# Data Driven Estimation of Likelihood Ratios Application to GANs

GEORGE MOUSTAKIDES

UNIVERSITY OF PATRAS, GREECE

# Outline

- Why likelihood ratios ?

- Data driven estimation of likelihood ratios

- Generative models

  - Design using likelihood ratio estimation

  - Generative models vs probability densities

  - Application to inverse problems

# Why Likelihood Ratios ?

Realizations

$X$: $X_1$ $X_2$ $\cdots$ $\cdots$ $\cdots$ $X_n$ $\mathsf{f}(X)$

$X_\mathrm{o}$  For $X_\mathrm{o}$, $\mathsf{f}(X_\mathrm{o})$ expresses the likelihood that $X_\mathrm{o}$ corresponds to a bird



$X$: $X_1$ $X_2$ $\cdots$ $X_n$ $\mathsf{f}(X)$ $\quad$ $Y$: $Y_1$ $Y_2$ $\cdots$ $Y_m$ $\mathsf{g}(Y)$

Statistical similarity ? $\qquad \mathsf{f}(X) \overset{?}{=} \mathsf{g}(X) \quad \Leftrightarrow \quad \dfrac{\mathsf{g}(X)}{\mathsf{f}(X)} \overset{?}{=} 1$

$X:$      $\cdots$    f$(X)$   Cat

$X_1$      $X_2$       $\cdots$     $X_n$

$Y:$      $\cdots$    g$(Y)$   Dog

$Y_1$      $Y_2$       $\cdots$     $Y_m$

New unlabeled data

$X_0$

Decide cat or dog

Decision making
Hypothesis testing
Classification

Optimum
Likelihood Ratio Test

$$\frac{\text{g}(X_0)}{\text{f}(X_0)} \underset{\text{Decide cat}}{\overset{\text{Decide dog}}{\gtrless}} \lambda$$

# Data Driven Estimation
# of Likelihood Ratios

**THEOREM:** Assume $X$ follows density f($X$), $Y$ follows density g($Y$).

Select strictly increasing function $\omega(\mathrm{r})$ and strictly positive function $\rho(z)$. Compute $\phi(z),\ \psi(z)$ from

$$\phi'(z) = \rho(z), \qquad \psi'(z) = -\omega^{-1}(z)\,\rho(z)$$

For D($X$) arbitrary function define cost

$$\mathrm{J}(\mathrm{D}) = \mathbb{E}_{\mathsf{f}}\big[\phi\big(\mathrm{D}(X)\big)\big] + \mathbb{E}_{\mathsf{g}}\big[\psi\big(\mathrm{D}(Y)\big)\big]$$

Then the solution to optimization

$$\max_{\mathrm{D}} \mathrm{J}(\mathrm{D}) = \max_{\mathrm{D}}\Big\{\mathbb{E}_{\mathsf{f}}\big[\phi\big(\mathrm{D}(X)\big)\big] + \mathbb{E}_{\mathsf{g}}\big[\psi\big(\mathrm{D}(Y)\big)\big]\Big\} \;\Rightarrow\; \mathrm{D}_{\mathsf{o}}(X) = \omega\left(\frac{\mathsf{g}(X)}{\mathsf{f}(X)}\right)$$

# Examples of functions

$\omega(\mathsf{r}) = \mathsf{r}$ (estimate likelihood ratio), if select $\rho(z) = 1$

$$\phi(z) = -\frac{z^2}{2}, \quad \psi(z) = z$$

Mean Square

$\omega(\mathsf{r}) = \log \mathsf{r}$ (estimate log-likelihood ratio), if select $\rho(z) = e^{-0.5z}$

$$\phi(z) = -2e^{0.5z}, \quad \psi(z) = -2e^{-0.5z}$$

Exponential

$\omega(\mathsf{r}) = \dfrac{\mathsf{r}}{1 + \mathsf{r}}$ (estimate posterior probability), if select $\rho(z) = \dfrac{1}{z}$

$$\phi(z) = \log(1 - z), \quad \psi(z) = \log(z)$$

Cross Entropy

# Data Driven Implementation

$$J(D) = \mathbb{E}_f\big[\phi\big(D(X)\big)\big] + \mathbb{E}_g\big[\psi\big(D(Y)\big)\big]$$

