

MagicDelta



開発駆動コース
仲山ゼミ 秋穂 正斗

自作VPLとレンダリングエンジンを用いた
AIセキュリティについて学習できるWebアプリ

X:Tebasaki_lab

MagicDeltaとは？

課題

AIが難しく挫折...



- ① AIセキュリティの認知度が低い！
- ② AI自体がめちゃくちゃ難しい！
- ③ 敷居が高いので誰もやらない！

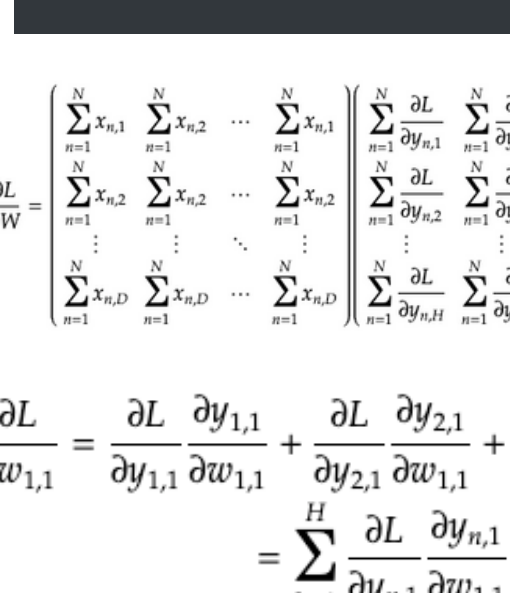
解決策

自作VPLとAnimationでAIセキュリティについて
