# Bird Transformers: Modifying Birds Using Natural Language

## Emma Qin, Will Schwarzer, Kyra Wilson, Orlando Zuniga

# "Make this bird green!"

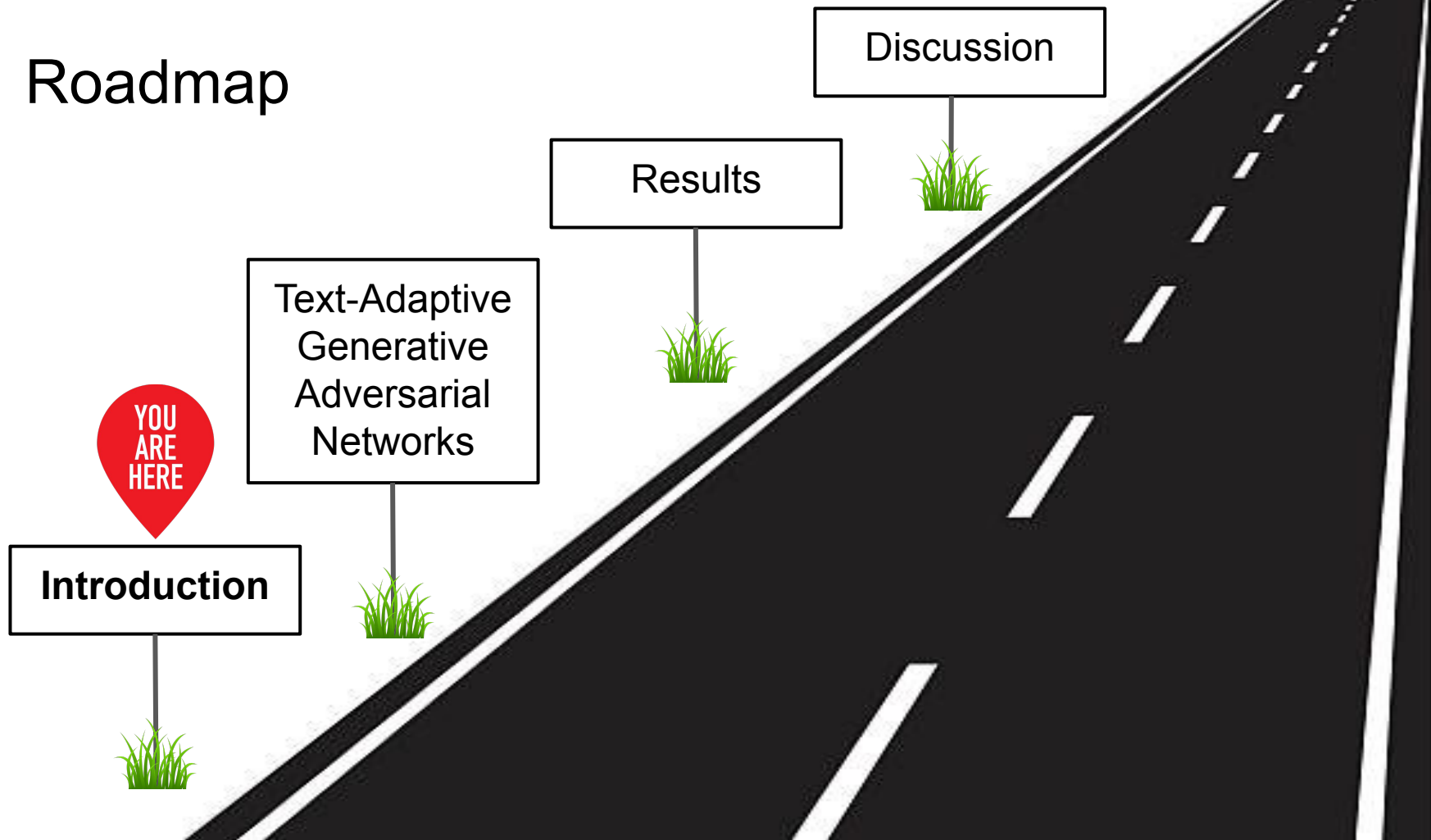Specialized Skills

Computers
