



Funded by
the European Union



From Attacks to Answers: Counterfactuals at the Intersection of Robustness and Explainability in AI

Davy Preuveneers

Industrial Research Manager, KU Leuven, Belgium

Adversarial Threats on Real Life Learning Systems

September 17th, 2025 – Esclangon building, Campus Pierre et Marie Curie, Paris

Overview

Counterfactuals at the Intersection of Robustness and Explainability in AI

- | Motivation and context
- | Problem landscape
- | Adversarial examples and counterfactual explanations
- | Commonalities and key differences
- | Future directions
- | Summary and takeaways



Support the Guardian

Fearless, independent, reader-funded

Support us →

The Guardian

News

Opinion

Sport

Culture

Lifestyle

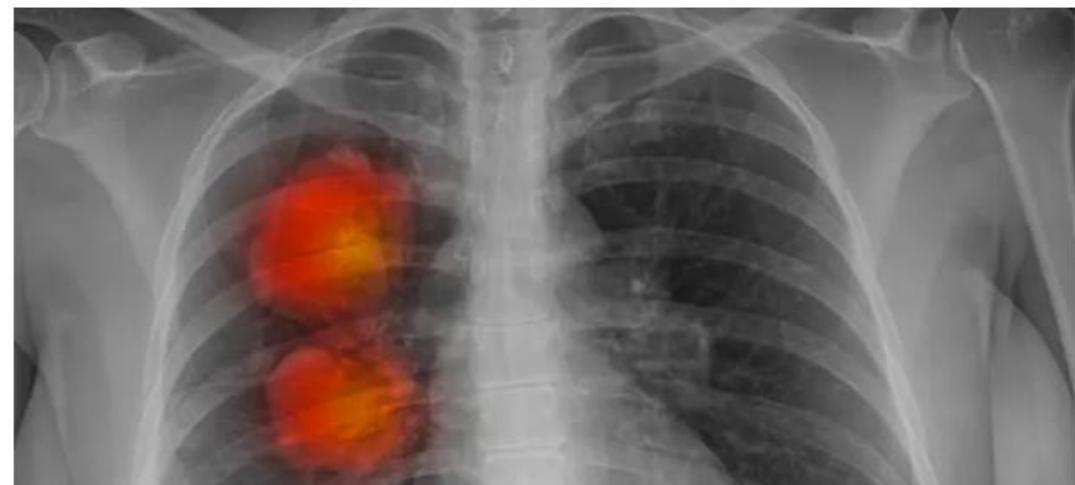
More ▾

UK World Climate crisis Newsletters Football Coronavirus Business Environment UK politics Education **Society**

Cancer

New artificial intelligence tool can accurately identify cancer

Exclusive: algorithm performs more efficiently and effectively than current methods, according to a study



“The team used CT scans of about 500 patients with large lung nodules to develop an AI algorithm using radiomics. The technique can extract vital information from medical images not easily spotted by the human eye.”

“The results showed the AI model could identify each nodule’s risk of cancer with an AUC of 0.87.”

“According to these initial results, our model appears to identify cancerous large lung nodules accurately”

<https://www.theguardian.com/society/2023/apr/30/artificial-intelligence-tool-identify-cancer-ai>



THE
EDGE

ANALYTICS

ATTACKS /
BREACHES

APP SEC

CLOUD

ENDPOINT

IoT

OPERATIONS

PERIMETER

PHYSICAL
SECURITY

RISK

THREAT
INTELLIGENCE

VULNS /
THREATS

VULNERABILITIES / THREATS // ADVANCED THREATS

4/27/2021
11:00 AM

Expect an Increase in Attacks on AI Systems



Robert Lemos
News

0 COMMENTS
[COMMENT NOW](#)

Companies are quickly adopting machine learning but not focusing on how to verify systems and produce trustworthy results, new report shows.

Research into methods of attacking machine-learning and artificial-intelligence systems has surged—with nearly 2,000 papers published on the topic in one repository over the last decade—but organizations have not adopted commensurate strategies to ensure that the decisions made by AI systems are trustworthy