$\{X_1, X_2, \dots, X_n\}$ following $f(X)$, $\qquad$ $\{Y_1, Y_2, \dots, Y_m\}$ following $g(Y)$

Approximate: $D(X)$ with neural network $D(X,\vartheta)$, Expectations with sample means

$$\max_{\vartheta} J(\vartheta) = \max_{\vartheta}\left\{ \frac{1}{n}\sum_{i=1}^{n}\phi\big(D(X_i,\vartheta)\big) + \frac{1}{m}\sum_{j=1}^{m}\psi\big(D(Y_j,\vartheta)\big)\right\} \;\Rightarrow\; \vartheta_o \Rightarrow D(X,\vartheta_o)$$

We expect $D(X,\vartheta_o) \approx D_o(X) = \omega\left(\dfrac{g(X)}{f(X)}\right) \;\Rightarrow\; \omega^{-1}\big(D(X,\vartheta_o)\big) \approx \dfrac{g(X)}{f(X)}$

Different $\omega(r),\ \phi(z), \psi(z)$ produce approximation of different quality

## Summary

$\{X_1, X_2, \dots, X_n\}$ following $f(X)$, $\qquad \{Y_1, Y_2, \dots, Y_m\}$ following $g(Y)$

Select $\omega(r), \rho(z)$, compute $\phi(z), \psi(z)$

$$\max_{\vartheta}\left\{\frac{1}{n}\sum_{i=1}^{n}\phi\big(D(X_i, \vartheta)\big) + \frac{1}{m}\sum_{j=1}^{m}\psi\big(D(Y_j, \vartheta)\big)\right\}$$

$$\Rightarrow\ \vartheta_o\ \Rightarrow\ D(X, \vartheta_o)\ \Rightarrow\ \omega^{-1}\big(D(X, \vartheta_o)\big) \approx \frac{g(X)}{f(X)}$$

Compare $\omega^{-1}\big(D(X, \vartheta_o)\big)$ to $1$ to assess whether the two datasets have the same statistical behavior or not
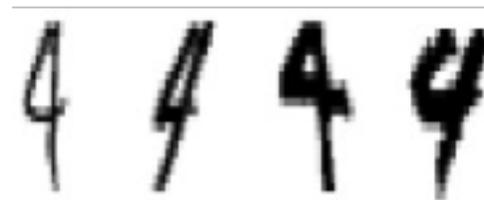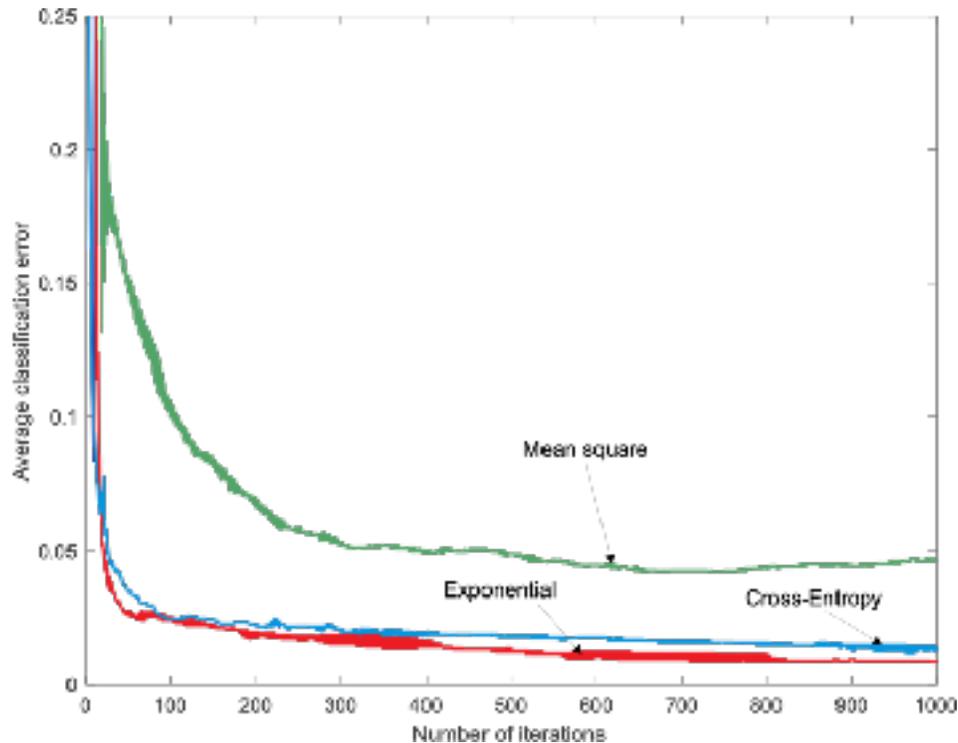
For new sample $X_0$ compare $\omega^{-1}\big(D(X_0, \vartheta_o)\big)$ to threshold $\lambda$ to decide whether $X_0$ statistically follows the first or the second set
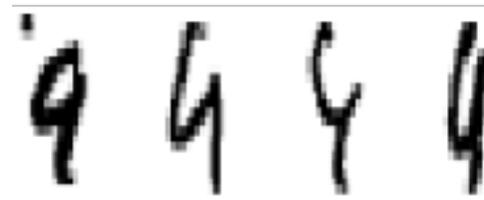
# Example: Classification Problem