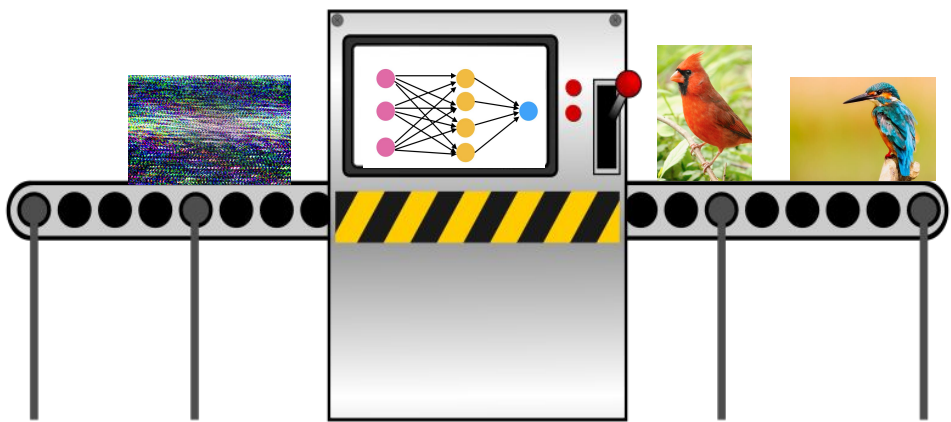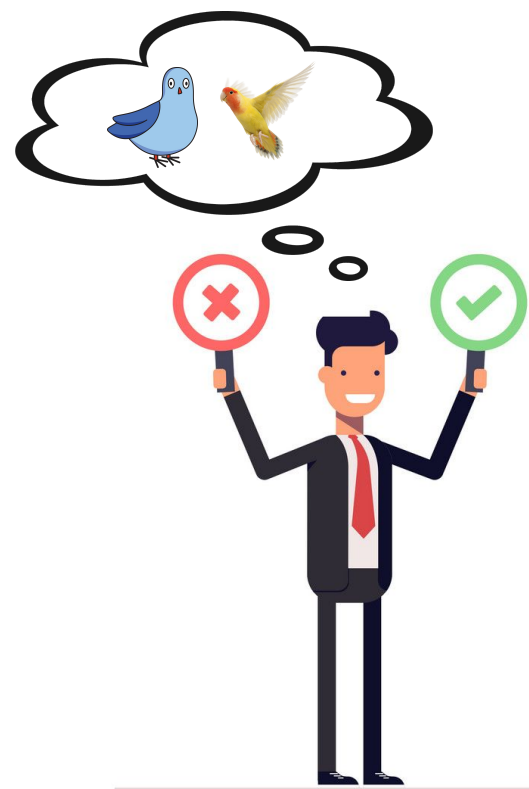
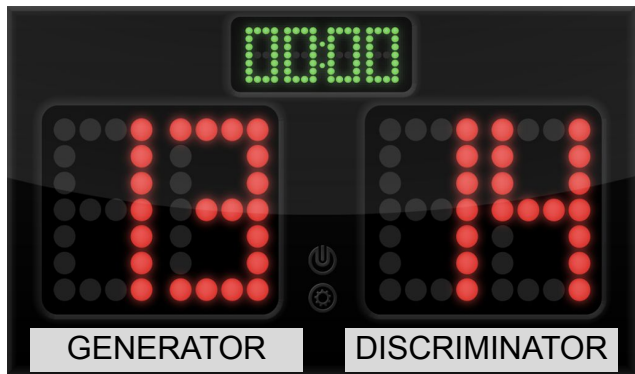# Image Modification with Natural Language

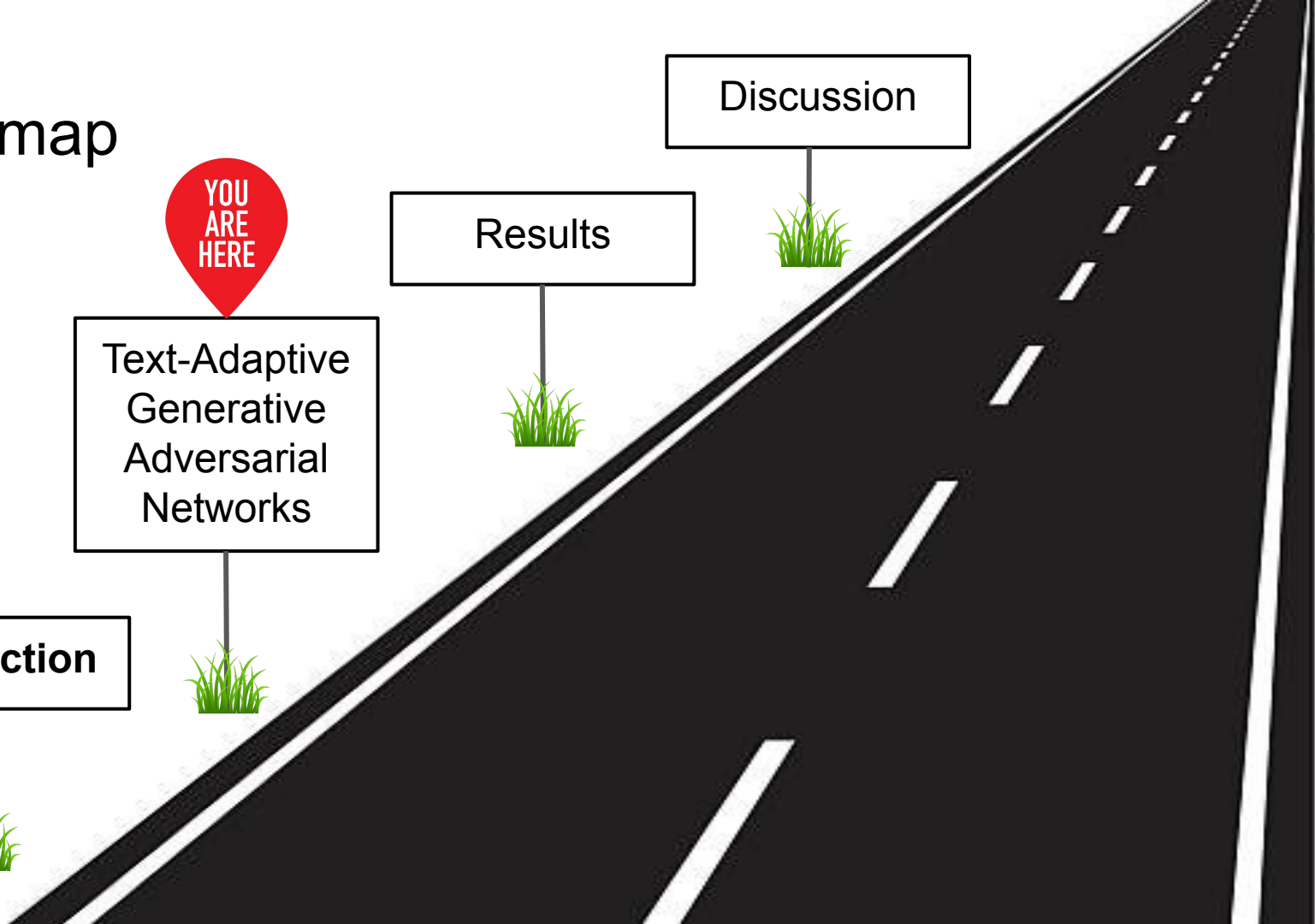# Generative Adversarial Networks (GANs)

