

Semi-Supervised Disentangling of Causal Factors

Sargur N. Srihari

srihari@cedar.buffalo.edu

Topics in Representation Learning

1. Greedy Layer-Wise Unsupervised Pretraining
2. Transfer Learning and Domain Adaptation
3. Semi-supervised Disentangling of Causal Factors
4. Distributed Representation
5. Exponential Gains from depth
6. Providing Clues to Discover Underlying Causes

Topics in Semi-supervised Disentangling of Causal Factors

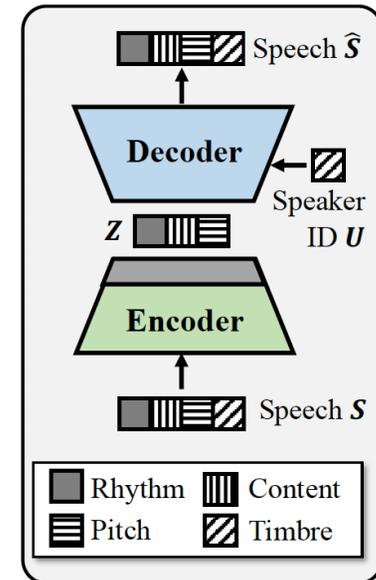
- Causal Representation
- Relationship to Semi-supervised learning
- Salient features: Autoencoders vs GANs
- Role of Cause and Effect in Deep Learning

Example of entangled factors

- Factors underlying speech

1. Language content (C):
transcribable to text
2. Timbre (V):
characterizes speaker identity
 - Prosody (speaker emotion)
3. Pitch variation (F):
aspects of speaker tone
4. Rhythm (R):
how fast word/syllable uttered

Autoencoder for speech



$S = \{S_t\}$ is the sound spectrogram
 S is generated through
 $S = g_s(C, R, F, V)$, $U = g_u(V)$

- Disentangling useful for

- Speech Analysis:
remove irrelevant interference
- Speech Generation:
voice conversion, prosody modification

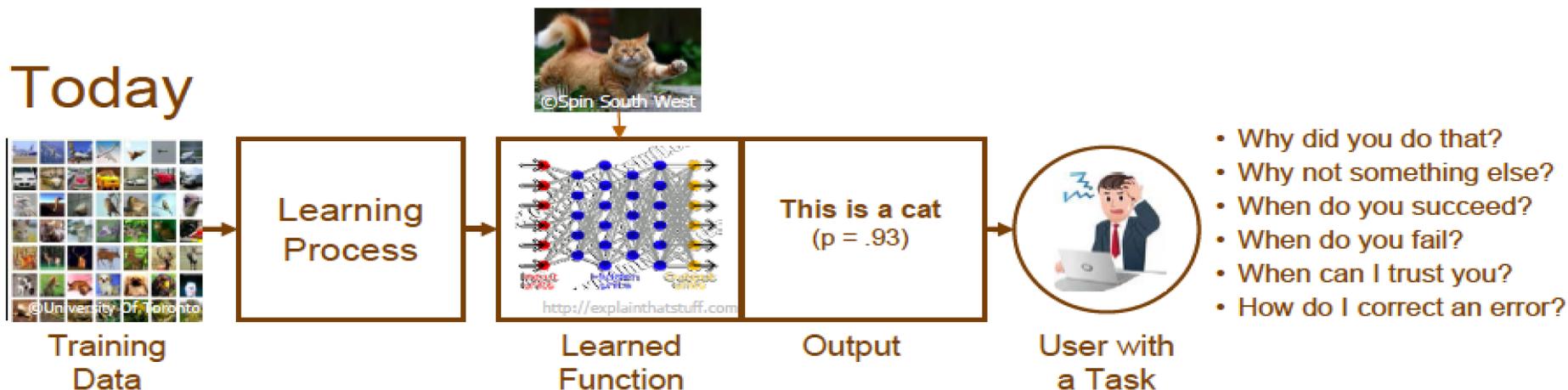
Since speaker ID is separately fed to the decoder, the decoder can still have access to all the information to perform self-reconstruction even if the encoder does not preserve the timbre information in its output Z .