効率よく学習できるコンテンツの開発

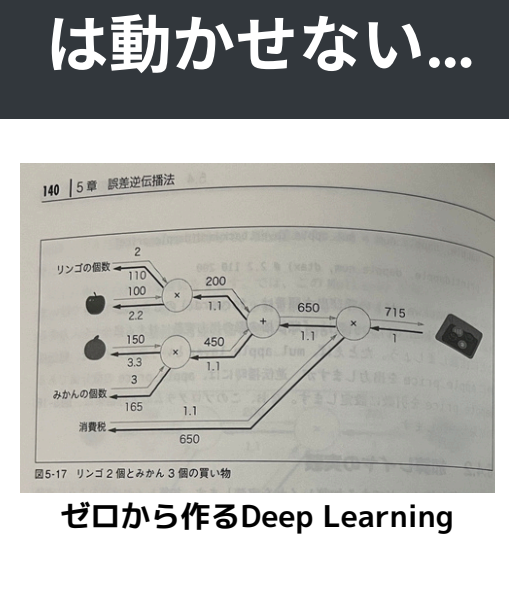
価値

- ① AIセキュリティの認知度が増える！
- ② AIを直感的に学習できる！
- ③ ゲーム感覚でみんなが取り組む！

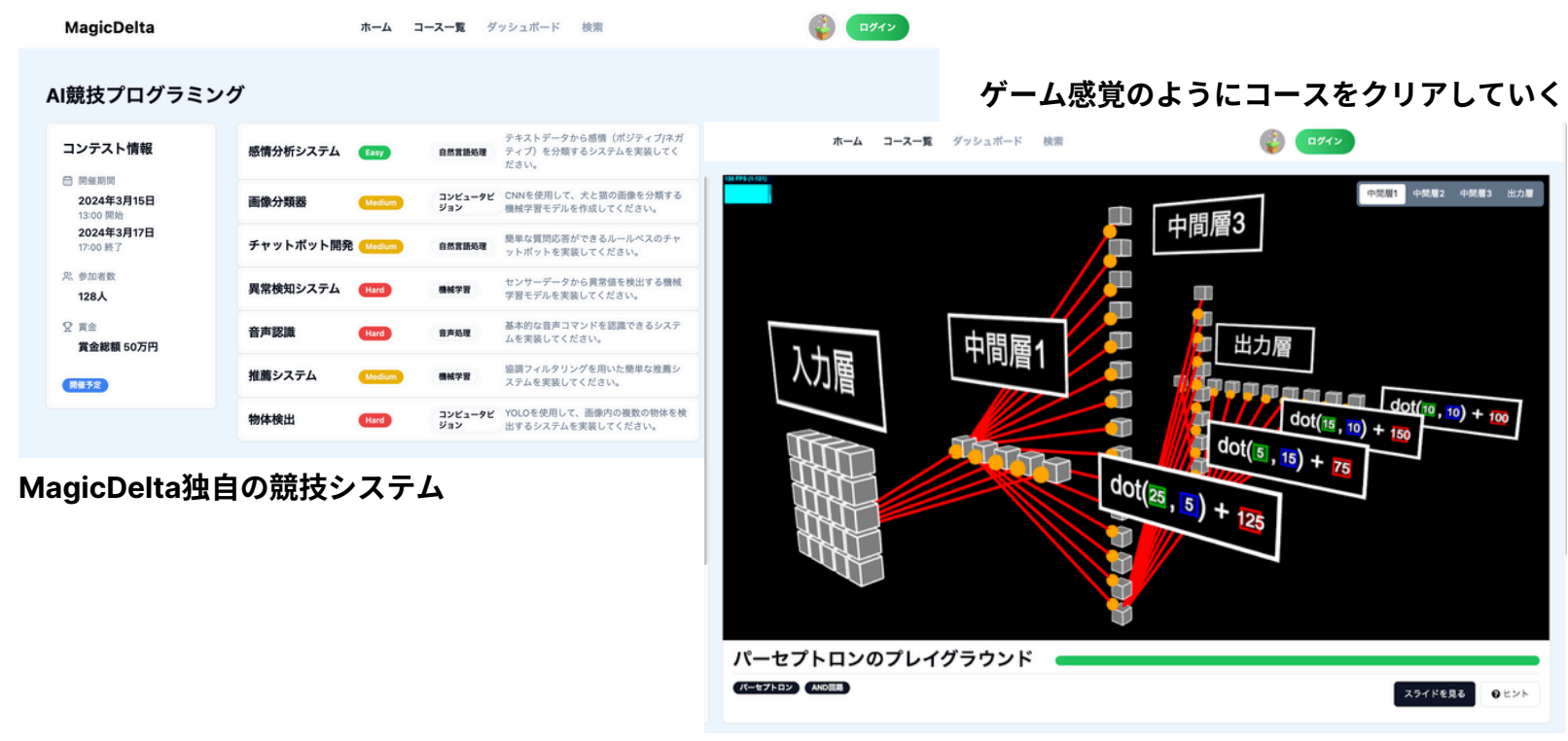
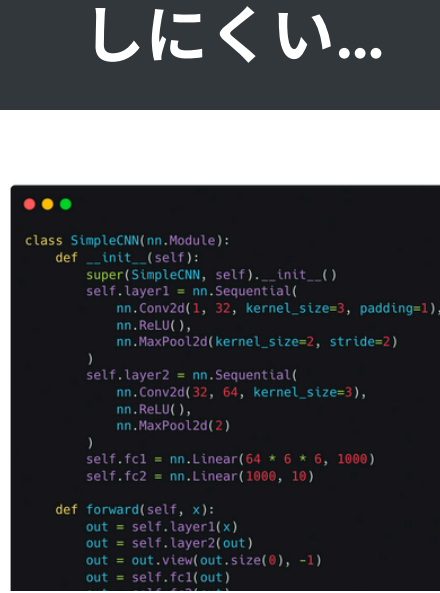
数学が難しい...



動画や本の内容は動かせない...



観測がしにくい...