Generator

Discriminator

# Roadmap

**Introduction**

Text-Adaptive Generative Adversarial Networks

YOU ARE HERE

Results

Discussion

# Image Encoding



| 0 | 0 | 0 | 0 | 0 | 30 | 0 |
|---|---|---|---|---|----|---|
| 0 | 0 | 0 | 0 | 30 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

# Image Encoding

(250, 24, 170)

30 25 27 2

| 0 | 0 | 0 | 0 | 0 | 30 | 0 |
|---|---|---|---|---|----|---|
| 0 | 0 | 0 | 0 | 30 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

# Convolutional Neural Network



INPUT     CONVOLUTION + RELU     POOLING     CONVOLUTION + RELU     POOLING

FEATURE LEARNING

# Intuition: modifying image with text



Black crown

Grey wings

Yellow
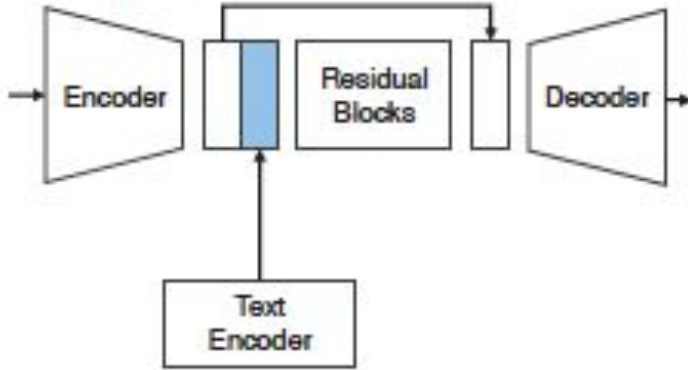
A yellow bird with grey wings and a black crown

