
臺轡之後,老師就叫學生回來上頭	鈴_A/完了/老師_N//叫/大家追來/念書	转_A/完了/老師_N/一起/念書	於監督了,老師和我們一起讀書。
16位事一定要考慮清楚。不要反反覆覆	你/你/事/什麼//你/鸡/清楚/再/一定//出爾反爾_B/不要	你/做/事/什麼/你/想/清楚/再/一定/出爾反爾_8/不要	你做什麼事,一定要先想清楚,不要再出爾反爾
音倒哲案正確	這/啓案/正確	這/答案/正確	這個答案是正確的。
表決定努力 學管手 語	手題/學/我_B/決定/努力	手頭/擊/我_8/決定/努力	我決定勞力學手語・
7一定要幫我	你/戴/教/一定	你/覧/教/一定	你一定要幫我。
也努力得到第一名 表現保護,人人有責	他_A/認真/努力/獲得/第一名_N 環境/保護/每/人/賽任	他_A/認真/努力/獲得/第一名_N 環境/保護/每/人/賽任	他認真努力以獲得第一名。 保護環境是每個人的責任。
R現時頃・八八月真 F天的模別比賽・結果我們難了	非天/極球/比賽/結果/我們/整	環境/神鏡/電/八/黄芒 昨天/極球/比賽/結果/我們/翻	我們昨天的棒球比賽輸了。
記述書籍で	我_8/比赛/勤	我 B/比賽/輔	我動了比賽。
以最好不要失約	你/失约_A/不要_s/配得	你/矢約_A/不要_5/記得	別配導你失約的事情
也時實的目的很明確	他_A/讀書/目的/清楚/很	他_A/讀書/目的/清楚/很	他讀書目的明確。
色爸僧促我立刻上床證實	爸爸/懂/立刻/去/睡覺	爸爸/備/立刻/去/睡覺	爸爸提供叫我去睡覺了。
包許明天會下面	明天/下雨/也許/會_N	明天/下南/也許/會_N	明天可能會下原呢。
的可以最致何一起去旅遊頃	教們/一起/去/旅行/玩/她/加入/可以	我們/一起/去/旅行/玩/她/加入/可以	
要求他是老師 #15年は第フ	他_A/老師_S/原來 徐珠/生病	他_A/老昕_s/原来	原來他是老師・
*妹生病了 *天的考試很難	昨天/考試 5/難	妹妹/生病 昨天/岩賊 5/難	妹妹生病了。 昨天的考試很難。
F大町等部(収無 也今天警我代班	作人/考証_5/辞 他 A/今天/幫我/代替/工作	昨大/考览_5/辞 他 A/今天/常装/代替/工作	作大助等的收耗。 他今天警我工作了。
1. 1. 住在台北	我 的 聚 的 台北	我」的家。例台北	我家在台北。
也們兩個人是我的朋友	他們兩個/我 8/朋友	他們兩個/我_8/朋友	他們兩個是我朋友。
1望世界和平	希望/世界/和平	粉罐/世界/和平	希望世界能夠和平。
5. 近来的行動程異	他_A/强去的最近/行動/奇径_s	他_A/退去的最近/行動/奇怪_S	他最近做的事情真的很奇怪。
己去森林很危險	自己_N/去_B/森林/危險/很	自己_N/去_B/森林/危險/很	自己去森林很危險。
整整面的工作 - 真叫我害怕	這/工作/復康/我_B/伯_N/很	道/工作/這些/我_B/怕_N/很	我對這基工作感到很害怕。 簽签完全不到這個件事情。
& 管完全不知道這件事 以則則已經回答了你的問題	道/事情/爸爸/完全/不知道_N 你/問題/我_B/剛剛/回答你/完了	值/事情/爸爸/完全/不知道_N 你/問題/我_B/剛剛/回答你/完了	我剛剛已經區著了你的問題。
N	明天/會議/我_B/安排/完了	明天/會議/我」B/安排/完了	明天的會議我已經安排好了。
「 病 昨 天 整理 好 了 客 籍	媽媽/昨天/整理/客廳/地方/整理/结束	妈妈/昨天/整理/客廳/地方/整理/纸束	媽媽昨天已經招客廳整理完了。
故每天騎部路車上學	我 B/每天/新腳踏車 A/去/讀書	我 B/每天/發腳踏車 A/去/讀書	我每天騎都路車去讀書・
0頭我是同學	我俩/同學_A/是	我师/同學_A/是	我俩是同學。
又服在桌子上	桌子/衣服/放	桌子/危險/放	桌子放得很不穩。
7.有收集製賣嗎	多票/張/緊票/你/收集/有沒有	郵票/限/郵票/你/收集/有沒有	你有沒有收集彭票呢?
T作環沒完成,繼續把它做完	工作/完/還沒/再_A/繼續/工作/完	工作/完/再_A/繼續/工作/完	工作完成後,再指續工作。
事安・你好 6展大以後想當老師	早/安靜/好_A	早/安靜/好_A	早安,但裡真的很難
0長人以情で異名的 0果你不出席會議・記得要頭表記	長大/以後_B/想/做/老師_S 如果/你/會議/参加/不要_S//你/記得/告訴我	長大/以後_B/想/做/老師_S 如果/你/會議/参加/不要_S/你/記得/告訴我	如果你不参加會議,配得告訴我。
的名字看起來很不錯	他 A/名字 5/看 N/투 A	他 A/名字 S/看 N/存 A	他看起來名字很棒。
電影時請保持安靜	電影/電/語/安藤/拉然_A	電影/書/碼/安藤/泣然	去靜看電影,依然不調。
也很有錢,到虞亂花	他」A/董嘉/有錢/花錢(到重)	他 A/豐富/有錢/花錢(到處)	也錢多,斯以花得很大方
e子長得很像爸爸	兒子/他_A/爸爸/他俩/猜/像	兒子/去/養養/他備/前/像	兒子像他爸爸,他俩的驗很像。
自動物皮製作產品很強忍	動物/皮/物/部/架_A/隆忍	動物/皮/做/那/哭_A/階忍	那種以動物皮為原料的行為實在太隨忍,讓人無法忍住