AIセキュリティとは

AIセキュリティについて

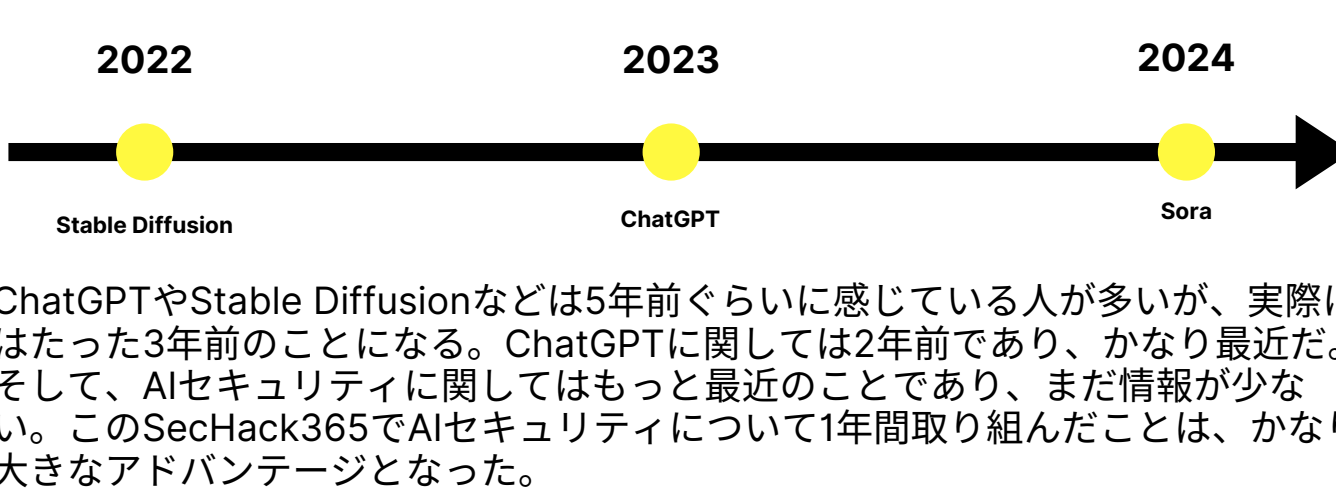
AIの脆弱性について定義する
脆弱性=AIへの不正な入力, AIによる間違えた出力

AIの脆弱性というものはかなり抽象的で、人によって意見が変わるため、今回のプロジェクトではこのように脆弱性を定義した。不正な入力というのは、ノイズや不正なデータが入った入力などを指し、間違えた出力というのは、画像の認識率が極端におかしい、ハルシネーション、正解ラベルとは違うものを出力するなどがある。これを元に、どのようにしたら脆弱性を潰すことができるのか、脆弱性の仕組みについて学べるのかを考え、コンテンツ制作に取り組んだ。

AIへの攻撃の例

- 敵対的攻撃 (サンプル)
 - 入力画像に人間が気づかないようなノイズを注入することにより、AIに間違えた判定をさせようもの
- 代表的な攻撃先
 - 自動運転
 - 顔認証
 - 文字認識

AIはつい最近の出来事



MagicDeltaの特徴

① 最高の学習体験

MagicDeltaの最大の特徴は、その学習体験にある。3Dで自由にAIモデルを観測できたり、AIに特化した自作VPLを搭載しているため、書籍に比べて圧倒的に分かりやすく、効率的にAIセキュリティの仕組みについて学習できる。

② コンテンツ

勾配降下法や敵対的攻撃などを学ぶために、それぞれコンテンツ(教材)を用意されている。このコンテンツは基礎から応用まで豊富に存在し、AIセキュリティについて網羅的に学習することができる。

③ 競技システム

これらを用いて構築されたのがMagicDelta独自の競技システムだ。ユーザーはCTF形式でAIセキュリティを学ぶことができる。例えば、うまく動かないAIモデルを修正したり、脆弱性のあるAIモデルに攻撃をし、間違えた推論家にするなど、既存のCTFとかなり近い形式をとっている。今後、この独自の競技システムは今後有名なものになる。

テーマの歴史

第1回 イベント	6月15日(土)	テーマ: AI脆弱性の可視化
第2回 イベント	7月19日(金)~21日(日)	テーマ: 敵対的サンプル等の可視化
第3回 イベント	9月27日(金)~29日(日)	
第4回 イベント	11月15日(金)~17日(日)	
第5回 イベント	2025年 2月1日(土)・2日(日)	テーマをAIセキュリティの学習に変更

最初に応募した時は、「AI脆弱性の可視化」だったが、イベント回で様々な方からアドバイスを受けた。AIセキュリティを学習するために技術を使った方が良いと判断し、最終的に「AIセキュリティの学習」というテーマに変更された。開発駆動はとにかく手を動かして開発してみるのがモットーのため、テーマ変更にも柔軟さを感じた。

VPLの思想

私は、Scratch (ビジュアルプログラミング言語) が初めて触った言語であり、つまりScratchである。もうScratchをやっているが、そこで得た知識は素晴らしいものだった。Scratchの良さは、ユーザーがうまく動かない学習でも、プログラムをマウスで動かして確認しながら学ぶことができる点にある。Scratchは、Scratchの良さを活かして、AIセキュリティの学習にも活用したい。Scratchの良さを活かして、AIセキュリティの学習にも活用したい。Scratchの良さを活かして、AIセキュリティの学習にも活用したい。