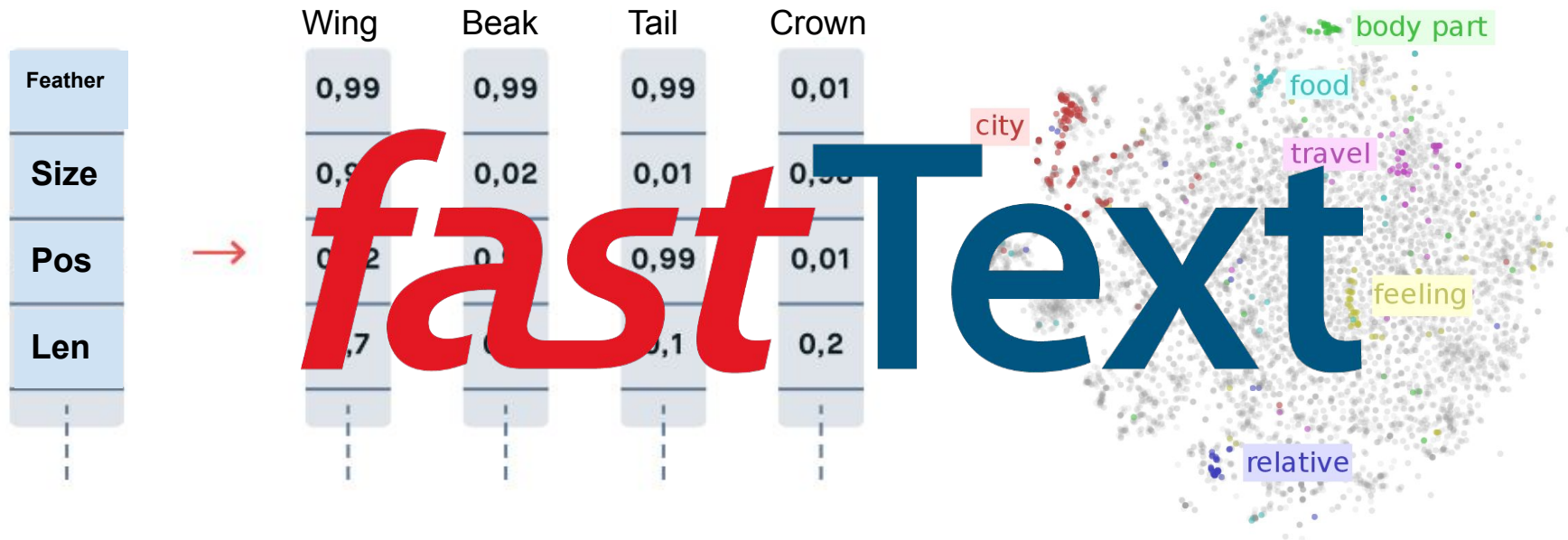A bird with a mix of black and white

Generator

Encoder

Residual Blocks

Decoder

Text Encoder

The bird is a mix of black and white..

# Text Encoding



| | Wing | Beak | Tail | Crown |
|---|---|---|---|---|
| Feather | 0,99 | 0,99 | 0,99 | 0,01 |
| Size | 0,9 | 0,02 | 0,01 | 0,0 |
| Pos | 0,2 | 0, | 0,99 | 0,01 |
| Len | ,7 | | 0,1 | 0,2 |

fastText

city
body part
food
travel
feeling
relative

# Bidirectional GRU

# Discriminator

# Loss

- Reconstruction loss
- Unconditional loss
- Conditional loss

# Recap

What's a GAN?



Generator

Discriminator

# Recap

The problem?

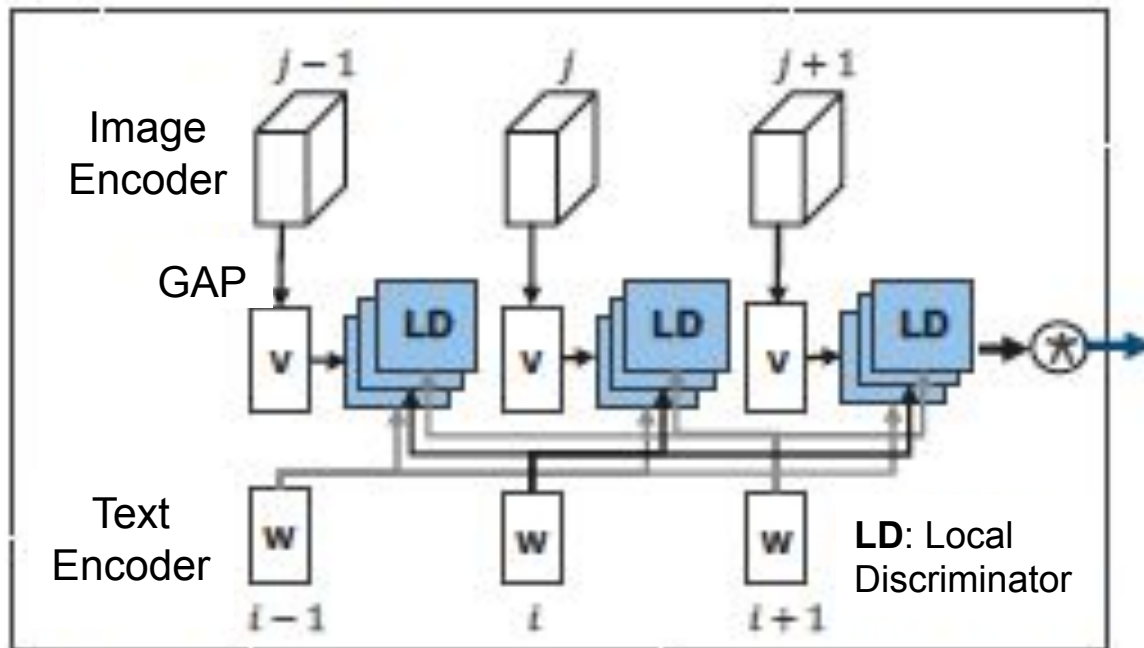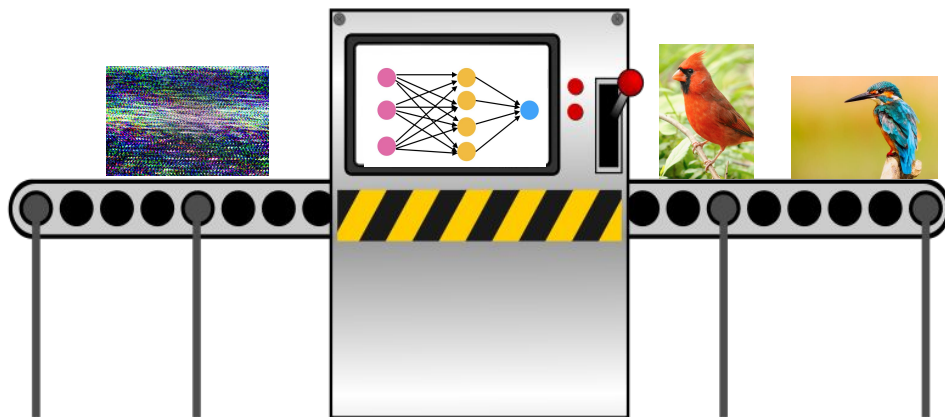- We want to manipulate images based on natural language descriptions

