

# Introduction to Counterfactual Explanations

Mateusz Krzyżiński

Bartek Sobieski



Warsaw, April 8th, 2024

# Agenda

1. Intuition & definition
2. Properties
3. CEs for tabular data
4. CEs for images

Note, much of this presentation is based on the article:  
R. Guidotti, ***Counterfactual explanations and how to find them: literature review and benchmarking***,  
Data Mining and Knowledge Discovery, 2022.

# **What are Counterfactual Explanations?**

# Counterfactual thinking



AMERICAN PSYCHOLOGICAL ASSOCIATION

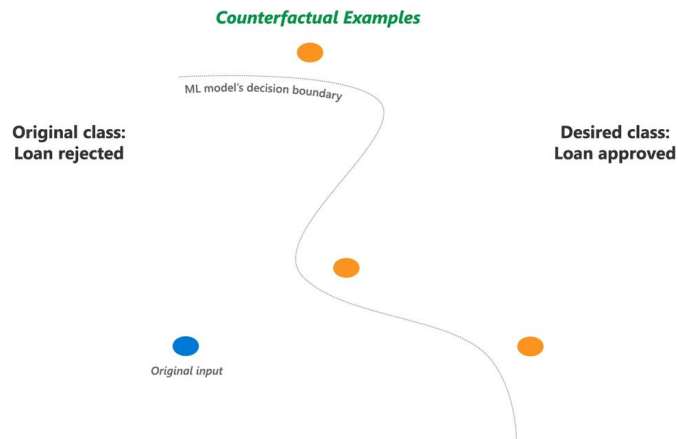
## counterfactual thinking

Updated on 04/19/2018

1. imagining ways in which events in one's life might have turned out differently. This often involves feelings of regret or disappointment (e.g., *If only I hadn't been so hasty*) but may also involve a sense of relief, as at a narrow escape (e.g., *If I had been standing three feet to the left...*).
2. any process of reasoning based on a conditional statement of the type "If X, then Y" where X is known to be contrary to fact, impossible, or incapable of empirical verification. Counterfactual thinking of the first sort is common in such historical speculations as *If Hitler had been killed in July 1944, then ...* . Counterfactual thinking of the second and third types can play a useful role in evaluating the implications of a theory or [heuristic](#) and in [thought experiments](#). See also [as-if hypothesis](#); [conditional reasoning](#).

# Counterfactual explanations - intuition

**Counterfactual explanations** suggest what should be different in the input instance to change the prediction of a model.



Amit Sharma, DiCE, Microsoft Research Blog

# Counterfactual explanations - formalization

Given a model  $f$  that outputs the prediction  $y = f(x)$  for an instance  $x$ , a counterfactual explanation consists of an instance  $x'$  such that the prediction for  $f$  on  $x'$  is different from  $y$ , i.e.,  $y \neq f(x')$ , and such that the difference between  $x$  and  $x'$  is *minimal*.

- ❑ post-hoc
- ❑ example-based
- ❑ local\*
- ❑ both model agnostic and model specific

\*sets of CEs and their aggregations can be considered global explanations

**We will discuss global approach later in this track.**

# Counterfactual explanations - first method

Counterfactual instance is to be found based on the optimization:

$$x_c = \arg \min_{x'} \{ \max_{\lambda} \lambda (\hat{f}(x') - y')^2 + d(x, x') \}$$



parameter to balance  
the trade-off

desired prediction

distance measure  
L1 weighted with the inverse  
median absolute deviation

**There are only(?) two properties considered here.**

S. Wachter, B. Mittelstadt, C. Russell, *Counterfactual Explanations without Opening the Black Box: Automated Decisions and the GDPR*, Harvard Journal of Law & Technology, 2018.

# **Properties of Counterfactual Explanations**



# Desirable properties

## Validity.

A counterfactual  $x'$  is valid iff it actually changes the prediction with respect to the original one, i.e.,  $f(x) \neq f(x')$

## Minimality (Sparsity).

$x'$  is minimal iff there is no other valid example  $x''$  with the smaller number of different attributes w.r.t.  $x$ , i.e.,  $\forall x'' \quad |\delta_{x,x''}| \geq |\delta_{x,x'}|$

## Similarity (Proximity).

$x'$  should be similar to  $x$ , i.e.,  $d(x, x') < \varepsilon$  given a distance function  $d$  and a predefined threshold  $\varepsilon$

# Desirable properties cont'd

## Plausibility.

A counterfactual  $x'$  is plausible iff it is coherent with a reference sample, i.e., it is not an outlier, out-of-distribution example.

