

AIによる誤認識・偽造・差別の問題

佐久間 淳

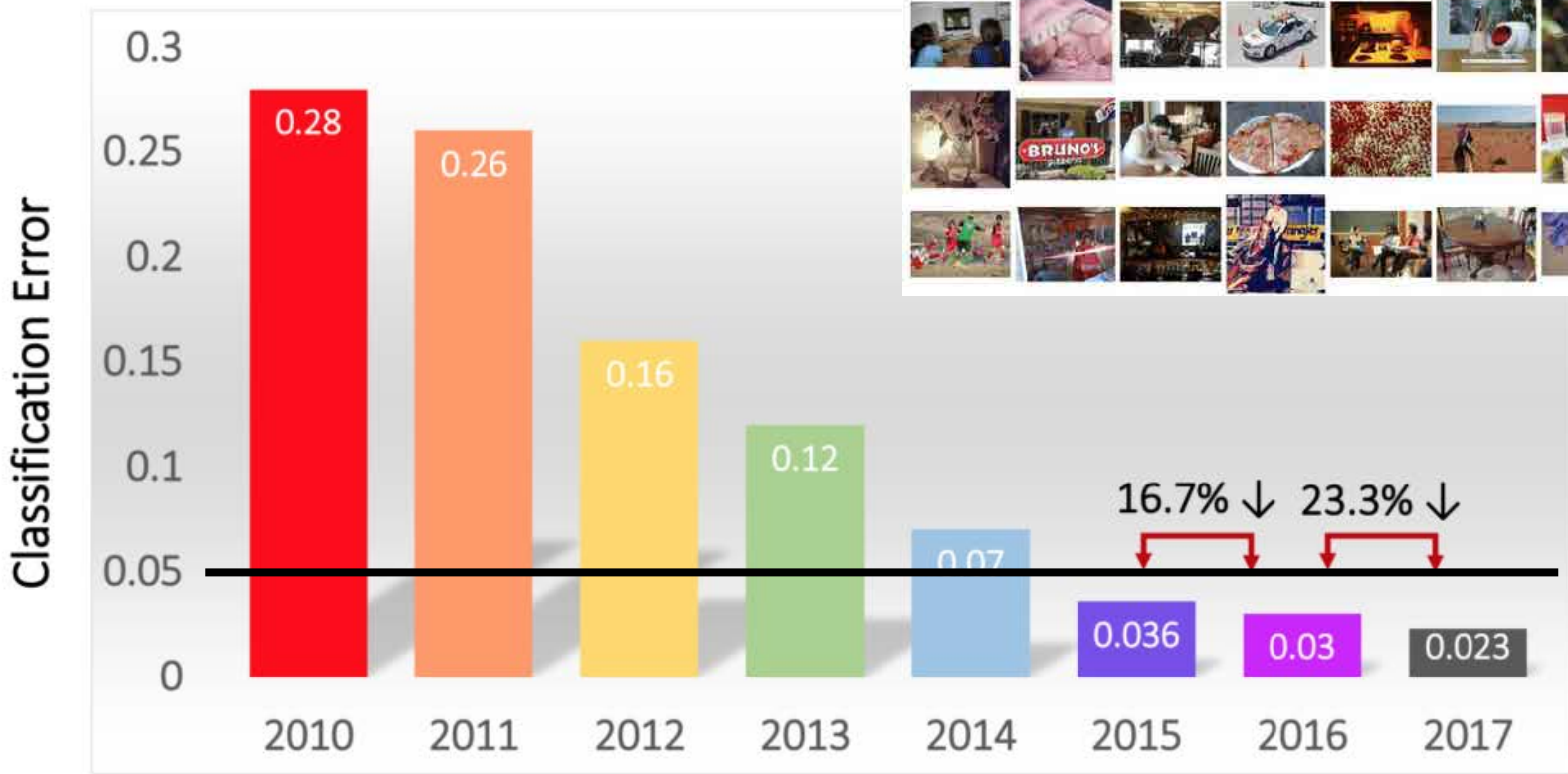
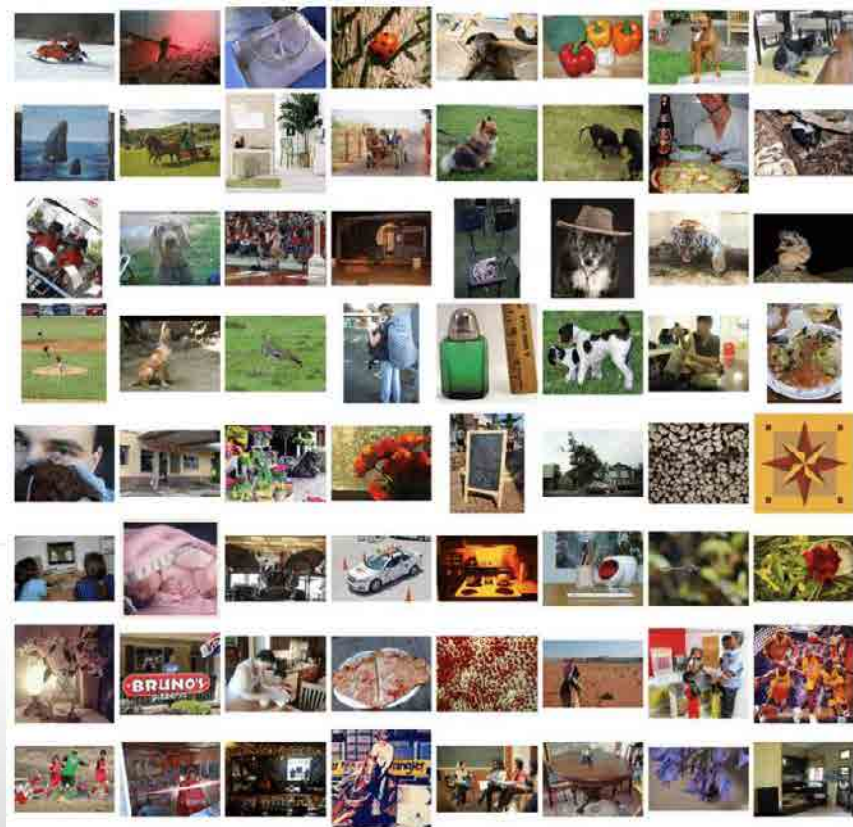
筑波大学 / 理研AIP