How does TAGAN differ?

- **Key difference:** the text-adaptive discriminator creates word-level local discriminators.

# Results (Generator)



5 epochs



20 epochs

# Results (Discriminator pt1.)

**Generator with no text**



**Network with small dataset**

Epoch 1

# Epoch 6

Epoch 10

Epoch 20

# 30 Epochs



Real Image

**Caption:** this particular bird has a **belly that is white and has black spots**



Fake Image

# TAGAN Qualitative Results

**TAGAN Results:**

**Our Results:**

Original



Original



The bird has **wings that are blue** and has a **white belly.**



The small bird is **yellow and brown speckled** with a **pointed beak**

# TAGAN Quantitative Results

| Method | Accuracy | Naturalness | $L_2$ Error |
|---|---|---|---|
| SISGAN | 2.33 | 2.34 | 0.30 |
| AttnGAN | 2.19 | 2.11 | 0.25 |
| TAGAN | 1.49 | 1.56 | 0.11 |
| Ours | N/A | N/A | ? |

Roadmap

Introduction

Text-Adaptive
Generative
Adversarial
Networks

Results

**Discussion**

YOU
ARE
HERE

# Discussion

- Why are our results different?
- Can we make them not different?
- If we made them not different, what would we then make different?

# Discussion

- **Why are our results different?**
- Can we make them not different?
- If we made them not different, what would we then make different?

Q: Why are our results different?
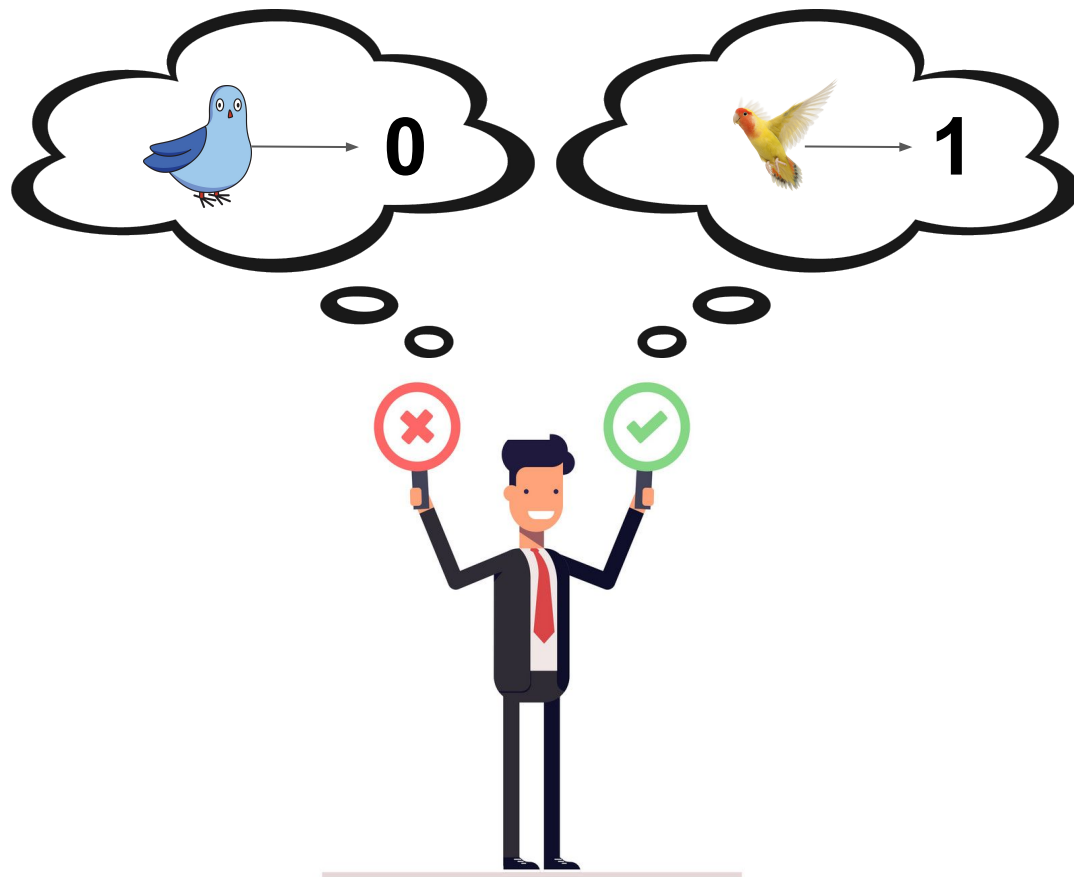
A: Unstable training