Related Content

Sponsored by



RESOURCES

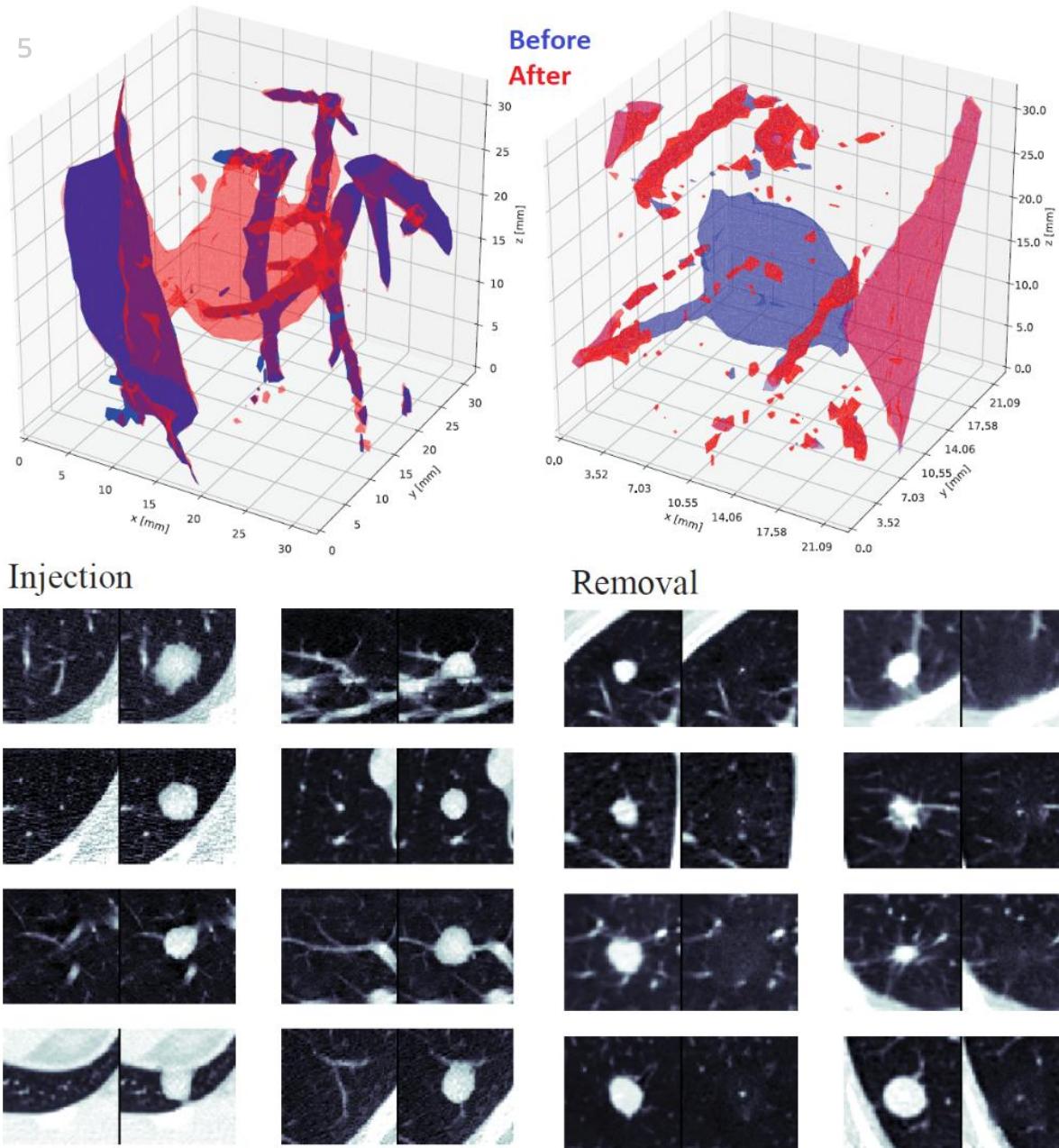


Zero Trust and the Power of Isolation for Threat Prevention

Cybersecurity is a top priority for enterprise networks, but as cyberattacks and data breaches become more prevalent, traditional cybersecurity approaches fall short. Here's how secure web gateways (SWGs) can help.

Login

<https://www.darkreading.com/vulnerabilities---threats/advanced-threats/expect-an-increase-in-attacks-on-ai-systems/d/d-id/1340833>



Injecting and removing cancerous pulmonary lung nodules with generative adversarial networks

Why? Insurance fraud, political motives, job theft, ...

Yisroel Mirsky, Tom Mahler, Ilan Shelef, and Yuval Elovici. CT-GAN: Malicious Tampering of 3D Medical Imagery using Deep Learning. 28th USENIX Security Symposium (USENIX Security 19)

https://www.usenix.org/system/files/sec19-mirsky_0.pdf

<https://github.com/ymirsky/CT-GAN>



6

BAD BOTS ON THE RISE: INTERNET TRAFFIC HITS RECORD LEVELS

| 09 JUN 2025 |

CYBERTHREATS CYBERATTACKS CYBERSECURITY DATA PRIVACY DATA BREACHES INNOVATION
DIGITAL TRANSFORMATION REVIEW DIGITAL IDENTITY AND SECURITY AI

“More than half of all online traffic is generated by automated software programs known as bots.”

<https://www.thalesgroup.com/en/worldwide/digital-identity-and-security/magazine/bad-bots-rise-internet-traffic-hits-record-levels>

2025

BAD BOT REPORT

The Rapid Rise of Bots and the Unseen Risk for Business

[2025 Bad Bot Report by Imperva](#)

The web bot arms race: detection and evasion

Motivation and key challenges

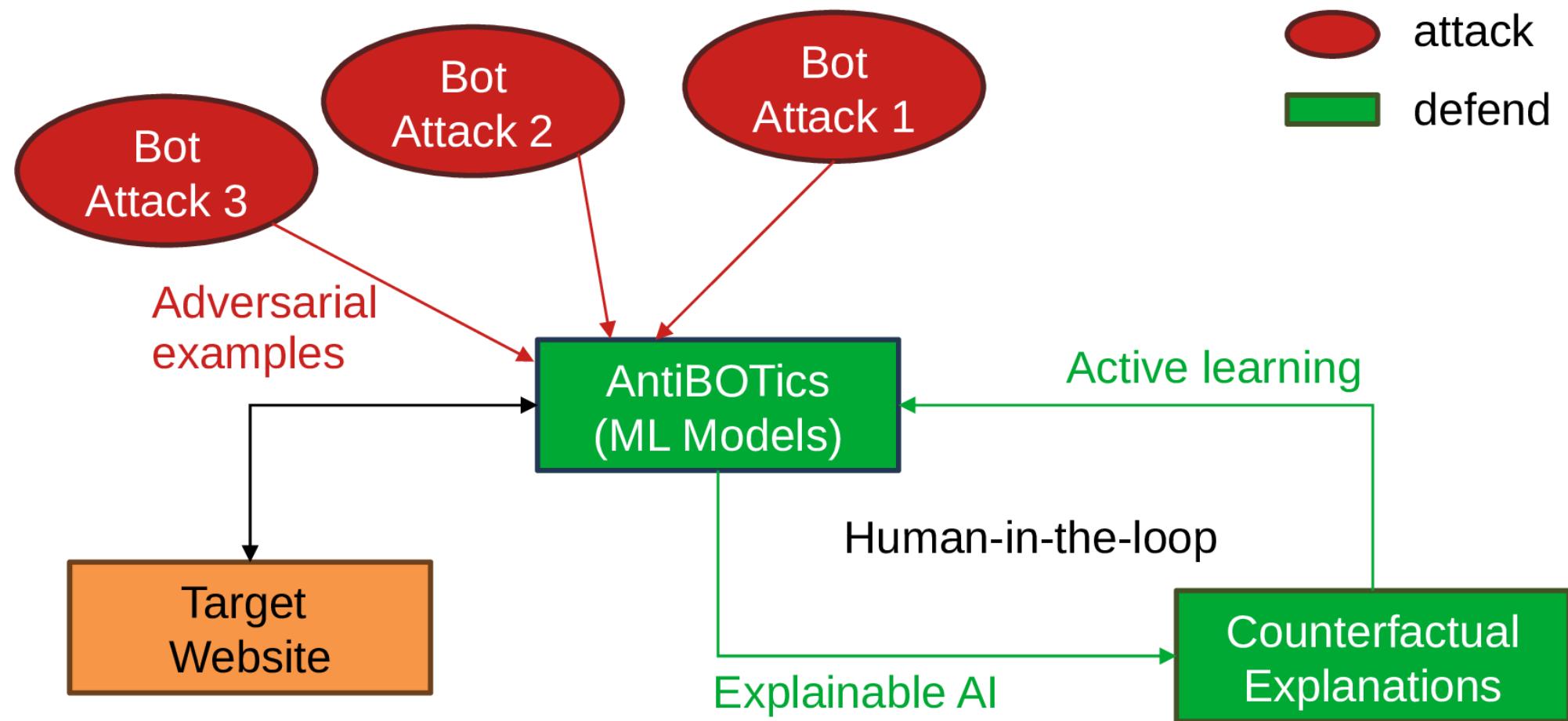
- | Modern web bots mimic human behavior to evade detection
- | They exploit residential IPs and proxies
- | Traditional detection methods (e.g., rule-based) ineffective
- | AI models misclassify traffic: false positives and negatives
- | AI models vulnerable to unseen behaviors and evasion attacks

→ Detection must evolve continuously!

