資料3-3

AI開発者の想定する脅威・対策

2025年10月17日 三井物産セキュアディレクション株式会社

アジェンダ



- AI開発者の位置付け
- AI開発者が被害者となる脅威の一覧
- AI開発者が実施者となる対策の一覧
- 画像識別AIに対する攻撃手法や対策のVLMへの転用可能性



AI開発者の位置付け

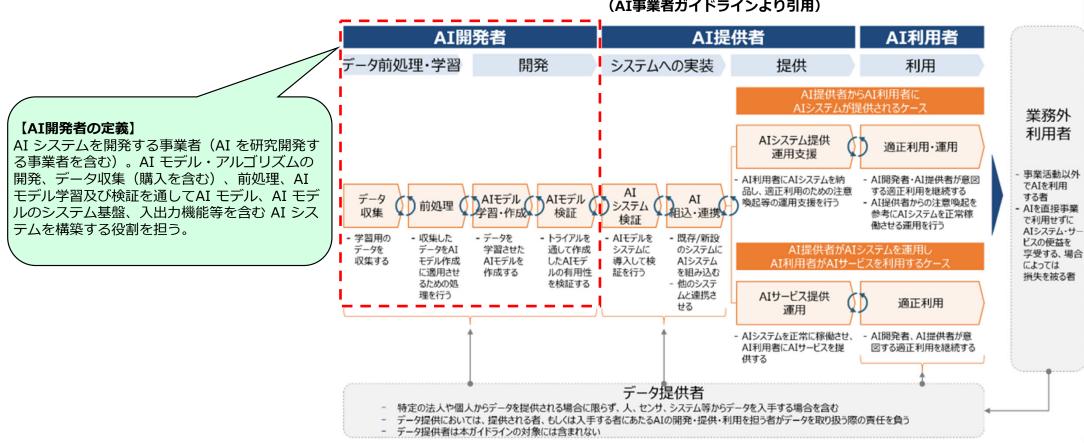
© 2025 Mitsui Bussan Secure Directions, Inc. All Rights Reserved.

3

AI開発者の位置付け

本資料が対象とする、AI開発者の定義およびライフサイクルとの対応関係は下図の通り。

一般的な AI 活用の流れにおける主体の対応 (AI事業者ガイドラインより引用)



本資料の対象範囲

- 本資料ではAI開発者が、「被害者となる脅威」および「実施者となる対策」を扱う。
 - 「被害者となる脅威」は、**「学習データ関連情報の窃取」**および「モデルの窃取」を目的とした攻撃を指す。
 - 「実施者となる対策」は、「データ前処理・学習」および「開発」のフェーズにおいて実施すべき対策を指す。

※必ずしも脅威と対策が一対一に対応していないことに留意が必要(開発者が被害者となる脅威に対して、開発者が対策を実施するとは限らない)



AI開発者が被害者となる脅威の一覧

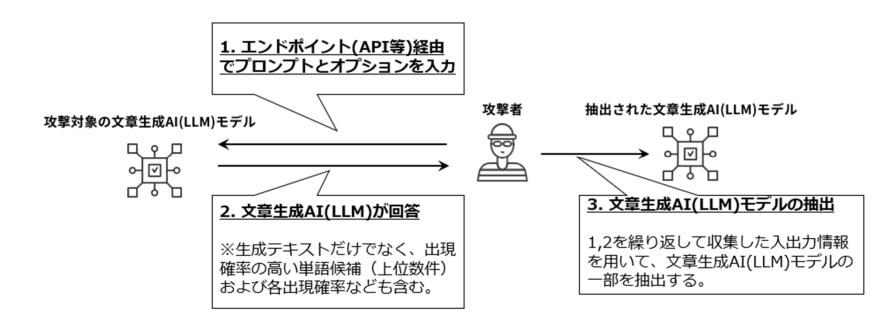
AI開発者が被害者となる脅威の一覧表(第一回分科会資料を加筆)