Discriminator

Discriminator

"this particular bird has a belly that is white and has black spots"

# Discussion

- Why are our results different? *A: Unstable training*
- **Can we make them not different?**
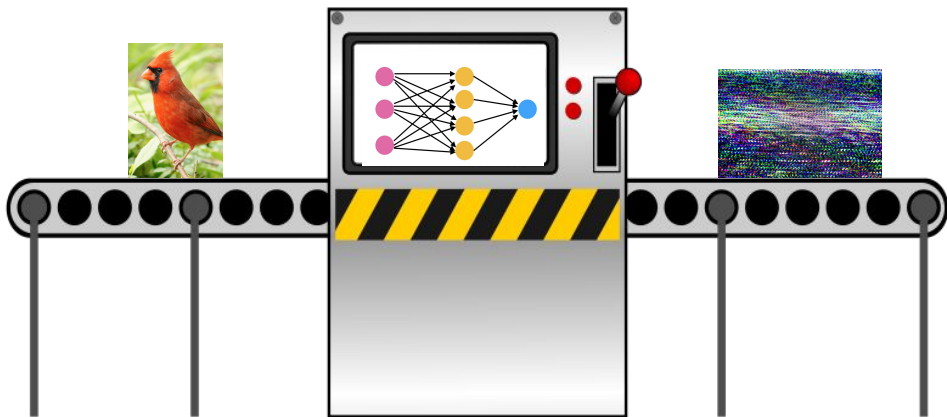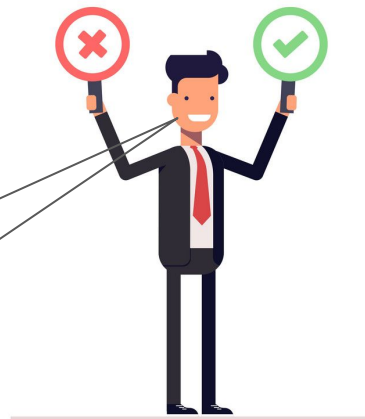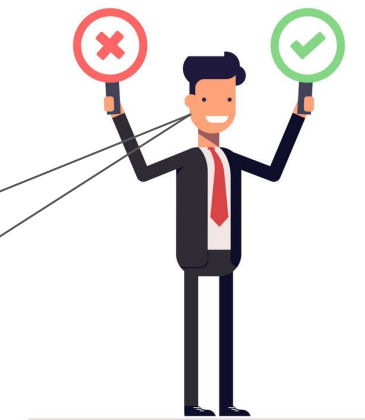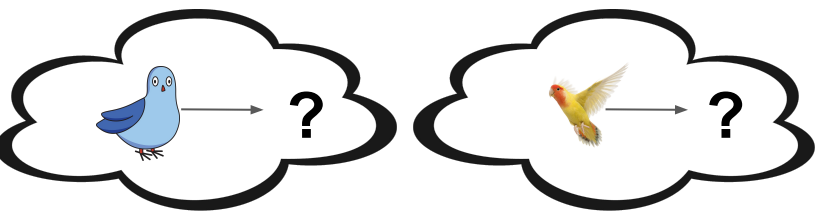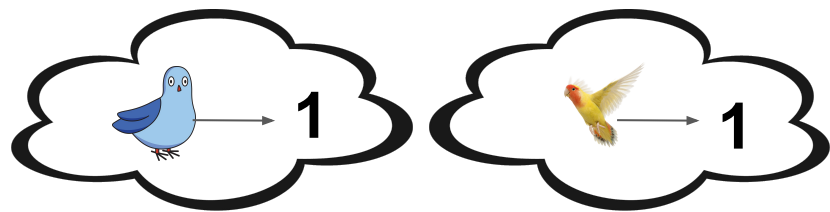- If we made them not different, what would we then make different?