From dataset MNIST isolate handwritten numerals 4 and 9

Training set: 5500 "4" and 5500 "9". Testing set: 982 "4" and 1009 "9"

4 mistaken for 9

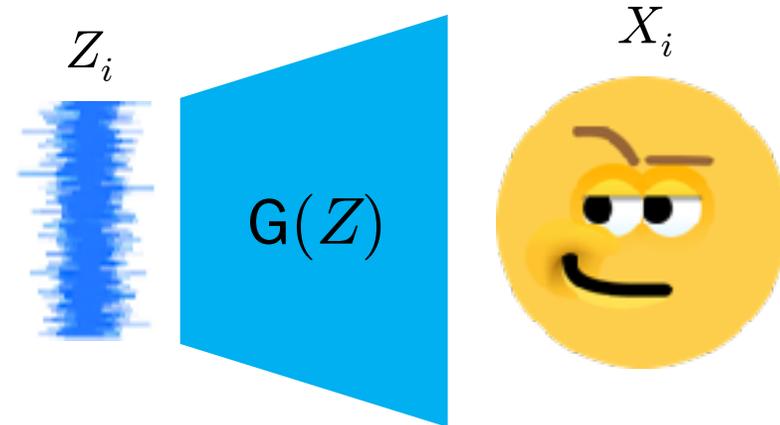9 mistaken for 4

# Generative Models

$X$:  ⋯  ⋯  f$(X)$

Is it possible to generate synthetic data (realizations $X_i$) that follow f$(X)$ ?
NOT an easy problem even if density f$(X)$ is known!

Begin with density h$(Z)$:  Simple to generate realizations $Z_i$
Find transformation G$(Z)$: Such that $X_i = $ G$(Z_i)$ follows f$(X)$

**THEOREM:** Under general conditions
a transformation  G  exists !!!

Pair $\{$G$(Z)$,h$(Z)\}$ Generative model
G$(Z)$ Generator

$Z_i$
$X_i$
G$(Z)$

$X$ follows $f(X)$ and $Y$ follows $g(Y)$

$$\max_{D} J(D) = \max_{D} \left\{ \mathbb{E}_f \big[ \phi\big(D(X)\big) \big] + \mathbb{E}_g \big[ \psi\big(D(Y)\big) \big] \right\} \quad \Rightarrow \quad D_o(X) = \omega \left( \frac{g(X)}{f(X)} \right)$$

$Z$ follows $h(Z)$, select $G(Z)$, define $Y = G(Z)$, check if likelihood ratio = 1