※「AISIの分類」は、AISI「AIシステムに対する既知の攻撃と影響」より引用。

AISIの分類	脅威	脅威の具体例	攻撃経路	侵害属性	対象となるAI
A:モデル抽出攻撃	① モデル抽出攻撃・ AIシステムのエンドポイント(API等)に対して繰り返しアクセスし、AIモデルへの入出力情報を観察することで、同程度の性能を持つ複製モデルを作成する攻撃手法。・ 競争上の地位低下や、モデルに含まれる機密情報の窃取などにつながる。	L-1-1 : LLMの入出力情 報を観察してAIモデルを 抽出する攻撃	AI 利用者が入力 するプロンプト	機密性	LLM
		I-1-1:画像識別AIを模 倣したAIを作成する攻撃	AI 利用者が入力 する画像	機密性	
B:学習データ情報 収集攻撃	• ハ ハラ 翼 テニタが田 TB O ろ ソ 製 主 注	学翌デークを容取するで	AI 利用者が入力 する画像	機密性	両偽⇒PU∧T
	③ メンバーシップ推論攻撃 ・ 入力画像に対する画像識別AIの応答を観察することで、入力した画像がAIの学習データに含まれているかどうかを推測する攻撃手法。 ・ 学習が不十分な画像識別AIは、学習した画像と学習していない画像では、明確に異なる応答を返すため、その性質を悪用するもの。	I-3-1:画像識別AIの挙 動を観察することで学習 データを推測する攻撃	AI 利用者が入力 する画像	機密性	画像識別AI
H: プロンプトイン ジェクション攻撃	①直接的プロンプトインジェクション攻撃 (学習データの窃取型)AIモデルの学習データを窃取する攻撃手法。攻撃者が細工したプロンプトをAIモデルに入力することで、AIモデルの回答に学習データの一部が出力される。	L-8-1 : プロンプトを細工 してLLMの学習データを 窃取する攻撃	AI 利用者が入力 するプロンプト	機密性	LLM

L-1-1:LLMの入出力情報を観察してAIモデルを抽出する攻撃

- 攻撃手法: LLMシステムのエンドポイント(API等)に対して繰り返しアクセスして、LLMが出力するトークンの候補(上位数件)と各出現確率を観察・分析することで、AIモデルを抽出する攻撃。
- **脅威程度**:本攻撃は、 AI利用者と同様に外部からエンドポイント(API等)にアクセスしてプロンプトを入力するだけで成立するが、 AIモデルの抽出には非常に多くの試行回数を要するため、攻撃実施のハードルは 比較的高い。



出典: https://arxiv.org/pdf/2403.06634

L-8-1:プロンプトを細工してLLMの学習データを窃取する攻撃

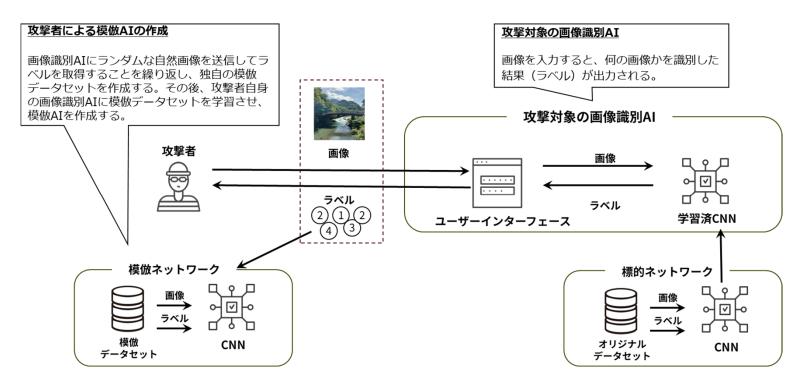
- **攻撃手法**: LLMシステムに悪意のあるプロンプトを入力することで、学習データに含まれる機密情報を窃取 する攻撃。
- **脅威程度**:本攻撃は、AI利用者と同様に外部からプロンプトを入力するだけで攻撃が成立するため、攻撃実施のハードルは低い。



出典: https://arxiv.org/abs/2012.07805

I-1-1:画像識別AIを模倣したAIを作成する攻撃

- 攻撃手法:画像識別AI(CNN)に多数の画像を入力し、入力した画像とその分類結果を紐付けることで、AI モデルを抽出する攻撃。
- **脅威程度**:本攻撃は、 AI利用者と同様に外部からエンドポイント(API等)にアクセスして画像を入力するだけで成立するが、 AIモデルの抽出には多くの試行回数を要するため、攻撃実施のハードルは比較的高い。