# Could it be our problem?



(Buggy) Discriminator

Training

(Buggy) Discriminator

# Could it be the paper's problem?

| Module | Layers | Input size | Output size |
|---|---|---|---|
| Image Encoder | Conv2d(4, 2), LeakyReLU(0.2) | $3\times128\times128$ | $64\times64\times64$ |
| | Conv2d(4, 2), BN, LeakyReLU(0.2) | $64\times64\times64$ | $128\times32\times32$ |
| conv3 | Conv2d(4, 2), BN, LeakyReLU(0.2) | $128\times32\times32$ | $256\times16\times16$ |
| conv4 | Conv2d(4, 2), BN, LeakyReLU(0.2) | $256\times16\times16$ | $512\times8\times8$ |
| conv5 | Conv2d(4, 2), BN, LeakyReLU(0.2) | $512\times8\times8$ | $512\times4\times4$ |
| Unconditional Discriminator | Conv2d(4, 0), Softmax | $512\times4\times4$ | $1\times1\times1$ |
| Text Encoder | Bidirectional GRU | # of words $\times$ 300 | # of words $\times$ 512 |
| $\beta_{ij}$ | Linear, Softmax | # of words $\times$ 512 | # of words $\times$ 3 |
| $\alpha_i$ | See Eq. (3) in the paper | # of words $\times$ 512 | # of words $\times$ 1 |
| $f_{\mathbf{w}_i,j}$ | Linear (See Eq. (2) in the paper) | N/A | N/A |
| From conv3 | Conv2d(3, 1), BN, LeakyReLU(0.2) | $256\times16\times16$ | $256\times16\times16$ |
| (a) | Global Average Pooling | $256\times16\times16$ | $256\times1\times1$ |
| From conv4 | Conv2d(3, 1), BN, LeakyReLU(0.2) | $512\times8\times8$ | $512\times8\times8$ |
| (b) | Global Average Pooling | $512\times8\times8$ | $512\times1\times1$ |
| From conv5 | Conv2d(3, 1), BN, LeakyReLU(0.2) | $512\times4\times4$ | $512\times4\times4$ |
| (c) | Global Average Pooling | $512\times4\times4$ | $512\times1\times1$ |
| Conditional Discriminator | See Eq. (5) in the paper with $(\alpha_i, \beta_{ij}, f_{\mathbf{w}_i,j}, (a), (b), (c))$ | N/A | $1\times1\times1$ |

# Could it be the model's problem?

Not only are GANs famous for training instability, the authors themselves mention having such problems:

- "Note that we do not penalize generated outputs using the conditional discriminator in Eq. (6) due to **instability of training**."
- "We set $\lambda_1$ and $\lambda_2$ to 10 and 2 respectively considering both the visual quality and the **training stability**."

Q: But then how did they transform birds???

(Possible) A: Trying over and over and over and over and over and over and over and over and over again
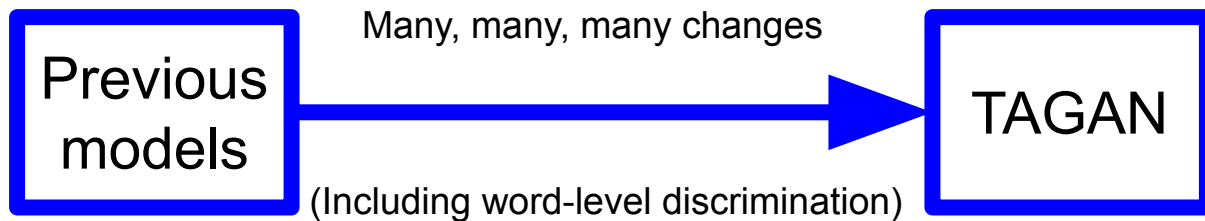
# Discussion

- Why are our results different? *A: Unstable training*
- Can we make them not different? *A: It depends, but probably*
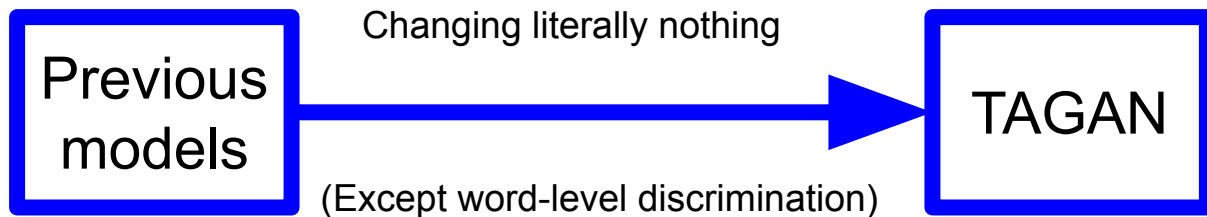- **If we made them not different, what would we then make different?**