## Actionability.

A counterfactual  $x'$  is actionable iff all the differences between  $x'$  and  $x$  are related to actionable features that can be mutated/changed.

## Diversity.

If a set of counterfactual examples is returned, it should be formed by diverse examples (similar to  $x$  but most different between each other).

## Desirable properties cont'd

**Discriminative power.** (Warning! Highly subjective)

A counterfactual  $x'$  must help in figuring out why different prediction can be obtained with it. Difficult to quantify without human-based evaluation.

**Causality.**

A counterfactual  $x'$  should maintain any known causal relationship between features (related to plausibility).

# **How to retrieve Counterfactual Explanations?**

# Strategies for finding counterfactuals

## Optimization

minimizing a loss function that accounts for desired properties

## Heuristic Search

local and heuristic choices performed iteratively

## Instance-Based

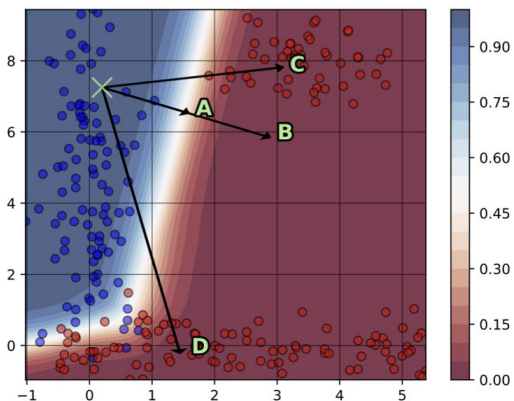
selecting the best example from a reference dataset

## Decision Tree

exploiting the structure of a decision tree approximating the black-box

# Strategies for finding counterfactuals

## Instance-Based



## FACE: Feasible and Actionable Counterfactual Explanations

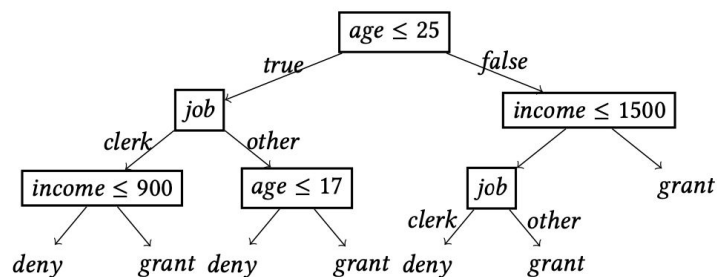
1. Construction of a weighted graph over data points (KDE, kNN,  $\epsilon$ -graph).
2. Preparing the list of candidate targets (based on additional constraints).
3. Finding shortest path in a constructed graph for each candidate target.
4. Selecting the candidate with the shortest path.

**We will discuss attacks  
on instance-based CEs next week.**

R. Poyiadzi et al., *FACE: Feasible and Actionable Counterfactual Explanations*, AIES, 2020.

# Strategies for finding counterfactuals

## Decision Tree



## LORE: Local Rule-Based Explanations

1. Generate a set of synthetic samples from the neighbourhood of  $x$  through a genetic algorithm.
2. Train a decision tree on this set labeled with the black-box predictions.
3. Extract the counterfactual rules from a trained decision tree (those with the smallest number of split conditions not satisfied by  $x$ ).

R. Guidotti et al., *Local Rule-Based Explanations of Black Box Decision Systems*, IEEE Intelligent Systems, 2019.

# Counterfactual explanations in computer vision



# Definition

Given a **classifier** and an **image**, what is the **minimal semantic change** that **flips the model's decision**.