Causal Representation

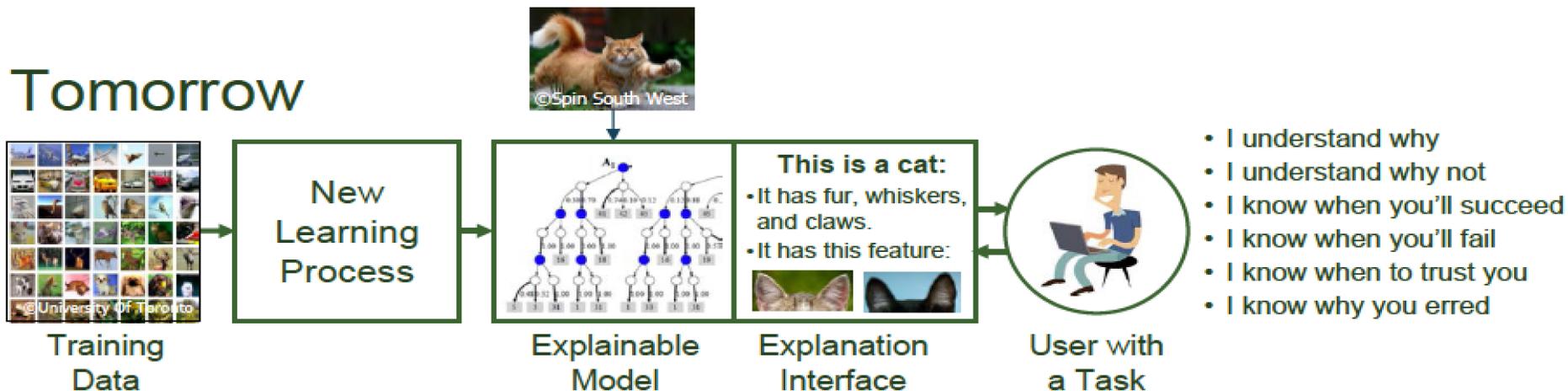
- One hypothesis of an ideal representation h is that it is a causal representation
- Representation is such that either
 1. It corresponds to underlying causes of x OR
 2. Directions in feature space correspond to causes
- So that representation disentangles causes
- Causal factors are useful for explanation
 - given next

Causal Explanation

Today



Tomorrow

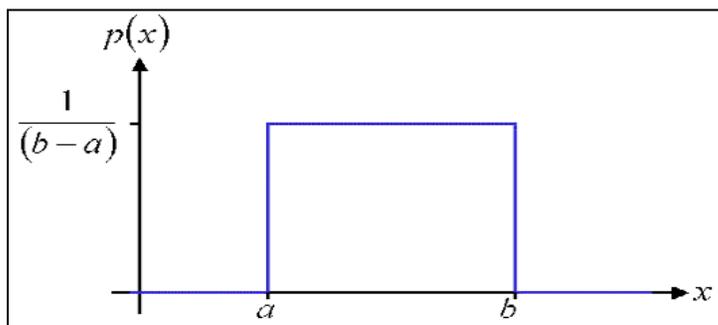


Two goals of representation learning

1. A representation that is easy to model, e.g.,
 1. Independent from each other
 2. Sparse
2. Representation that separates causal factors
 - May not be easy to model
- For some tasks the two coincide

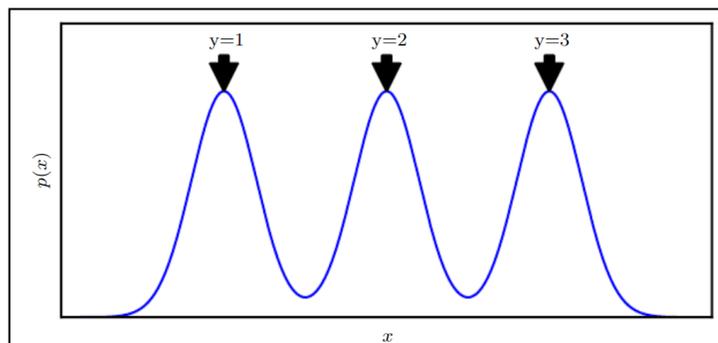
Where semi-supervised can fail/succeed

- When is $p(x)$ is of no help to learning $p(y|x)$?
 - $p(x)$ is uniform and we want to learn $f(x)=E[y|x]$



Clearly observing the training set of x values alone gives us no information about $p(y|x)$