# Putting the science back in CS: the value of controls

# Putting the science back in CS: the value of controls

# Putting the science back in CS: the value of controls

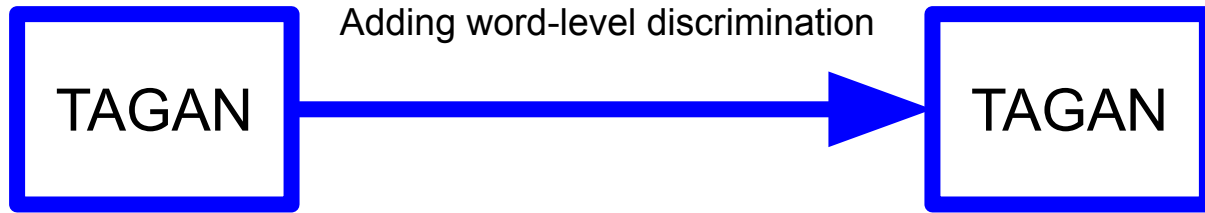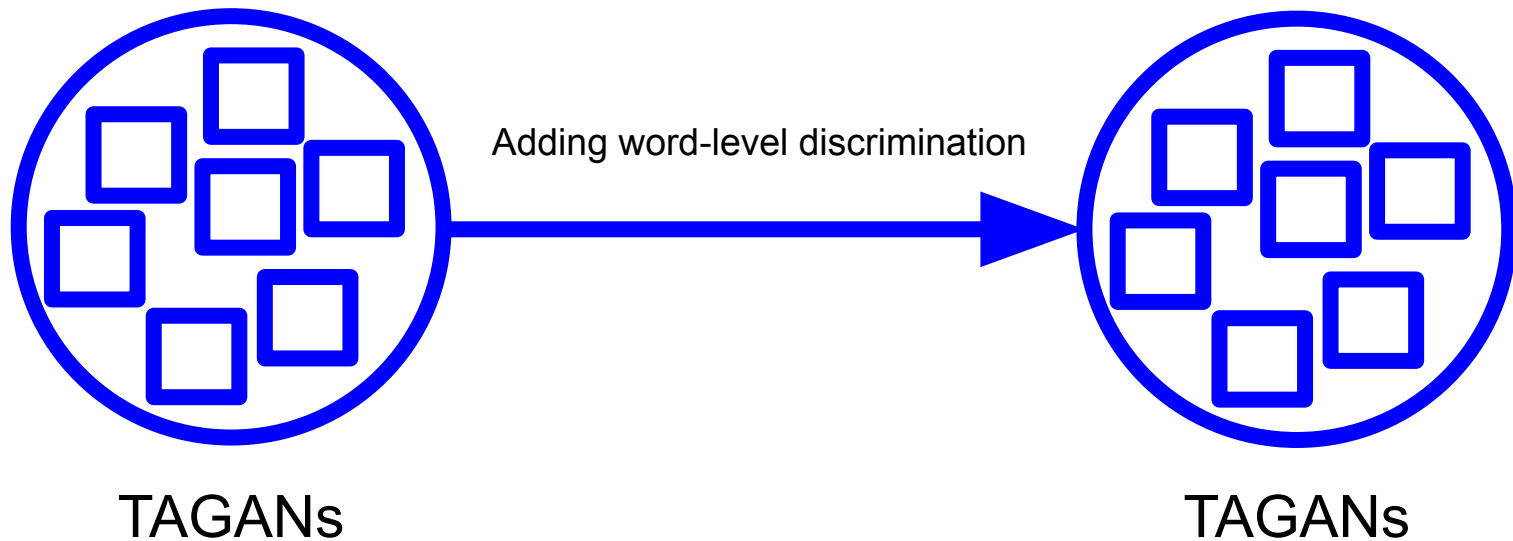| Image Encoder | Conv2d(4, 2), LeakyReLU(0.2) | $3 \times 128 \times 128$ | $64 \times 64 \times 64$ |
|---|---|---|---|
| | Conv2d(4, 2), BN, LeakyReLU(0.2) | $64 \times 64 \times 64$ | $128 \times 32 \times 32$ |
| conv3 | Conv2d(4, 2), BN, LeakyReLU(0.2) | $128 \times 32 \times 32$ | $256 \times 16 \times 16$ |

```python
self.conv123 = nn.Sequential(
    nn.Conv2d(3, 64, 4, 2, padding=1, bias=False),
    nn.LeakyReLU(negative_slope=0.2, inplace=True),
    nn.Conv2d(64, 128, 4, 2, padding=1, bias=False),
    nn.BatchNorm2d(128),
    nn.LeakyReLU(negative_slope=0.2, inplace=True),
    nn.Conv2d(128, 256, 4, 2, padding=1, bias=False),
    nn.BatchNorm2d(256),
    nn.LeakyReLU(negative_slope=0.2, inplace=True)
)
```

# Putting the science back in CS: DIY replication

# Discussion

- Why are our results different? *A: Unstable training*
- Can we make them not different? *A: It depends, but probably*
- If we made them not different, what would we then make different? *A: Science*

# Acknowledgements





Seonghyeon Nam, Yunji Kim, and Seon Joo Kim

# Thank you!

Questions?