**THEOREM** (Goodfellow et al. 2014): $Z$ follows $h(Z)$, define $Y = G(Z)$ and cost

$$J(G, D) = \mathbb{E}_f \big[ \phi\big(D(X)\big) \big] + \mathbb{E}_h \left[ \psi\Big( D\big(G(Z)\big) \Big) \right]$$

then the optimum solution to the adversarial problem

$$\min_{G} \max_{D} J(G, D) = \min_{G} \max_{D} \left\{ \mathbb{E}_f \big[ \phi\big(D(X)\big) \big] + \mathbb{E}_h \left[ \psi\Big( D\big(G(Z)\big) \Big) \right] \right\}$$

is such that $Y = G_o(Z)$ follows $f(Y)$

$D(X)$ Discriminator          $G(Z)$ Generator

# Data Driven Implementation

$\{X_1, X_2, \ldots, X_n\}$ following $f(X)$, $\qquad \{Z_1, Z_2, \ldots, Z_m\}$ following $h(Z)$

Likelihood ratio $D(X)$ approximated by neural network $D(X,\vartheta)$ (Discriminator)

Generator function $G(Z)$ approximated by neural network $G(Z,\theta)$ (Generator)

$$J(\theta, \vartheta) = \frac{1}{n} \sum_{i=1}^{n} \phi\big(D(X_i, \vartheta)\big) + \frac{1}{m} \sum_{j=1}^{m} \psi\Big(D\big(G(Z_j, \theta), \vartheta\big)\Big)$$

Adversarial optimization becomes

$$\min_{\theta} \max_{\vartheta} J(\theta, \vartheta) = \min_{\theta} \max_{\vartheta} \left\{ \frac{1}{n} \sum_{i=1}^{n} \phi\big(D(X_i, \vartheta)\big) + \frac{1}{m} \sum_{j=1}^{m} \psi\Big(D\big(G(Z_j, \theta), \vartheta\big)\Big) \right\}$$

$$\Rightarrow \{\theta_o, \vartheta_o\} \Rightarrow \theta_o \Rightarrow G(Z, \theta_o)$$

Generative Adversarial Networks

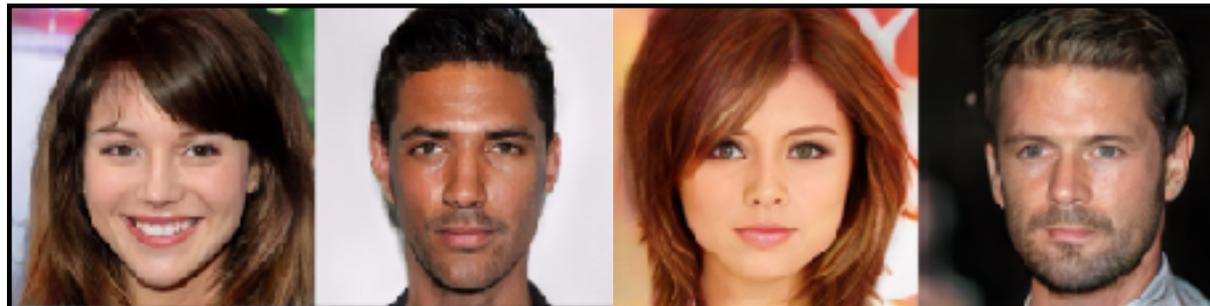IF $Z$ follows $h(Z)$   THEN   $Y = G(Z, \theta_o)$ follows $f(Y)$

**Example** (NVIDIA)

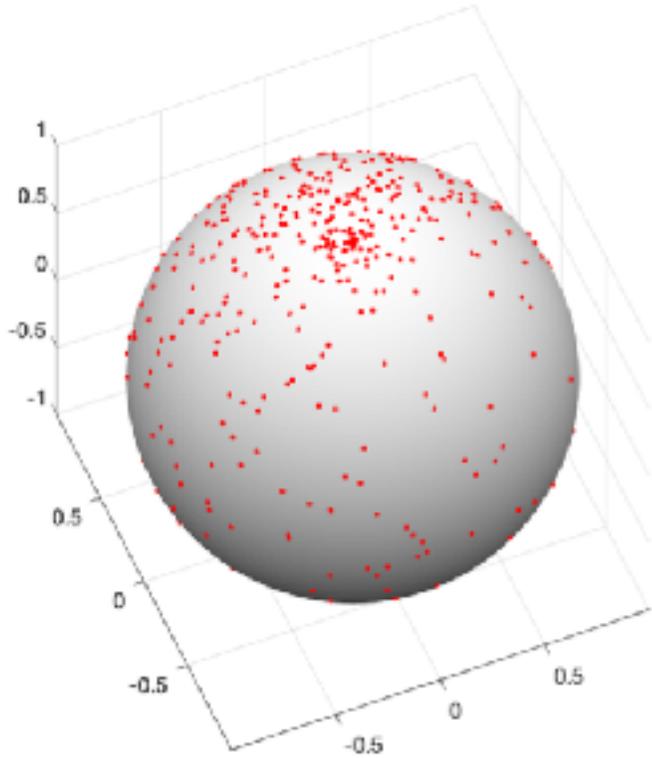HD-CelebA (30 000 high definition images 1024 X 1024 of celebrities)