→ But manual verification and sample labeling is burdensome!

The web bot arms race: detection and evasion

AI-based attacks and defenses



AntiBOTics: AI-based web bot attack detection

Robust ML models for a hybrid AI defense system

| Feature extraction from HTTP logs over 10, 100 and 1000-second intervals

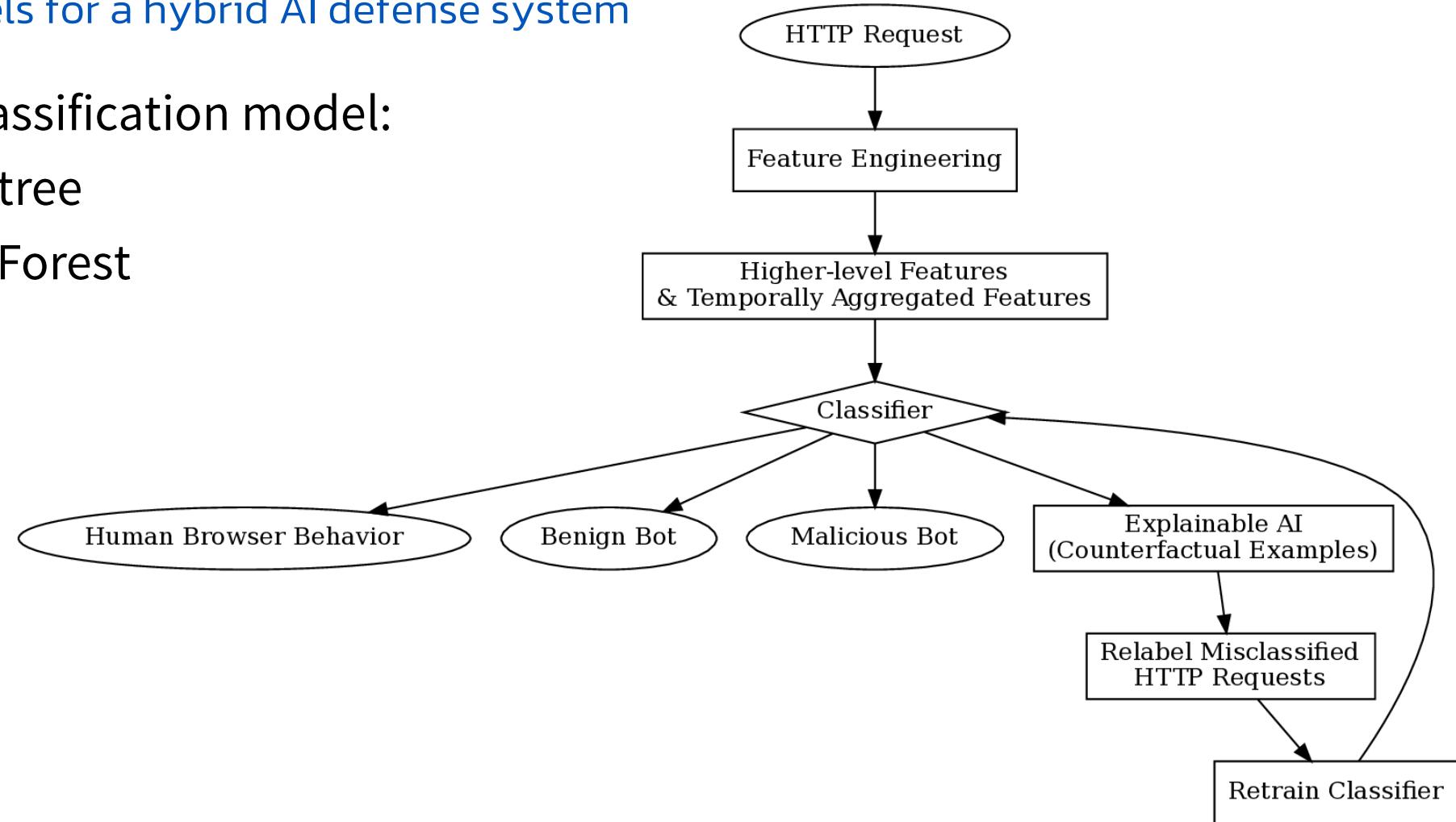
Feature Group	Description
<i>Total Number of Requests</i> *	The cumulative count of HTTP requests made.
<i>HTTP Request Methods</i> *	Distribution of HTTP methods used (e.g., GET, POST, PUT, HEAD).
<i>HTTP Response Codes</i> *	Distribution of HTTP response status codes received (e.g., 200, 301, 404).
<i>Session Duration</i>	Total time elapsed from the first to the last request in a session.
<i>Time Between Requests</i> *	Average interval between consecutive HTTP requests.
<i>Unique URLs Accessed</i> *	Count of distinct URLs visited.
<i>User-Agent String</i>	Analysis of the User-Agent header to identify the browser or bot type.
<i>Referrer Analysis</i>	Examination of the HTTP Referrer header to understand navigation paths.
<i>Request Patterns</i> *	Detection of repetitive or predictable request sequences.
<i>Resource Types Requested</i> *	Types of resources requested (e.g., HTML, images, scripts) and their frequencies.

C. Iliou et al., *Detection of advanced web bots by combining web logs with mouse behavioural biometrics*, Digital threats: research and practice 2 (2021), 1–26.

AntiBOTics: AI-based web bot attack detection

Robust ML models for a hybrid AI defense system

- | Training a classification model:
 - | Decision tree
 - | Random Forest
 - | XGBoost
 - | CatBoost
 - | ...



