

Rise of Causality in Computer Vision

The repository contains lists of papers on causality and how relevant techniques are being used to further enhance deep learning era computer vision solutions.

The repository is organized by [Maheep Chaudhary](#) and [Haohan Wang](#) as an effort to collect and read relevant papers and to hopefully serve the public as a collection of relevant resources.

Causality

- [The Seven Tools of Causal Inference with Reflections on Machine Learning](#)
 - ► Maheep's notes
- [On Pearl's Hierarchy and the Foundations of Causal Inference](#)
 - ► Maheep's notes
- [Unit selection based on counterfactual logic](#)
 - ► Maheep's notes

$\text{argmax } c \beta P(\text{complier}|c) + \gamma P(\text{always-taker}|c) + \theta P(\text{never-taker}|c) + \delta P(\text{defier}|c)$

where benefit of selecting a complier is β , the benefit of selecting an always-taker is γ , the benefit of selecting a never-taker is θ , and the benefit of selecting a defier is δ . Our objective, then, should be to find c .

Theorem 4 says that if Y is monotonic and satisfies gain equality then the benefit function may be defined as: -

$(\beta - \theta)P(y, x|z) + (\gamma - \beta)P(y, x'|z) + \theta$

"Third, the proposed approach could be used to evaluate machine learning models as well as to generate labels for machine learning models. The accuracy of such a machine learning model would be higher because it would consider the counterfactual scenarios."

[Sundar et al., 1998; Blumenthal et al., 2001; Winer,

2001; Resnick et al., 2006; Lewis and Reiley, 2014]

- [Unit Selection with Causal Diagram](#)
 - ► Maheep's notes

$W + \sigma U \leq f \leq W + \sigma L$ if $\sigma < 0$, $W + \sigma L \leq f \leq W + \sigma U$ if $\sigma > 0$, where " f " is the objective function. Previously in normal case the objective function is bounded by the equation:

$\max\{p_1, p_2, p_3, p_4\} \leq f \leq \min\{p_5, p_6, p_7, p_8\}$ if $\sigma < 0$,
 $\max\{p_5, p_6, p_7, p_8\} \leq f \leq \min\{p_1, p_2, p_3, p_4\}$ if $\sigma > 0$,

In the extension of the same the author proposes the new situations which arise such as the when " z " is partially observable. and if " z " is a pure mediator.

The author then discusses about the availability of the observational and experimental data. If we only have experimental data then we can simply remove the observational terms in the theorem

$$\max\{p_1, p_2\} \leq f \leq \min\{p_3, p_4\} \text{ if } \sigma < 0,$$

$$\max\{p_3, p_4\} \leq f \leq \min\{p_1, p_2\} \text{ if } \sigma > 0,$$

but if we have only observational data then we can take use of the observed back-door and front-door variables to generate the experimental data, but if we have partially observable back-door and front-door variables then we can use the equation:

$$LB \leq P(y|do(x)) \leq UB$$

The last topic which author discusses about is the reduction of the dimensionality of the variable "z" which satisfies the back-door and front-door variable by substituting the causal graph by substituting "z" by "W" and "U" which satisfies the condition that "no_of_states_of_W * no_of_states_of_U = no_of_states_of_z".

- [The Causal-Neural Connection: Expressiveness, Learnability, and Inference](#)

- ► Maheep's notes

The author proposes Neural Causal Models, that are a type of SCM but are capable of amending Gradient Descent. The author proposes the network to solve two kinds of problems, i.e. "causal effect identification" and "estimation" simultaneously in a Neural Network as generative model acting as a proxy for SCM.

"causal estimation" is the process of identifying the effect of different variables

"Identification" is obtained when we apply backdoor criterion or any other step to get a better insight. The power of identification has been seen by us as seen in the papers of Hanwang Zhang.

Theorem 1: There exists a NCM that is in sync with the SCM on ladder 3

- [The Causal Loss: Driving Correlation to Imply Causation\(autonomous\)](#)

- ► Maheep's notes

- [Double Machine Learning Density Estimation for Local Treatment Effects with Instruments](#)

- ► Maheep's notes

The author argues that by obtaining the PDF may give very valuable information as compared to only estimating the Cumulative Distribution Function.

Causality & Computer Vision

- [Counterfactual Samples Synthesizing and Training for Robust Visual Question Answering](#)

- ► Maheep's notes

- [How Should Pre-Trained Language Models Be Fine-Tuned Towards Adversarial Robustness?](#)

- ► Maheep's notes

$$I(S; Y, T) = I(S; Y) + I(S; T|Y),$$

The author proposes by this equation that the two models overlap, i.e. the objective model

and the pretrained model. S represents the features extracted the model by the objective model and T is the features extracted by the pretrained model.

- [Counterfactual Zero-Shot and Open-Set Visual Recognition](#)

- ► Maheep's notes

ZSL is usually provided with an auxiliary set of class attributes to describe each seen- and unseen-class whereas the OSR has open environment setting with no information on the unseen-classes [51, 52], and the goal is to build a classifier for seen-classes. The author describes previous works in which the generated samples from the class attribute of an unseen-class, do not lie in the sample domain between the ground truth seen and unseen, i.e., they resemble neither the seen nor the unseen. As a result, the seen/unseen boundary learned from the generated unseen and the true seen samples is imbalanced.

The author proposes a technique using counterfactual, i.e. to generate samples using the class attributes, i.e. Y and sample attribute Z by the counterfactual, i.e. X would be \tilde{x} , had Y been y , given the fact that $Z = z(X = x)$ and the consistency rule defines that if the ground truth is Y then \tilde{x} would be x . The proposed generative causal model $P(X|Z, Y)$ generate examples for ZSL and OSR.

- [Counterfactual VQA: A Cause-Effect Look at Language Bias](#)

- ► Maheep's notes

- [CONTERFACTUAL GENERATIVE ZERO-SHOT SEMANTIC SEGMENTATION](#)

- ► Maheep's Notes

The model will contain a total of 4 variables R, W, F and L . The generator will to generate the fake features using the word embeddings and real features of the seen class and will generate fake images using word embeddings after learning. However, this traditional model cannot capture the pure effect of real features on the label because the real features R not only determine the label L by the link $R \rightarrow L$ but also indirectly influence the label by path $R \rightarrow F \rightarrow L$. This structure, a.k.a. confounder. Therefore they remove the $R \rightarrow F \rightarrow L$ and let it be $W \rightarrow F \rightarrow L$, removing the confounding effect of F . Also they use GCN to generate the image or fake features from word embeddings using the GCN which also provides to let the generator learn from similar classes.

- [Adversarial Visual Robustness by Causal Intervention](#)

- ► Maheep's notes

1. Augments the image with multiple retinotopic centres

2. Encourage the model to learn causal features, rather than local confounding patterns.

They propose the model to be such that $\max P(Y = \hat{y}|X = x + \delta) - P(Y = \hat{y}|\text{do}(X = x + \delta))$, subject to $P(Y = \hat{y}|\text{do}(X = x + \delta)) = P(Y = \hat{y}|\text{do}(X = x))$, in other words they focus on annihilating the confounders using the retinotopic centres as the instrumental variable.

- [What If We Could Not See? Counterfactual Analysis for Egocentric Action Anticipation](#)

- ► Maheep's notes

They ask this question so as to only get the effect of semantic label. As the visual feature is the main feature the semantic label can act as a confounder due to some situations occurring frequently. Therefore the author proposes to get the logits "A" from the pipeline without making any changes to the model and then also getting the logits "B" when they provide a random value to visual feature denoting the question of counterfactual, i.e. "what action would be predicted if we had not observed any visual representation?" getting the unbiased logit by:

Unbiased logit = A - B

- [Transporting Causal Mechanisms for Unsupervised Domain Adaptation](#)

- ► Maheep's notes

- [WHEN CAUSAL INTERVENTION MEETS ADVERSARIAL EXAMPLES AND IMAGE MASKING FOR DEEP NEURAL NETWORKS](#)

- ► Maheep's notes

$\text{Effect}(x_i \text{ on } x_j, Z) = P(x_j | \text{do}(x_i_dash), Z_{-}x_i) - P(x_j | Z_{-}x_i) \dots\dots\dots(1)$ The expected casual effect has been defined as: $E_{x_i}[\text{Effect}(x_i \text{ on } x_j, Z)] = (P(X_i = x_i | Z) * (\text{equation}_1))$

The author proposes three losses to get the above equaitons, i.e. the effect of pixels. The losses are interpretability loss, shallow reconstruction loss, and deep reconstruction loss. Shallow reconstruction loss is simply the L 1 norm of the difference between the input and output of autoencoder to represent the activations of the network. For the second equation they applied the deep reconstruction loss in the form of the KL-divergence between the output probability distribution of original and autoencoder-inserted network.

These losses are produced afer perturbtaing the images by maksing the images and inserting adversarial noise.

- [Interventional Few-Shot Learning](#)

- ► Maheep's notes

The author proposes the solution by proposing 4 variables, i.e. "D", "X", "C", "Y" where D is the pretrained model, X is the feature representaiton of the image, C is the low dimesion representation of X and Y are the logits. The author says the D affects both the X and C, also X affects C, X and C affects the logit Y. The autho removes the affect of D on X using backdoor.