Prediction: **smiling**



Prediction: **not smiling**

# Tabular vs visual

# Curse of dimensionality

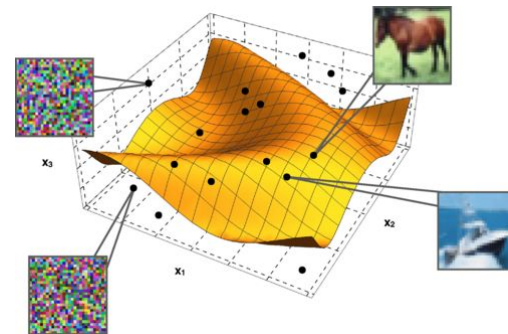
	Tabular	Visual
Dimensionality	~ 10 - 1000 features	~ 50 000 - 200 000 pixels

Visual data is **much more sparse** than tabular data

# Manifolds differ

Given an observation of each modality,  
move along a random direction:

$$\mathbf{x} + \alpha \cdot \mathbf{d}, \text{ where } \|\mathbf{d}\| = 1, \alpha = \beta \cdot n_{\text{features}}$$

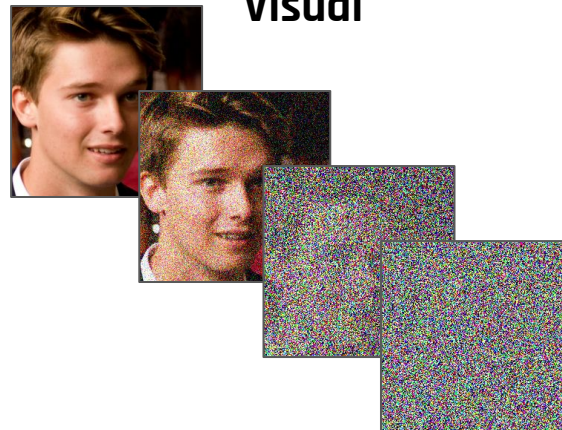


## Tabular

**X**

	Age: 44	Salary (k): 15	Flat size (m <sup>2</sup> ): 96	Savings (k): 12
$\beta = 0.001$	Age: 43.99	Salary (k): 14.9994	Flat size (m <sup>2</sup> ): 95.999	Savings (k): 11.999
$\beta = 0.01$	Age: 43.96	Salary (k): 14.994	Flat size (m <sup>2</sup> ): 95.99	Savings (k): 11.99
$\beta = 0.1$	Age: 43.65	Salary (k): 14.94	Flat size (m <sup>2</sup> ): 95.9	Savings (k): 11.88

## Visual




**X**

$$\beta = 0.000001$$

$$\beta = 0.00001$$

$$\beta = 0.0001$$

# Modality-specific differences

$$x_c = \arg \min_{x'} \{ \max_{\lambda} \lambda (\hat{f}(x') - y')^2 + d(x, x') \}$$


The diagram shows two blue arrows pointing upwards from question marks. The first arrow points to the term  $\hat{f}(x')$  in the equation, and the second arrow points to the term  $d(x, x')$ . This indicates that the model's output and the distance metric are the focus of the discussion on modality-specific differences.

How can the differences between visual and tabular data be accounted for in the optimization objective?