筑波大学
University of Tsukuba

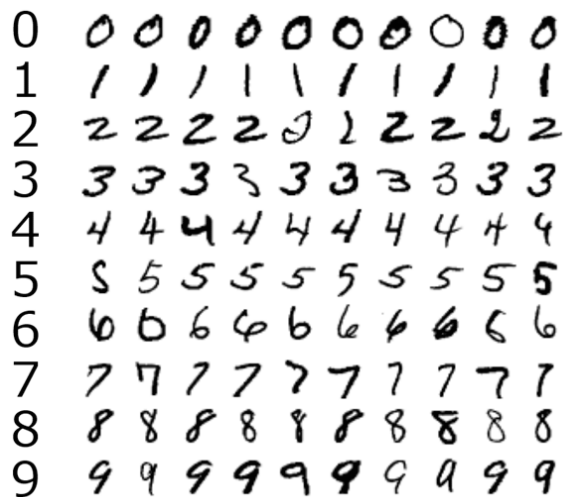


人間を超える AIの画像認識能力

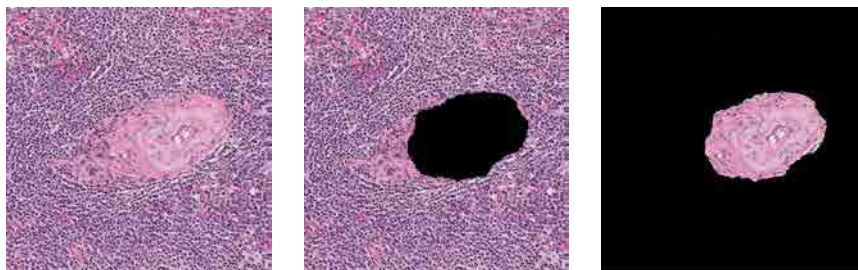


<https://www.kaggle.com/getting-started/149448>

AIに何を望む?