- [CLEVRER: COLLISION EVENTS FOR VIDEO REPRESENTATION AND REASONING](#)

- ► Maheep's notes

The dataset is build on CLEVR dataset and has predictive both predictive and counterfactual questions, i.e. done by, Predictive questions test a model's capability of predicting possible occurrences of future events after the video ends. Counterfactual questions query the outcome of the video under certain hypothetical conditions (e.g. removing one of the objects). Models need to select the events that would or would not happen under the designated

condition. There are at most four options for each question. The numbers of correct and incorrect options are balanced. Both predictive and counterfactual questions require knowledge of object dynamics underlying the videos and the ability to imagine and reason about unobserved events.

The dataset is being prepared by using the physics simulation engine.

- [Towards Robust Classification Model by Counterfactual and Invariant Data Generation](#)

- ► Maheep's notes

They augment using the augmentations as: None, CF(Grey), CF(Random), CF(Shuffle), CF(Tile), CF(CAGAN) and the augmentations which alter the invariant features using: F(Random) F(Shuffle) F(Mixed-Rand) F(FGSM)

- [Unbiased Scene Graph Generation from Biased Training](#)

- ► Maheep's Notes

1. To take remove the context bias, the author compares it with the counterfactual scene, where visual features are wiped out(containing no objects).

The author argues that the true label is influenced by Image(whole content of the image) and context(individual objects, the model make a bias that the object is only to sit or stand for and make a bias for it) as confounders, whereas we only need the Content(object pairs) to make the true prediction. The author proposes the $TDE = y_e - y_e(x_{bar}, z_e)$, the first term denote the logits of the image when there is no intervention, the latter term signifies the logit when content(object pairs) are removed from the image, therefore giving the total effect of content and removing other effect of confounders.

- [Counterfactual Visual Explanations](#)

- ► Maheep's Notes

1. He selects 'distractor' image I' that the system predicts as class c' and identify spatial regions in I and I' such that replacing the identified region in I with the identified region in I' would push the system towards classifying I as c' .

The author proposes the implementation by the equation:

$$f(I^*) = (1-a)*f(I) + a*P(f(I'))$$

where

I^* represents the image made using the I and I' * represents the Hamdard product.

$f(.)$ represents the spatial feature extractor

$P(f(.))$ represents a permutation matrix that rearranges the spatial cells of $f(I')$ to align with spatial cells of $f(I)$

The author implements it using the two greedy sequential relaxations – first, an exhaustive search approach keeping a and P binary and second, a continuous relaxation of a and P that replaces search with an optimization.

- [Counterfactual Vision and Language Learning](#)

- ► Maheep's Notes

1. The author replaces the embedding of the question or image using another question or image so as to predict the correct answer and minimize counterfactual loss.

- [Counterfactual Vision-and-Language Navigation via Adversarial Path Sampler](#)

- ► Maheep's Notes

1. The author APS, i.e. adversarial path sampler which samples batch of paths P after augmenting them and reconstruct instructions I using Speaker. With the pairs of (P, I) , so as to maximize the navigation loss L_{NAV} .
2. The NAV, i.e. navigation model trains so as to minimize the L_{Nav} making the whole process more robust and increasing the performance.

The APS samples the path based on the visual features v_t which are obtained using the attention on the feature space f_t and history h_{t-1} and previous action taken a_{t-1} to output the path using the predicted a_t and the features f_t .

- [Beyond Trivial Counterfactual Explanations with Diverse Valuable Explanations](#)

- ► Maheep's Notes

1. DiVE uses an encoder, a decoder, and a fixed-weight ML model.
2. Encoder and Decoder are trained in an unsupervised manner to approximate the data distribution on which the ML model was trained.
3. They optimize a set of vectors E_i to perturb the latent representation z generated by the trained encoder.

The author proposes 3 main losses: **Counterfactual loss** : It identifies a change of latent attributes that will cause the ML model f to change it's prediction. **Proximity loss** : The goal of this loss function is to constrain the reconstruction produced by the decoder to be similar in appearance and attributes as the input, therefore making the model sparse. **Diversity loss** : This loss prevents the multiple explanations of the model from being identical and reconstructs different images modifying the different spurious correlations and explaining through them. The model uses the beta-TCVAE to obtain a disentangled latent representation which leads to more proximal and sparse explanations and also Fisher information matrix of its latent space to focus its search on the less influential factors of variation of the ML model as it defines the scores of the influential latent factors of Z . This mechanism enables the discovery of spurious correlations learned by the ML model.

- [SCOUT: Self-aware Discriminant Counterfactual Explanations](#)

- ► Maheep's Notes

The author implements using a network by giving a query image x of class y , a user-chosen counter class $y' \neq y$, a predictor $h(x)$, and a confidence predictor $s(x)$, x is then forwarded to get the $F_h(x)$ and $F_s(x)$. From $F_h(x)$ we predict $h_y(x)$ and $h_{y'}(x)$ which are then combined with the original $F_h(x)$ to produce the $A(x, y)$ and $A(x, y')$ to get the activation tensors and they are then combined with $A(x, s(x))$ to get the segmented region of the image which is discriminative of the counter class.

- [Born Identity Network: Multi-way Counterfactual Map Generation to Explain a Classifier's Decision](#)

- ► Maheep's Notes

1. The author proposes Counterfactual Map Generator (CMG), which consists of an encoder E , a generator G , and a discriminator D . First, the network design of the encoder E and the generator G is a variation of U-Net with a tiled target label concatenated to the skip connections. This generator design enables the generation to synthesize target conditioned maps such that multi-way counterfactual reasoning is possible.
2. The another main technique proposes is the Target Attribution Network (TAN) the objective of the TAN is to guide the generator to produce counterfactual maps that transform an input sample to be classified as a target class. It is a complementary to CMG.

The author proposes 3 main losses:

Counterfactual Map loss: The counterfactual map loss limits the values of the counterfactual map to grow as done by proximity loss in DiVE.

Adversarial loss: It is an objective function retained due to its stability during adversarial training.

Cycle Consistency loss: The cycle consistency loss is used for producing better multi-way counterfactual maps. However, since the discriminator only classifies the real or fake samples, it does not have the ability to guide the generator to produce multi-way counterfactual maps.

- [Introspective Distillation for Robust Question Answering](#)

- ► Maheep's Notes

1. The author proposes to have a causal feature to teach the model both about the OOD and ID data points and take into account the P_{OOD} and P_{ID} , i.e. the predictions of ID and OOD.
2. Based on the above predictions the it can be easily introspected that which one of the distributions is the model exploiting more and based on it they produce the second branch of the model that scores for S_{ID} and S_{OOD} that are based on the equation $S_{ID} = 1/XE(P_{GT}, P_{ID})$, where XE is the cross entropy loss. further these scores are used to compute weights W_{ID} and W_{OOD} , i.e. $W_{OOD} = S_{OOD}/(S_{OOD} + S_{ID})$ to train the model to blend the knowledge from both the OOD and ID data points.
3. The model is then distilled using the knowledge distillation manner, i.e. $L = KL(P_T, P_S)$, where P_T is the prediction of the teacher model and the P_S is the prediction of the student model.

- [Counterfactual Explanation and Causal Inference In Service of Robustness in Robot Control](#)

- ► Maheep's Notes

The additional component in the model is the predictor takes the modified image and produces real-world output. The implementation of it in mathematics looks like:

$\min d_g(x, x') + d_c(C(x'), t_c)$, where d_g is the distance b/w the modified and original image, d_c is the distance b/w the class space and C is the predictor that x' belongs to t_c class.

The loss defines as: $total_loss = (1-\alpha)*L_g(x, x') + (\alpha)*L_c(x, t_c)$, where L_c is the loss x belongs to t_c class

- Counterfactual Explanation Based on Gradual Construction for Deep Networks

- Maheep's Notes

The proposed also focuses on 2 things, i.e. Explainability and Minimality. while implementing the techniuie the authors observe the target class which were being generated were getting much perturbed so as to come under asverserial attack and therefore they propose the logit space of x' to belong to the space of training data as follows:

$$\operatorname{argmin}(\sigma(f_k'(x')) - (1/N) * \sigma(f_k'(X_{i,c_t}))) + \lambda(x' - x)$$

where f' gives the logits for class k , X_{i,c_t} represents the i -th training data that is classified into c_k class and the N is the number of modifications.

- CoCoX: Generating Conceptual and Counterfactual Explanations via Fault-Lines

- Maheep's Notes

The proposed model is implemented by taking the CNN captured richer semantic aspect and construct xconcepts by making use of feature maps from the last convolution layer. Every feature map is treated as an instance of an xconcept and obtain its localization map using the Grad-CAM and are spatially pooled to get important weights, based on that top p pixels are selected and are clustered using K-means. The selection is done using the TCAV technique.

Algorithm 1: Generating Fault-Line Explanations

Input: input image I , classification model M , predicted class label c_{pred} , alternate class label c_{alt} and training dataset χ

1. Find semantically meaningful superpixels in χ ,

$$\alpha_{m,L}^c = \frac{1}{Z} \sum_i \sum_j \frac{\partial y^c}{\partial A_{ij}^{m,L}}$$

2. Apply K-means clustering on superpixels and obtain xconcepts (Σ).

3. Identify class specific xconcepts (Σ_{pred} and Σ_{alt}) using TCAV,

$$S_{c,X} = \nabla g_c(f(I)) \cdot v_X$$

4. Solve Equation 4 to obtain fault-line Ψ ,

$$\Psi \leftarrow \min_{\delta_{pred}, \delta_{alt}} \alpha D(\delta_{pred}, \delta_{alt}) + \beta \|\delta_{pred}\|_1 + \lambda \|\delta_{alt}\|_1$$

return Ψ .