# Measuring perceptual distance

# Semantic similarity

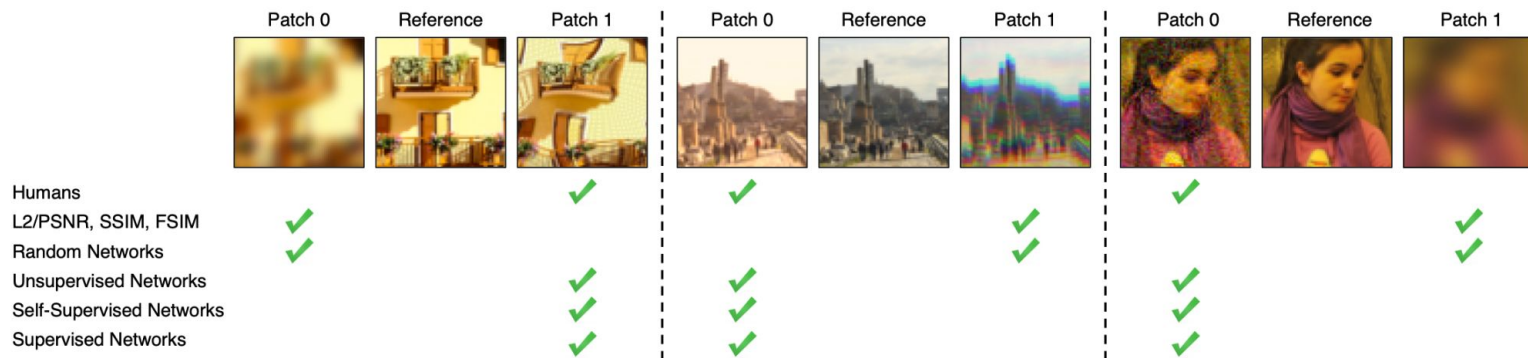
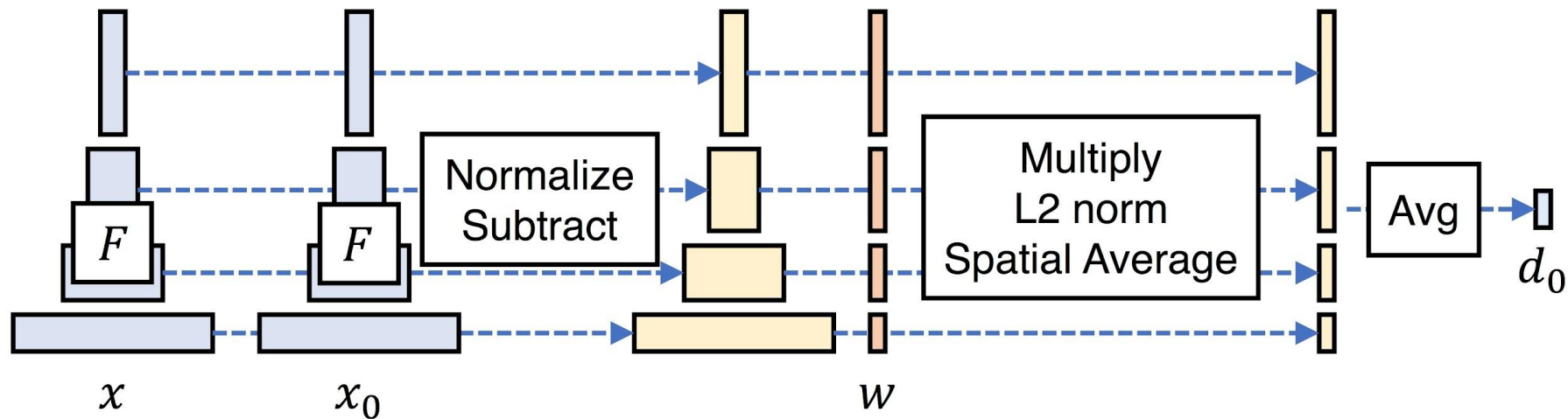


Figure 1: **Which patch (left or right) is “closer” to the middle patch in these examples?** In each case, the traditional metrics (L2/PSNR, SSIM, FSIM) disagree with human judgments. But deep networks, even across architectures (Squeezenet [20], AlexNet [27], VGG [52]) and supervision type (supervised [47], self-supervised [13, 40, 43, 64], and even unsupervised [26]), provide an *emergent embedding* which agrees surprisingly well with humans. We further calibrate existing deep embeddings on a large-scale database of perceptual judgments; models and data can be found at <https://www.github.com/richzhang/PerceptualSimilarity>.

R. Zhang, P. Isola, A. Efros, E. Shechtman, O. Wang, *The Unreasonable Effectiveness of Deep Features as a Perceptual Metric*, CVPR, 2018.