NVIDIA used progressive growing of GANs (4X4), (8X8),...,(1024X1024)

$Y$ of size 3 X 10$^6$, $Z$ Gaussian vector of length 500

# Generative Models
## vs
# Probability Densities

Points in N-D space can be random and lie on a lower dimensional surface (manifold)



Example red points on sphere (2-D in 3-D space)

Points are random with coordinates $Y = [y_1, y_2, y_3]$ satisfying the deterministic equation

$$y_1^2 + y_2^2 + y_3^2 = r^2$$

Then density has the form

$$\mathsf{f}(y_1, y_2, y_3) = \delta(y_1^2 + y_2^2 + y_3^2 - r^2)\mathsf{h}(y_1, y_2)$$

Dirac $\delta(x)$ generalized function is defined as

$$\delta(x) = \left\{ \begin{array}{ll} 0 & x \neq 0 \\ \infty & x = 0 \end{array} \right. , \quad \int_{-\epsilon}^{\epsilon} \delta(x)\,dx = 1$$

Generative model would describe the random data with input density $h(z_1,z_2)$ and generator vector function $G(z_1,z_2)$

$$Y = G(z_1, z_2) \Rightarrow \begin{bmatrix} y_1 = G_1(z_1, z_2) \\ y_2 = G_2(z_1, z_2) \\ y_3 = G_3(z_1, z_2) \end{bmatrix} \Rightarrow \begin{bmatrix} y_1 = r\cos(2\pi z_1)\sin(\pi z_2) \\ y_2 = r\sin(2\pi z_1)\sin(\pi z_2) \\ y_3 = r\cos(\pi z_2) \end{bmatrix}$$

$h(z_1,z_2)$ defined on $[0,1]\times[0,1]$ and $G(z_1,z_2)$ is an ordinary function

Data are representable as $Y = G(Z)$, $Z$ follows $h(Z)$. Many datasets satisfy

$$\dim(Z) \ll \dim(Y)$$

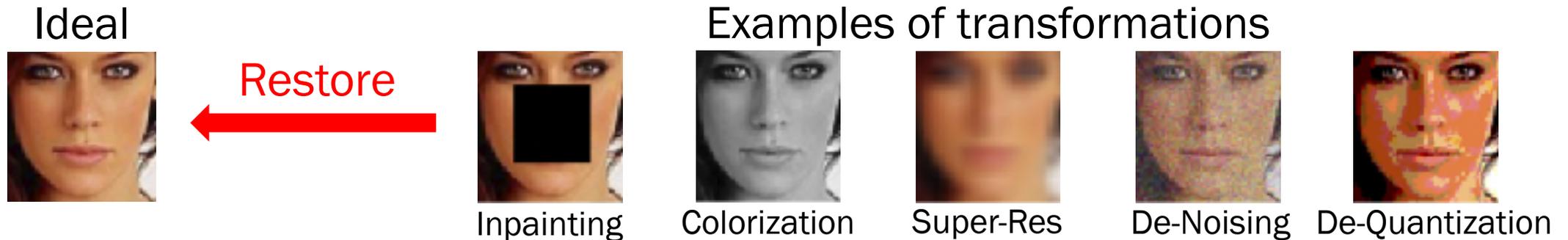In HD CelebA:      $\dim(Y) = 3\ \text{X}\ 1024\ \text{X}\ 1024 = 3\ \text{X}\ 10^6$
Input to Generator $G(Z)$:      $\dim(Z) = 500$ (independent Gaussians)

# Application to
# Inverse Problems

Several image restoration problems in Computer Vision can be formulated as follows