文字認識



医療診断、自動運転、etc.

- 単純な認識から、重大な判断へ
- 正確だけでなく、安定した判断が重要
 - 正確 = 精度が高い
 - 安定 = 多少の変化に影響されない(天候、ノイズ、攻撃、etc.)

敵対的サンプル(AE) [Szegedy+14]

The diagram illustrates the concept of adversarial examples (AE) for a neural network. It shows three stages:

- Original Image (x):** A photograph of a panda, classified as "panda" with 57.7% confidence.
- Adversarial Perturbation:** A square of random noise, representing the perturbation $\text{sign}(\nabla_x J(\theta, x, y))$, which is added to the original image. A red figure with a white envelope icon is positioned above this noise, indicating it is a crafted perturbation. The classification for this noise alone is "nematode" with 8.2% confidence.
- Adversarial Example:** The original image plus the perturbation, resulting in a new classification: "gibbon" with 99.3% confidence.

Two callout boxes provide context for the final classification:

- 人間はパンダと認識 (Humans recognize it as a panda)
- AIはテナガザルと認識 (AI recognizes it as a gibbon)

現実世界における敵対的サンプル [Etimov et al.]



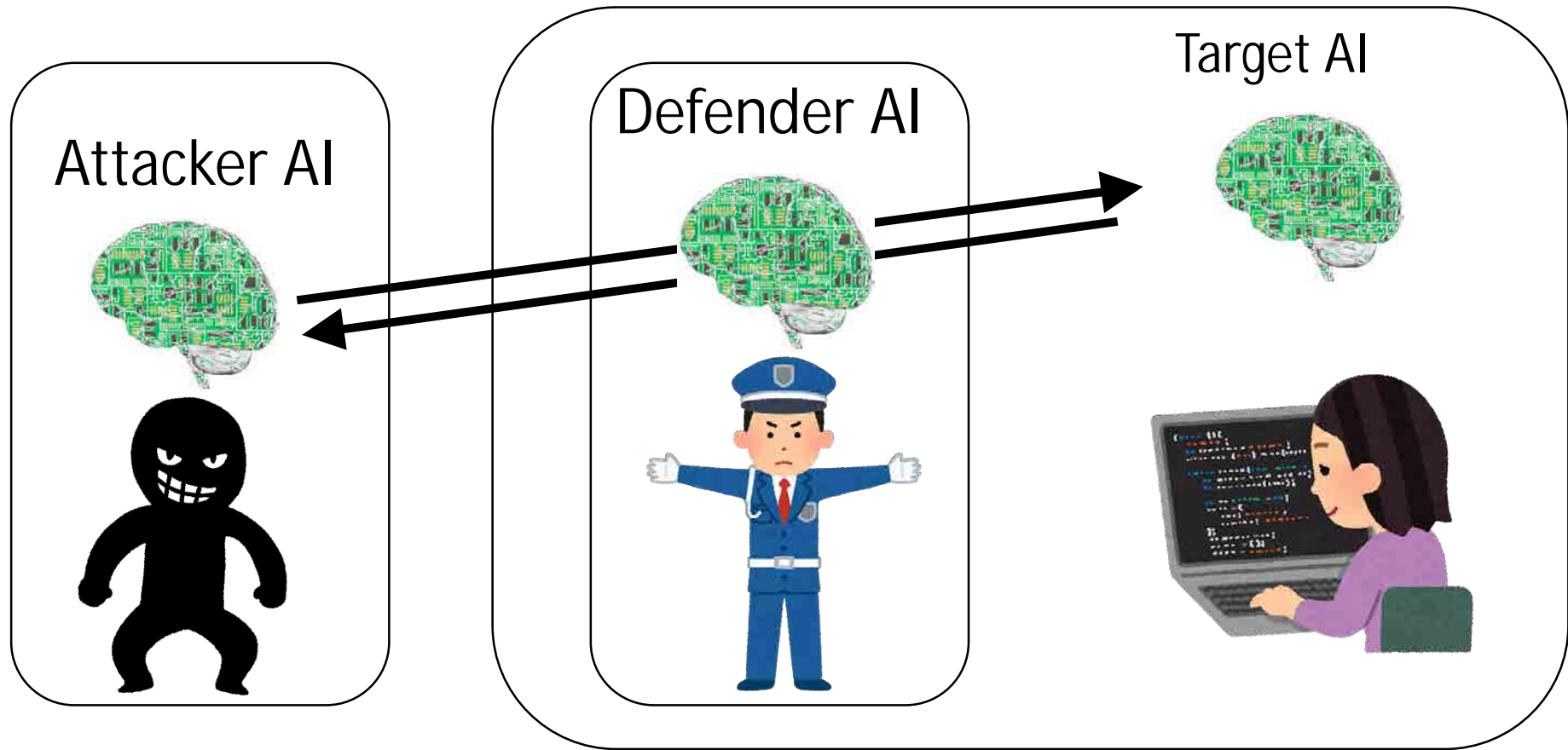
Images: [Etimov et al](#)