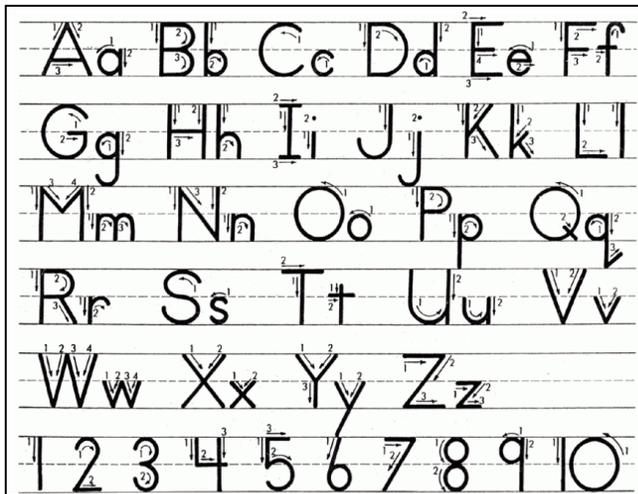
- Compare with case when $p(x)$ is helpful
 - $p(x)$ is multimodal Gaussian



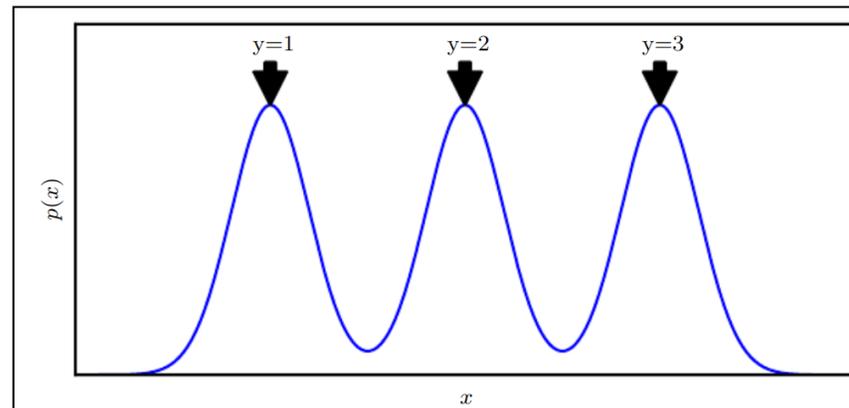
Modeling $p(x)$ reveals $p(y|x)$ where a single labeled example per class enough to learn $p(y|x)$

How semi-supervised can succeed

- Density over x is a mixture over three components, one per value of y
- Components well-separated by feature x
 - modeling $p(x)$ reveals where each component is
 - A single labeled example per class enough to learn $p(y|x)$



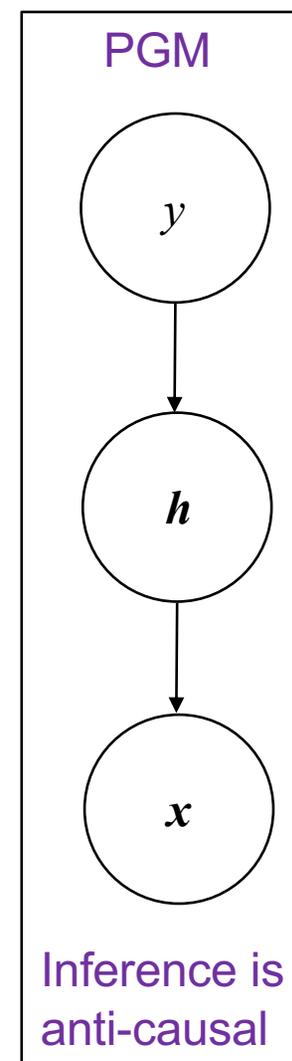
x = no. of black pixels



In this case $p(y|x)$ is a univariate Gaussian for $y=1,2,3$

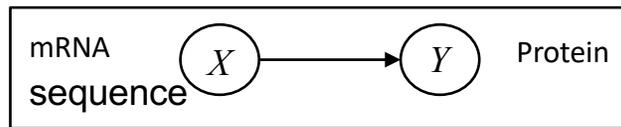
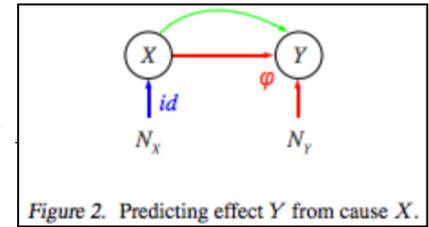
Causal h useful to predict y