- [CX-ToM: Counterfactual Explanations with Theory-of-Mind for Enhancing Human Trust in Image Recognition Models](#)

- ► Maheep's Notes

Algorithm 1 Generating Fault-Line Explanations

input image I , classification model M , predicted class label c_{pred} , alternate class label c_{alt} and training dataset \mathcal{X}

1. Find semantically meaningful superpixels in \mathcal{X} ,

$$\alpha_{m,L}^c = \frac{1}{Z} \sum_i \sum_j \frac{\partial y^c}{\partial A_{ij}^{m,L}}$$

2. Apply K-means clustering on superpixels and obtain xconcepts (Σ).

3. Identify class specific xconcepts (Σ_{pred} and Σ_{alt}) using TCAV,

$$S_{c,X} = \nabla g_c(f(I)) \cdot v_X$$

4. Solve Equation 5 to obtain fault-line Ψ ,

$$\Psi \leftarrow \min_{\delta_{pred}, \delta_{alt}} \alpha D(\delta_{pred}, \delta_{alt}) + \beta \|\delta_{pred}\|_1 + \lambda \|\delta_{alt}\|_1$$

return Ψ .

- [DeDUCE: Generating Counterfactual Explanations At Scale](#)

- ► Maheep's Notes

Algorithm 1 DeDUCE

Inputs: original input $\mathbf{x} \in \mathcal{X}$, target class t , trained model f with feature extractor f_Z , training data $X_{tr} \in \mathcal{X}^N$, coefficient λ , step size δ , max iterations max_iter, max pixel changes p , number of pixels m , target confidence γ .

Output: counterfactual $\mathbf{x}' \in \mathcal{X}$

```

1: (before deployment) apply  $f$  to  $X_{tr}$  and fit the GMM  $p(\mathbf{z}) = \frac{1}{|C|} \sum_{c \in C} p_c(\mathbf{z})$ 
2:  $\mathbf{x}' \leftarrow \mathbf{x}$ 
3:  $k \leftarrow 0$ 
4:  $\mathbf{P} \leftarrow \mathbf{0}_{dim(\mathcal{X})}$ 
5: while  $f(\mathbf{x}')_t \leq \gamma$  and  $k < \text{max\_iter}$  do
6:   compute gradient  $\mathbf{g}$  in input space
7:   select  $m$  most salient pixels:  $I = \text{select\_q\_largest\_masked}(|\mathbf{g}|, \mathbf{P} < p)$ 
8:   update these pixels:  $\forall i \in I : \mathbf{x}'[i] \leftarrow \mathbf{x}'[i] + \text{sign}(\mathbf{g}[i]) \cdot \delta$ 
9:   clip to input domain:  $\mathbf{x}' \leftarrow \text{clip}(\mathbf{x}')$ 
10:   $\forall i \in I : \mathbf{P}[i] \leftarrow \mathbf{P}[i] + 1$ 
11:   $k \leftarrow k + 1$ 
12: return  $\mathbf{x}'$ 

```

- [Designing Counterfactual Generators using Deep Model Inversion](#)

- ► Maheep Notes

```

argmin( lambda_1*sigma_on_l(layer_l(x'), layer_l(x)) +
lambda_2*L_mc(x';F) + lambda_3*L_cf(F(x'), y'))

```

where,

layer_l: The differentiable layer "l" of the neural network, it is basically used for semantic preservation.

L_mc: It penalizes \mathbf{x}' which do not lie near the manifold. L_mc can be Deterministic Uncertainty

Quantification (DUQ).

L_fc: It ensures that the prediction for the counterfactual matches the desired target

- **ECINN: Efficient Counterfactuals from Invertible Neural Networks**

- ► Maheep's Notes

$$x' = f_inv(f(x) + \alpha * \delta_x)$$

where,

x' :Counterfactual image.

f : INN and therefore f_inv is the inverse of f .

δ_x : the information to be added to convert the latent space of image to that of counterfactual image.

$||z + \alpha_0 * \delta_x - \mu_p|| = ||z + \alpha_0 * \delta_x - \mu_q||$ where the $z + \alpha_0 * \delta_x$ is the line separating the two classes and μ_p and μ_q are the mean distance from line. Therefore

$$\alpha = \alpha_0 + 4/5 * (1 - \alpha_0)$$

- **EXPLAINABLE IMAGE CLASSIFICATION WITH EVIDENCE COUNTERFACTUAL**

- ► Maheep's Notes

Algorithm 1 SEDC

Inputs:

I % Image to classify

$C_M : I \rightarrow \{1, 2, \dots, k\}$ % Trained classifier with k classes

$S = \{s_i, i = 1, 2, \dots, l\}$ % Segmentation of the image with l segments

Procedure:

$E = \{\}$ % List of explanations

for s_i in S **do**

if class change after removing s_i from I **then**

$E = E \cup \{s_i\}$

end if

end for

while $E = \emptyset$ **do**

 Select *best* % Best-first: segment set with highest reduction in predicted class score

$best_set = \text{expansions of } best \text{ with one segment}$

for C_0 in $best_set$ **do**

if class change after removing C_0 from I **then**

$E = E \cup \{C_0\}$

end if

end for

end while

Output:

Explanations in E

- **Explaining Visual Models by Causal Attribution**

- ► Maheep Notes

- **Explaining the Black-box Smoothly-A Counterfactual Approach**

- ► Maheep's Notes

$$L_{cgan} = \log(P_{data}(x)/q(x)) + \log(P_{data}(c|x)/q(c|x))$$

where $P_{data}(x)$ is the data distribution and learned distribution $q(x)$, whereas $P_{data}(c|x)/q(c|x) = r(c|x)$ is the ratio of the generated image and the condition.

2.) **Classification model consistency:** The generated image should give desired output.

Therefore the condition-aware loss is introduced, i.e. $L := r(c|x) + D_{KL}(f(x') || f(x) + \text{delta})$, where $f(x')$ is the output of classifier of the counterfactual image is varied only by delta amount when added to original image logit. They take delta as a knob to regularize the generation of counterfactual image.

3.) **Context-aware self-consistency:** To be self-consistent, the explanation function should satisfy three criteria

- (a) Reconstructing the input image by setting $\delta = 0$ should return the input image, i.e., $G(x, 0) = x$.
- (b) Applying a reverse perturbation on the explanation image x should recover x .

To mitigate this conditions the author propose an identity loss. The author argues that there is a chance that the GAN may ignore small or uncommon details therefore the images are compared using semantic segmentation with object detection combined in identity loss. The identity loss is : $L_{identity} = L_{rec}(x, G(x, 0)) + L_{rec}(x, G(G(x, \text{delta}), -\text{delta}))$

- [Explaining the Behavior of Black-Box Prediction Algorithms with Causal Learning](#)

- ► Maheep's Notes

$V = (Z, Y')$

$Y' = g(z_1, \dots, z_s, \text{epsilon})$ On the basis of possible edge types, they find out which high level causes, possible causes or non-causes of the balck-box output Y' .

- [Explaining Classifiers with Causal Concept Effect \(CaCE\)](#)

- ► Maheep's Notes

$\text{Effect} = E(F(I) | \text{do}(C = 1)) - E(F(I) | \text{do}(C = 0))$ where F gives output on image I and C is the concept. This can be done at scale by intervening for a lot of values in a concept and find the spurious correlation. But due to the insufficient knowlegde of the Causal Graph teh author porposes a VAE which can calculate the precise CaCE by by generating counterfactual image by just changing a concept and hence computing the difference between the prediction score.

- [Fast Real-time Counterfactual Explanations](#)

- ► Maheep's Notes

1.) **Adversarial loss:** It measures whether the generated image is indistinguishable from the real world images

2.) **Domain classification loss:** It is used to render the generate image $x + G(x, y')$ conditional on y' . $L = E[-\log(D(y' | x + G(x, y')))]$ where $G(x, y')$ is the perterbuation introduced by generator to convert image from x to x'

3.) **Reconstruction loss:** The Loss focuses to have generator work properly so as to produce the image need to be produced as defined by the loss. $L = E[x - (x + G(x, y')) + G(x + G(x, y'), y)]$ 4.) **Explanation loss:** This is to gurantee that the generated fake image produced belongs to the distribution of H . $L = E[-\log H(y' | x + G(x, y'))]$

5.) **Perturbation loss:** To have the perturbation as small as possible it is introduced. $L =$

$$E[G(x, y') + G(x + G(x, y'), y)]$$

All these 5 losses are added to make the final loss with different weights.

- [GENERATIVE_COUNTERFACTUAL_INTROSPECTION_FOR_EXPLAINABLE_DEEP_LEARNING](#)

- ► Maheep's Notes

$\min(\lambda * \text{loss}(I(A')) + ||I - I(A')||)$, where loss is cross-entropy for predicting image $I(A')$ to label c' .

