

OH YEAH DITTO, I CHOOSE YOU.

OH YEAH
百變怪
就決定是你了！



一、前言

本專題聚焦於當代生成式 AI 與影像合成技術對 EKYC (電子身分驗證) 流程所帶來的挑戰。隨著靜態人臉合成、臉部動畫與證件影像偽造技術日益成熟，傳統僅以相似度比對或單一模態檢測的驗證流程可能出現漏洞。本計畫希望在主辦提供的受控 **SANDBOX** 環境中，建立一套模組化且可量化的攻測框架，系統化地重現常見攻擊向量，衡量不同情境下系統的反應與弱點，並以清晰的數據支持可行的防護建議。整體工作將嚴守倫理與法規，所有測試僅使用合成或授權資料，並配合負責任揭露流程以確保研究過程的透明與安全。

二、創意描述

本專題的核心創意在於把複雜的攻擊鏈條拆解成彼此可比較、可重複的單元，藉由一致的測試矩陣來量化不同攻擊類型在 **EKYC** 流程中造成的影響。與其描述一個高度自動化的攻擊代理，我們強調「模組化思維」：將人臉合成、臉部動畫、證件偽造與影像前處理視為可參數化的模組，並在同一套測試規範下，依序或交互驗證系統的防護能力。透過這樣的設計，我們能夠把抽象的威脅具體化為可測量的實驗條件，讓每一次測試都能回溯其來源與設定，進而為後續的防護策略提供有力的實證依據。

三、系統功能簡介

本系統提供一組用於沙盒測試的工具與流程，主要負責樣本製作、批次測試與結果彙整。樣本製作部分包含靜態臉像的選取或合成、簡短的臉部動畫生成，以及證件影像的欄位修改與合成，所有樣本會依測試矩陣被標記並分類。測試執行時，系統會以既定格式將樣本逐一提交至受測的 **EKYC** 環境，紀錄系統回傳的相似度分數、活體檢測回應、**OCR** 結果與相關錯誤碼，並保留每一筆測試的原始輸入與系統回應以供事後查核。分析模組則負責對回收資料進行清洗、統計與視覺化，將成功率、分數分佈、常見失敗情境等資訊輸出為可閱讀的報告，最後彙整為對應的弱點類型與建議修補方向。

四、系統特色

此框架的特色在於它的可重現性與比較性。每一個攻擊模組均以明確的參數定義，例如合成影像的品質等級、動畫的長度與互動性、證件欄位被修改的類型，以及影像經過不同壓縮或噪聲處理後的情況，這些參數構成了測試矩陣，使得不同實驗之間的結果具可比性。系統同時強調資料治理，所有測試資料來源、處理方式與結果使用限制都被記錄與管理，以保障倫理與法遵要求。最後，結果呈現側重於實務可操作的輸出，將技術性的發現轉換成企業或開發團隊能理解且能直接採納的建議。

五、系統開發工具與技術

在技術選擇上，本計畫採用現有成熟工具與可公開取得的資源以確保穩定性與可驗證性。樣本製作會利用公開合成資料或現成模型生成的影像作為基底，影像處理與簡單合成工作以常用影像處理函式庫完成，證件影像的 OCR 與欄位比對則採用可被驗證的 OCR 工具來讀取與分析。測試執行與資料擷取以標準化的資料格式進行，自動化腳本負責批次上傳與回收結果，而分析採用常見的統計與視覺化工具來呈現指標。整體技術堆疊以透明、可重現為原則，並在系統設計上留下足夠的紀錄與日誌，方便審查與延伸研究。