- If h represents causes of observed x , and outputs y are among most salient causes, then it is easy to predict y from h
 - Representation for $p(x)$ may also be good for $p(y|x)$ if output y is among salient causes of x
 - E.g., cat (y) influences shape of ear h_i
- Motivates semi-supervised learning

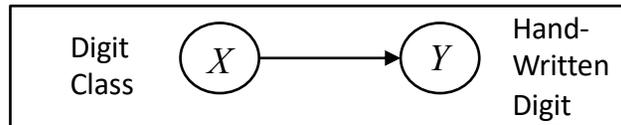


Semi-supervised useful in anti-causal problem

- **Generative process: Effect Y , Cause X**
 - **Ex1: Predict protein Y from mRNA sequence**
 - It is a *causal problem*



- **Ex 2: Predict class X from handwritten digit Y**
 - it is an *anti-causal problem*



- **Modeling $p(X)$ with extra data does not help in Ex 1**
 - We assume that $p(X)$ is independent of $p(Y|X)$, from Bayes Net
- **But in Ex 2 modeling $p(Y)$ is helpful**
 - because $p(X|Y)$ is dependent on $p(Y)$
- **Problems like Ex 2 benefit from semi-supervised learning**

$$p(X | Y) = \frac{p(Y | X)p(X)}{p(Y)}$$

Semi-supervised possible when factors are disentangled

- What could tie $p(y|\mathbf{x})$ and $p(\mathbf{x})$ together?
 - If y is closely associated with one of the causal factors of \mathbf{x} , then $p(\mathbf{x})$ and $p(y|\mathbf{x})$ will be strongly tied
- Unsupervised learning that disentangles underlying factors is useful for semi-supervised learning

Formalizing best possible model

- Assume y is one of the causal factors of x
- Let h represent all those factors
- The true generative process can be conceived as structured according to this directed model with h as the parent of x : $p(h, x) = p(h)p(x|h)$
 - Thus data has marginal probability $p(x) = E_h p(x|h)$
- Thus we conclude that the best possible model of x is one that uncovers the above true structure with h as a latent variable that explains the observed variations in x

Ideal representation learning

- It should recover the latent factors
- If \mathbf{y} is one of these then it will be easy to predict \mathbf{y} from such a representation
- We also see from Bayes rule:
$$p(\mathbf{y} | \mathbf{x}) = \frac{p(\mathbf{x} | \mathbf{y})p(\mathbf{y})}{p(\mathbf{x})}$$
- Thus the marginal $p(\mathbf{x})$ is intimately tied to the conditional $p(\mathbf{y}|\mathbf{x})$
 - Knowledge of the structure of $p(\mathbf{x})$ should help learn $p(\mathbf{y}|\mathbf{x})$
 - Therefore in situations respecting these assumptions, semi-supervised learning should improve performance

Brute force for large no of causes

- Most observations are formed by an extremely large no of causes
- Suppose $\mathbf{y} = h_i$, but the unsupervised learner does not know which h_i
- *Brute-force solution*
 - Unsupervised learning: a representation that captures *all* the reasonably salient generative factors h_j
 - Disentangle them from each other thus making it easy to predict \mathbf{y} from \mathbf{h} regardless of which h_i is associated with \mathbf{y}

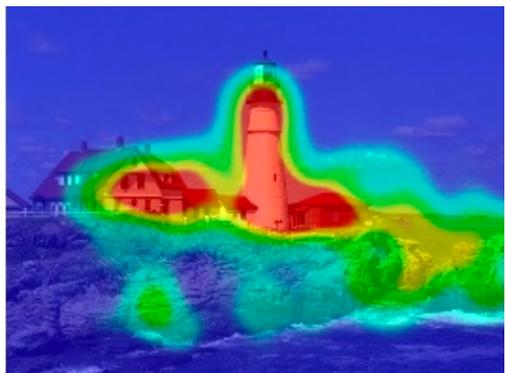
Brute force is infeasible

- It is not possible to capture all factors of variation that influence the observation
 - Ex: in a visual scene, should the representation encode all the smallest objects in the background?
 - Humans fail to perceive changes in environment not relevant to task they are performing
- Research frontier in semi-supervised learning:
What to encode in each situation

Saliency Detection



Question: What have you seen?