- [Generative_Counterfactuals_for_Neural_Networks_via_Attribute_Informed_Perturbations](#)

- ► Maheep's Notes

$\min(E[\sigma_{\text{diff_attributes}} * (-a * \log(D(x')) - (1-a) * (1-D(x)))]) + E[||x - x'||]$ where $D(x')$ generates attributes for counterfactual image.

To generate the counterfactual 2 losses are produced, one ensures that the perturbed image has the desired label and the second one ensures that the perturbation is minimal as possible, i.e.

$$L_{\text{gen}} = \text{Cross_entropy}(F(G(z, a)), y) + \alpha * L(z, a, z_0, a_0)$$

The $L(z, a, z_0, a_0)$ is the l2 norm b/w the attribute and the latent space.

- [Question-Conditioned Counterfactual Image Generation for VQA](#)

- ► Maheep's Notes

- [FINDING AND FIXING SPURIOUS PATTERNS WITH EXPLANATIONS](#)

- ► Maheep's Notes

$$P(\text{Spurious} | \text{Main}) = P(\text{Spurious} | \text{not Main}) = 0.5$$

The second step consist of minimizing the potential for new SPs by setting the

$$P(\text{Main} | \text{Artifact}) = 0.5).$$

SPIRE moves images from {Both, Neither} to {Just Main, Just Spurious} if $p > 0.5$, i.e. $p =$

$P(\text{Main} | \text{Spurious})$ but if $p < 0.5$ then SPIRE moves images from {Just Main, Just Spurious} to {Both, Neither}.

- [Contrastive_Counterfactual_Visual_Explanations_With_Overdetermination](#)

- ► Maheep's Notes

- [Training_calibration-based_counterfactual_explainers_for_deep_learning](#)

- ► Maheep's Notes

Validity: ratio of the counterfactuals that actually have the desired target attribute to the total number of counterfactuals

The confidence of the **image** and **sparsity**, i.e. ratio of number of pixels altered to total no of pixels. The other 2 metrics are **proximity**, i.e. average l2 distance of each counterfactual to the K-nearest training samples in the latent space and **Realism score** so as to have the generated image is close to the true data manifold.

TraCE reveals attribute relationships by generating counterfactual image using the different

attribute like age "A" and diagnosis predictor "D".

$\Delta_{A_x} = x - x_{a'}; \Delta_{D_x} = x - x_{d'}$

The $x_{a'}$ is the counterfactual image on the basis for age and same for $x_{d'}$.

$x' = x + \Delta_{A_x} + \Delta_{D_x}$ and hence atleast we evaluate the sensitivity of a feature by $F_d(x') - F_d(x_{d'})$, i.e. F_d is the classifier of diagnosis.

- [Generating Natural Counterfactual Visual Explanations](#)

- ► Maheep's Notes

- [On Causally Disentangled Representations](#)

- ► Maheep's Notes

- [INTERPRETABILITY_THROUGH_INVERTIBILITY_A_DEEP_CONVOLUTIONAL_NETWORK](#)

- ► Maheep's Notes

$z' = z + \alpha * w$ where $x' = \text{phi_inverse}(z + \alpha * w)$. Any change orthogonal to w will create an "isofactual". To show that their counterfactuals are ideal, therefore they verify that no property unrelated to the prediction is changed. Unrealized properties = $e(x)$, $e(x) = v^T * z$, where v is orthogonal to w . $e(x') = v^T * (z + \alpha * w) = v^T * z = e(x)$. To measure the difference between the counterfactual and image intermediate feature map h , i.e. $m = |\Delta_h| * \cos(\text{angle}(\Delta_h, h))$ for every location of intermediate feature map.

- [Model-Based Counterfactual Synthesizer for Interpretation](#)

- ► Maheep's Notes

Algorithm 1: Building Pipeline of MCS for Interpretation

Setup Phase

- Prepare f for explanation, and select l for measurement;
- Design the generator G with domain \mathcal{D} based on Theo. 2;

Training Phase

- Data Modeling with Eq. 9;
- **for training batch k do**
 - Utilize umbrella sampling to prepare a set of queries q_k ;
 - Weigh batch k with w over q_k based on Theo. 1;
 - Update G and D in the min-max game of Eq. 5;

Interpretation Phase

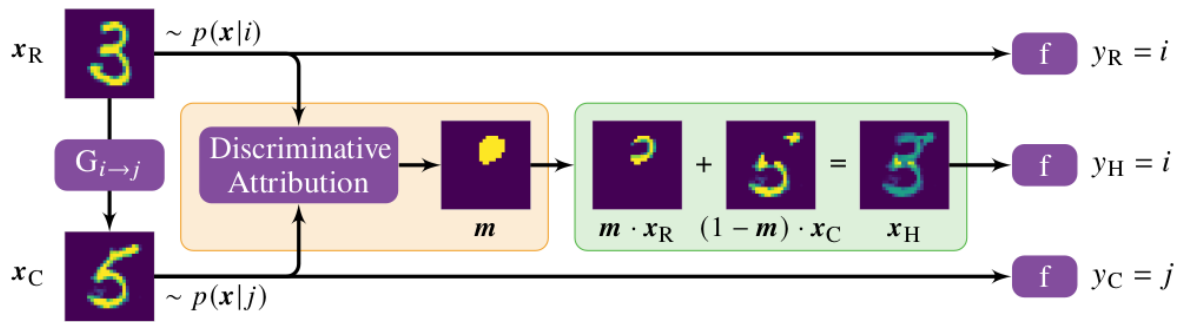
- Feed the user query q to G for counterfactual generation.
-

- [The Intriguing Relation Between Counterfactual Explanations and Adversarial Examples](#)

- ► Maheep's Notes

- Discriminative Attribution from Counterfactuals

- Maheep's Notes

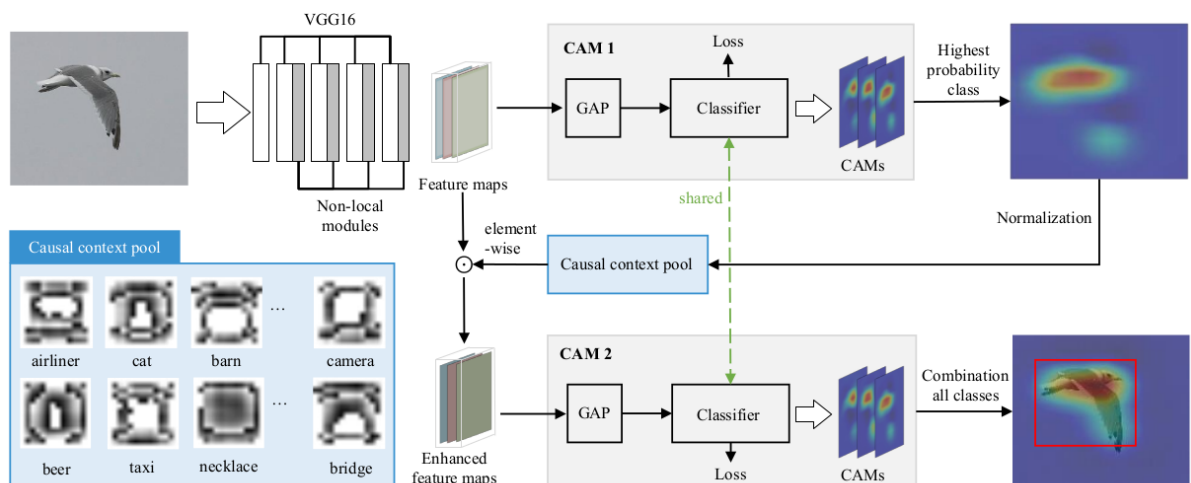


- Causal Interventional Training for Image Recognition

- Maheep's Notes

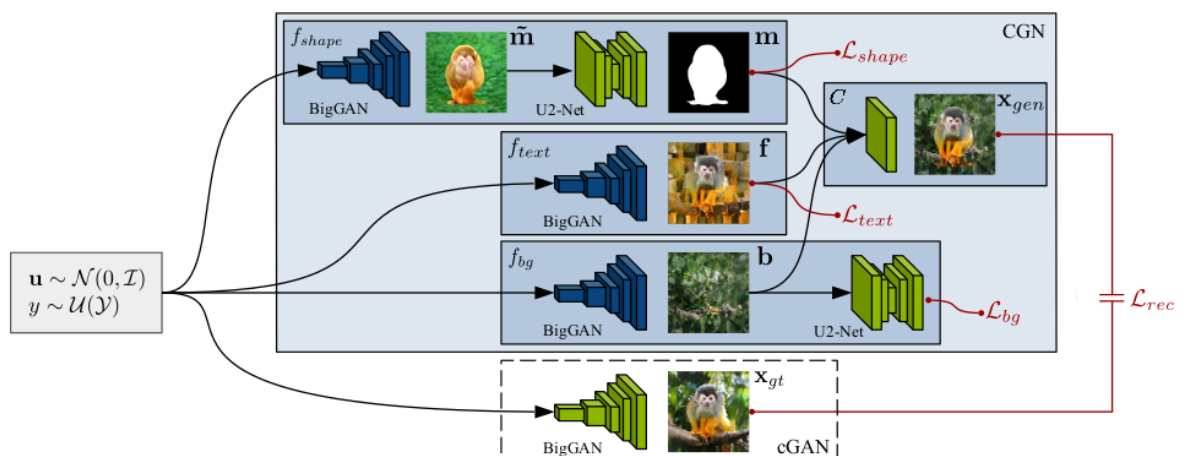
- Improving Weakly-supervised Object Localization via Causal Intervention

- Maheep's Notes



- COUNTERFACTUAL GENERATIVE NETWORKS

- Maheep's Notes



- Discovering Causal Signals in Images

- Maheep's Notes

X causes Y , the cause, noise and mechanism are independent but we can identify the footprints of causality when we try to Y causes X as the noise and Y will not be independent.

