

深度多模態生成式 AI 模型綜合評估報告：架構解析、翻譯應用與跨領域創意工作流

1. 生成式 AI 生態系統的演進與多模態典範轉移

1.1 從大型語言模型到大型多模態模型的技術跨越

2025 年至 2026 年初標誌著人工智慧發展史上的一個關鍵轉折點，即從單純處理文本的大型語言模型（Large Language Models, LLMs）全面轉向原生的大型多模態模型（Large Multimodal Models, LMMs）。這一轉變不僅是功能的疊加，更是底層神經網絡架構的根本性重構。過去，處理音訊、視訊和圖像通常需要依賴獨立的專用模型系統——例如使用 Whisper 進行自動語音識別（ASR），使用 ResNet 進行圖像分類，再將這些非結構化數據轉換為文本後輸入 LLM 進行處理。然而，新一代的模型如 Google 的 Gemini 1.5 Pro、OpenAI 的 GPT-4o、Anthropic 的 Claude 3.5 Sonnet、Microsoft 的 Copilot 以及 xAI 的 Grok 2/3，已經在不同程度上實現了對視覺、聽覺和文本資訊的「端到端」原生理解¹。

這種「原生多模態」（Native Multimodality）意味著模型不再將視訊視為一連串靜態圖片的集合，或者將音訊視為單純的波形轉錄文本，而是能夠捕捉時間序列上的動態變化、語氣中的情感色彩以及畫面空間中的相對關係。例如，Gemini 1.5 Pro 引入了長達 100 萬甚至 200 萬 token 的上下文視窗，這使其能夠在單次推理中處理長達數小時的視訊或音訊，從根本上改變了資訊檢索與分析的模式²。與此同時，Claude 3.5 Sonnet 透過「Artifacts」功能重新定義了人機協作的介面，將 AI 的輸出從線性的對話流擴展為可互動的應用程式視窗⁴。Grok 則利用其與 X（前 Twitter）平台的深度整合，展示了即時數據流與生成式 AI 結合的可能性，儘管這也帶來了新的倫理挑戰⁵。

1.2 五大核心模型的架構定位與適用場景分析

在探討具體的應用操作之前，必須對這五個主流模型的架構哲學與市場定位進行精確的剖析。這些底層差異直接決定了它們在處理「英文語音轉繁體中文」或「MP4 視訊轉文字」等複雜任務時的效能表現。

- **ChatGPT (OpenAI):** 被視為「全能型生產力基準」。其核心優勢在於生態系統的極度成熟，整合了 DALL-E 3 的圖像生成能力與 Whisper 的語音識別能力。GPT-4o 的推出進一步強化了其多模態的即時反應速度，使其在處理短語音翻譯和日常對話時表現出極高的流暢度與準確性。其「Advanced Voice Mode」模擬了人類的對話節奏，提供了極具沈浸感的互動體驗¹。
- **Gemini (Google):** 定位為「長文本與原生視訊理解的霸主」。Google DeepMind 的策略在於極致的上下文窗口（Context Window）與跨模態整合。Gemini 是目前唯一能夠在不需第三方插件的情況下，原生理解長達 1 小時以上高解析度視訊內容的模型。其能夠同時處理音訊軌與視覺軌的資訊，使其在複雜的視訊分析任務中具有壟斷性的優勢³。
- **Claude (Anthropic):** 以「安全性、可解釋性與程式碼推理」著稱。Claude 3.5 Sonnet 在程式碼生成與邏輯推理的基準測試中經常超越競爭對手。雖然其原生音訊與視訊上傳功能相對

受限（通常需要依賴 API 或轉錄後的文本），但其強大的文本處理能力使其成為翻譯文學作品、法律文件或進行深度語義分析的首選工具。其「Artifacts」功能為開發者與設計師提供了前所未有的原型製作效率²。

- **Copilot (Microsoft):** 定位為「企業工作流的智慧引擎」。不同於其他面向消費者的聊天機器人，Copilot 深度嵌入於 Microsoft 365 生態系統（Word, Excel, PowerPoint, Teams, Outlook）。它利用 Microsoft Graph 來存取使用者的郵件、行事曆和文件，這使其在處理企業內部的會議錄影、郵件摘要與跨應用程式協作上具有不可替代的地位⁵。
- **Grok (xAI):** 標榜為「即時、反叛與不受限的創意工具」。Grok 的最大護城河在於其對 X 平台全量數據的即時存取權，這使其在新聞分析、趨勢預測與社交情緒分析上具有絕對的時效性優勢。其內建的 Flux 模型在圖像生成上展現了極高的寫實度與風格化能力，且其「Spicy Mode」提供了一種區別於其他安全過濾嚴格模型的互動體驗¹⁵。

2. 深度任務解析：英文語音檔翻譯成繁體中文的操作模式與策略

將英文語音檔案轉換為繁體中文文本是一項涉及自動語音識別（ASR）、自然語言理解（NLU）與機器翻譯（MT）的複合型任務。針對使用者特別提到的繁體中文需求，本節將詳細拆解五大模型的具體操作流程、提示工程（Prompt Engineering）技巧以及潛在的技術限制。