出典: https://ojs.aaai.org/index.php/AAAI/article/view/32734/34889



I-2-1:画像識別AIの内部構造を観察することで学習データを窃取する攻撃

- 攻撃手法:画像識別AI(※Softmax、MLP、DAE)にランダムに初期化した画像を入力し、応答として得られるAIモデルの内部情報(勾配情報)を基に、学習データ(画像)を窃取する攻撃。
- **脅威程度**:本攻撃は、画像識別AIの勾配情報を基に推測する仕組み上、攻撃対象の画像識別AIに直接アクセスする必要があるため、攻撃実施のハードルは比較的高い。また、実際に学習に使われた画像データと完全に同一のデータを復元することは難しい。

画像識別AIの内部情報(勾配情報)から復元されたデータのイメージ

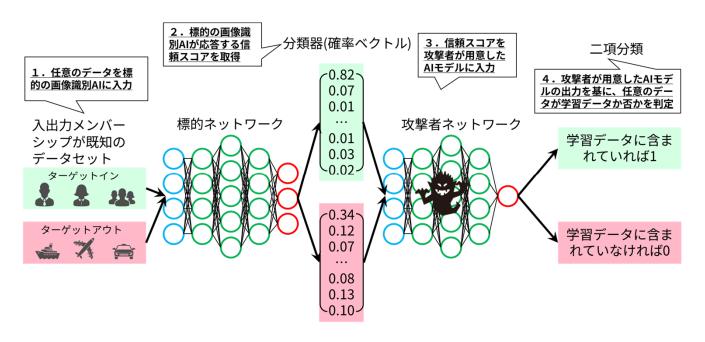


※ソフトマックス回帰、多層パーセプトロン、 デノイジングオートエンコーダーの略

出典: https://dl.acm.org/doi/10.1145/2810103.2813677

I-3-1:画像識別AIの挙動を観察することで学習データを推測する攻撃

- 攻撃手法:画像識別AI(CNN)にデータを入力し、これに対するAIの応答(分類結果=ラベル、信頼スコア)を分析することで、入力したデータがAIの学習データに含まれるか否か(メンバーシップか否か)を推論する攻撃。
- **脅威程度**:本攻撃は、 AI利用者と同様に外部から画像を入力するだけで攻撃が成立するため、攻撃実施の ハードルは低い。ただし画像1枚ごとに分類結果を収集し精密に分析する必要があるため、効率的に大量の 機微情報を収集するには限界がある。



出典: https://ojs.aaai.org/index.php/AAAI/article/view/32734/34889



AI開発者が実施者となる対策の一覧

AI開発者が実施者となる対策の一覧(LLM特有のもの)

AIシステムに対する対策のうち、**LLM特有なものについては、**具体例として以下が挙げられる。

対策の具体例	対応する脅威の具体例		LLM/ 画像識別AI
 ファインチューニングデータの精査 ファインチューニングデータを精査し、細工されたデータが可能な限り含まれないようにする対策手法。 AIのファインチューニングに使用するデータを精査し、安全性を低下させる可能性のあるデータ(悪意のある指示や嘘を含む情報等)を可能な限り除外する。 	LLMの挙動を操作する攻撃 L-4-1:ファインチューニングデータを細工して	開発	LLM
LLM自体の頑健性向上 ・ 悪意あるプロンプトにLLMが応答しないように頑健性を高める対策手法。 ・ 開発段階において、攻撃者からの悪意のあるプロンプト(不法行為や差別、偽情報、機密情報等の文章を作成させることを意図した指示)に応答しないよう、LLMの頑健性を高める。	を窃取する攻撃 L-8-1:プロンプトを細工してLLMの学習デー タを窃取する攻撃		LLM