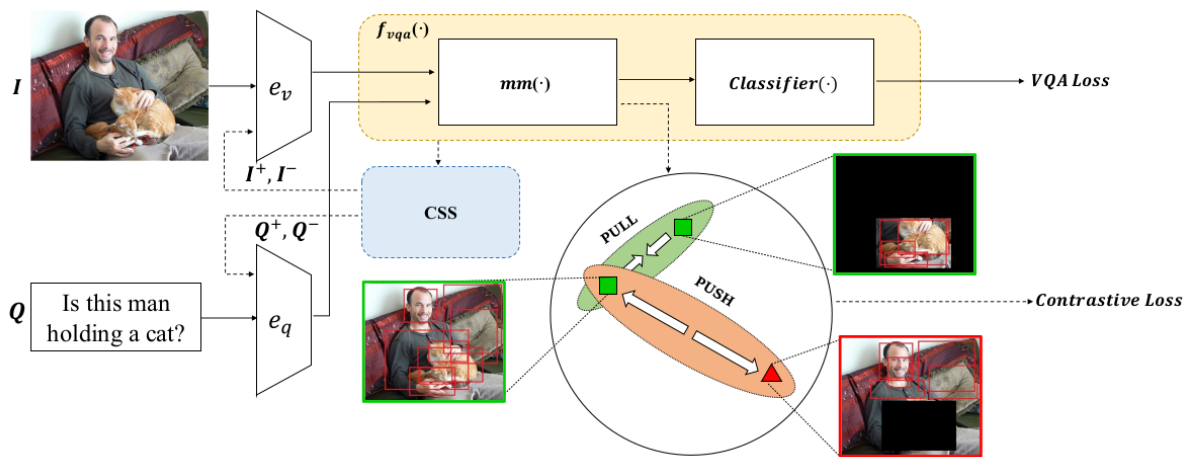
- Learning to Contrast the Counterfactual Samples for Robust Visual Question Answering

- Maheep's Notes

Q and V ($\text{mm}(Q, V) = a$), and joint embedding of Q and V_+ (factual) ($\text{mm}(Q, V_+) = p$) by taking a cosine similarity b/w them. They also aim to decrease mutual information b/w $\text{mm}(Q, V_-) = n$ and a by taking cosine similarity ($s(a, n)$). The final formula becomes:

$$L_c = E[-\log(e^{s(a, p)} / (e^{s(a, p)} + e^{s(a, n)}))]$$

The total loss becomes $L = \lambda_1 L_c + \lambda_2 L_{vqa}$



- Latent Causal Invariant Model

- Maheep's Notes

S and (b) others that are spuriously correlated Z from V (latent variable).

There exists a spurious correlation b/w S and Z . The author argues that we will get a

$$p(y|\text{do}(s^*)) = p(y|s^*)$$

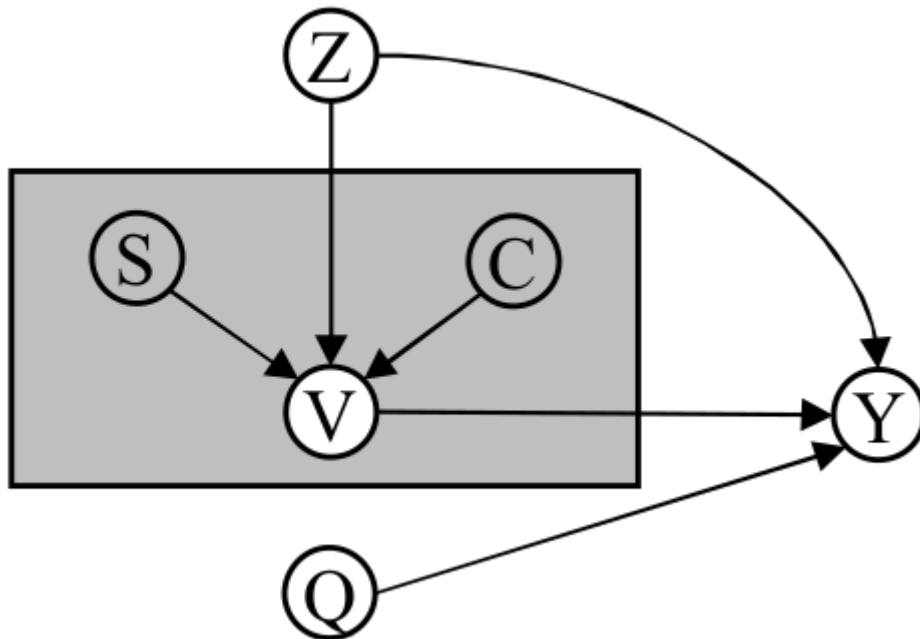
- Two Causal Principles for Improving Visual Dialog

- Maheep's Notes

a_i (answer) is a sentence observed from the "mind" of user u during dataset collection. Then, $\sigma(P(A) * P(u|H))$, H is history and A is answer can be approximated as $\sigma(P(A)P(a_i|H))$. They further use $p(a_i|QT)$, where QT is Question Type to approximate $P(a_i|H)$ because of two reasons: First, $P(a_i|H)$ essentially describes a prior knowledge about a_i without comprehending the whole $\{Q, H, I\}$ triplet.

- Weakly-Supervised Video Object Grounding via Causal Intervention

- Maheep's Notes



Z that occurs due to some specific objects occurring frequently. The style confounder is replaced by using the contrastive learning, where the counterfactual examples are created by taking the vectors from a memory bank by taking the top selected top regions for described object and then the selected regions and frames are grouped together into frame-level content(H_c) and region-level content(U_c), and the rest of the regions are grouped as U_s and H_s . These regions are converted to counterfactual using these memory vectors which were created by taking the randomly selected regions in training set. The most similar one and replaces the original one, to generate examples to have them hard to distinguish from real ones contrastive learning is used. The equation looks like:

$$IE(p|do(U_s = U_{s_generated})) < IE(p|do(U_c = U_{c_generated}))$$

$$IE(p|do(H_s = H_{s_generated})) < IE(p|do(H_c = H_{c_generated}))$$

where the IE is Interventional Effect. As for the next confounder they use the textual embedding of o_k (object) essentially provides the stable cluster center in common embedding space for its vague and diverse visual region embeddings in different videos. Therefore, by taking the textual embedding of the object as the substitute of every possible object z and apply backdoor adjustment.

- **Towards Unbiased Visual Emotion Recognition via Causal Intervention**

- ► Maheep's Notes

IERN, which is composed of four parts:

- 1.) **Backbone**

> It extracts the feature embedding of the image.

- 2.) **Feature Disentanglement**

> It disentangles the emotions and context from the image, having emotion discriminator(d_e) and context discriminator(d_c) which ensures that the extracted feature are separated and has the desired feature. The loss comprises as :

$L = CE(d_e(g_e(f_b(x))), y_e) + MSE(d_c(g_e(f_b(x))), 1/n)$ where g_e is emotion generator and y_e is the emotion label and n is the number of confounder and the same loss is for context replacing d_e , g_e and d_c by d_c , g_c and d_e , here n represents

number of emotions. To ensure that the separated features fall within reason-able domains, IERN should be capable of reconstructing the base feature $f_b(x)$, i.e. L

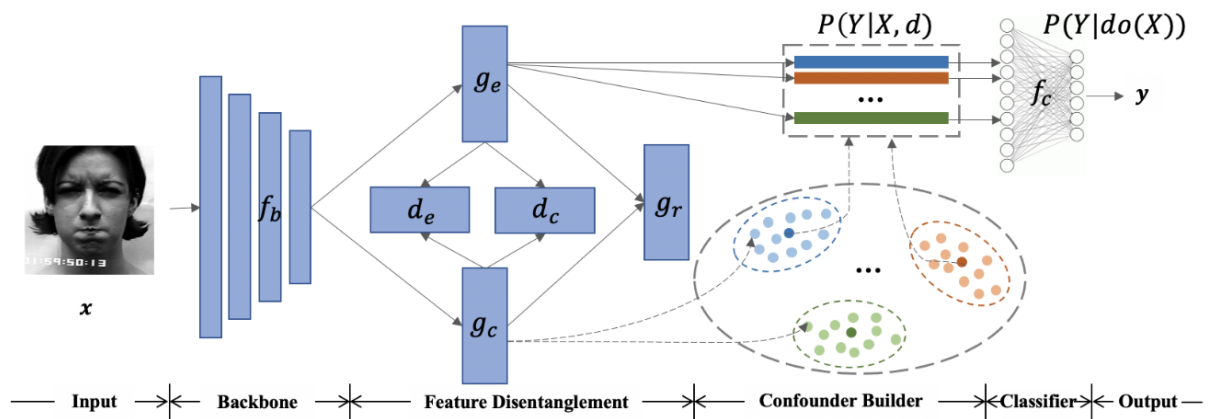
$$= \text{MSE}(g_r(g_e(f_b(x))), g_c(f_b(x))), f_b(x))$$

3.) Confounder Builder

> The purpose of the confounder builder is to combine each emotion feature with different context features so as to avoid the bias towards the observed context strata.