AntiBOTics: AI-based web bot attack detection

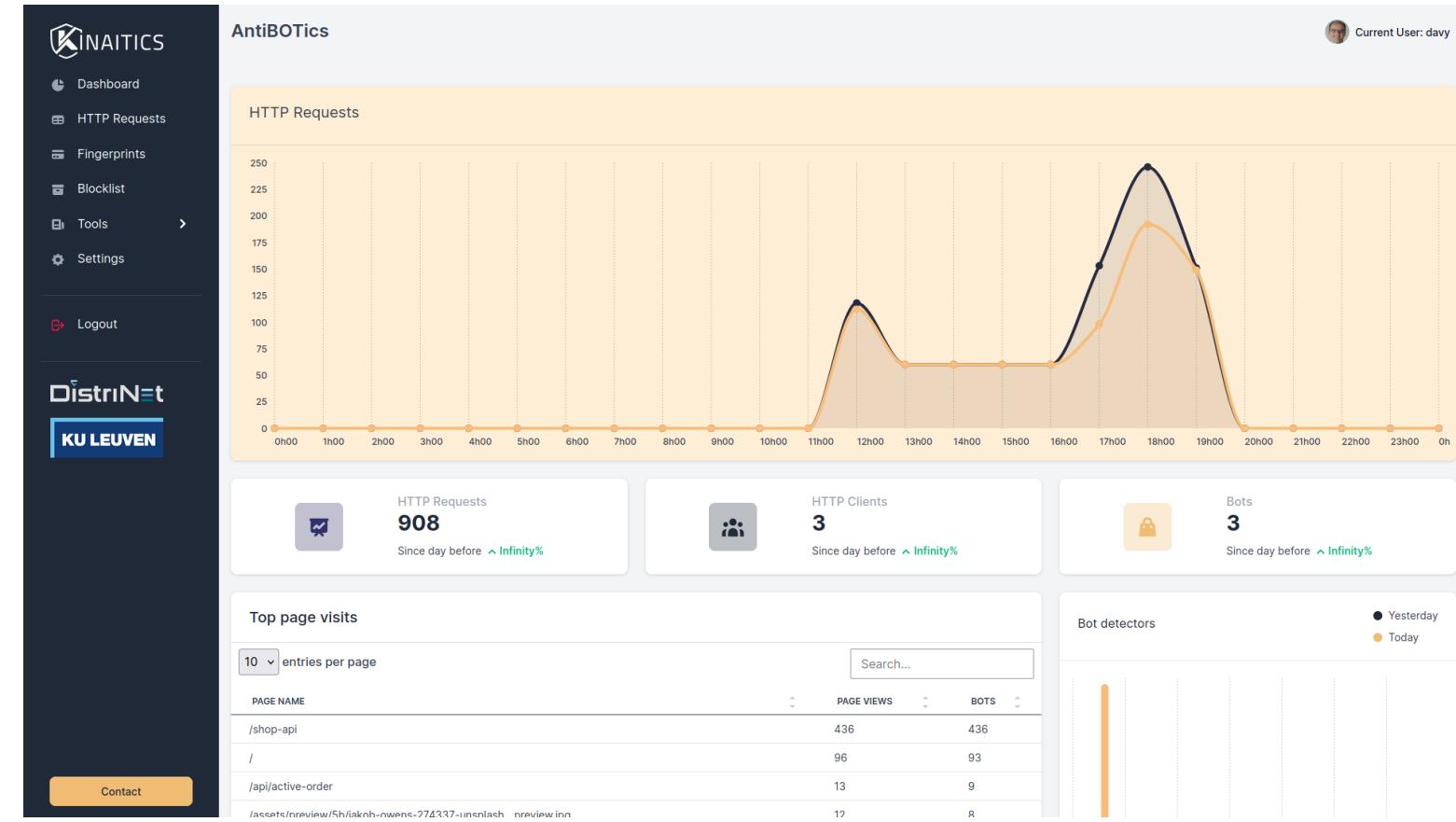
Robust ML models for a hybrid AI defense system