2.1 ChatGPT：依賴 Whisper 的高精度轉錄與語境翻譯

ChatGPT 在處理語音翻譯任務時，主要依賴後端整合的 Whisper 模型。Whisper 是一個在 68 萬小時多語言與多任務監督數據上訓練而成的強大 ASR 模型，對於口音、背景噪音與技術術語具有極強的魯棒性。

具體操作流程：

1. **移動端應用（iOS/Android）：**這是目前最直觀的操作方式。使用者點擊對話框旁的「+」號或波形圖示，選擇上傳音訊檔案（支援 MP3, WAV, M4A 等格式）。上傳後，ChatGPT 會自動調用 Whisper 進行轉錄，並將轉錄後的文本顯示在對話中。隨後，使用者可輸入指令進行翻譯¹⁸。
2. **桌面端與網頁版：**直接將音訊檔案拖曳至對話視窗。需要注意的是，免費版與 Plus 版在檔案大小與數量上可能存在限制（通常單檔不超過 512MB，但在對話模式下建議更小以確保穩定性）。

提示工程策略（Prompt Engineering）：

為了獲得最佳的「繁體中文」翻譯，單純的「Translate this」指令往往不夠，容易產出帶有中國大陸用語（如「視頻」、「軟件」）的文本。建議使用以下結構化提示 20：

- **角色設定（Persona）：**"You are a professional translator specializing in localization for the Taiwanese market."（你是一位專精於台灣市場在地化的專業翻譯。）
- **上下文賦予（Context）：**"The attached audio is a snippet from a technical engineering podcast."（附檔是一段關於技術工程的 Podcast 片段。）

- **具體指令（Instruction）**： "First, transcribe the English audio verbatim. Then, translate it into Traditional Chinese (Taiwan standard). Ensure that technical terms are accurate and colloquialisms are adapted to local usage. Avoid Mainland Chinese terminology." (首先，逐字轉錄英文音訊。然後，將其翻譯成台灣標準的繁體中文。確保技術術語準確，口語表達符合在地習慣。避免使用中國大陸術語。)

創意應用變體：

- **雙語對照學習模式**：要求 ChatGPT 輸出一個 Markdown 表格，第一欄為英文原文，第二欄為繁體中文翻譯，第三欄為「關鍵詞彙解析」，這對於語言學習者或需要校對的專業人士極具價值⁷。
- **風格轉換翻譯**：如果音訊是輕鬆的訪談，可以指令：「請用台灣 PTT 或 Dcard 論壇的輕鬆口吻進行翻譯」；如果是學術演講，則指令：「請使用正式、學術性的繁體中文風格」²⁰。

2.2 Gemini：長音訊的原生理解與跨模態分析

Gemini 1.5 Pro 在處理音訊翻譯任務上展現了與 ChatGPT 截然不同的技術路徑。它不需要先將音訊轉錄為純文本再進行處理，而是能夠將音訊波形編碼為 tokens 直接輸入模型。這種「原生音訊理解」使得它能捕捉到語氣、停頓甚至背景聲音中的訊息。

具體操作流程：

1. **Google AI Studio 或 Gemini Advanced**：使用者可直接點擊「上傳」按鈕選擇音訊檔案。Gemini 支援的檔案大小上限極高（可達數小時的錄音），這得益於其百萬級別的 token 窗口²¹。
2. **Google Drive 整合**：對於超大型檔案，可先上傳至 Google Drive，然後在 Gemini 對話框中透過 @Google Drive 指令直接調用該檔案進行分析，無需消耗本地頻寬重新上傳。

提示工程策略：

- **時間戳記索引**：由於 Gemini 能理解時間軸，提示詞可以包含：「請翻譯這段錄音，並在每個主要段落前標註原始音訊的時間戳記（Timestamp），以便我回溯收聽。」⁹。
- **多語者識別（Speaker Diarization）**："Identify different speakers (Speaker A, Speaker B) and translate their dialogue into Traditional Chinese format:: 翻譯內容"²¹。

創意應用變體：

- **會議情感分析與翻譯**：除了翻譯文字，還可以要求 Gemini 分析講者的情緒：「請翻譯這段談判錄音，並在括號中註記講者當時的語氣（例如：憤怒、猶豫、諷刺），這有助於我理解談判的潛台詞。」這是原生多模態模型的獨特優勢。

2.3 Claude：高精度文本重組與風格化翻譯（需前置轉錄）

截至 2026 年初，Claude 的網頁介面尚未全面開放直接上傳音訊檔案進行原生分析的功能（雖然 API 已支援部分多模態輸入，但主要集中在視覺）。因此，使用 Claude 進行音訊翻譯通常需要一個「轉錄中介」步驟。儘管如此，Claude 3.5 Sonnet 在文本生成的細膩度與邏輯性上往往優於其

他模型，使其成為高品質翻譯的首選。

操作方式 (Workaround) :

1. **前置轉錄**：使用 OpenAI Whisper API、Otter.ai 或其他語音轉文字工具生成 SRT 或 TXT 檔。
2. **文本上傳與分析**：將轉錄好的文本檔案上傳至 Claude。Claude 的 200k token 上下文視窗足以處理整本書籍長度的對話錄。