Measurement $\quad X = \mathsf{T}(Y) + W \qquad$ more general $\ X = \mathsf{T}(Y, \alpha) + W$

Ideal

Known transformation

Noise

Unknown parameters

**Problem:** Recover (restore) ideal $Y$ from measurements $X$

Ideal

Examples of transformations

Restore

Inpainting · Colorization · Super-Res · De-Noising · De-Quantization

Recovering $Y$ from measurements $X$ is an ill posed problem

$Y$

$X$

Inpainting

More unknowns than equations

Classical approach: Impose "smoothness" constraints to obtain a (unique) solution



… …

Available generative model $\{\mathsf{G}(Z),\mathsf{h}(Z)\}$: $\quad Y = \mathsf{G}(Z)$

Since $Y = \mathsf{G}(Z)$, instead of estimating $Y$, estimate input to generator $Z$ then recover $Y$ as the output of the generator

Because $\dim(Z) \ll \dim(Y)$, significant computational gain and stable processing

# Ad-Hoc Approaches

Select $Z$ so that measurement $X$ and $\mathsf{T}(\mathsf{G}(Z))$ are "close"

$$\min_{Z} \|X - \mathsf{T}(\mathsf{G}(Z))\|^2 \;\Rightarrow\; Z_{\mathsf{o}} \;\Rightarrow\; Y_{\mathsf{o}} = \mathsf{G}(Z_{\mathsf{o}})$$

Well defined optimization, computationally stable



$Y$　　　$X$　　　$\mathsf{G}(Z_{\mathsf{o}})$

Failure

Generative model is a pair $\{\mathsf{G}(Z), \mathsf{h}(Z)\}$

Even for $\mathsf{T}(\mathsf{G}(Z_{\mathsf{o}}))$ "close" to $X$, if likelihood $\mathsf{h}(Z_{\mathsf{o}})$ is very small
then $Y_{\mathsf{o}} = \mathsf{G}(Z_{\mathsf{o}})$ is a bad solution

Must take into account <span style="color:red">input density $\mathsf{h}(Z)$</span>

Yeh et al. (2017), (2018)

$$J(Z) = \|X - \mathsf{T}(\mathsf{G}(Z))\|$$

Regularizer

$\log \mathsf{h}(Z)$

$$+ \lambda \left\{ \log\big(1 - \mathsf{D}(\mathsf{G}(Z))\big) - \log\big(\mathsf{D}(\mathsf{G}(Z))\big) - \frac{1}{2}\|Z\|^2 \right\}$$

$$\min_Z \mathsf{J}(Z) \;\Rightarrow\; Z_\mathsf{o} \;\Rightarrow\; Y_\mathsf{o} = \mathsf{G}(Z_\mathsf{o})$$

Parameter needs tuning
Complicated



Success ?

Asim et al. (2019)

$$\mathsf{J}(Z) = \|X - \mathsf{T}(\mathsf{G}(Z))\|^2 + \lambda\|Z\|^2 \qquad \min_Z \mathsf{J}(Z) \;\Rightarrow\; Z_\mathsf{o} \;\Rightarrow\; Y_\mathsf{o} = \mathsf{G}(Z_\mathsf{o})$$

Both methods require exact knowledge of $\mathsf{T}(Y)$

## Statistical Estimation

Following classical optimal Statistical estimation theory, in particular the Maximum Aposteriori Probability (MAP) method we obtain

$$\mathsf{J}(Z, \alpha) = \log\left(\|X - \mathsf{T}(\mathsf{G}(Z), \alpha)\|^2\right) + \frac{1}{N}\|Z\|^2, \quad N = \dim(X)$$

$$\min_{Z, \alpha} \mathsf{J}(Z, \alpha) \; \Rightarrow \; \{Z_{\mathsf{o}}, \alpha_{\mathsf{o}}\} \; \Rightarrow \; Y_{\mathsf{o}} = \mathsf{G}(Z_{\mathsf{o}})$$