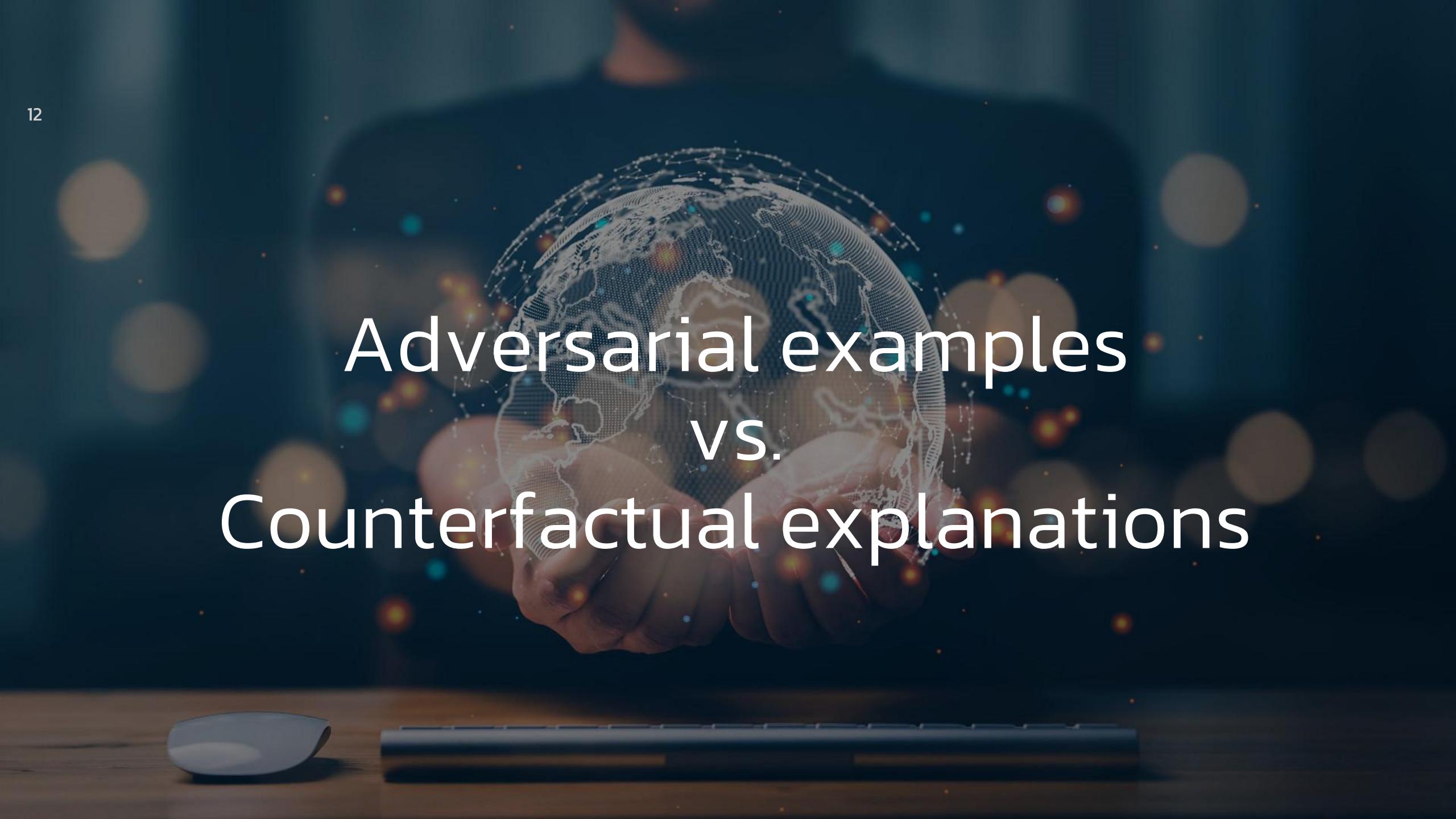
Key features

- | Checklists (e.g. IP address)
- | Browser Fingerprinting
- | Heuristics / Statistics
- | Mix of ML Classifiers

Human-in-the-loop

- | Explainable AI with Counter-factual Explanations
- | ML Model Updates Through Active Learning



A hand is shown interacting with a large, semi-transparent digital globe of the Earth. The globe is overlaid with a network of white lines and small colored dots (blue, red, green) representing data or connections. The background is dark blue with blurred lights, suggesting a futuristic or high-tech environment. In the foreground, a portion of a keyboard and a computer mouse are visible on a light-colored desk.

Adversarial examples vs. Counterfactual explanations

Adversarial examples

Modify inputs slightly to change model predictions



$$+ .007 \times$$



=



x

“panda”

57.7% confidence

$$\text{sign}(\nabla_x J(\theta, x, y))$$

“nematode”

8.2% confidence

$$\epsilon \text{sign}(\nabla_x J(\theta, x, y))$$

“gibbon”

99.3 % confidence

- | Goodfellow, Ian J., Jonathon Shlens, and Christian Szegedy. "Explaining and harnessing adversarial examples." *arXiv preprint arXiv:1412.6572* (2014). <https://arxiv.org/abs/1412.6572>

14

Adversarial examples

Modify inputs slightly to change model predictions

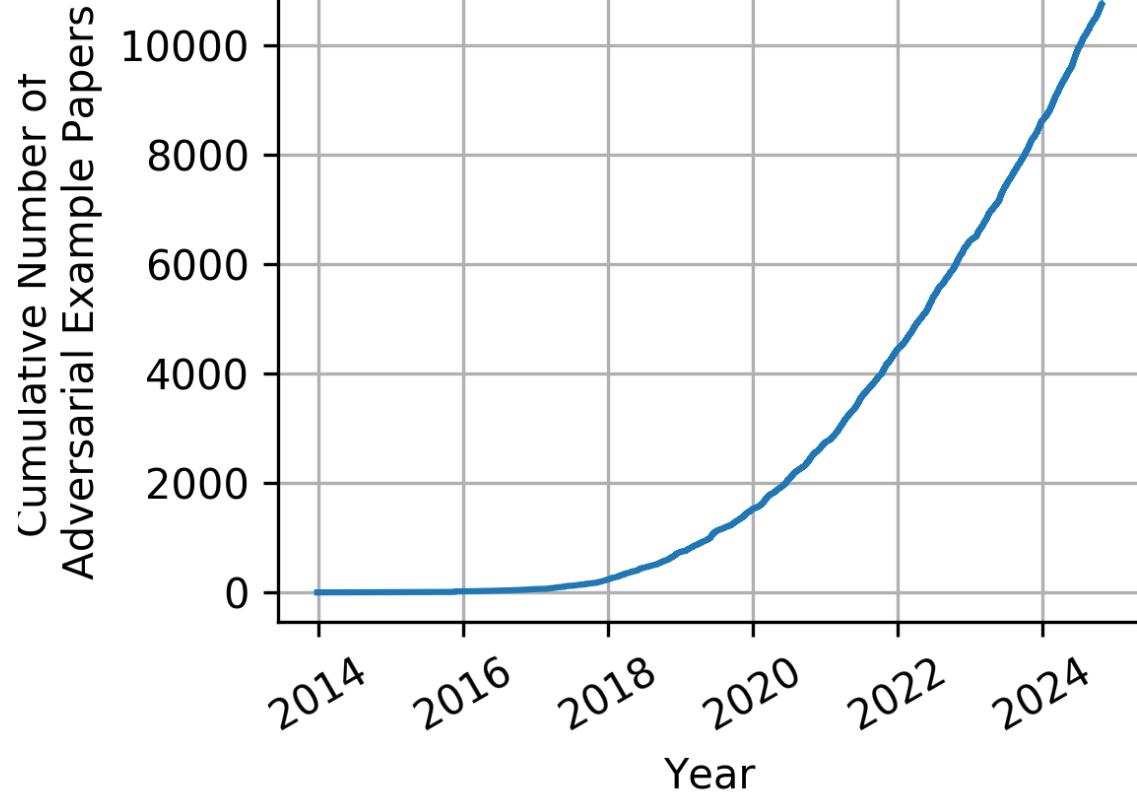


Attackers provide unusual inputs in the physical world

| Eykholt et al., *Robust Physical-World Attacks on Deep Learning Models*, CVPR, 2017.