$M^{\dagger}B_{\dagger}S^{\dagger}D_{*}$

ファインチューニングデータの精査

- · 想定脅威:
 - データポイズニング攻撃
 - DoS攻撃
- **対策内容**:ファインチューニングデータをあらかじめ精査し、可能な限り不適切なデータを除外する。不適切なデータの具体例としては、以下が挙げられる。
 - 内容が極端に偏っていることや、事実と異なるもの
 - 明らかに悪意のある指示に誘導するような構造やプロンプトインジェクション攻撃の兆候が見られるもの
 - 差別的、暴力的、またはその他倫理的に問題のある表現を含むもの
 - 過剰に繰り返された構文や、不自然に同一出力を誘導するもの

第二回分科会にて説明済み

- 想定脅威:
 - プロンプトインジェクション攻撃(全類型)
 - ※プロンプトインジェクション攻撃は、直接プロンプトインジェクション攻撃および間接プロンプトインジェクション攻撃の両方を指す。
- 対策内容: AIモデル自体の頑健性を高めることで、LLMが悪意あるプロンプトに従わないようにする。頑健性向上の具体的手法としては、以下が挙げられる。
 - **安全性アラインメント**: LLMが不法行為や差別的表現、偽情報、機密情報の漏えいなどを引き起こすような悪意あるプロンプトに応答しないように、人間のフィードバックに基づく強化学習等を用いて、安全基準を事後学習させる。
 - **指示の階層化**: LLMが従うべき指示の優先度を定義し、高優先度(例:システムプロンプト)を常に優先的に処理するようにLLMの内部ロジックを調整する。

※頑健性向上は、一般には信頼性等の観点からも実施すべき対策になるが、本資料ではセキュリティの観点からの対策として記載している



AI開発者が実施者となる対策の一覧(LLMと画像識別AIで共通のもの)

AIシステムに対する対策のうち、原理上、LLMと画像識別AIで共通し得るものについては、具体例として以下が挙げられる。

対策の具体例	対応する脅威の具体例	対策段階	LLM/ 画像識別AI
 学習データの精査 ・ 収集した学習データを精査し、可能な限り機密情報が含まれないようにする対策手法。 ・ データ前処理・学習工程において、収集したデータを精査し、外部に出力することを意図しない機密情報が含まれている場合、これらを可能な限り除外する。 	L-8-1:プロンプトを細工してLLMの学習データを窃取する攻撃 I-2-1:画像識別AIの内部構造を観察する ことで学習データを窃取する攻撃 I-3-1:画像識別AIの挙動を観察することで 学習データを推測する攻撃	ナークBIIが t中・ラジ	LLM 画像識別AI
信頼できる事前学習データの使用 ・ 不正行為を意図した学習データの利用を防止する対策手法。 ・ AIの学習データの出所等が公開されている場合はそれらを参考に入手する事前学習データの安全性を確認する	L-3-1:学習データを細工してLLMの挙動を 操作する攻撃 I-5-1:学習データを細工して画像識別AIに バックドアを設置する攻撃	ナークロ W t甲・	LLM 画像識別AI
 ファインチューニングデータによる補正 ファインチューニングを行うことでLLMおよび画像識別AIが悪意のある入力データに応答しないようにする対策手法 細工されたデータを学習したAIに対し、安全性が確認されたデータでファインチューニングを行うことでバックドアを無効化し、AIが正しく応答できるようにする。 	I-5-1:学習データを細工して画像識別AIに バックドアを設置する攻撃	開発	LLM 画像識別AI

- 想定脅威:
 - プロンプトインジェクション攻撃(学習データの窃取型)
 - ・ モデル反転攻撃
 - メンバーシップ推論攻撃
 - ※プロンプトインジェクション攻撃は、直接プロンプトインジェクション攻撃および間接プロンプトインジェクション攻撃の両方を指す。
- **対策内容**:学習用に収集したデータを精査し、外部に出力することを意図しない機密情報を可能 な限り除外する。対策の具体例としては、以下が挙げられる。
 - **具体例**:機密情報の検出には自動化されたツールやAIベースの検出器を活用しつつも、検出の網をすり 抜ける情報に対応するためにセキュリティ教育を受けたレビュアーによる人手の確認を併用することに より、機密情報の漏えいリスクを軽減することができる。