- 道路標識上に特殊なパターンを持つステッカーを添付
- AIは“STOP”標識を“speed limit 45 mph”と認識
- このステッカーが貼られた標識では、自動運転車はブレーキを働かせない可能性

音声の敵対的サンプル [YS, IJCAI'19]



AI同士の攻防



攻撃者は常に防御AIを上回る攻撃AIを作成しようとする
いたちごっこを続けること自体がAIの安全性を守ることに繋がる

AIによる人工メディア生成・操作



<https://deepfakesweb.com/?locale=ja>

実在しない人物画像生成

実在する人物の動画偽造

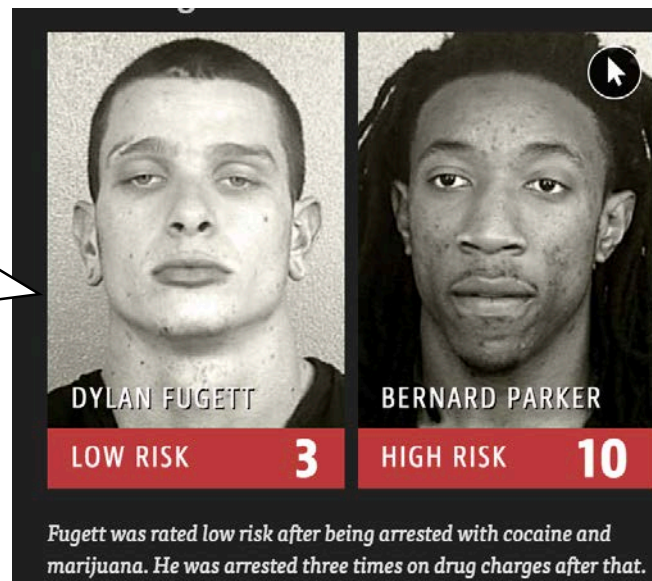
8

- 画像/音声/動画生成・操作
 - 偽造ポルノ、リベンジポルノ、証拠の捏造
- 人間の文章と見分けがつかない長文テキストの生成
 - デマ、フェイクニュース、虚偽方法、世論の誘導
- 従来より、巧妙かつ大量の自動生成

AIによる差別の実例 [Angwin+16]

- COMPASSは再犯リスクを予測（米国のいくつかの州）
 - スコアは被告人の保釈・判決・仮釈放等に影響
- COMPASSによる再犯リスクスコアリングは人種の影響を受けている [Angwin+16]
 - アルゴリズムによる決定が差別を生み出す

再犯リスク低と判定、
実際には三回再犯



再犯リスク高と判定、
実際には再犯なし

AIが社会に信頼されるためには

- エキスパートレベルの品質
 - 法令遵守
 - 環境変化・攻撃に対する安定性
 - 人権の尊重
 - 自己決定権の尊重
 - プライバシー保護
 - 公平性への配慮
 - 判断根拠の説明可能性
 - 判断過程の追跡可能性
-
- 必要条件
- 制約
- 事後検証