Adversarial examples

A Complete List of All (arXiv) Adversarial Example Papers



- | Nicolas Carlini (Google DeepMind)
- | <https://nicholas.carlini.com/writing/2019/all-adversarial-example-papers.html>

Adversarial examples with constrained domains

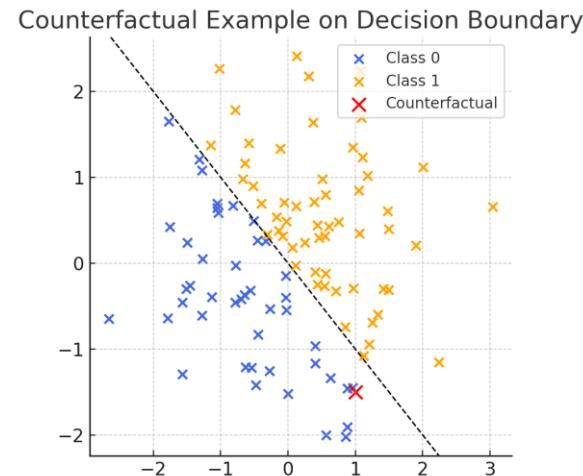
Modify inputs slightly to change model predictions

- | Restrictions on Adversarial Examples:
 - | **Malware:** “*feature changes can cause the application in question to lose its malware functionality in parts or completely*”
 - | Grosse, Kathrin, et al. "Adversarial examples for malware detection." *European symposium on research in computer security*. Cham: Springer International Publishing, 2017.
 - | **Network packets:** “*network packets that do not obey the TCP/IP protocol would not be permissible*”
 - | Sheatsley, Ryan, et al. "Adversarial examples in constrained domains." *arXiv preprint arXiv:2011.01183* (2020).

Counterfactual explanations

Modify inputs slightly to change model predictions

- | Counterfactual example: type of explanation that shows how a small change to the input would lead to a different prediction by the model.



Answer the 'what-if' question for stakeholders
Example: 'If income $\uparrow \$5k \rightarrow$ loan approved'

- | Wachter, Sandra, Brent Mittelstadt, and Chris Russell. "Counterfactual explanations without opening the black box: Automated decisions and the GDPR." *Harv. JL & Tech.* 31 (2017): 841.
- | Guidotti, Riccardo. "Counterfactual explanations and how to find them: literature review and benchmarking." *Data Mining and Knowledge Discovery* 38.5 (2024): 2770-2824.

Counterfactual examples as an “explainable AI” technique

When misclassifications happen

- | DiCE [4,5]: a model-agnostic approach to implementing counterfactual examples
- | Find them as an optimization problem, like constructing adversarial examples:
 - | Feature perturbations alter the outcome of the classification model
 - | Feature perturbations must be diverse and **realistically feasible**
- | From high-level features to HTTP requests?
 **Temporal constraints and dependencies across features**

[4] <https://interpret.ml/DiCE/>

[5] R. K. Mothilal, A. Sharma, C. Tan, *Explaining machine learning classifiers through diverse counterfactual explanations*, in: Proceedings of the 2020 conference on fairness, accountability, and transparency, 2020, pp. 607–617.

Causally and temporally constrained counterfactual explanations

Extending DiCE with a new sampling technique

- | Web logs (input space) to about 100 high-level features (feature space) (web logs)
- | Examples of dependencies and constraints across features:
 1. **Causal constraint:** Increasing the “number of GET requests” feature:
 - | Increases “total number of requests” feature
 - | Alters “GET/POST/PUT ratio” feature
 - | Potentially affect “HTML-to-image ratio” feature
 - | ...
 2. **Temporal constraint:** Aggregating smaller intervals (e.g. 10 s) into larger ones (e.g. 100 s) means that changes in the former feature propagate to the larger

Causally and temporally constrained counterfactual explanations

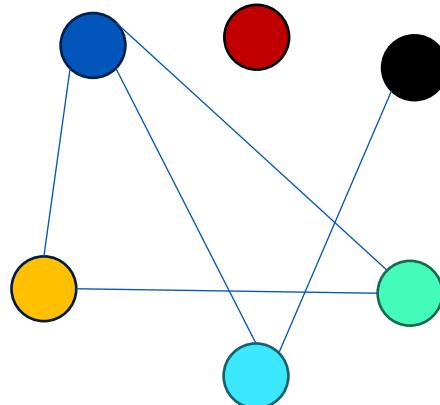
Extending DiCE with a new sampling technique

- | **Solution:**
 1. Extending DiCE with a constrained random sampling method by modeling dependencies across features as a graph
 2. Validity checking and projection
- | **Validity checking:** Is updated feature vector still a counterfactual example?
- | **Projection:** Improve the quality of the counterfactual example by minimizing deviations from the original instance
- | Details: see paper later on

Causally and temporally constrained counterfactual explanations

Extending DiCE with a new sampling technique

- | Extending DiCE with a constrained random sampling method by modeling dependencies across features as a graph



Algorithm 1 Iterative Adjustment of Causally Dependent and Time-Aggregated Features

Require: Feature list F , dependency graph G , original feature vector X_{orig} , perturbed feature vector X_{pert}

Ensure: Updated feature vector X_{final} satisfying dependencies and constraints

```

1:  $X_{final} \leftarrow X_{pert}$ 
2:  $updateQueue \leftarrow$  all perturbed features in  $X_{pert}$ 
3:  $updateCount \leftarrow$  dictionary mapping each feature to 0
4:  $maxUpdates \leftarrow$  predefined upper limit for feature updates
5: while  $updateQueue$  is not empty do
6:    $f \leftarrow$  dequeue from  $updateQueue$ 
7:   if  $updateCount[f] \geq maxUpdates$  then
8:     continue ▷ Avoid infinite loops due to cyclic dependencies
9:   end if
10:   $updateCount[f] \leftarrow updateCount[f] + 1$ 
11:  for each dependent feature  $g$  in  $G$  connected to  $f$  do
12:    Adjust  $X_{final}[g]$  based on  $X_{final}[f]$  and dependencies
13:    if  $g$  not in  $updateQueue$  then
14:      Enqueue  $g$  into  $updateQueue$ 
15:    end if
16:  end for
17:  for each time-aggregated feature  $h$  linked to  $f$  do
18:    Adjust  $X_{final}[h]$  to maintain consistency over time scales
19:    if  $h$  not in  $updateQueue$  then
20:      Enqueue  $h$  into  $updateQueue$ 
21:    end if
22:  end for
23: end while
24: return  $X_{final}$ 
  