# LPIPS

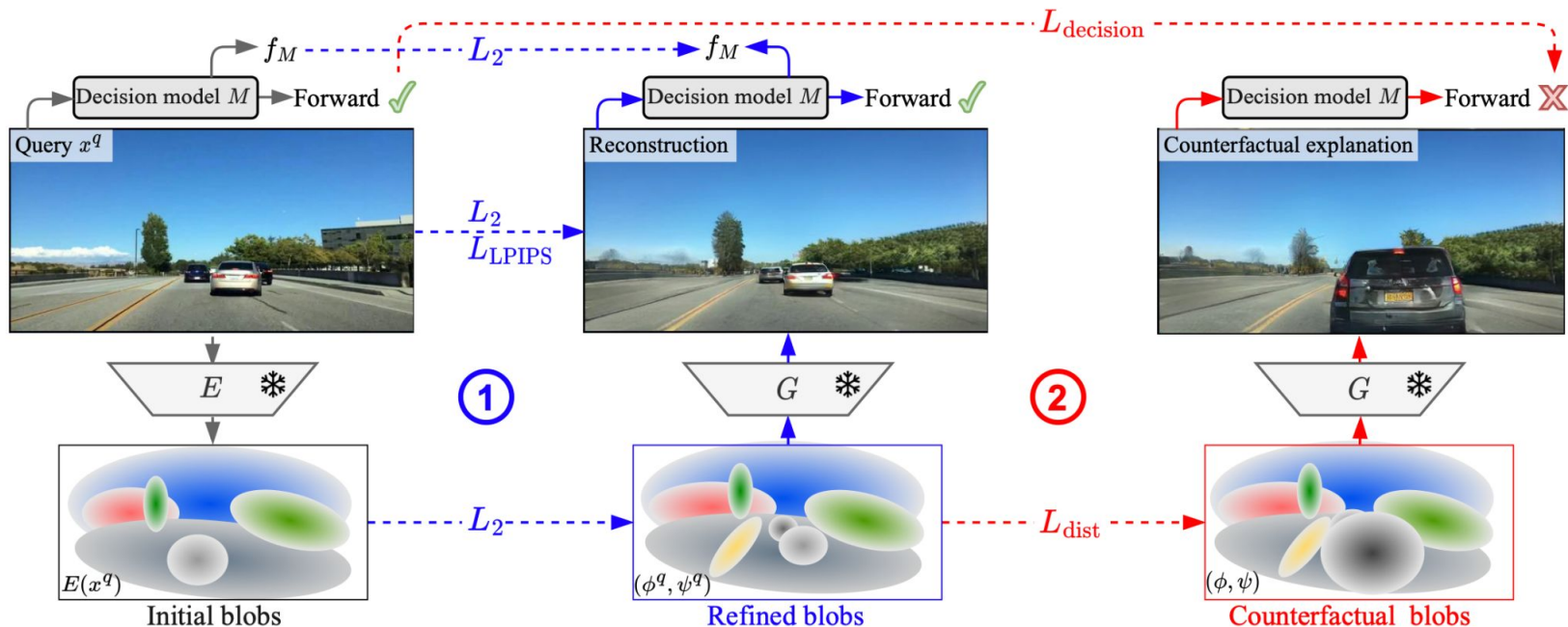


R. Zhang, P. Isola, A. Efros, E. Shechtman, O. Wang, *The Unreasonable Effectiveness of Deep Features as a Perceptual Metric*, CVPR, 2018.



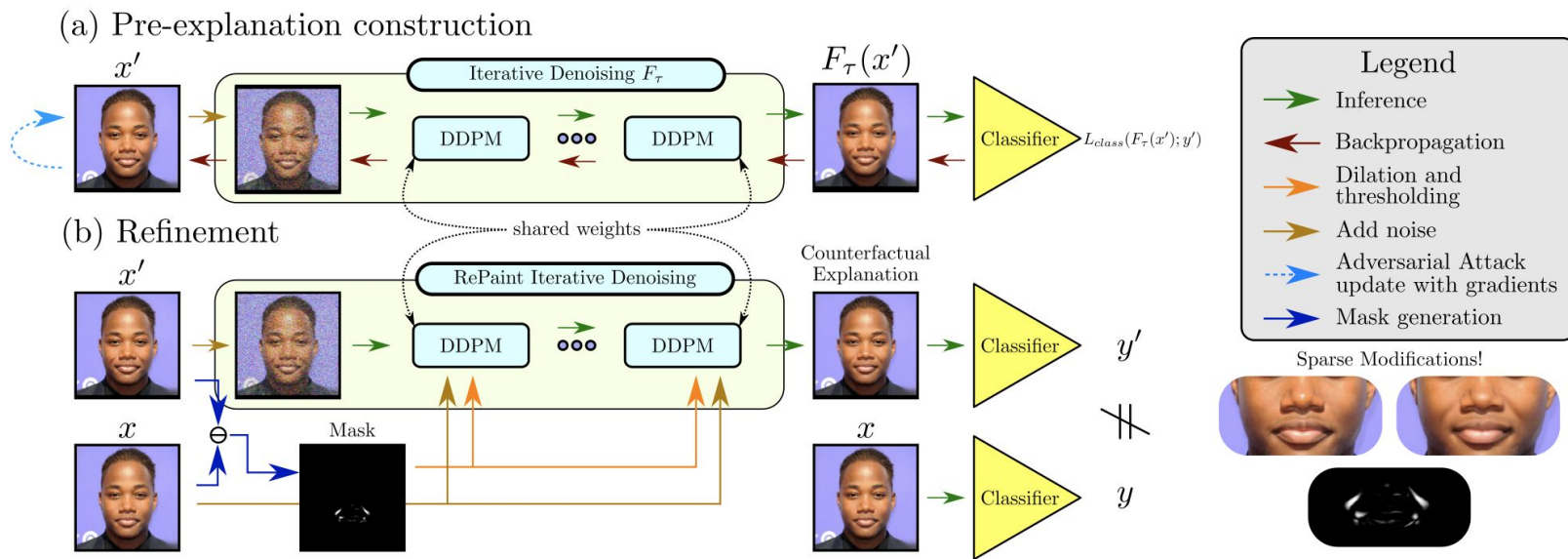
# **State-of-the-art in generating explanations**