開発技術 / こだわりなど

技術構成図

今回のWebアプリは、フロントエンドとバックエンドに分かれており、主にフロントエンドが大きな役割を担っている。全体の技術要素としては、Next.js (App Router), React, TypeScript, Tailwind CSS, Three.js, shadcn, Canvas APIなどを用いた。

自作VPL

LLMがプログラムや図解を下書きのコンテンツに生成するVPLは、全て自分で作り、ライブラリを使わずに、レンダリングエンジンから動作確認することができた。これは、レンダリングエンジンから動作確認することができた。これは、レンダリングエンジンから動作確認することができた。これは、レンダリングエンジンから動作確認することができた。

自作エンジン

2023年まで開発していたアニメーションエンジン、これは全て自作エンジンによって構築されている。現在はThree.jsを用いたが、そのコア部分は自作エンジンである。アニメーションを生成するロジックや自作のアニメーションエンジンを実装する部分、良い例を作るように工夫した。よく、Webアプリケーションやゲームのアニメーションの作り方は一切無かった。自分自身の経験やアニメーションの経験、Scratchのアニメーションの経験から構築した。Scratchのアニメーションの経験から構築した。Scratchのアニメーションの経験から構築した。

開発でのこだわり

リソース配分
・コスト削減に努める。開発目的にならないようにした。
・動的なテーマ変更
・当時は脆弱性の可視化だったが、AIセキュリティの学習というテーマに変更した。

苦労した点

自作VPLとアニメーションエンジンに関しては私の専門的な知識が不足していた。最初は、AIセキュリティに関する知識がなかった。AIセキュリティに関する知識がなかった。AIセキュリティに関する知識がなかった。AIセキュリティに関する知識がなかった。

外部での発表

SECCON 13 Contest of Contest 入賞

2024年12月、こちらのコンテストに応募した。以前のオフラインイベント回でトレーナーに応募してみたら？と誘われたのがきっかけだった。結果入賞し、3/1の電脳会議でワークショップを行うことに。

福岡未踏 2024

経済産業省 福岡未踏に2024年度に採択。これは、実は同じWebサイトの開発を行っていたが、ちゃんと申し込みをしていた。福岡未踏では自作VPLやアニメーションエンジンの基礎技術を開発し、SecHack365ではそれを用いてコンテンツ開発や足りない機能や新しい設計を導入した。つまり、福岡未踏でFWを作り、それを元にWebサイトをSecHack365で作っていたと表現するのが分かりやすい。実際、双方で見つけた課題を並行して解決したことも多かった。

今後について

社会実装

- toB
 - 大学の授業 - 10件
 - 企業の研修 - 5件
- toC
 - ユーザー数 5万人
 - MAU 1000人

機能開発

- ログイン
 - Google, emailログイン
- 学習ログ
 - 日々の勉強内容やクリア数など
- コンテンツ開発
 - +50コース
- OSS化
 - VPL
 - アニメーションエンジン

3年以内に実現

この1年で学んだこと

SecHack365の1年間で、私は沢山のことを学んだ。もともと、AIの基礎技術となる線形代数、誤差逆伝播、RNNやCNNなどはすでに知っていたが、AIセキュリティや、そもそもセキュリティについては素人だった。去年は100冊ほど専門書を購入したが、そのうち、SecHack365での活動の後押ししてくれた本たちを紹介する。基本的に機械学習、AI周りの本がもちろん多いが、SecHackのメンバーがOAuthをやっていたり、Web系のセキュリティについて詳しいのに憧れたため、認証認可やWebセキュリティについても個人で勉強した。SecHack365に参加する人はこれらのセキュリティについて詳しい人が多いが、私のようにAIに詳しい人はあまりおらず、AI×セキュリティなシステムを作れるようになったのは強みになると実感した。



謝辞

本活動を始めるにあたり、多くの方々に支えられました。まず、ゼミの指導員である仲山先生に、貴重なご指導を賜り、心より感謝申し上げます。また、メンバー、トレーナー、アシスタント、外部講師の皆様にも、多大なるご支援をいただきました。SecHack365に取り組み環境を整えてくれた家族に深く感謝いたします。皆様のご協力と応援のおかげで、ここまで成長することができました。心より御礼申し上げます。