```

Counterfactual explanations

Modify inputs slightly to change model predictions

```
1 X_test = pd.read_pickle("x_test.pkl.gz")
2 Y_test = pd.read_pickle("y_test.pkl.gz")
3 test_dataset = pd.concat([X_test, Y_test], axis=1)
4 column_names = list(X_test.columns)
5
6 num_features = [f for f in X_test.columns if "_10" in f]
7 num_features.extend(["response_size", "request_size"])
8
9 # initiate DiCE
10 d = dice_ml.Data(dataframe=train_dataset, continuous_features=num_features, outcome_name="label")
11 model = pickle.load(open("catboost.pkl.gz", "rb"))
12 m = dice_ml.Model(model=model, backend="sklearn")
13 exp = dice_ml.Dice(d, m, method="constrained_random") # custom sampling method
14
15 query_instance = test_dataset[test_dataset["label"] == 1].sample(n=1)
16 query_instance.drop(["label"], axis=1, inplace=True)
17
18 # generate counterfactuals
19 dice_exp = exp.generate_counterfactuals(query_instance, total_CFs=1, desired_class="opposite", proximity_weight=1.5,
20                                         diversity_weight=1.0, features_to_vary=column_names, posthoc_sparsity_param=0)
21
22 print(dice_exp.visualize_as_list(show_only_changes=True))
```

A hand is shown from the side, holding a glowing, translucent globe. The globe is covered in a network of blue and white lines, representing latitude and longitude or a similar coordinate system. The background is dark, with several out-of-focus, glowing circular lights in shades of orange, yellow, and blue, creating a bokeh effect. The overall atmosphere is futuristic and technological.

Are adversarial and counterfactual
examples the same?

Commonalities

Modify inputs slightly to change model predictions

Aspect	Adversarial Examples	Counterfactual Explanations
Minimal Perturbation	Both focus on small changes to the input	Both find the smallest change possible
Boundary Crossing	Crosses the model's decision boundary	Also crosses the decision boundary, but for explanation
Optimization Problem	Often solved by optimizing a loss function	Similarly uses optimization to find minimal modifications
Input Similarity	Changes are typically imperceptible or minimal	Changes are designed to be close to the original input

Key differences

Modify inputs slightly to change model predictions

Aspect	Adversarial Examples	Counterfactual Explanations
Primary Goal	Fool the model	Explain the model's decision
Perspective	Attacker-centric	User-centric
Optimization Objective	Minimize input perturbation while maximizing misclassification	Minimize input perturbation while achieving a specific target outcome
Plausibility of Change	Perturbations may be unrealistic (e.g., random pixel noise)	Changes must be plausible and actionable (e.g., higher income, lower debt)
Output Usefulness	Rarely interpretable for humans	Designed for human interpretability
Applications	Model robustness, security testing, adversarial training	Explainable AI, fairness auditing, regulatory compliance

Bridging explainability and adversarial-generation methods

Modify inputs slightly to change model predictions

- | Towards a unified view on perturbations:
 - | **Adversarial training** improves model resilience by exposing it to worst-case perturbations. However, adversarial examples are often unnatural and hard for humans to interpret (e.g., tiny imperceptible pixel changes).
 - | **Counterfactual explanations**, on the other hand, generate perturbations aligned with human reasoning (e.g., “If this feature changed slightly, the prediction would flip”).
 - | A hybrid approach ensures the model is robust to malicious, imperceptible changes while also being transparent in how decisions could change under meaningful, human-aligned modifications.

Bridging explainability and adversarial-generation methods

[Modify inputs slightly to change model predictions](#)

| Unified objective (minimize):

$$\mathcal{L}(x') = \underbrace{\lambda_1 d(x, x')}_{\text{proximity}} + \underbrace{\lambda_2 \ell_{\text{task}}(f(x'), y_{\text{target}})}_{\text{target attainment}} + \underbrace{\lambda_3 \phi(C(x'))}_{\text{plausibility penalty}} + \underbrace{\lambda_4 \psi(A(x'))}_{\text{actionability penalty}} + \underbrace{\lambda_5 H(x')}_{\text{human preference / safety}}$$

- | x = factual input, $y = f(x)$ predicted label or score.
- | x' = candidate modification (adversarial or counterfactual).
- | $d(\cdot, \cdot)$ = distance on input or latent manifold (e.g., L2 in latent space or edit distance in feature space).
- | $C(\cdot)$ = plausibility / causal constraints (learned density model, generative prior, causal graph).
- | $A(\cdot)$ = actionability constraints (which features a user can change).
- | $H(\cdot)$ = human feedback / preference penalty (from HITL)

- | Goldwasser, Jeremy, and Giles Hooker. "Unifying Image Counterfactuals and Feature Attributions with Latent-Space Adversarial Attacks." *arXiv preprint arXiv:2504.15479* (2025).

Bridging explainability and adversarial-generation methods

Modify inputs slightly to change model predictions

| Unified objective (minimize):

$$\mathcal{L}(x') = \underbrace{\lambda_1 d(x, x')}_{\text{proximity}} + \underbrace{\lambda_2 \ell_{\text{task}}(f(x'), y_{\text{target}})}_{\text{target attainment}} + \underbrace{\lambda_3 \phi(C(x'))}_{\text{plausibility penalty}} + \underbrace{\lambda_4 \psi(A(x'))}_{\text{actionability penalty}} + \underbrace{\lambda_5 H(x')}_{\text{human preference / safety}}$$

- | For **adversarial examples**, you set λ_2 to encourage any misclassification and usually omit C, ψ, H (or set them weak).
- | For **counterfactuals**, you set λ_2 to target a specific outcome and **strongly enforce** C, ψ, H so outputs are plausible, causal, and actionable.
- | Training and evaluation can share the same solver (gradient attack in latent space, mixed-integer solver for discrete features, or search over generative-model latent codes).

Experimental setup and evaluation

A human-in-the-loop analysis

- | **Generation and retraining with counterfactual examples:** A HTTP request initially classified as malicious (label = 1) needs to be modified to benign (label = 0)