提示工程策略：

- **文學性與修辭**：Claude 擅長處理複雜的句式與修辭。提示詞："Here is a transcript of a literary speech. Translate it into elegant Traditional Chinese, preserving the rhetorical devices and the emotional resonance of the original speaker."（這是一份文學演講的逐字稿。請將其翻譯成優雅的繁體中文，保留原講者的修辭手法與情感共鳴。）¹²。

創意應用變體：

- **方言與特定語境模擬**：雖然 Claude 不能聽，但在文字翻譯上，你可以要求它：「請將這段英文對話翻譯成帶有台灣眷村氣息的繁體中文小說對白」，Claude 對於文化語境的模擬能力極強²³。

2.4 Copilot：企業生態系中的自動化翻譯

Copilot 的優勢在於它不是一個孤立的工具，而是微軟生產力軟體的「副駕駛」。

操作流程：

1. **Teams 會議整合**：這是 Copilot 最強大的場景。會議結束後，Copilot 會自動基於 Teams 的即時轉錄 (Transcript) 生成內容。使用者只需指令：「Generate a Traditional Chinese summary of this meeting and list all action items.」（生成這場會議的繁體中文摘要並列出所有待辦事項。）¹⁴。
2. **Word 聽寫整合**：使用者可在 Word 中開啟「聽寫」功能讀入英文，然後側邊欄呼叫 Copilot：「Rewrite this text into Traditional Chinese.」

創意應用變體：

- **跨語言郵件輔助**：在 Outlook 中收到一封含英文語音附件的郵件，使用 Copilot 摘要語音內容並直接草擬繁體中文的回覆，實現無縫的跨語言商務溝通²⁴。

2.5 Grok：即時性與次文化翻譯

Grok 目前在音訊處理功能上相對較新，但其背後的 xAI 正在快速迭代。

操作預期：未來 Grok 將能直接分析 X 貼文中的音訊片段。目前主要依靠使用者輸入轉錄文本。

創意應用變體：

- **俚語與迷因（Meme）翻譯：**由於 Grok 訓練數據包含大量 Twitter/X 的即時對話，它對於最新的網路俚語（Gen Z slang）、加密貨幣術語或流行迷因的理解遠超其他模型。指令：「Translate this transcript into Traditional Chinese using current Taiwanese internet slang.」（用台灣當前的網路流行語翻譯這段文字。）Grok 能產出非常「接地氣」甚至帶有諷刺幽默的翻譯，適合社群媒體內容的地化⁵。

3. 深度任務解析：傳送 MP4 檔並將內容轉成文字的技術路徑

「視訊理解」（Video Understanding）比單純的語音轉錄更為複雜，因為它涉及視覺幀（Visual Frames）與音訊軌（Audio Track）的同步分析。使用者希望「將內容轉成文字」，這可能包含語音逐字稿、畫面場景描述、或是畫面中出現的文字（OCR）。

3.1 Gemini：原生視訊理解的技術護城河

Gemini 是目前市場上處理此任務最強大的工具，因為它具備「原生多模態」能力。它不是透過每秒截圖再用圖像模型分析（雖然早期版本是這樣），而是能夠處理連續的視覺編碼。Gemini 1.5 Pro 能夠以每秒 1 幀（1 FPS）的速度對視訊進行採樣與標記（Tokenization），並結合音訊進行綜合理解³。

具體操作流程：

1. **直接上傳：**支援上傳 MP4, MPEG, MOV 等格式。在 Google AI Studio 中，影片會被自動分解為 tokens。
2. **提示詞策略：**
 - **全面轉錄：** "Analyze this video. 1. Provide a verbatim transcript of the audio. 2. Describe the visual actions occurring at each key timestamp. 3. Extract any on-screen text (OCR). Output everything in Traditional Chinese."（分析這段影片。1. 提供語音逐字稿。2. 描述每個關鍵時間點的視覺動作。3. 提取螢幕上的文字。全部以繁體中文輸出。）
 - **細節檢索：** "At what timestamp does the red car appear, and what is the text on its license plate?"（紅車在什麼時間點出現？車牌上的文字是什麼？）Gemini 能精確回答這類視聽結合的問題²⁶。

比較優勢：相比於其他模型需要將影片「看」作一連串圖片，Gemini 能理解動作的連續性（例如：「這個人正在做什麼手勢？」），這對於體育分析或手語翻譯至關重要。

3.2 ChatGPT：幀採樣與代碼解釋器的結合

ChatGPT (GPT-4o) 處理視訊通常採取兩種路徑：一是透過 Vision 模型分析使用者上傳的關鍵幀，二是利用 Code Interpreter（現稱 Advanced Data Analysis）編寫 Python 腳本來處理視訊檔。

具體操作流程：

1. **上傳與分析**：上傳 MP4 後，GPT-4o 會抽取影片中的關鍵畫面進行視覺識別，同時利用 Whisper 轉錄聲音。
2. **利用 Python 進行精確提取**：對於需要精確時間軸或畫面 OCR 的任務，可以提示：「Write a Python script using OpenCV to extract frames every 5 seconds and recognize text in them. Then summarize the content.」（寫一個 Python 腳本，利用 OpenCV 每 5 秒截取一幀並識別其中的文字。然後總結內容。）這種方式雖然不是「原生」觀看，但對於特定任務（如提取簡報中的文字）非常精確²⁷。