六、系統使用對象

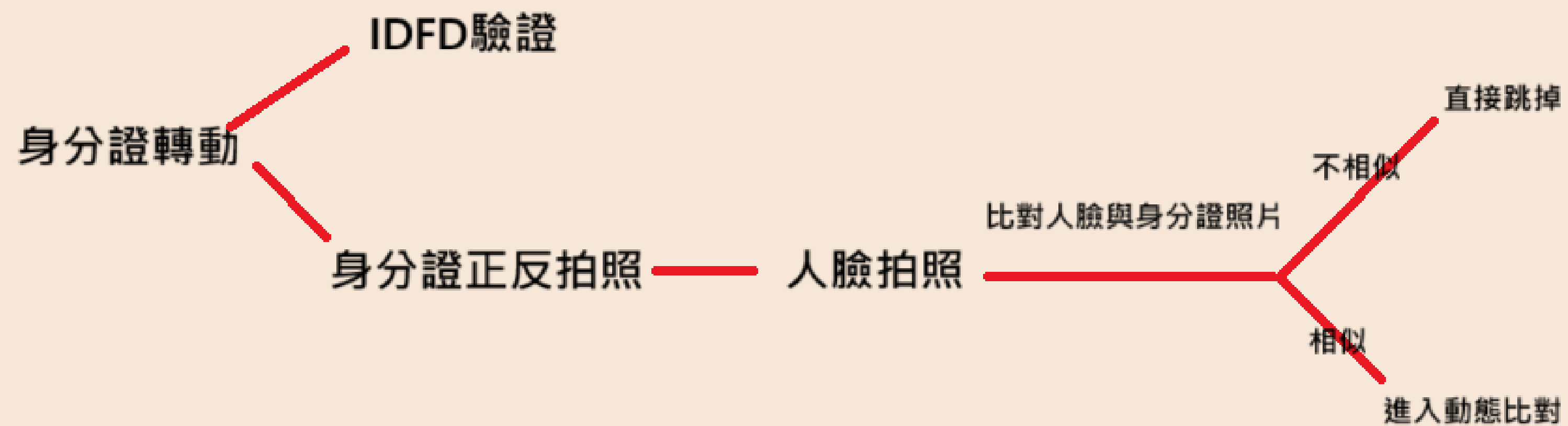
本專題的成果面向多元的受眾。對於比賽主辦與評審，我們提供一套能展示攻測方法論與量化結果的技術方案，說明如何在受控環境下驗證EKYC的多模態防護；對於企業的藍隊或資安工程師，我們的報告能作為弱點驗證與修補優先順序的參考，幫助決策者在資源有限的情況下合理分配防護投入；對於學術研究者與教學單位，這套框架能作為研究或實驗教學的樣板，讓後續研究在相同標準下進行比較與擴展。

七、系統使用環境

所有的測試工作將在主辦方提供的隔離 **SANDBOX** 環境中進行，以避免對外部系統或真實使用者造成風險。系統部署上，我們將使用可移植的執行單元與標準化的資料接口，確保在不同測試環境間能快速復現實驗結果。資料來源會事先經過審查，僅採用合成或授權的測試樣本，並落實存取權限與使用紀錄。測試期間的所有操作、異常與回收資料都會被完整記錄，以便於結果驗證與負責任揭露。

八、結語

本專題提出一個結合模組化思維與量化指標的 **EKYC** 攻測框架，旨在以系統化、可重現的方式揭示生成式 **AI** 所帶來的驗證風險，並將技術發現轉化為具體的防護建議。我們希望透過清晰的實驗設計、嚴謹的資料治理與實務導向的分析，為業界與學界提供一套既能說服評審、又能在實務上產生效益的成果。若主辦與評審認同此願景，我們將依授權與倫理規範，在 **SANDBOX** 中執行實驗並回報具體的量化結果與改善建議。