4.) Classifier

> It is simply used for prediction.



- **Human Trajectory Prediction via Counterfactual Analysis**

- ► Maheep's Notes

They Y_{causal} is defined as $Y_{\text{causal}} = Y_i - Y_i(\text{do}(X_i = x_i))$

They define a generative model which generates trajectory by a noise latent variable Z indicated by Y^*_i . Finally the loss is defined as:

$$Y_{\text{causal}} = Y^*_i - Y^*_i(\text{do}(X_i = x_i))$$

$$L_{\text{causalGAN}} = L_2(Y_i, Y_{\text{causal}}) + \log(D(Y_i)) + \log(1 - D(Y_{\text{causal}})),$$

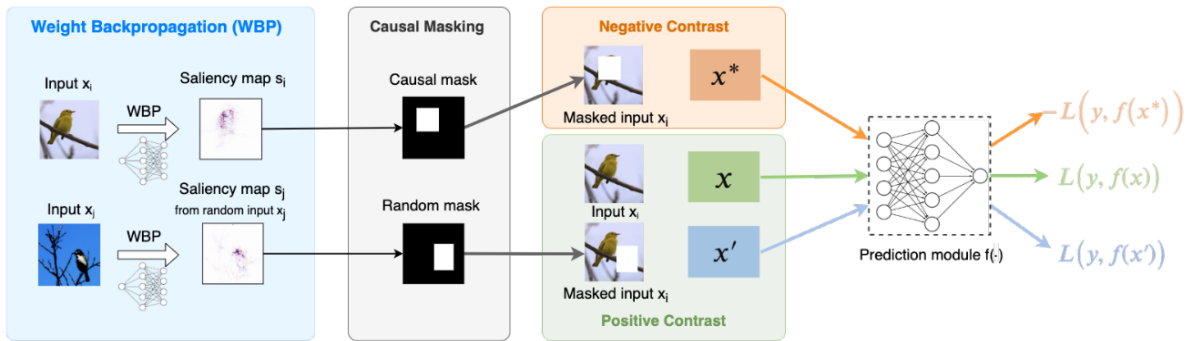
where D is the discriminator.

- **Proactive Pseudo-Intervention: Contrastive Learning For Interpretable Vision Models**

- ► Maheep's Notes

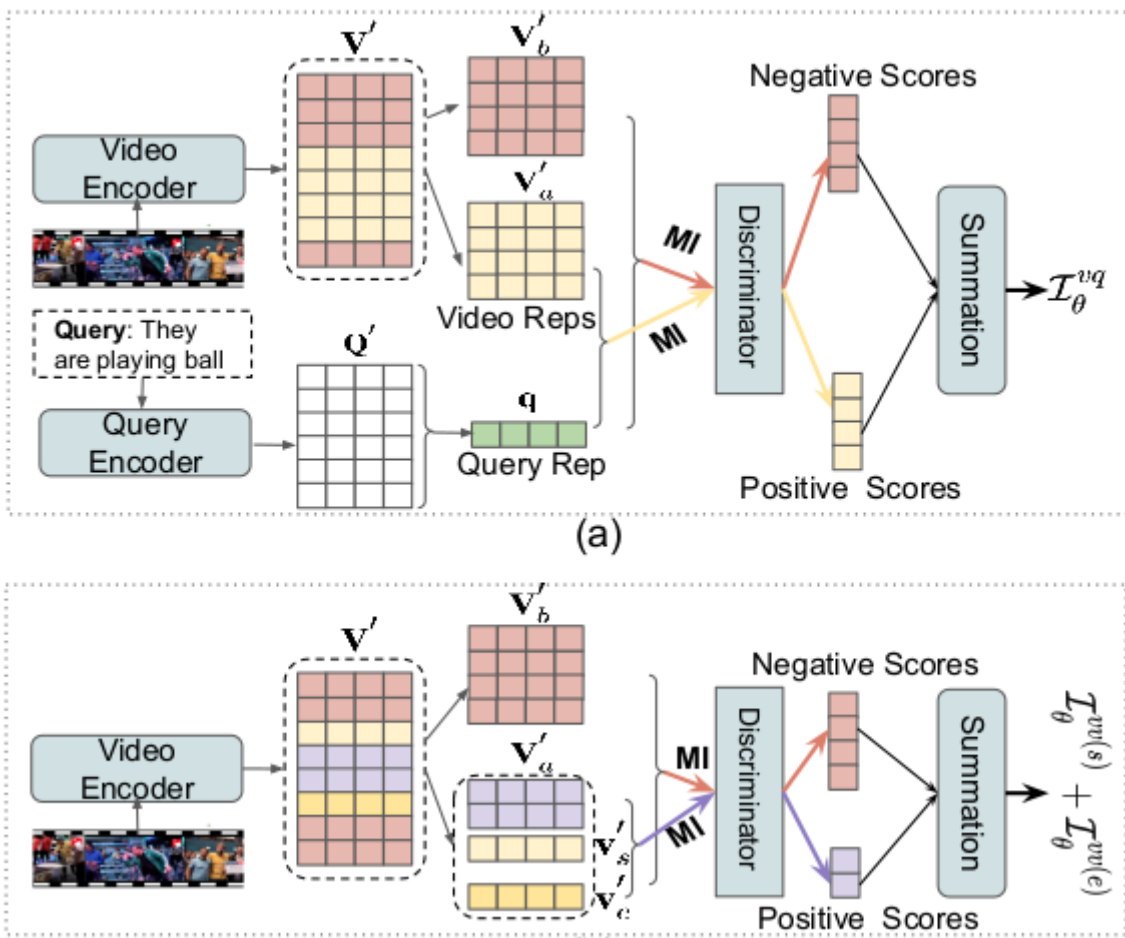
$f(\theta)$ are produced and the main features are masked out of the image giving us x^* . Now the loss becomes $L = \text{sigma}(l(x^*, \text{not_y}; f(\theta)))$

A trivial solution can be the saliency maps covers the whole image therefore L1-norm of saliency map is used to encourage succinct (sparse) representations. The another problem that now arises is that the model can learn a shortcut that when it get a masked image then it has to always give not_y as prediction, so as to counter it the author proposes to send images with random masks on them, making the loss $L = \text{sigma}(l(x', y; f(\theta)))$

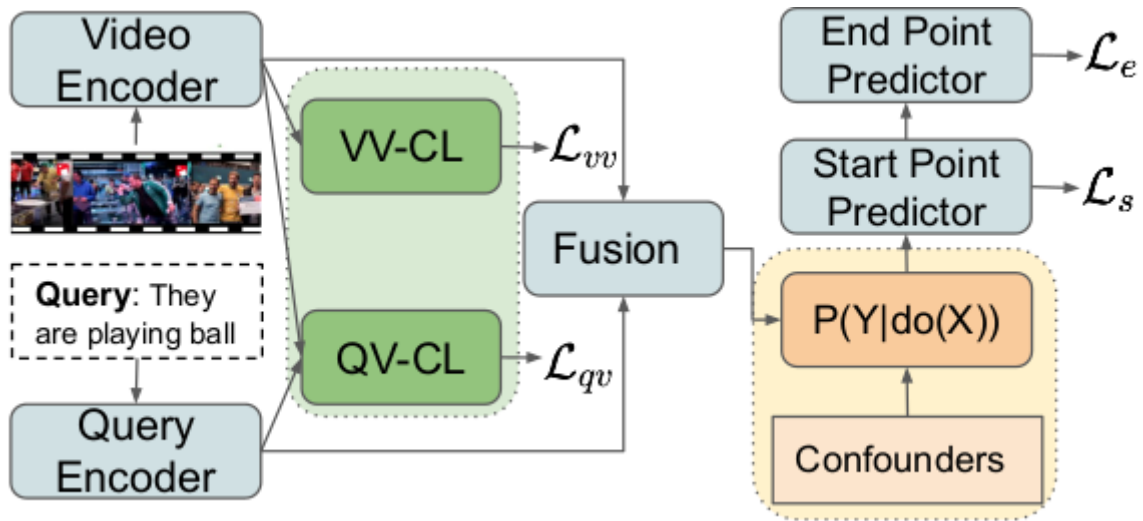


- Interventional Video Grounding with Dual Contrastive Learning

- Maheep's Notes



3. The output of two feature encoders are fed to a fusion module with a context-query attention mechanism to capture the cross-modal interactions between visual and textual features.
4. Next, to mitigate the spurious correlations between textual and visual features, they use causal interventions $P(Y | \text{do}(X))$ with event as surrogate confounders to learn representations.
5. Finally, two losses L_s and L_e for the start and end boundaries are introduced.