3.3 Copilot：結合 Stream 與 Clipchamp 的視覺化摘要

Copilot 在處理 MP4 時，不僅僅是轉文字，更強調「導航」與「編輯」。

具體操作流程：

1. **Clipchamp 智慧剪輯**：在微軟的 Clipchamp 視訊編輯器中，Copilot 可以自動生成字幕，甚至根據使用者的文字指令（例如：「Create a video for social media based on this footage」）自動剪輯精彩片段。
2. **Stream 錄影深度搜尋**：對於企業內部的培訓影片或會議錄影，Copilot 能在 Stream 中建立「智慧章節」。使用者可以問：「影片中哪裡提到了『第三季預算』？」Copilot 會直接給出時間連結，點擊即跳轉²⁸。

3.4 Claude：圖像序列分析的替代方案

由於 Claude 缺乏原生的視訊上傳介面（截至 2025 年底主要透過 API 支援），一般使用者需要採取變通方法。

操作方式 (Workaround) :

將影片的關鍵影格（Keyframes）截圖，組合成 PDF 或直接作為多張圖片上傳，並配合音訊轉錄稿。提示詞：「Based on these screenshots and the transcript, reconstruct the narrative of the video in Traditional Chinese.」（根據這些截圖和逐字稿，用繁體中文重構影片的敘事。）這適合靜態較多的影片（如講座），不適合動作片³⁰。

3.5 Grok：發展中的視覺能力與 X 生態

Grok 的視覺模型（基於 Grok-2 Vision）可以理解圖片內容，對於短影片的理解能力正在整合中。其優勢在於結合 X 平台上的熱門影片趨勢，能解釋影片為何在社群上爆紅，或是識別影片中的迷因元素。

4. 五大 AI 的創意使用方式：基於人物誌（Persona）的深度場景應用

本節將跳脫基礎功能，根據網路上不同領域專家的實踐，整理出極具創意的進階使用方式。

4.1 程式設計師與全端開發者 (The Full-Stack Developer)

- **Claude 3.5 Sonnet + Artifacts 的「微型應用工廠」：**
 - 創意場景：開發者不再只是要求 AI 寫程式碼片段，而是利用 **Artifacts** 功能生成完整的互動式 Web App。例如，輸入指令：「製作一個 React 應用程式，模擬太陽系行星的重力軌道，並允許我調整各行星的質量。」Claude 會在側邊欄直接渲染出可執行的遊戲或模擬器，開發者可以即時與其互動並迭代修改（如「把背景改成深藍色」、「增加一顆彗星」）。這將 Claude 從「聊天機器人」變成了「即時軟體原型機」⁴。
 - **SVG 向量圖生成**：「畫一張解釋神經網絡反向傳播 (Backpropagation) 的 SVG 流程圖，使用扁平化設計風格。」Claude 能直接生成代碼並預覽圖形，這對於撰寫技術部落格的開發者極為實用³¹。
- **Grok 的「代碼吐槽大會」 (Code Roast) :**
 - 創意場景：開發者將一段寫得很爛或充滿 Bug 的代碼貼給 Grok，並開啟「Fun Mode」或「Spicy Mode」，指令：「Roast this code as if you are a senior engineer who hasn't had coffee yet.」（像一個還沒喝咖啡的資深工程師一樣吐槽這段代碼。）Grok 會用刻薄、幽默但技術精準的語言指出錯誤，這在開發社群 Reddit 和 X 上成為一種流行的減壓與學習方式¹⁷。

4.2 內容創作者與社群行銷專家 (The Creator & Marketer)

- **Grok 的「即時新聞劫持」 (Newsjacking) :**
 - 創意場景：這是 Grok 的殺手級應用。行銷人員利用其存取 X 平台即時數據的能力，指令：「現在 X 上關於『奧斯卡頒獎典禮』的熱門討論趨勢是什麼？有沒有什麼突發的迷因 (Meme)？請根據這些迷因，為我的披薩品牌生成 5 條帶有幽默感的行銷推文。」其他模型因為數據延遲，無法做到這種分秒必爭的熱點跟進³⁵。
- **Gemini 的「長影片內容煉金術」 (Content Repurposing) :**
 - 創意場景：創作者上傳一段 1 小時的 YouTube 直播錄影給 Gemini，指令：「請分析這段影片，並將其重製為：1. 一個包含 10 條推文的 Twitter Thread，總結精華。2. 一篇適合發布在 LinkedIn 的 800 字專業文章。3. 5 個適合製作 TikTok/Reels 短影音的腳本，並標註對應的原始影片時間戳記。」Gemini 的長上下文能力保證了它不會遺漏細節，且能精準定位高光時刻³⁸。
- **ChatGPT Advanced Voice Mode 的「模擬對手戲」：**
 - 創意場景：劇本創作者或 Podcaster 在通勤開車時，開啟 ChatGPT 的語音模式，設定 Prompt：「你現在是世界上最挑剔的電影製片人，我要向你推銷我的劇本，請你不斷反駁我、挑戰我的邏輯漏洞，直到我說服你為止。」這種高強度的語音角色扮演 (Roleplay) 能有效訓練口語表達與邏輯思維，且完全釋放雙手⁸。