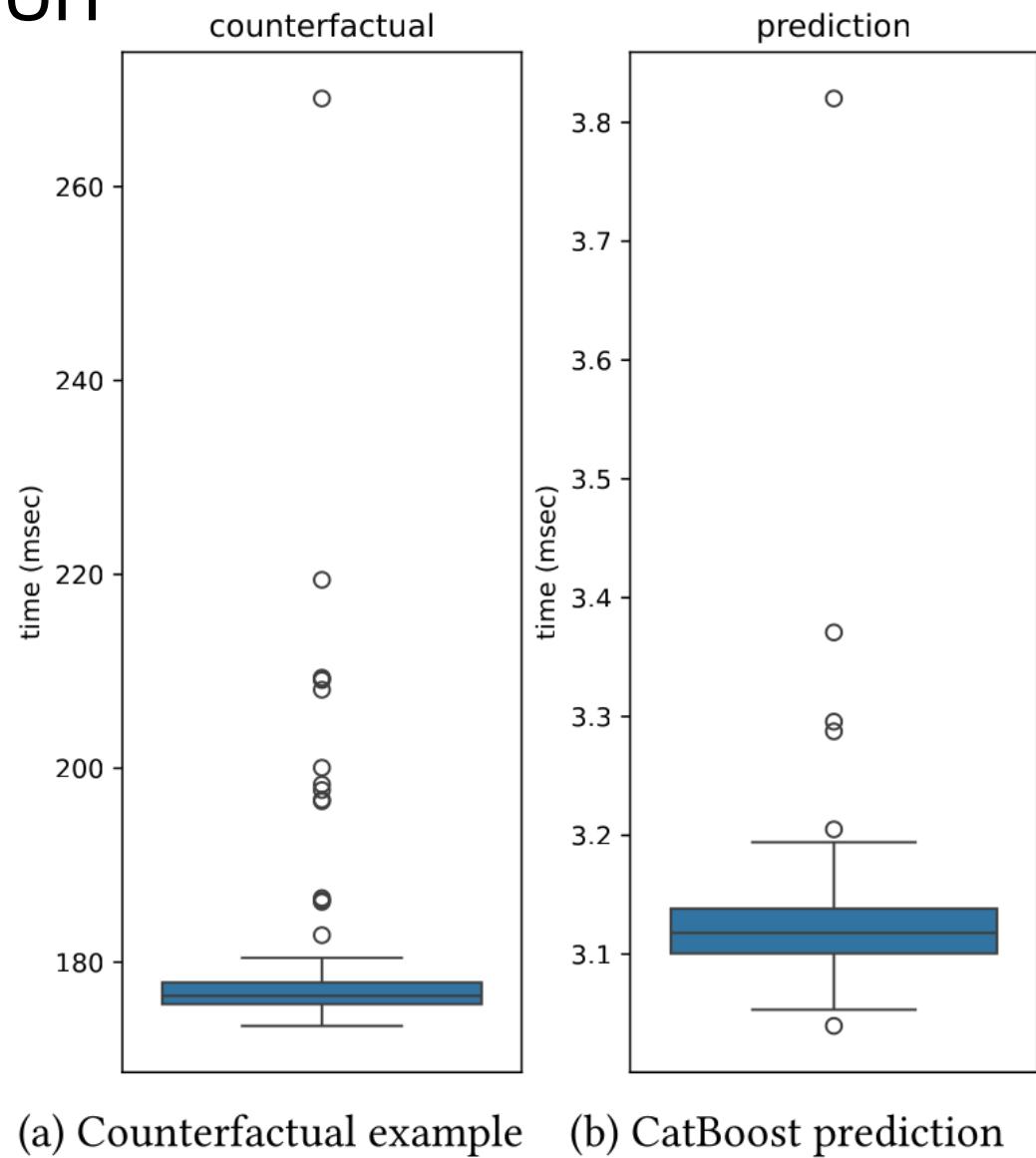
Feature	Original Value	Target Value
html_pages_10	21.0	1.5
label	1	0

- | A malicious bot sent 21 HTTP requests for HTML pages within a 10-second window, significantly exceeding the browsing speed of a typical user
- | Such a perturbation would not be realistic, but a perturbation to 20.0 might
- | Multiple attempts required to generate counterfactual examples!

Experimental setup and evaluation

A human-in-the-loop analysis

- | Benchmarking time: counterfactual example vs. CatBoost prediction on a system with:
 - | 12th Gen Intel Core i9-12900 CPU
 - | 32 GB of memory
- | Counterfactual example generation is almost **two orders of magnitude** slower
- | Why? Costly iterative approach to optimizing counterfactual examples
- | Only effective under selective generation



Ongoing work

Adversarial and counterfactual examples

- | Practical HITL loop (iterative): simultaneously improves robustness and the quality of counterfactual explanations
 - | Generate candidate x' (adversarial or counterfactual) using latent-space perturbation or constrained search.
 - | Present top-k candidates via UI with explanations of changed features.
 - | Collect human feedback: accept/reject, plausibility score, improvement suggestions, or alternative target choices.
 - | Update models: (i) retrain plausibility/actionability scorer, (ii) adversarially-train the predictor using accepted adversarial seeds, (iii) refine generator to prefer human-validated CFs.

Future directions

Unifying Adversarial Examples, Counterfactual Explanations & HITL

| Causality-Aware Explanations

- | Integrate causal models into counterfactual and adversarial generation.
- | Ensure perturbations respect real-world causal dependencies.
- | Study counterfactual fairness — guaranteeing equitable outcomes across groups.

| Robustness + Explainability Co-Training

- | Combine adversarial training with counterfactual recourse to improve both robustness and interpretability.
- | Develop multi-objective training pipelines balancing:
 - | Prediction accuracy
 - | Robustness against attacks
 - | Quality of explanations

Summary & Takeaways

Adversarial examples and counterfactual explanations

- | Adversarial robustness and explainability are interlinked.
 - | Both AEs and CFs involve minimal input perturbations
 - | Both operate near decision boundaries
 - | Both can be formulated as optimization problems
- | Why unification matters
 - | For Research: Shared frameworks improve understanding of decision boundaries.
 - | For Practice: Combining adversarial training + counterfactual recourse → robust, interpretable, and user-friendly models.
 - | With HITL: Humans guide plausibility, fairness, and actionability, bridging security and interpretability goals.
- | Unifying them yields more robust and interpretable AI systems.

References

Next steps

- | Wachter, Sandra, Brent Mittelstadt, and Chris Russell. "Counterfactual explanations without opening the black box: Automated decisions and the GDPR." *Harv. JL & Tech.* 31 (2017): 841.
- | Pawelczyk, Martin, et al. "Exploring counterfactual explanations through the lens of adversarial examples: A theoretical and empirical analysis." *International Conference on Artificial Intelligence and Statistics*. PMLR, 2022.
- | Guidotti, Riccardo. "Counterfactual explanations and how to find them: literature review and benchmarking." *Data Mining and Knowledge Discovery* 38.5 (2024): 2770-2824.
- | Goldwasser, Jeremy, and Giles Hooker. "Unifying Image Counterfactuals and Feature Attributions with Latent-Space Adversarial Attacks." *arXiv preprint arXiv:2504.15479* (2025).

35

Thanks!

Any questions?

Davy Preuveneers

Davy.Preuveneers@kuleuven.be