信頼できる事前学習データの使用



- 想定脅威:
 - データポイズニング攻撃
- 対策内容: LLMおよび画像識別AIに対する不正行為を意図した学習データの混入を防止する。信頼できる事前学習データの使用の具体的手法としては、以下が挙げられる。
 - **具体例**:入手した事前学習データの出所、加工履歴などが公開されている場合は事前学習データの信頼性を確認する。それらの情報の提供がない場合は、実績ある提供元、あるいは著名な研究機関・公的機関から提供されるデータの使用を優先する。

ファインチューニングデータによる補正

- 想定脅威:
 - データポイズニング攻撃
 - MLaaSの悪用(バックドアを設置する攻撃)
- **対策内容**: バックドアが仕込まれたLLMおよび画像識別AIに対して、安全性が確認された良質なデータセットを用いたファインチューニングを実施する。ファインチューニングの具体的手法としては、以下が挙げられる。
 - **具体例**:信頼できるプロンプトや画像とそれらに対する正しい応答例を多く含むデータでモデルをファインチューニングすることで、以前に学習した悪意あるパターンの影響を希釈・無効化する。それにより、事前の学習データの完全性を保証できない場合でも、データ汚染による挙動の異常を緩和することができる。

AI開発者が実施者となる対策の一覧(画像識別AI特有のもの)

AIシステムに対する対策のうち、**画像識別AI特有なものについては、**具体例として以下が挙げられる。

対策の具体例	対応する脅威の具体例	対策段階	LLM/ 画像識別AI
 画像識別AI自体の頑健性向上 画像識別AIが敵対的サンプルを誤識別しないように学習する対策手法。 画像識別AIの学習時に、通常のデータに加えて画像識別AIの誤識別を誘発するデータや様々な変化を加えたデータを学習させることで、画像識別AI自体の頑健性を向上させ、敵対的サンプルによる誤識別やデータ窃取を抑制する。 	I-3-1:画像識別AIの挙動を観察することで 学習データを推測する攻撃 I-6-1:入力画像を細工することで画像識別 AIの誤識別を誘発する攻撃	開発	画像識別AI
DoS攻撃に対するAIの頑健性向上 ・ 画像識別AIのボトルネックを解消し、DoS攻撃耐性を高める対策手法。 ・ 画像識別AIの設計段階で最悪ケースのエネルギー消費や処理遅延が発生する箇所を特定し、アーキテクチャを最適化することでDoS攻撃への耐性を高める。	 I-7-1 : 入力画像を細工してDoSを行う攻撃 	開発	画像識別AI
 信頼できるAI作成サービスの利用 ・ 不正行為を意図したMLaaSの利用を防止する対策手法。 ・ 信 頼 で き る MLaaS(Machine Learning as a Service)を利用し、AIが細工されないようにする。例えば、実績のある提供元のサービスを利用することを優先する。 	I-8-2: MLaaSサービスを利用して画像識別 AIにバックドアを設置する攻撃	開発	画像識別AI

画像識別AI自体の頑健性(※)向上

- 想定脅威:
 - 回避攻撃
 - メンバーシップ推論攻撃
- **対策内容**: AIモデル自体の頑健性を高めることで、敵対的サンプルによる誤識別やデータ窃取を抑制する。頑健性向上の具体的手法としては、以下が挙げられる。
 - **データ拡張**:通常の画像データだけでなく、敵対的サンプルや、拡大・縮小・回転・コントラスト変更など多様な変換処理を加えたデータを組み込むことで、頑健性を高める。
 - **複数モデルの併用:** アーキテクチャの異なる複数のモデルを併用する手法(アンサンブル学習やMixture of Experts等)を用いることで、頑健性を高める。

※頑健性向上は、一般には信頼性等の観点からも実施すべき対策になるが、本資料ではセキュリティの観点からの対策として記載している

$M^{I}B_{I}S^{I}D_{\bullet}$

DoS攻撃に対するAIの頑健性(※)向上

- 想定脅威:
 - DoS攻撃
- 対策内容: 画像識別AIの設計段階で、想定される最悪ケースの入力(画像の大きさ・色情報の複雑さ・歪みなど)に対する処理経路を分析し、モデルのアーキテクチャ最適化(例: 軽量なモデル設計、前処理による入力制限、バッチ処理やタイムアウトの設定など)を行う。DoS攻撃に対するAIの頑健性向上の具体的手法としては、以下が挙げられる。
 - **前処理による入力制限**:入力段階でのフィルタリング機構を設け、明らかに処理コストの高い画像を遮断またはリサイズ・間引き処理することで、モデル本体への過剰負荷を防止する。
 - **バッチ処理やタイムアウトの設定**:処理負荷を測定するためのテストスイート を導入し、事前に過負荷 状態を再現・評価し、モデルが過負荷とならないようにタイムアウト等を設定する。