4.3 學術研究者與學生 (The Academic Researcher)

- **Gemini 的「跨文檔綜合分析師」：**
 - 創意場景：研究生將 20 篇相關領域的 PDF 論文（總計數千頁）一次性上傳給 Gemini 1.5 Pro，指令：「這 20 篇論文中，關於『氣候變遷對台灣高山農業影響』的觀點有哪

些矛盾之處？請建立一個比較表格，列出每篇論文的作者、年份、主要論點以及與其他論文的衝突點。」這利用了 Gemini 驚人的「大海撈針」（Needle In A Haystack）檢索能力，大幅縮短文獻回顧的時間²。

- **Copilot 的「Excel 數據視覺化魔術師」：**

- **創意場景：**學生將實驗數據貼入 Excel，直接呼叫 Copilot：「請分析這些數據的趨勢，找出異常值（Outliers），並幫我生成三張不同類型的圖表（散佈圖、熱力圖）來展示變數之間的關係，並將圖表配色調整為學術期刊常用的灰階風格。」Copilot 能直接操作 Excel 物件，實現「對話即作圖」²⁵。

4.4 企業經理人與行政人員 (The Executive)

- **Copilot 的「會議時光機」：**

- **創意場景：**經理人因為行程衝突，遲到了 20 分鐘才進入 Teams 線上會議。他不需要打斷會議詢問進度，而是直接問側邊欄的 Copilot：「我錯過了什麼？剛剛誰發言最積極？有沒有提到我的名字或我的專案？」Copilot 會即時生成「至今為止的會議摘要」，讓經理人能無縫接軌，被稱為「非同步參會」的神器⁴³。

- **ChatGPT 的「跨文化商務禮儀顧問」：**

- **創意場景：**在準備接待來自中東或日本的客戶前，詢問 ChatGPT：「我即將接待一組沙烏地阿拉伯的商務代表團，請模擬他們的對話風格與禁忌，並幫我檢查這份繁體中文的歡迎致詞是否得體，有無觸犯文化禁忌的詞彙。」

4.5 視覺藝術家與前衛設計師 (The Artist)

- **Grok (Flux) 的「無過濾藝術探索」：**

- **創意場景：**雖然主流模型（DALL-E 3）對圖像生成的審查極為嚴格，但 Grok（整合 Flux 模型）在某些藝術風格（如超現實主義、諷刺漫畫、甚至是帶有恐怖美學的圖像）上提供了較大的創作自由度。藝術家利用它來生成一些可能被其他平台拒絕的「前衛（Edgy）」概念草圖，用於激發靈感（需注意，近期因 Deepfake 爭議，Grok 也收緊了對真實人物的生成限制，但對風格化藝術仍較開放）⁴⁵。

- **Claude Artifacts 的「色彩配置生成器」：**

- **創意場景：**設計師指令：「寫一個 React 組件，讓我上傳一張照片，然後自動提取其中的 5 個主色調，並生成對應的 Hex Code、RGB 值以及推薦的互補色搭配建議。」Claude 不僅給出代碼，還直接生成一個可用的色彩提取工具介面³²。

5. 比較分析與數據化評估

為了讓讀者更直觀地選擇適合的工具，以下透過表格呈現各模型在核心任務上的表現數據與特性比較。

5.1 英文語音轉繁體中文能力比較

特性比較	ChatGPT (Plus)	Gemini (1.5 Pro)	Claude (3.5 Sonnet)	Copilot (M365)	Grok (2/3)
語音識別核心	OpenAI Whisper (業界頂尖)	Google Universal Speech Model	無原生音訊 (需轉錄)	Microsoft Speech Services	xAI Audio (發展中)
操作便捷性	★★★ ★★ (App 直接錄音/上傳)	★★★ ★ (網頁上傳)	★★ (需轉錄文本)	★★★ ★ (Teams/Word整合)	★★ (主要靠文本)
長音訊處理	受限 (需分割檔案)	極強 (1M+ tokens, 數小時)	強 (但需先轉成文字檔)	強 (依賴會議錄影長度)	未知/受限
繁中翻譯品質	優 (需 Prompt 引導)	優 (流暢度高)	極優 (文學性與修辭最佳)	優 (商務風格為主)	普通 (偏口語/俚語)
最佳適用場景	短錄音、即時對話、學習	長演講、會議記錄、多檔案	文學翻譯、深度文本分析	企業內部會議、郵件	社群媒體、迷因翻譯

5.2 MP4 視訊轉文字能力比較

特性比較	Gemini (1.5 Pro)	ChatGPT (GPT-4o)	Copilot (M365)	Claude (3.5 Sonnet)
原生視訊理解	是 (Native Multimodal)	否 (抽幀 + 語音)	是 (結合 Stream 索引)	否 (需轉成截圖)
時間軸精確度	極高 (可精確到秒)	中 (依賴抽幀頻率)	高 (依賴字幕軌)	無
畫面細節識別	極強 (動作、微	強 (主要畫面)	中 (主要依賴語)	強 (單張圖片分

	小文字)		音)	析)
處理速度	快 (批次處理)	中 (需上傳處理)	快 (已預處理)	慢 (需手動截圖)
最佳適用場景	全能視訊分析、體育、食譜	數學解題、短片分析	會議回顧、教育訓練	靜態簡報分析