Answer 1: Lighthouse

Answer 2: Lighthouse and Houses

Answer 3: Lighthouse, Houses and Rocks

Two ways to deal with many causes

- Two main strategies to deal with a large no of underlying causes:
 1. Use both supervised and unsupervised learning
 - Use a supervised signal at the same time as the unsupervised learning signal so that the model will choose to capture the most relevant factors of variation
 2. Use much larger representations if using purely unsupervised learning

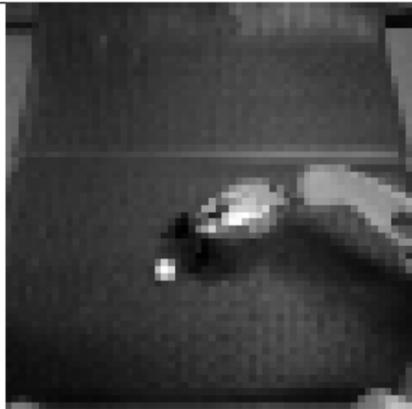
Modifying definition of saliency

- Emerging strategy for unsupervised learning is to modify the definition of which underlying causes are most salient
- Autoencoders and generative models usually optimize a fixed criterion, say MSE
- These fixed criteria determine which causes are considered salient
 - Ex: MSE applied to pixels implies that an underlying cause is salient only if it significantly changes the brightness of a large no of pixels
 - Problematic if task involves interacting with small objects
 - Example next

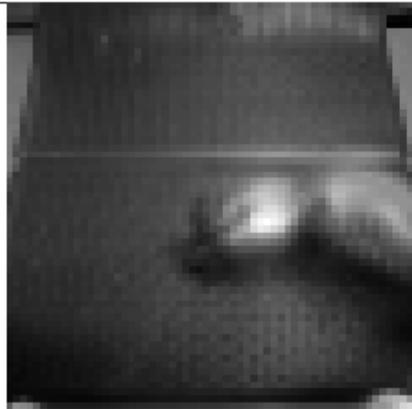
Failure of saliency detection

- Autoencoder trained with MSE for a robotics task fails to reconstruct a ping pong ball

Input



Reconstruction



The existence of the ping-pong ball and all its spatial coordinates are important underlying causal factors that generate the image and are relevant to the robotics task

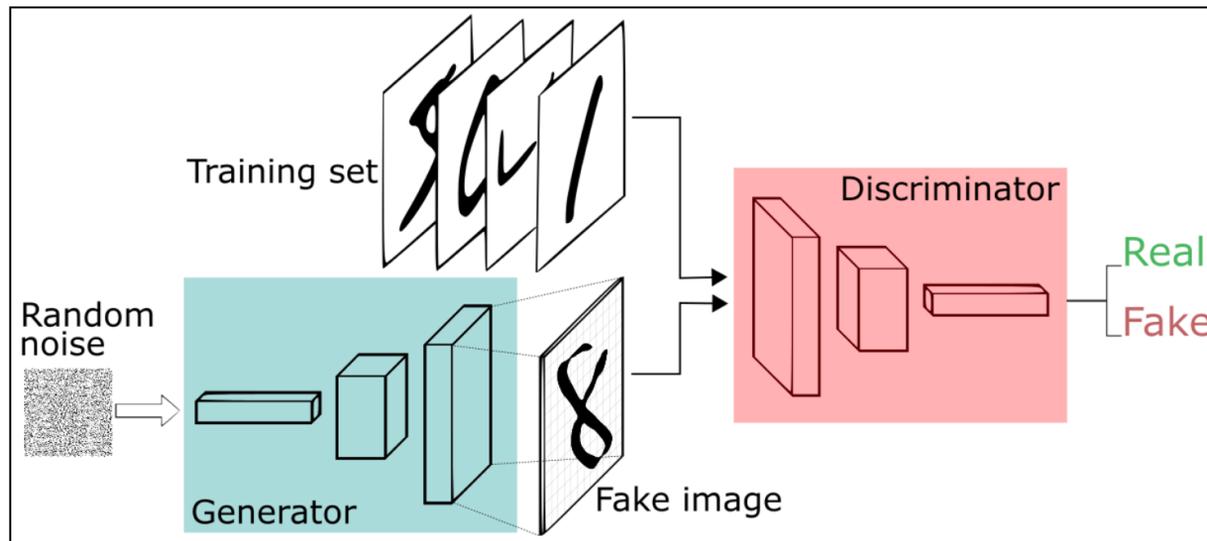
- The autoencoder has limited capacity and training with MSE did not identify ball as salient enough
- Same robot succeeds with larger objects
 - Such as baseballs which are more salient according to MSE