※頑健性向上は、一般には信頼性等の観点からも実施すべき対策になるが、本資料ではセキュリティの観点からの対策として記載している

信頼できるAI作成サービスの利用



- 想定脅威:
 - MLaaSの悪用
- 対策内容:信頼できるMLaaS(Machine Learning as a Service)を利用し、AIが細工されないようにする。信頼できるAI作成サービスの利用の具体的手法としては、以下が挙げられる。
 - **具体例**:実績と透明性のあるMLaaSプロバイダのサービスを選定する。例えば、大手クラウド事業者が 提供するMLaaSは、一定程度、AI開発環境に関する技術仕様も開示されており、運用体制やセキュリ ティ監査体制が整っているので、サプライチェーン・レベルでの信頼性確保が期待できる。



画像識別AIに対する攻撃手法や対策のVLMへの転用可能性

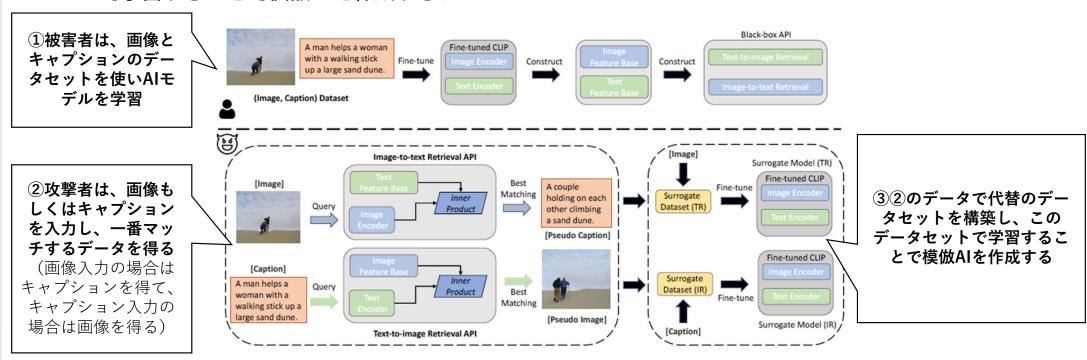
画像識別AIに対する攻撃手法や対策のVLMへの転用可能性

- ・ 画像識別AIに対する攻撃手法や対策については、原理的にVLM※に転用可能であると考えられる。
- 画像識別AIに対する以下の脅威については、同様なアプローチでVLMを攻撃する事例がある。
 - 「I-1-1:画像識別AIを模倣したAIを作成する攻撃」
 - 「I-2-1:画像識別AIの内部構造を観察することで学習データを窃取する攻撃」
 - 「I-3-1:画像識別AIの挙動を観察することで学習データを推測する攻撃」

※Vision Language Modelの略。画像等の視覚情報と、テキスト等の言語情報を統合的に処理するAI技術。

I-1-1:画像識別AIを模倣したAIを作成する攻撃の転用事例1

- VLMにおいても、「I-1-1:画像識別AIを模倣したAIを作成する攻撃」と同様のアプローチにより模倣AIを作成する研究がある。
 - 画像とテキストを扱うVLMに対し、入出力の挙動を観察してデータセットを構築し、そのデータセットで学習することで模倣AIを作成する。