6. 深度學習模式與潛在風險分析

在享受這些創意應用與高效工作流的同時，使用者必須理解這些模型在「深度學習模式」下的運作機制及其帶來的潛在風險。

6.1 幻覺（Hallucination）與事實查核機制

- **機制：**所有 LLM/LMM 本質上都是機率預測模型，它們是在「預測下一個字」而非「陳述事實」。
- **風險：**在視訊分析中，如果畫面模糊，Gemini 或 ChatGPT 可能會「腦補」出不存在的細節（例如看錯車牌號碼）。
- **緩解策略：**
 - **Gemini:** 利用「Grounding with Google Search」功能，讓模型在生成答案後主動搜尋 Google 確認事實⁴⁷。
 - **ChatGPT:** 要求模型引用來源，或使用 Web Browsing 功能進行雙重驗證。

6.2 上下文視窗（Context Window）的戰略意義

- **機制：**上下文視窗決定了模型能「記住」多少資訊。
- **趨勢：**Gemini 的 100 萬 token 代表了一種「暴力美學」的深度學習應用模式。使用者不再需要像過去那樣精心修剪 Prompt 或使用 RAG（檢索增強生成）來切分檔案，而是可以直接將整個程式碼庫、整本法律法條或整部電影丟進去。這改變了「微調（Fine-tuning）」的需求，許多過去需要微調才能做到的任務，現在透過長上下文的 In-context Learning 就能達成³。

6.3 隱私、安全性與倫理爭議

- **Copilot:** 企業版 Copilot 強調「商業數據保護」（Commercial Data Protection），保證使用者的數據不會被用來訓練模型，這對於處理包含機密資訊的會議錄音至關重要。
- **Grok 的倫理挑戰：**報告必須指出，Grok 由於其寬鬆的內容審查政策，曾被濫用於生成未經同意的 Deepfake 圖像（如 Nudification），這引發了全球監管機構的關注。使用者在使用 Grok 進行圖像生成時，必須嚴格遵守法律與道德規範，避免觸犯法律風險（如生成真實人物

的虛假圖像) ⁴⁸。

7. 結論與操作建議

綜合以上分析，針對您的兩大核心需求，本報告提出以下結論性的操作建議：

1. 對於「英文語音轉繁體中文」：
 - 首選：若檔案較大或追求極致的理解深度，請使用 **Google Gemini 1.5 Pro**。其原生長音訊支援與時間戳記功能是目前的最佳解決方案。
 - 次選：若追求手機操作的便捷性與即時性，**ChatGPT App** 是不二之選。
2. 對於「傳送 MP4 檔並將內容轉成文字」：
 - 唯一推薦：**Google Gemini 1.5 Pro**。目前市場上只有它能真正做到「看懂」整部影片的視覺與聽覺流，並進行結構化的繁體中文輸出。
 - 企業場景：若是在公司內部處理會議錄影，則應優先使用 **Microsoft Copilot** 以確保資安與整合性。
3. 對於「創意使用方式」：
 - 探索者：請務必嘗試 **Claude 的 Artifacts** 來製作小遊戲或視覺化圖表，這會徹底改變您對 AI 「只能聊天」的刻板印象。
 - 觀察者：利用 **Grok** 來觀察即時的社群輿論與迷因趨勢，這是其他模型無法提供的視角。

隨著 2026 年的到來，AI 的使用已經從單一的問答工具進化為多模態的協作夥伴。掌握這些不同模型的特性與創意用法，將是未來數位生產力的關鍵分水嶺。

Works cited

1. ChatGPT vs Gemini vs Copilot vs Claude vs Perplexity vs Grok | AI Assistants - Gmelius, accessed January 10, 2026,
<https://gmelius.com/blog/best-ai-assistants-comparison>
2. The Best AI in October 2025? We Compared ChatGPT, Claude, Grok, Gemini & Others, accessed January 10, 2026,
<https://felloai.com/the-best-ai-in-october-2025-we-compared-chatgpt-claude-grok-gemini-others/>
3. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context - Googleapis.com, accessed January 10, 2026,
https://storage.googleapis.com/deepmind-media/gemini/gemini_v1_5_report.pdf
4. Everything I built with Claude Artifacts this week - Simon Willison's Weblog, accessed January 10, 2026,
<https://simonwillison.net/2024/Oct/21/claude-artifacts/>
5. The Ultimate AI Showdown September 2025: A deep dive into ChatGPT vs. Copilot vs. Gemini and when to use Claude, Grok or Perplexity instead. The Hidden Strengths and Weaknesses of Every Major AI - Reddit, accessed January 10, 2026,
https://www.reddit.com/r/ThinkingDeeplyAI/comments/1njyr71/the_ultimate_ai_sh

[owdown september 2025 a deep/](#)