# OCTET: Object-aware Counterfactual Explanations



M. Zemni, M. Chen, E. Zablocki, H. Ben-Younes, P. Perez, M. Cord, *OCTET: Object-aware Counterfactual Explanations*, CVPR, 2023.

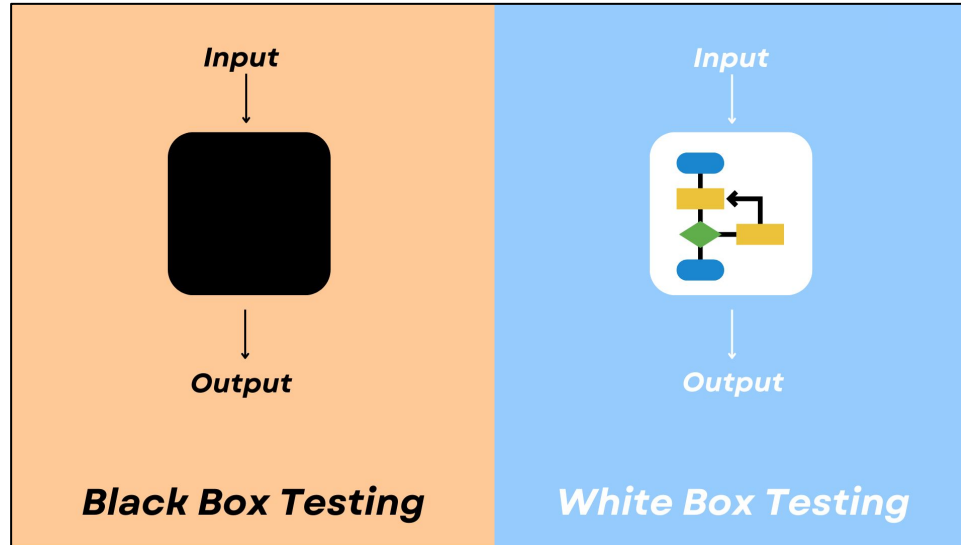
# Adversarial Counterfactual Visual Explanations