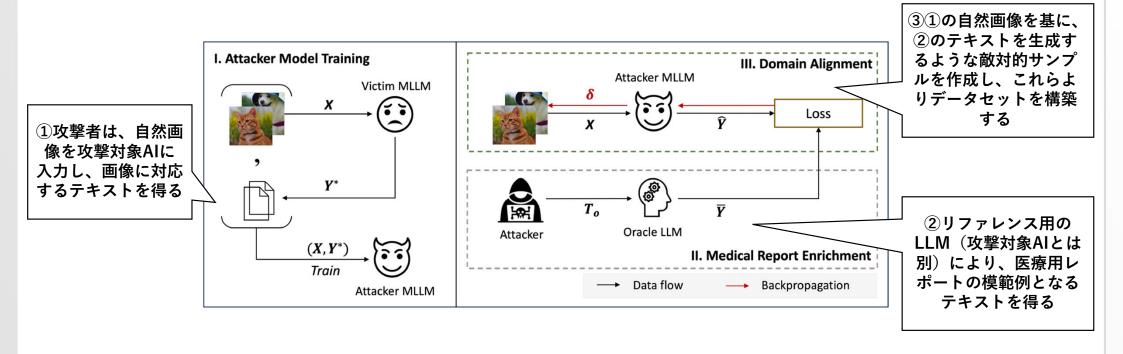


出典: https://openreview.net/forum?id=v-rx235RlfI

28

I-1-1:画像識別AIを模倣したAIを作成する攻撃の転用事例2

- 医療用VLMにおいても、「I-1-1:画像識別AIを模倣したAIを作成する攻撃」と同様のアプローチにより模倣AIを作成する研究がある。
 - 医療用VLMが対象で、AIモデル学習のための医療用データセットが不足している(プライバシーの観点で非公開の場合が多い)ため、自然画像を基にデータセットを構築する。

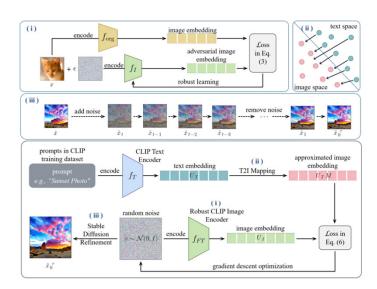


出典: https://ojs.aaai.org/index.php/AAAI/article/view/32734/34889

I-2-1:画像識別AIの内部構造を観察することで学習データを窃取する攻撃の転用事例

- VLMにおいても、「I-2-1:画像識別AIの内部構造を観察することで学習データを窃取する攻撃」と同様のアプローチにより学習データを窃取する研究がある。
 - CLIPモデル(OpenAIが2021年に発表した画像と自然言語を扱うマルチモーダル基盤モデル)に対し、 敵対的ファインチューニング、埋め込みアラインメント(テキストと画像を共通の特徴空間で扱えるよ うにする)、Stable Diffusionによる画像の再構築により、学習データを抽出する。



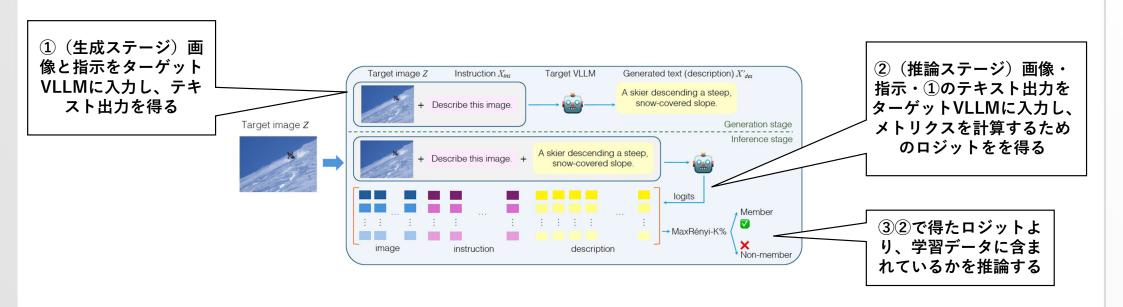


出典: https://arxiv.org/pdf/2508.00756

30

I-1-3:画像識別AIの挙動を観察することで学習データを推測する攻撃の転用事例

- VLMにおいても、「I-3-1:画像識別AIの挙動を観察することで学習データを推測する攻撃」と同様のアプローチにより学習データの一部を間接的に入手する研究がある。
 - ロジットをもとにメトリクス(MaxRényi-K%)を計測することで、入力が学習データに含まれているかを推論する。



出典: https://proceedings.neurips.cc/paper_files/paper/2024/file/b2c892312af07f8a77afbeed188391f4-Paper-Conference.pdf