6. The best use of Twitter is to use Grok to expose Elon Musk's paper-thin ego, make fun of him, and post the results. - Reddit, accessed January 10, 2026, https://www.reddit.com/r/TrueUnpopularOpinion/comments/1iush4w/the_best_use_of_twitter_is_to_use_grok_to_expose/
7. Can ChatGPT Translate Languages? A Detailed Guide for 2026 - Upwork, accessed January 10, 2026, <https://www.upwork.com/resources/chatgpt-for-translation>
8. Voice - fun ideas? : r/ChatGPT - Reddit, accessed January 10, 2026, https://www.reddit.com/r/ChatGPT/comments/1775nc6/voice_fun_ideas/
9. Video understanding | Gemini API - Google AI for Developers, accessed January 10, 2026, <https://ai.google.dev/gemini-api/docs/video-understanding>
10. Unlocking Multimodal Video Transcription with Gemini — Part 2:  Setup | by Laurent Picard | Google Cloud - Medium, accessed January 10, 2026, <https://medium.com/google-cloud/unlocking-multimodal-video-transcription-with-gemini-part2-43c491a0c4f1>
11. What kinds of documents can I upload to Claude?, accessed January 10, 2026, <https://support.claude.com/en/articles/8241126-what-kinds-of-documents-can-i-upload-to-claude>
12. Claude 3.5 Sonnet ranks #1 in the new creative story-writing benchmark. Claude 3.5 Haiku is #2 : r/ClaudeAI - Reddit, accessed January 10, 2026, https://www.reddit.com/r/ClaudeAI/comments/1hv3ido/clause_35_sonnet_ranks_1_in_the_new_creative/
13. Create a video with the Microsoft 365 Copilot app, accessed January 10, 2026, <https://support.microsoft.com/en-us/topic/create-a-video-with-the-microsoft-365-copilot-app-4edd41f6-a7ad-47d5-9a55-3fd25622c9f8>
14. What is Microsoft 365 Copilot?, accessed January 10, 2026, <https://learn.microsoft.com/en-us/copilot/microsoft-365/microsoft-365-copilot-overview>
15. How Grok AI Is Trending With Its New Features: Image-to-Video Generation, Voice Mode & Workspace Projects - Digidrome, accessed January 10, 2026, <https://www.digidrome.com/blog/grok-ai-image-to-video-voice-mode-workspace-projects>
16. Grok multimodal capabilities: using images, audio, and video in AI workflows for 2025., accessed January 10, 2026, <https://www.datastudios.org/post/grok-multimodal-capabilities-using-images-audio-and-video-in-ai-workflows-for-2025>
17. Okay, Maybe Grok-2 is Decent. : r/LocalLLaMA - Reddit, accessed January 10, 2026, https://www.reddit.com/r/LocalLLaMA/comments/1etl028/okay_maybe_grok2_is_decent/
18. Can You Upload Audio Files to ChatGPT? - VOMO.ai, accessed January 10, 2026, <https://vomo.ai/blog/can-you-upload-audio-files-to-chatgpt>
19. Can ChatGPT transcribe audio? Everything you need to know - Techpoint Africa, accessed January 10, 2026,

- <https://techpoint.africa/guide/can-chatgpt-transcribe-audio/>
- 20. 35 ChatGPT Prompts for High-Quality Translation [2026] - Pairaphrase, accessed January 10, 2026,
<https://www.pairaphrase.com/blog/chatgpt-prompts-translation>
 - 21. Audio understanding | Gemini API - Google AI for Developers, accessed January 10, 2026, <https://ai.google.dev/gemini-api/docs/audio>
 - 22. Upload & analyze files in Gemini Apps - Computer - Google Help, accessed January 10, 2026,
<https://support.google.com/gemini/answer/14903178?hl=en&co=GENIE.Platform%3DDesktop>
 - 23. My custom instructions prompt for being a translator into Mandarin Chinese. : r/ChatGPT, accessed January 10, 2026,
https://www.reddit.com/r/ChatGPT/comments/1l97sdu/my_custom_instructions_prompt_for_being_a/
 - 24. Standard versus priority access to features in Microsoft 365 Copilot Chat, accessed January 10, 2026,
<https://support.microsoft.com/en-us/topic/standard-versus-priority-access-to-features-in-microsoft-365-copilot-chat-12c8d9f8-db32-4f99-8ebe-d8d85879137f>
 - 25. How I Use Copilot with Microsoft Office to Be Super Productive and Save 2-3 Hours Every Day - Reddit, accessed January 10, 2026,
https://www.reddit.com/r/GPTHackers/comments/1i0br61/how_i_use_copilot_with_microsoft_office_to_be/
 - 26. Has Anyone Tried Gemini 1.5 Pro with the December Video? Seeking Feedback! - Reddit, accessed January 10, 2026,
https://www.reddit.com/r/Bard/comments/1b2tml5/has_anyone_tried_gemini_15_pro_with_the_december/
 - 27. Unlocking Multimodal Video Transcription with Gemini — Part 5: Finalization | by Laurent Picard | Google Cloud - Medium, accessed January 10, 2026,
<https://medium.com/@PicardParis/unlocking-multimodal-video-transcription-with-gemini-part5-488b357b53b1>
 - 28. Frequently asked questions about Copilot in the Clipchamp video player - Microsoft Support, accessed January 10, 2026,
<https://support.microsoft.com/en-us/topic/frequently-asked-questions-about-copilot-in-the-clipchamp-video-player-20ef6f0e-10f8-47aa-8bb7-697db7445fae>
 - 29. Ask questions & get summaries of any video with Microsoft Copilot in the Clipchamp player, accessed January 10, 2026,
<https://support.microsoft.com/en-us/office/ask-questions-get-summaries-of-any-video-with-microsoft-copilot-in-the-clipchamp-player-0b531ea9-2d9d-4830-97e4-2c1b2b8ca31d>
 - 30. Claude: Using images, audio, and video in practical workflows - Data Studios, accessed January 10, 2026,
<https://www.datastudios.org/post/clause-using-images-audio-and-video-in-practical-workflows>
 - 31. What is Artifacts in Claude 3.5 Sonnet? - Xuyun Zeng, accessed January 10, 2026,
<https://xyzcreativeworks.com/what-is-artifacts-in-claude-3-5-sonnet/>