Other definitions of salience

- If a group of pixels follows a highly recognizable pattern even if that pattern does not involve extreme brightness or darkness then that pattern could be considered salient
- One way to implement such a definition of salience is called generative adversarial networks (GANs)

Generative Adversarial Network

- *A generative model (G-network)*
 - is trained to fool a feedforward classifier
 - A feedforward network that generates images from noise
- *A discriminative model (D-network)*
 - A feedforward classifier that attempts to recognize samples from G as fake and samples from training set as real



GANs can determine saliency

- Any structured pattern that the feedforward network (D-network) can recognize is highly salient
 - GANs learn how to determine what is salient

Ex: MSE vs GANs

- Models trained to generate human heads neglect to generate the ears when trained with MSE
- But generate ears when trained with GANs
- Because the ears are not especially bright or dark compared to surrounding skin
- But their highly recognizable shape and consistent position means the feedforward network can easily learn to detect them

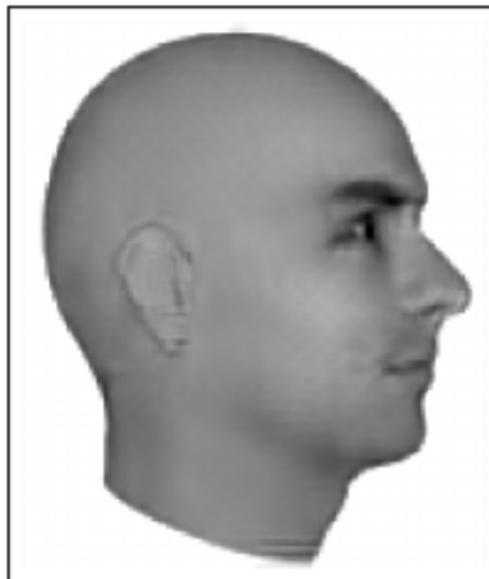
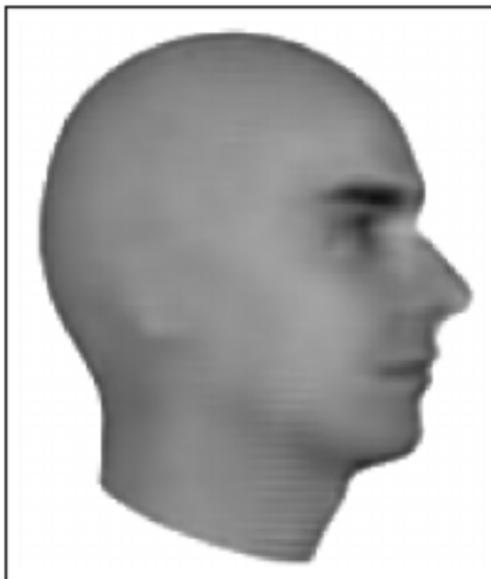
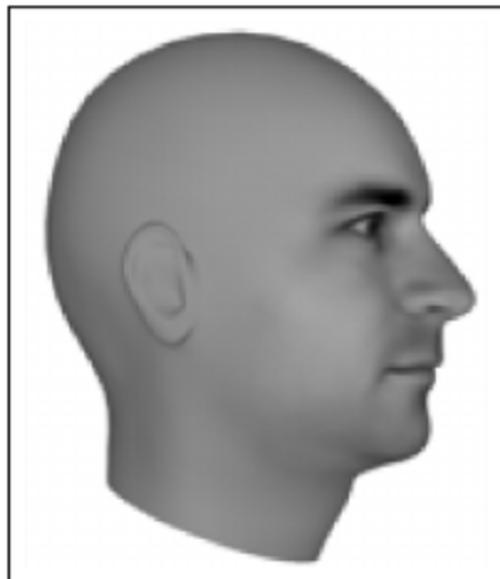
Predictive generative networks

- Models have been trained to predict the appearance of a 3-D model at a view angle

Ground Truth: Correct image that network should emit

MSE: Network trained with MSE alone. Considers ears to be not salient to learn to generate them