(本來是書師	我_B/厚來/老師_5	我_6/原宋/老師_5	我原來是老朝唯
天要考英文	今天/英文/考試5/要_5	今天/英文/考試s/要_s	今天我們要考察文了
(們一起走吧 [天的時候我們會吃西瓜	我們兩個/來/走 夏天_B/夏天_A/西瓜/吃	我們兩個/來/走 夏天_B/夏天_A/西瓜/吃	我們兩個一起走吧! 廣天就該的西瓜期!
的眼睛美了起来	他 A/眼睛/德/亮	他 A/眼睛/蓮/亮	他的眼睛真的好亮啊
試不可作弊	看試。5/作弊/不能	看試_5/作弊/不能	者試時不能作學
到困難時 - 我可以解你	磁見/困難/有//我 8/報你/可以	避見/函難/有/我_e/難你/可以	如果你碰見困難・我可以離你・
(各種性的避時	他_A/密博技//我_D/答覆	他_A/邀請我/我_B/菩薩	他邀請我,我答應了+
力學習是為了有更好的將來	學習/努力/目的/未來_B/更/釋_A	學習/努力/目的/更/博_A	勞力學習是為了更棒的目標。
(本)	每/出/錢/安包/放進/要_S	每/出/钱/贞包/放進/要_s	每次出門,要把錢放進支包。
5天游泳動奏體很好 5〒00日本	每天/游泳_A/解我/身種/好_A	每天/指示_A/報我/與體/好_A	母天游泳到身體有好處。
(可以做完 :是一種資產金費	做/絕吏/教_B/可以 後/金/種寢/其中/一/金/妻	街/結束/我 b/可以 撮/全/捶膊/其中/一/全/費	我可以做到結束。 金是其中一種貴金屬的種類。
(現一種典里哲學 (祖快樂	現/並/理規/其中/一/並/責 我 B/快樂	類/並/推理/共中/一/並/資 班 B/快機	並是其中一種實施學可種類。
○10 (大) ●20 明天是猶天	知。 知識/明天/嫡天_A	新堂/明天/寶天_A	希望明天墨邁天
是學生	我 8/學主	我_B/學生	我是學生。
1 選了我一顿	他_A/見舞	他_A/見我	他見到我了 -
1胡我很漂亮	他_a/告訴/我_s/漂亮_s	也_A/告訴/我_B/清亮_S	他告訴我我很漂亮。
挑曲永遠愛我	他_A/告訴/我/愛/我_B/藝天_B	他_A/告訴/我/爱/我_B/整天_B	他一點天都在對我說變我。
(信是好朋友、從小一起長大 ## - #19949	我俩/朋友/表大 我俩/	我俩/朋友/長大	我偏是從小一起長大的朋友・
(何一起照相 (() 本人) 本在2000年	表價/一起/照相	我债/一起/照相	我們一起來拍單吧
好多小孩在跑來跑去 對話语上傳	小孩/他們/她來說去 他 A/道/香煙/酒/上等 B	小孩/他們/雞來跑去 他 A/香煙/酒/上轉 B	他們的小孩在到盧跑來應去 他對香煙和志上傷。
(野於塔上時 2重展止抽煙	他_4/理/香/厚/上時_8 學校/描/抽煙/禁止	他_女者塔/诗/上榜_B 學校/塔/抽得/新止	地野客運和資上職。 這裡的學校是不允許檢摸的彈
《集宗正描程 · 媽抱抱小孩	小孩/媽媽/抱	小孩/妈妈/抱	場場地響小孫
本的櫻花全都開花了	日本/劉/司載/開花	日本/堂/開花	日本的當化就像花開一樣
(真的是我朋友 (他是一個真誠的朋友)	他_A/我_B/朋友/真/是	他_A/我_B/朋友/真/是	他真的是我朋友
的朋友在日本	我_8/朋友/那/日本/那	我_8/朋友/都/日本/都	那是我在日本的一位朋友
1进我一张卡片	卡片/他_A/挂载	卡片/他_A/结我	他送我平片。
知道我的眼鏡在哪裡嗎	我_8/银鏡/哪裡/你/知道_5	我_8/银廣/哪裡/你/知道_5	你知不知道我眼鏡放在睇裡?
· 克曼嘎奶	雙兒/奶_s/要_s	嬰兒/奶_s/要_s	他拉了我
門記得要用越 (4000) 見分記 (12000) (4	時/出/記得/錢/放口線 三七/第)/第三/[(古本/中漢/東京/第	等/出/記標/強/放口袋	時,出門記傳把錢放口鎖
2. 相関心是婚姻幸福的關鍵 2. 禁事機(大好)	互相/第心/這兩個/夫妻/率福/重要/這 程 A/對/集體 n/体表/不任 n	互相/關心/趙兩個/夫妻/幸福/重要/建 ※ 人為時 の/妹弟/エペス	這對夫妻互相關心,這對他們的幸福很重要 身體太胖對健康不好。
B群對身體不好 D菜時要用油	肝_A/討/身體_B/健康/不好_B 菜_N/炒/油_S/例/要_S	并_A/身體_B/健康/不好_B 菜 N/炒/油 S/倒/要 S	身體不許對健康小好。 炒菜的時候,要先倒油碗
(把我拉锡去	来_N/2////////	東 N/2/車3/東/東 3 他 A/拉/我	以来的時候,要先到海峡 我們兩個的地方很近。
1500 (1200 A) (百下南了	外面/下南	外面/下南	外面正在下南。
(約兩個組得很近	我們兩個/地方/近	我們兩個/地方/班	我們兩個的地方很近。
能體提供舒服的床	舒服_A/末/舒服_A/那/旅館/那/負責_B/提供_B	舒服_A/床/種類/那/旅館/那/負責_B/提供_B	那家旅館負責提供各種類型的舒服床。
	報紙/批 B/放賞/要 5	郵纸/我 B/欣賞/要 S	我想着看到纸。
我要看取纸		THE SAW THE BY THE PARTY AND IN	Sarrad TE TE + STILL