sensity-ai/dot

The Deepfake Offensive Toolkit



7

Contributors



22

Issues



4k

Stars



465

Forks

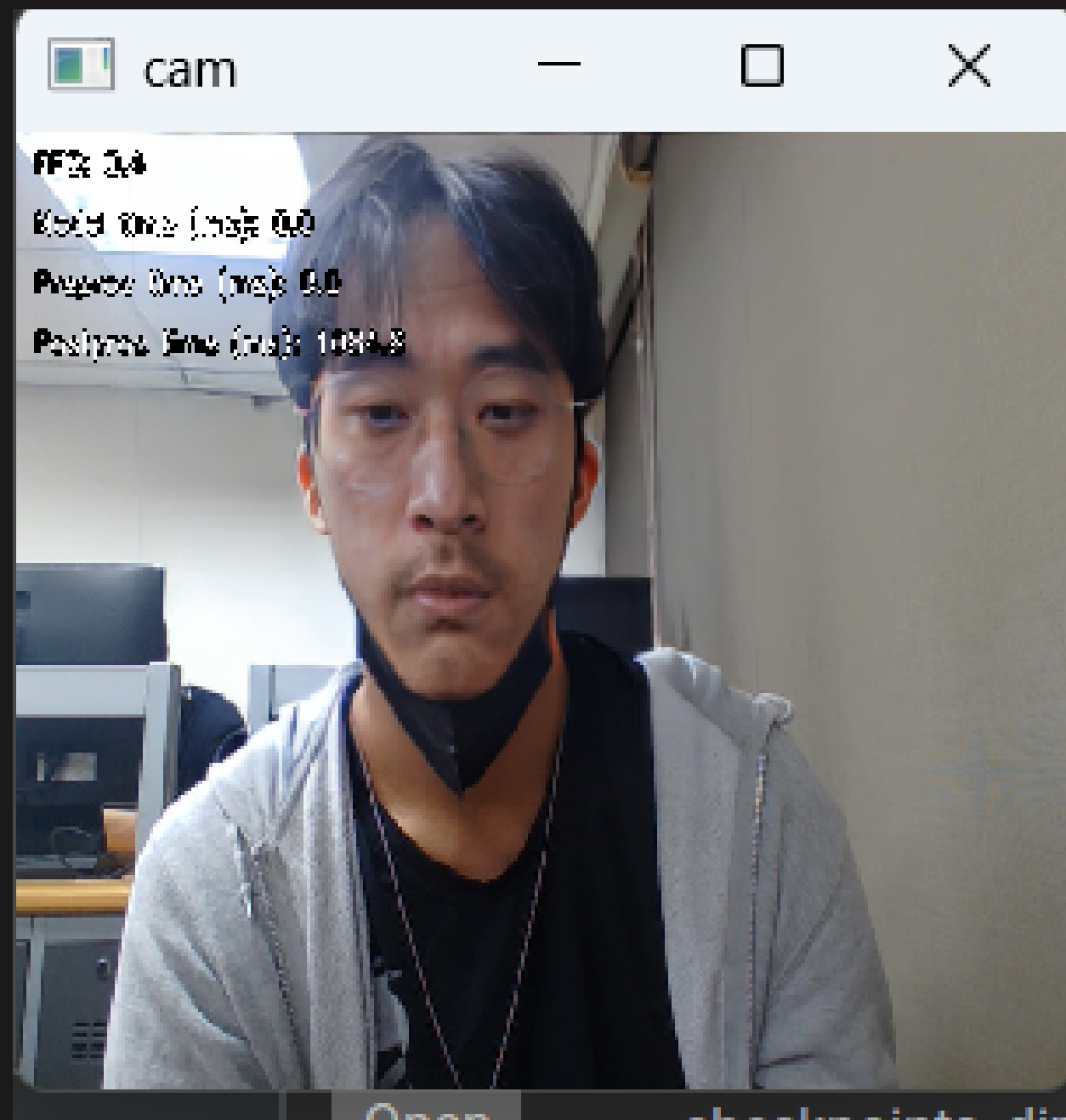


sensity-ai/dot: The Deepfake Offensive Toolkit

The Deepfake Offensive Toolkit. Contribute to sensity-ai/dot development by creating an account on GitHub.



GitHub



THANKS FOR LISTENING