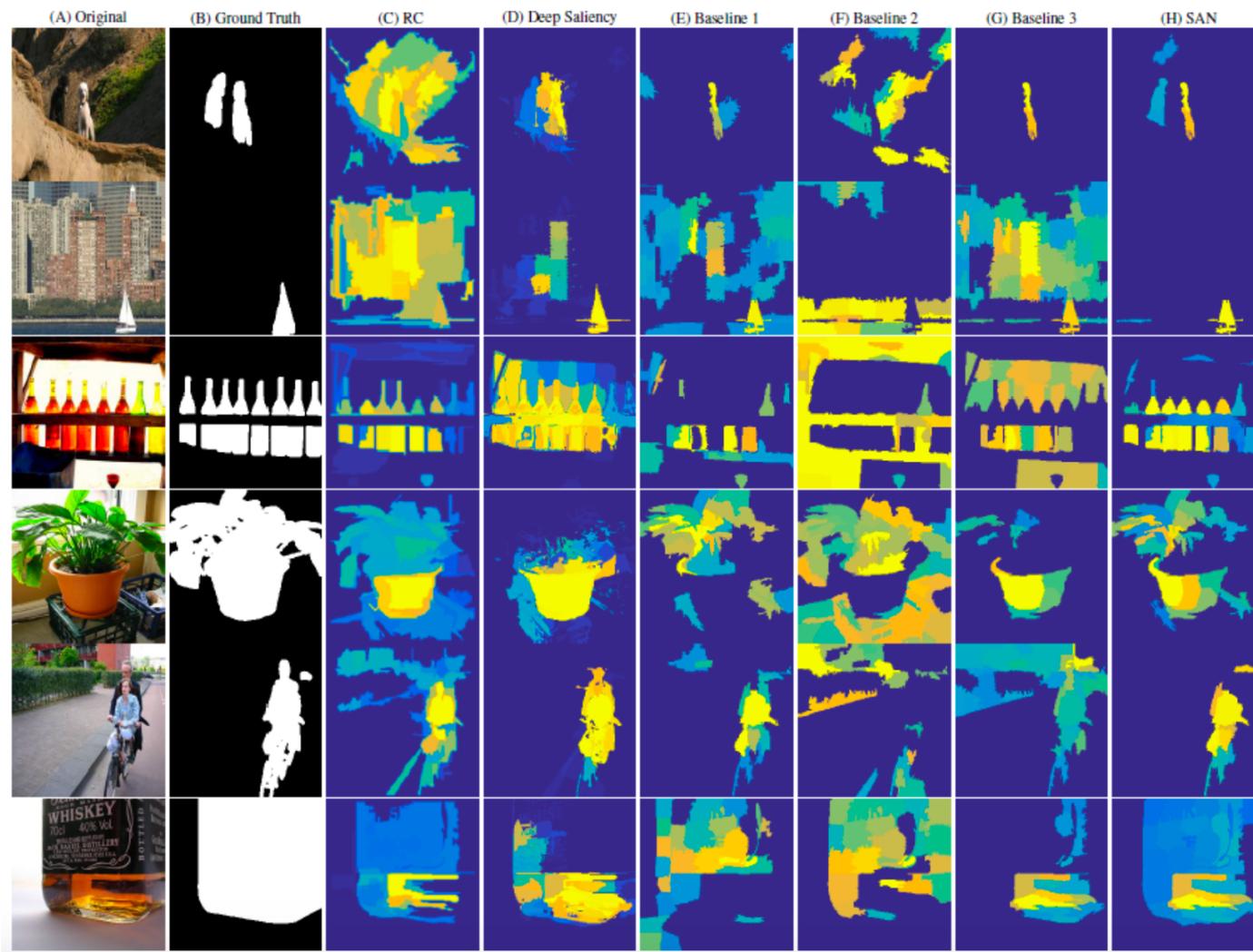
Adversarial: Trained with MSE and adversarial loss. Ears are salient since they follow a predictable pattern



Research on determining salient features

- GANs are only one step toward determining which factors should be represented
- Ongoing research is on
 - ways of determining which factors to represent
 - Develop mechanisms for representing different factors depending on the task

Ex: Saliency Detection using SANs



Invariance of Causal Factors

- Causal factors remain invariant
 - Hence learn a generative model that attempts to recover the causal factors \mathbf{h} and $p(\mathbf{x}|\mathbf{h})$