表 5-3: 手語翻譯實驗成果圖(紅色底格子為錯誤輸出,正確率計算以最右欄計算)

最後,我們整合實驗用的代碼,開發了簡易的操作介面。如圖 5-8



3.按下模型翻譯連接GPT,最後翻譯結果顯示在畫面上

圖 5-8:翻譯系統操作介面

六、問題與討論

本研究在手語翻譯上達到了 **88%**的準確率,本研究者詳細調查剩下 12%的錯誤,並歸納出以下幾點:

(一) Mediapipe 的不準確

在本研究的手語詞彙辨識模型實驗中,大部分的辨認錯誤都來自於 Mediapipe 的辨認 失誤,而導致模型輸入資料完全錯誤。可惜的是,現在市面上仍無比 Mediapipe 更準確的手 部點座標辨認模型。



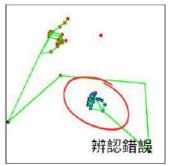


圖 5-9: Mediapipe 辨認錯誤(來源自行製作)

(二) LLM 翻譯問題

雖然 LLM 可以很好的協助我們重組中文句子,但是還是有其翻譯失誤的時候。 比如在"你最好不要失約"中,其拆解單字為"你/失約_A/不要_S/記得",但輸入進 LLM 輸出的結果卻是"別記得你失約的事情"。

然而,本研究仍採用 LLM 最主要原因來自於他的推理能力。像是"肥胖對身體不好"中,原拆解文字為"胖_A/對/身體_B/健康/不好_B",模型辨認錯誤導致輸出變成"胖_A/身體_B/健康/不好_B",但是 LLM 有完整理解語意,翻譯為"身體太胖對健康不好"。