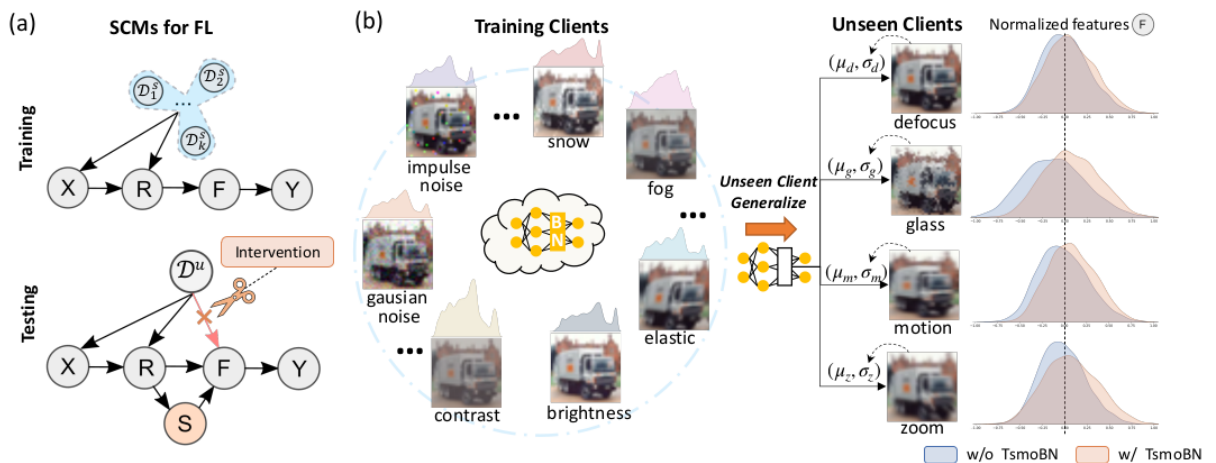
- Causality matters in medical imaging

- ► Maheep's Notes

- TSMOBN GENERALIZATION FOR UNSEEN CLIENTS IN FEDERATED LEARNING

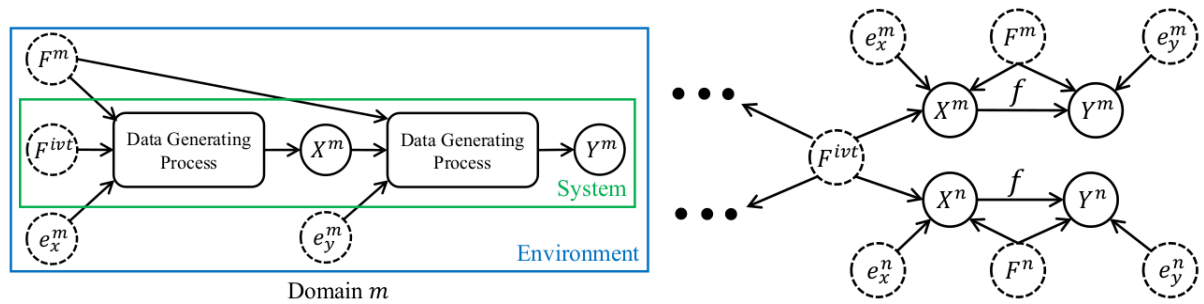
- ► Maheep's Notes

\mathcal{D}_{s_i} for datasets of different domain, i.e. coming from different users but used in training, \mathbf{X} are the samples, \mathbf{R} are the raw extracted features of \mathbf{X} , \mathbf{F} is the normalized feature representation of \mathbf{R} and \mathbf{Y} is the classifier. To remove the confounding effects brought by \mathcal{D}_u , a direct way is using causal intervention on normalized features (i.e., $\text{do}(\mathbf{F})$) to let the feature distribution similar to training distributions. This intervention by introducing the surrogate variable \mathbf{S} , which is test-specific statistics of raw features \mathbf{R} during testing by obtaining the test normalized features that have similar distributions as the training normalized features. More specifically by calculating the mean and variance pair at test time in BN to normalize features. Additionally they further propose to use momentum to integrate relations among different batches, thus reducing the variances. Precisely by giving the unseen client with M batches of data to be tested in sequential manner.



- Learning Domain Invariant Relationship with Instrumental Variable for Domain Generalization

- ► Maheep's Notes



- Latent Space Explanation by Intervention

- Maheep's Notes

$Z = \{z_1, z_2, \dots, z_n\}$ which is intervened to flip the output of the model and are finally visualized from the hidden representation using the loss $L = l(g(\phi'(x), x))$, where $\phi'(x)$ is the counterfactual model. The next goal is to ensure that the generated concepts follow the same concepts the discriminator employ. They achieve this by maximizing the amount of information that the explanatory learner (i.e., g) extracted from the latent representation with respect to the discriminative learner's (i.e., f_K) information.

- The Blessings of Unlabeled Background in Untrimmed Videos

- Maheep's Notes

$P(x_1, x_2, \dots, x_n | Z = z) = \prod P(x_t | Z = z)$, i.e. the features will become independent if we are able to observe Z but if there exists an unobserved confounder c , which affects multiple input video features within x and segment-level labels A . Then, x would be dependent, even conditional on z , due to the impact of c , in this case with the blessings of weak ignorability we can replace the expectation over C with a single z in

$$E[E[A|X = x, C = c]] = A$$

- Selecting Data Augmentation for Simulating Interventions

- Maheep's Notes

- Counterfactual Attention Learning for Fine-Grained Visual Categorization and Re-identification

- Maheep's Notes

$$Y_{\text{effect}} = E[Y(A = A, X = X)] - E[Y(A = A', X = X)]$$

The loss comprises as:

$L = L_{\text{crossentropy}}(Y_{\text{effect}}, y) + L_{\text{others}}$, where L_{others} represents the original objective such as standard classification loss.

- Meaningful Explanations of Black Box AI Decision Systems

- Maheep's Notes

- Are VQA Systems RAD? Measuring Robustness to Augmented Data with Focused Interventions

- Maheep's Notes

$$RAD = |J(D;F) \text{ and } J(D';F)| / |J(D;F)|$$

, where $J(D;F)$ as the set of example indices for which a model f correctly predicts y . D' represents the augmented example which is prepared as *VQA dataset* there are three answer types: "yes/no", "number" and "other", and 65 question types. In augmentations, they generate "yes/no" questions from "number" and "other" questions, i.e. What color is the ? is changed to Is the color of is ?

RAD is in $[0, 1]$ and the higher the RAD of f is, the more robust f is.

- [Adversarial Robustness through the Lens of Causality](#)

- ► Maheep's Notes

$$P_{\theta}(X, Y) = \sigma(P_{\theta}(Y, s|X) * P_{\theta}(X)),$$

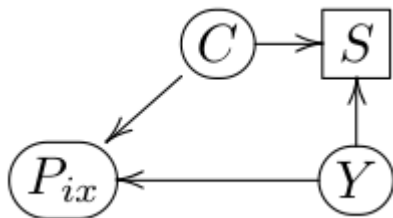
where s is the spurious correlation. As we know that the distribution of X can be hardly changed therefore $P_{\theta}(X) = P(X)$. Therefore it can be assumed that the difference b/w $P_{\theta}(Y, s|X)$ and $P(Y, s|X)$ is the main reason of the adversarial inrobustness. Therefore they define the loss as:

$$\min CE(h(X + E_{adv}; \theta), Y) + CE(h(X; \theta), Y) + CE[P(Y|g(X, s)), P(Y|g(X + E_{adv}, s))]$$

where E_{adv} adversarial perturbation, θ are parameters of the model, and g represents the parameter optimized to minimize the CE , i.e. Cross Entropy loss.

- [Causality-aware counterfactual confounding adjustment for feature representations learned by deep models](#)

- ► Maheep's Notes



In order to remove/reduce the influence of C on the predictive performance of the classifier, they apply the causality-aware adjustment proposed to generate counterfactual features, X' . These counterfactual examples are used to train a logistic regression classifier, and then use the same algorithm to generate counterfactual in test set X'_{test} to generate predictions that are no longer biased by the confounder.

- [Domain Generalization using Causal Matching](#)

- ► Maheep's Notes

Algorithm 1 MatchDG

In: Dataset $(d_i, x_i, y_i)_{i=1}^n$ from m domains, τ, t

Out: Function $f : \mathcal{X} \rightarrow \mathcal{Y}$

Create random match pairs Ω_Y .

Build a $p * q$ data matrix \mathcal{M} .