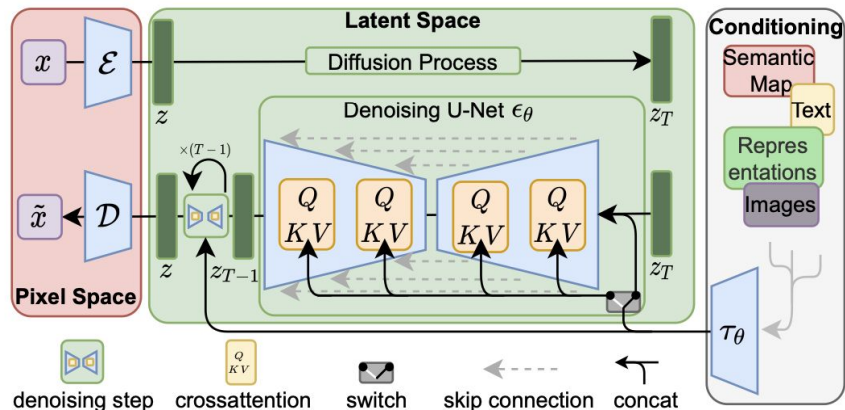


G. Jeanneret, L. Simon, F. Jurie, *Adversarial Counterfactual Visual Explanations*, CVPR, 2023.

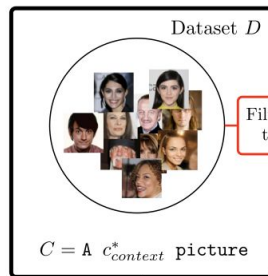
# Current challenges

# Lack of black-box methods



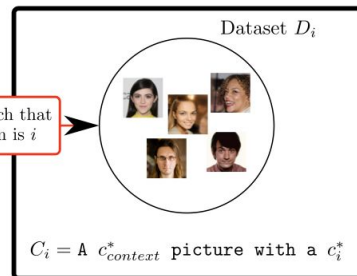


(a) Context Bias Learning



Textual Inversion for  $c_{context}^*$  using  $C$  and the dataset  $D$

(b) Class Bias Learning



Textual Inversion for  $c_i^*$  using  $C_i$  and the filtered dataset  $D_i$

(c) Counterfactual Generation from  $i$  to  $j$

Input (Classified as  $i$ )



EDICT Inversion



Output (Classified as  $j$ )



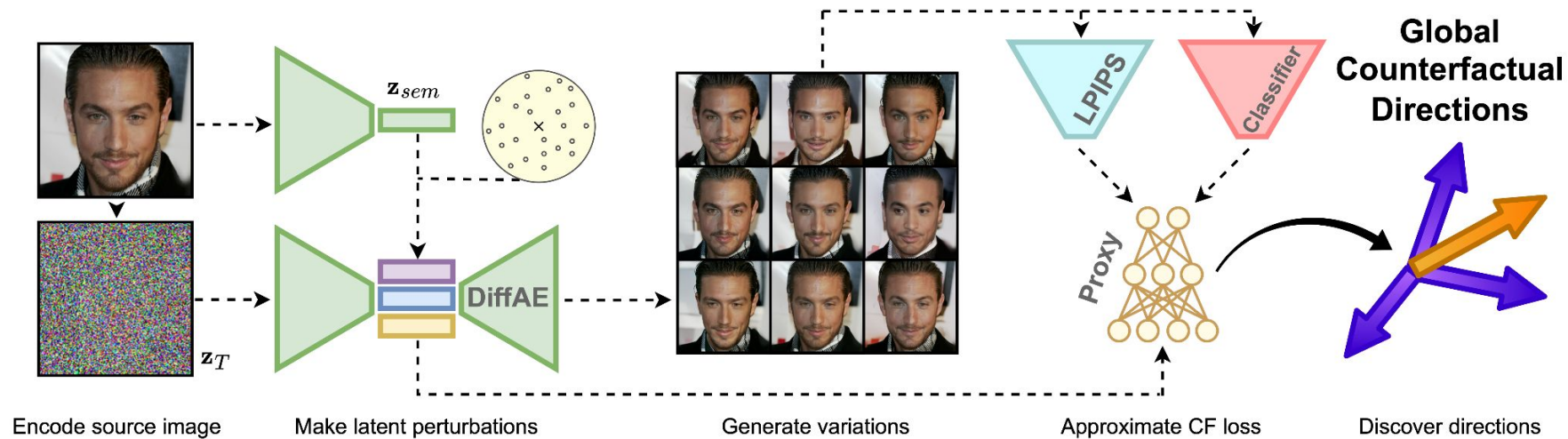
EDICT Denoising

Use  $C_i$  as the positive prompt and  $C_j$  as the negative

Flip  $C_i$  and  $C_j$

G. Jeanneret, L. Simon, F. Jurie, *Text-to-image Models for Counterfactual Explanations: a Black-Box approach*, WACV, 2024.

# GCD



# Problematic evaluation

**There is no perfect metric to evaluate visual counterfactual explanations and many of the existing ones are just noise**



**Thank you for attention**