本研究者推斷原因為在中文及手語的文法架構上的出入,導致 LLM 會錯意。日後期待更進提示詞,讓 LLM 的表現更精準。

原句	原詞	模型運行輸出	大型語言模型重組輸出	
你最好不要失約	你/失約_A/不要_S/記得	你/失約_A/不要_S/記得	別記得你失約的事情	
肥胖對身體不好	胖 A/對/身體 B/健康/不好 B	胖 A/身體 B/健康/不好 B	身體太胖對健康不好。	

圖 5-9:LLM 翻譯問題, (來源自行製作)

陸、結論

本研究貢獻在於,第一次將自監督學習應用在台灣手語辨識,擺脫了過去的研究需要 大量標記樣本的困境。成為台灣第一個手語詞彙量突破百位數,達到了 242 個可辨識詞彙以 及 94.8%的辨識準確率,而且本研究之作法僅需 5 個標記樣本即可訓練模型辨認詞彙。

本研究證明遮蔽一定資訊並使模型預測遮蔽內容的作法可適用於手語辨識任務。實驗結果顯示,結合手型模型與骨架模型的融合模型表現最佳,比沒有採用預訓練模型的對照組高出 82.5%的辨識準確率。

本研究也基於此自監督的預訓練模型,開發了首個可實際應用的手語翻譯的系統,在手語句子翻譯的表現達到優秀的 88%的準確率,且 BLEU-4 分數取得 20.98,證明了本研究的手語翻譯系統可真實應用在日常使用上。本研究者期待此技術在更妥善的完善後,可以投入實際應用的場合,幫助聾人與聽人的交流、溝通,增進弱勢族群的福祉,同時也可為手語教育帶來貢獻,增進社會的共榮和諧。

柒、參考文獻資料

- [1] Vaswani, A. (2017, June 12). Attention Is All You Need. https://arxiv.org/pdf/1706.03762.pdf
- [2] Devlin, J. (2018, October 11). *BERT: Pre-Training of Deep Bidirectional Transformers for Language Understanding*. https://arxiv.org/pdf/1810.04805.pdf
- [3] Alexey, D. (2020, October 22). *An Image Is Worth 16x16 Words: Transformers for Image Recognition at Scale*. https://arxiv.org/pdf/2010.11929.pdf
- [4] Kaiming, H. (2021, November 11). *Masked Autoencoders Are Scalable Vision Learners*. https://arxiv.org/pdf/2111.06377.pdf
- [5] Hvilshøj, F. (2023, March 3). What Is One-Shot Learning in Computer Vision. https://encord.com/blog/one-shot-learning-guide/
- [6] Huang, C. chiu. (2020, December 8). 論文閱讀筆記 3D 人體姿態辨識 Camera Distance-Aware Top-down Approach for 3D Multi-Person Pose Estimation from a Single RGB Image. https://williamchiu0127.medium.com/%E8%AB%96%E6%96%87%E9%96%B1%E8%AE%80%E7%AD%86

3d%E4%BA%BA%E9%AB%94%E5%A7%BF%E6%85%8B%E8%BE%A8%E8%AD%98-camera-distance-aware-top-down-approach-for-3d-multi-person-pose-estimation-from-3d89a33eeb33

[7] Li, M. (2021, October 29). *Transformer 论文逐段精读*. https://youtu.be/nzqlFIcCSWQ?si=5bXdhqd8Q3S_zff

[8] Mu, L. (2021, December 10). *MAE 论文逐段精读【论文精读】*. https://youtu.be/mYlX2dpdHHM?si=JzMmuL3Y1bi6-15L

[9] Mu, L. (2021, November 30). ViT 论文逐段精读【论文精读】. https://youtu.be/FRFt3x0bO94?si=8Xe34URtNDwvd5H9

- [10] Lee, H. (2019, June 1). Transformer. https://youtu.be/ugWDIIOHtPA?si=udow_2gw22RXRB5a
- [11] 中正大學. (n.d.). 臺灣手語線上辭典. https://twtsl.ccu.edu.tw/TSL/index.php
- [12] 教育部. (n.d.). 常用手語辭典. https://special.moe.gov.tw/signlanguage
- [13] 劉秀丹、曾進興(2007)。文法手語構詞語句法特性對聾生詞義與句義理解的影響。特殊教育研究學刊。 http://bse.spe.ntnu.edu.tw/upload/journal/prog/6O6_21SL_209R_35CM518.pdf

【評語】190027

手語辨識能力建議可以擴充到其他國家、其他語言以增加本系統 的可用性及貢獻度。

手語如何從單字轉換集結成句子,其流暢度、表達清晰度可以使 用 LLM 做進一步改善。

成果評估廣泛度應該持續加強。