Phase I

while notconverged **do**

for $batch \sim \mathcal{M}$ **do**

 Minimize contrastive loss (4).

end for

if epoch % $t == 0$ **then**

 Update match pairs using Φ_{epoch} .

end if

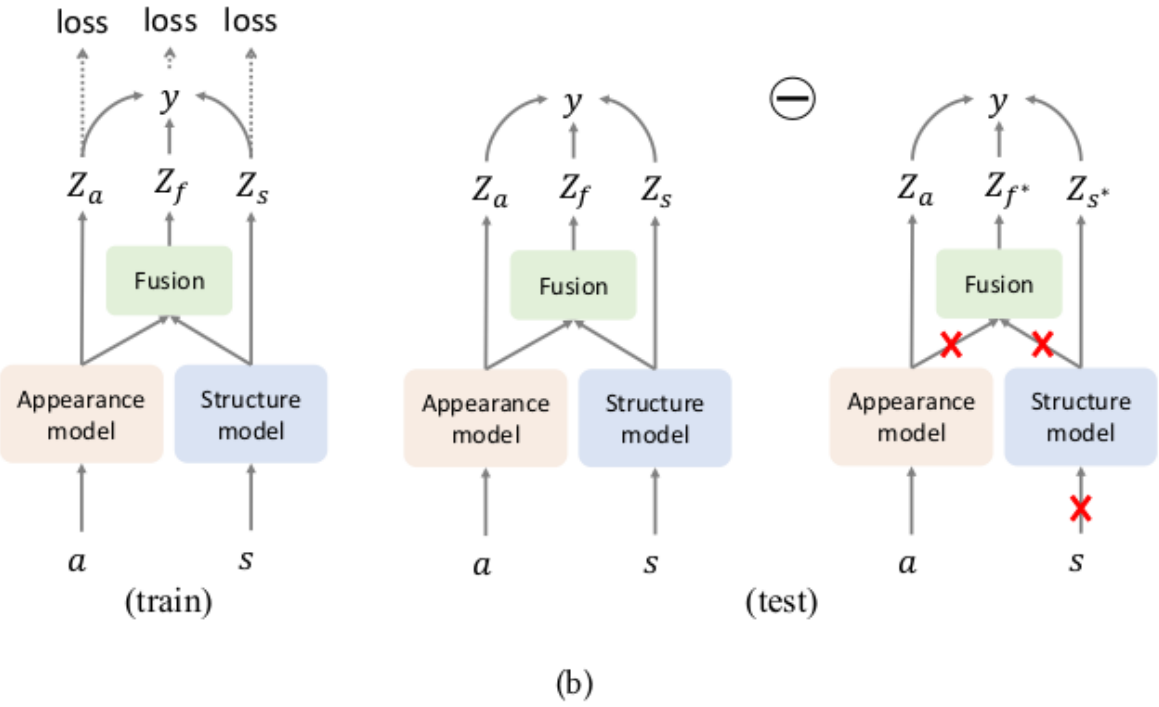
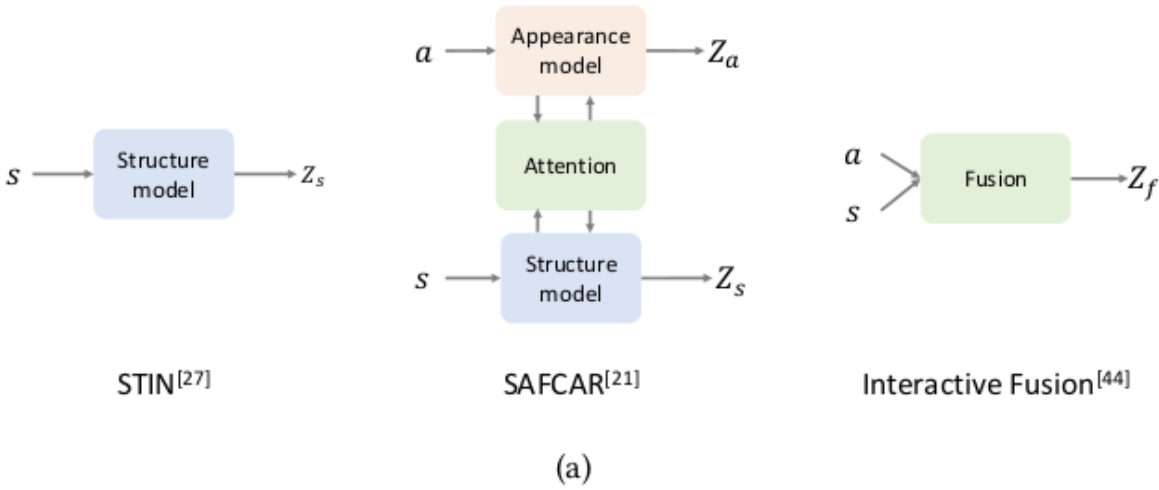
end while

Phase II

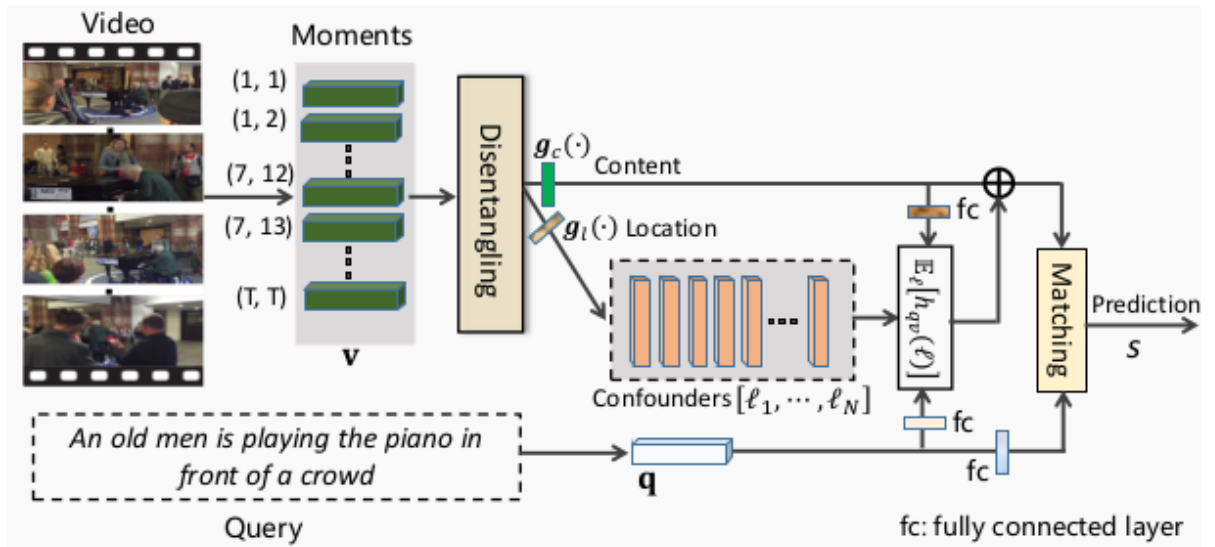
Compute matching based on Φ .

Minimize the loss (3) with learnt match function Φ to obtain f .

- [Counterfactual Debiasing Inference for Compositional Action Recognition](#)
 - ► Maheep's Notes



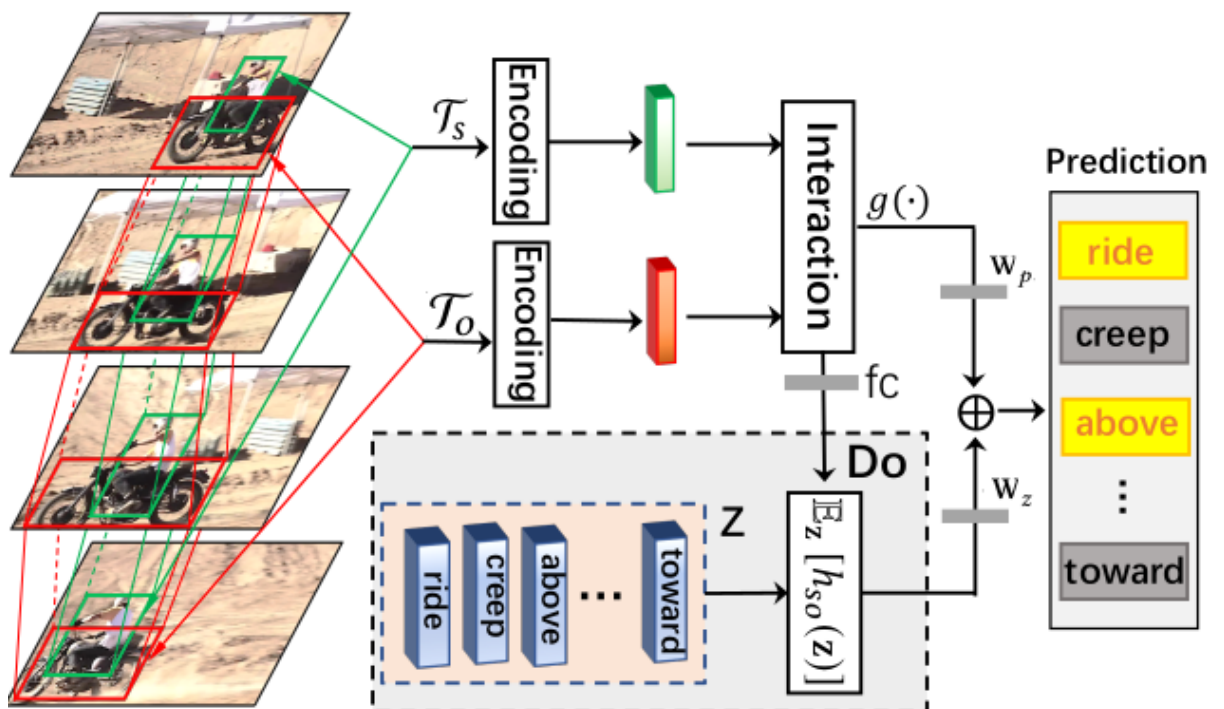
- [Deconfounded Video Moment Retrieval with Causal Intervention](#)
 - ► Maheep's Notes



- Intervention Video Relation Detection

- Maheep's Notes

$L = L_{obj} + \lambda * L_{pred}$, where L_{obj} is the cross entropy loss function to calculate the loss of classifying video object trajectories and L_{pred} is binary cross entropy loss used for predicate prediction.



- Visual Commonsense R-CNN

- Maheep's Notes

$E[g(z)]$ to get the top confounders from the dictionary. Now there is a complexity arise where the confounder is the doctionary act as the colliders therefore they are eridacted through the use of Neural Causation coefficient(NCC).