No parameters to tune

Can accommodate unknown parameters in transformation $\mathsf{T}(Y, \alpha)$

# Examples

## Blurring with 3 X 3 mask

| $Y$ | $X$ | Yeh | Asim | Known | Unknown |
|---|---|---|---|---|---|



## Colorization (green channel)

| $Y$ | $X$ | Yeh | Asim | Known | Unknown |
|---|---|---|---|---|---|

# De-Quantization

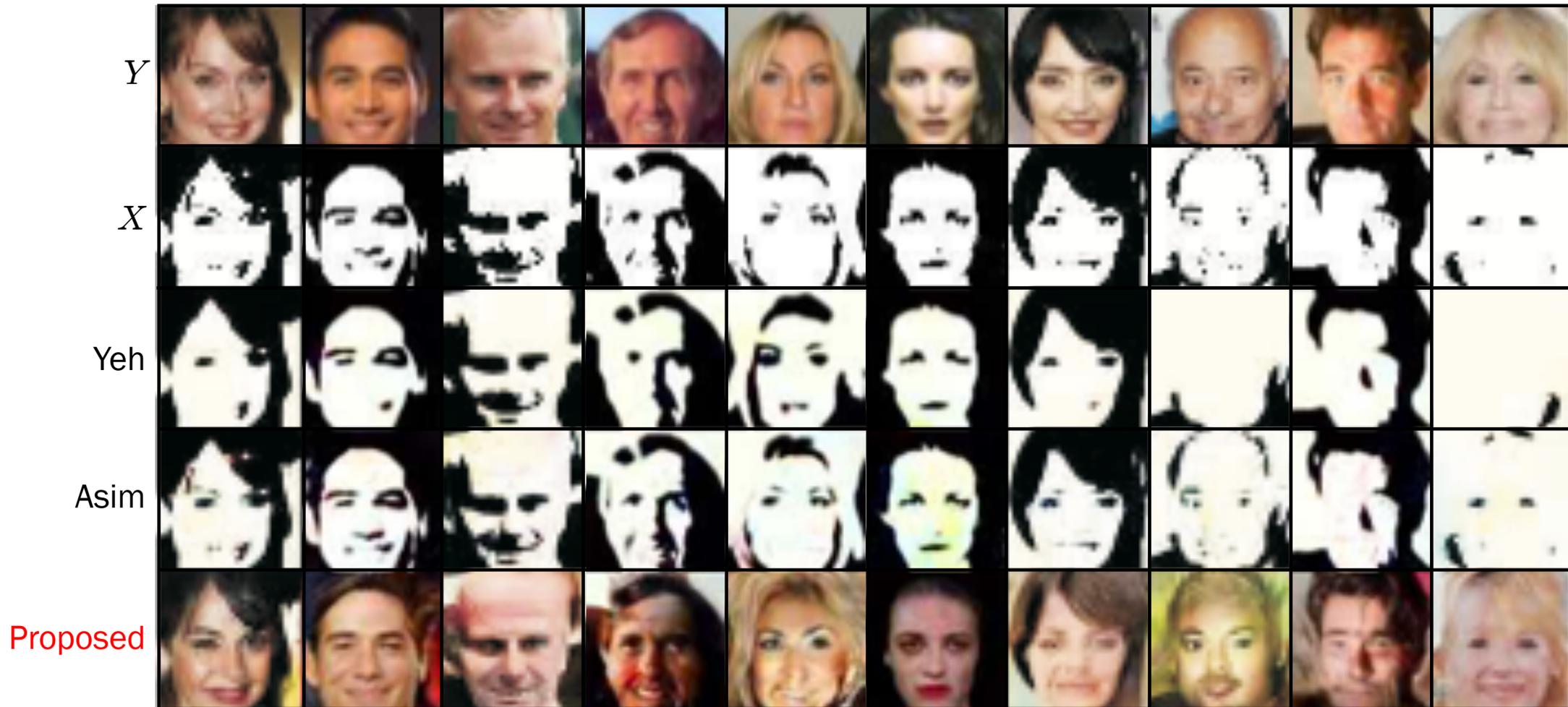2 levels per RGB channel, 8 colors



$Y$

$X$

Yeh

Asim

Proposed

# De-Quantization

3 levels per RGB channel, 27 colors

# De-Quantization and Colorization

RGB → Gray → BW (2 levels)

# Data Mixtures

$$X = aY + a'Y'$$

$Y \qquad Y' \qquad X \qquad$ Known $\qquad$ Unknown