32. Use artifacts to visualize and create AI apps, without ever writing a line of code, accessed January 10, 2026,
<https://support.claude.com/en/articles/11649427-use-artifacts-to-visualize-and-create-ai-apps-without-ever-writing-a-line-of-code>
33. Interactive Fiction Game Design with Claude Artifacts - TCEA Blog, accessed January 10, 2026,
<https://blog.tcea.org/interactive-fiction-game-design-with-claude-artifacts/>
34. Claude Artifacts - Build Interactive Apps and Dashboards - LLMindset.co.uk, accessed January 10, 2026,
<https://lilmindset.co.uk/posts/2024/10/claude-amazing-artifacts/>
35. Top 15+ Grok AI Prompts & How to Use Them - ClickUp, accessed January 10, 2026, <https://clickup.com/blog/grok-ai-prompts/>
36. 200+ Best Grok Prompts for Every Use Case in 2026, accessed January 10, 2026, <https://chatsmith.io/blogs/prompt/grok-prompts-00107>
37. How to Use Grok: A Simple Guide for Beginners - AI Tools - God of Prompt, accessed January 10, 2026, <https://www.godofprompt.ai/blog/how-to-use-grok>
38. 20 Ways Gemini can watch and analyze videos for you : r/Bard - Reddit, accessed January 10, 2026,
https://www.reddit.com/r/Bard/comments/1pl5l21/20_ways_gemini_can_watch_and_analyze_videos_for/
39. LMMS & Google's Gemini 1.5 Pro Watching Television News: Converting Videos To Text For Universal RAG & Summarization - The GDELT Project, accessed January 10, 2026,
<https://blog.gdeltproject.org/lmms-googles-gemini-1-5-pro-watching-television-news-converting-videos-to-text-for-universal-rag-summarization/>
40. Voice Mode Productivity Hack : r/ChatGPTPro - Reddit, accessed January 10, 2026,
https://www.reddit.com/r/ChatGPTPro/comments/1g3t1rj/voice_mode_productivity_hack/
41. 10 Creative Ways To Leverage ChatGPT's Advanced Voice Mode for Work and Play, accessed January 10, 2026,
<https://hotelemarketer.com/2024/09/25/10-creative-ways-to-leverage-chatgpts-advanced-voice-mode-for-work-and-play/>
42. Everyday Productivity Hacks for Work and Home | Microsoft 365, accessed January 10, 2026,
<https://www.microsoft.com/en-us/microsoft-365-life-hacks/everyday-ai/time-saving-tips/everyday-productivity-hacks-for-work-and-home-with-microsoft-365>
43. Unlock your productivity: Here are our Top 10 tips for using Microsoft 365 Copilot every day, accessed January 10, 2026,
<https://www.microsoft.com/insidetrack/blog/unlock-your-productivity-here-are-our-top-10-tips-for-using-microsoft-365-copilot-every-day/>
44. Microsoft 365 Copilot Power User Tips - YouTube, accessed January 10, 2026,
<https://www.youtube.com/watch?v=MuocazjjZmc>
45. 10 Best Grok Imagine Prompts: Golden Formula + Examples+ - HitPaw Edimakor, accessed January 10, 2026,

<https://edimakor.hitpaw.com/video-editing-tips/best-grok-imagine-prompts.html>

46. Grok can anyone send me some 18+ prompts that work please - Reddit, accessed January 10, 2026,
https://www.reddit.com/r/grok/comments/1putl2m/grok_can_anyone_send_me_s/ome 18 prompts that work/
47. People think ChatGPT, Claude, Gemini, Grok are just "different brands" of the same tool. : r/ChatGPTPromptGenius - Reddit, accessed January 10, 2026,
https://www.reddit.com/r/ChatGPTPromptGenius/comments/1p29ihh/people_think_chatgpt_claude_gemini_grok_are_just/
48. How Grok pushed deepfake “nudification” mainstream - The Hindu, accessed January 10, 2026,
<https://www.thehindu.com/sci-tech/technology/how-grok-pushes-deepfake-nudification-mainstream/article70485551.ece>
49. Musk's AI chatbot faces global backlash over sexualized images of women and children, accessed January 10, 2026,
https://www.shootonline.com/shoot_column/musks-ai-chatbot-faces-global-backlash-over-sexualized-images-of-women-and-children/
50. Grok AI still being used to digitally undress women and children despite suspension pledge, accessed January 10, 2026,
<https://www.theguardian.com/technology/2026/jan/05/elon-musk-grok-ai-digitally-undress-images-of-women-children>