



**Semi-supervised learning with GANs : putting together the pieces together**

**OLX Berlin Data Science Team**

[illegible]

**OLX** Free classifieds

**Find what  
you're looking for,  
at prices you want**

What are you looking for?

## Categories

Browse through some of our most popular categories



## Vehicles

- Cars & Bakkies
- Car Parts & Accessories
- Motorcycles & Scooters
- Trucks & Commercial Vehicles

All ads in Vehicles >



### Property

Houses & Flats for rent  
Rooms for rent & Shared  
Houses & Flats for sale  
Land

[All ads in Property >](#)



**Electronics & Computers**

- Cell Phones
- TV, Audio & Visual
- Computers & Laptops
- Gaming & Consoles

All ads in Electronics & Computers >

# OLX

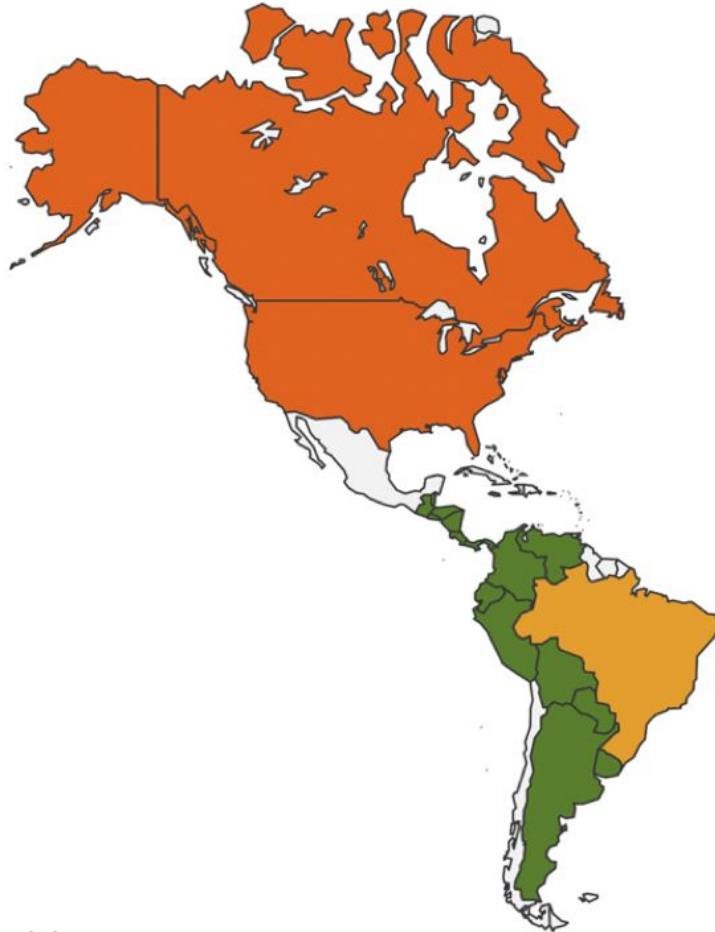


OLX GROUP

# OLX in a glance

**350+ million  
active users**

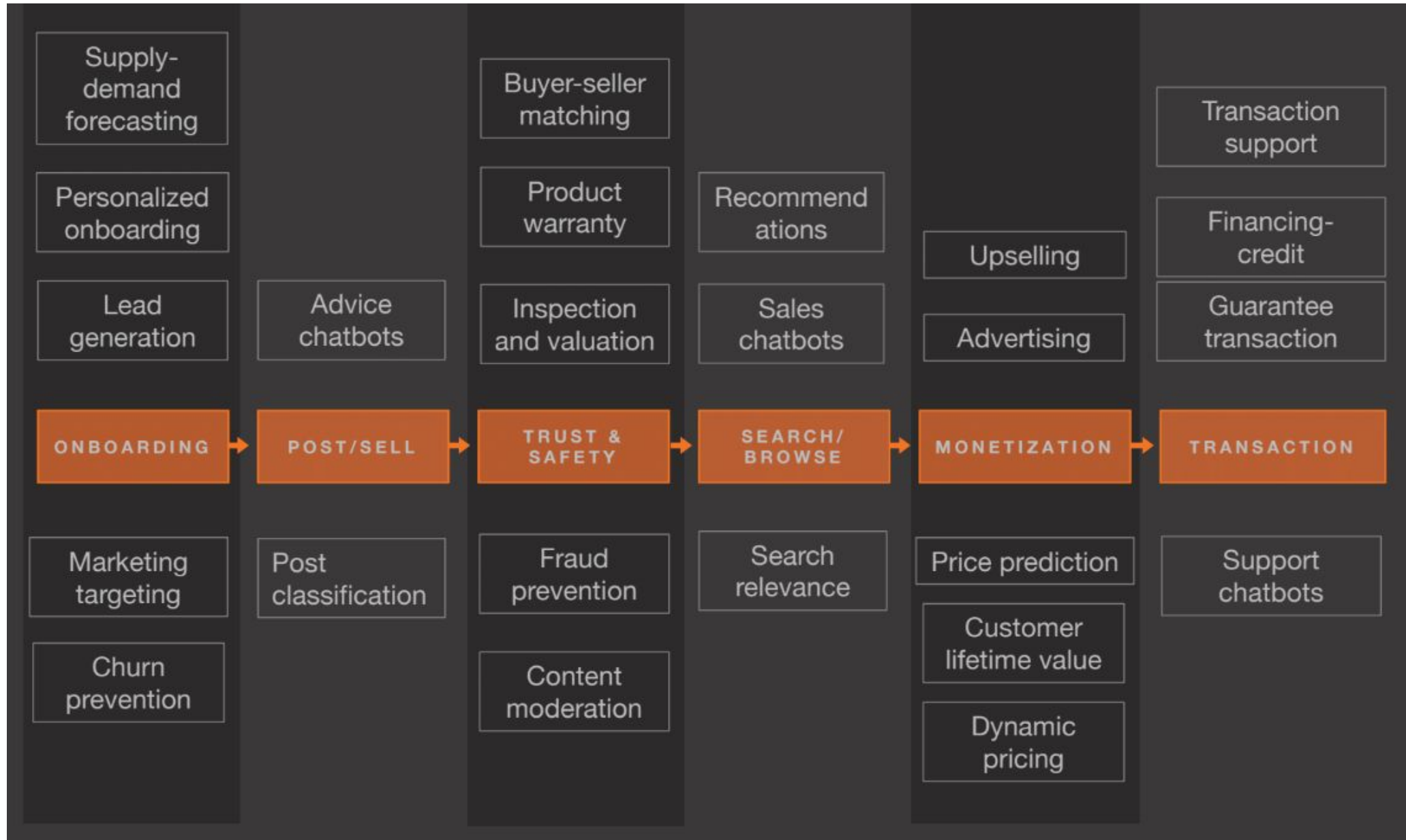
**40+ countries**



**60+ million new  
listings per month**

**Market leader  
in 35 countries**

# Classifieds - Data Science areas of Application



# Outline

Short Intro

What are GANs

How GANs help in a semi-supervised setup

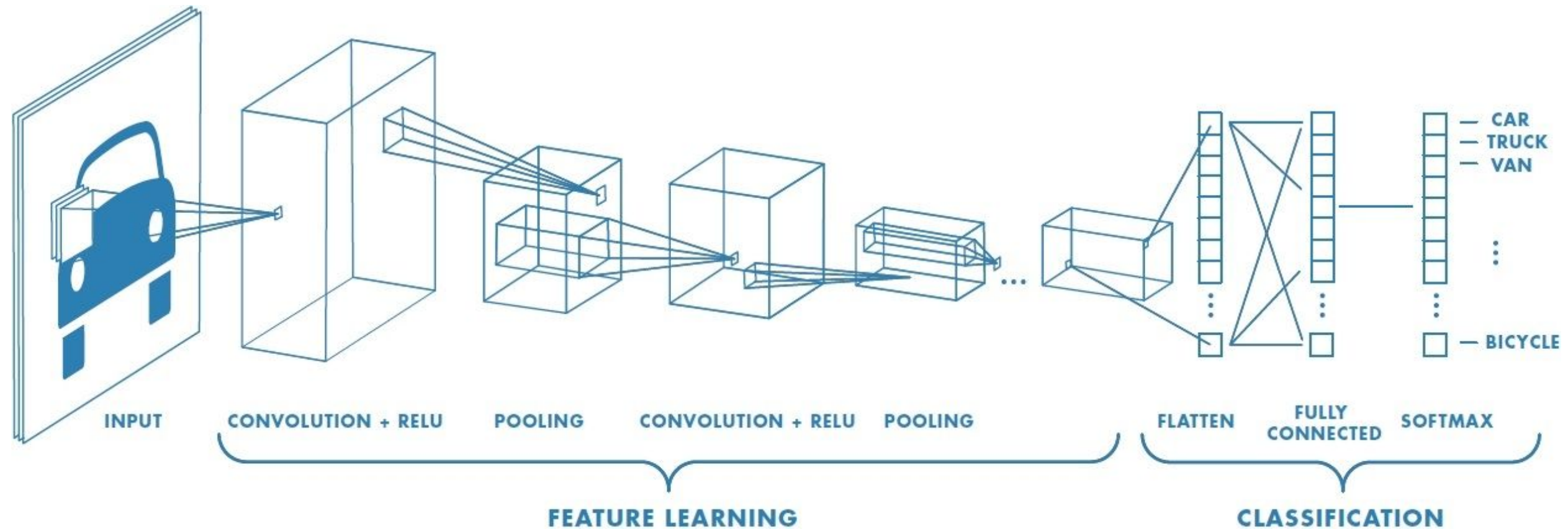
Using Colaboratory to train a GAN for semi-supervised learning

Options for putting the model in production

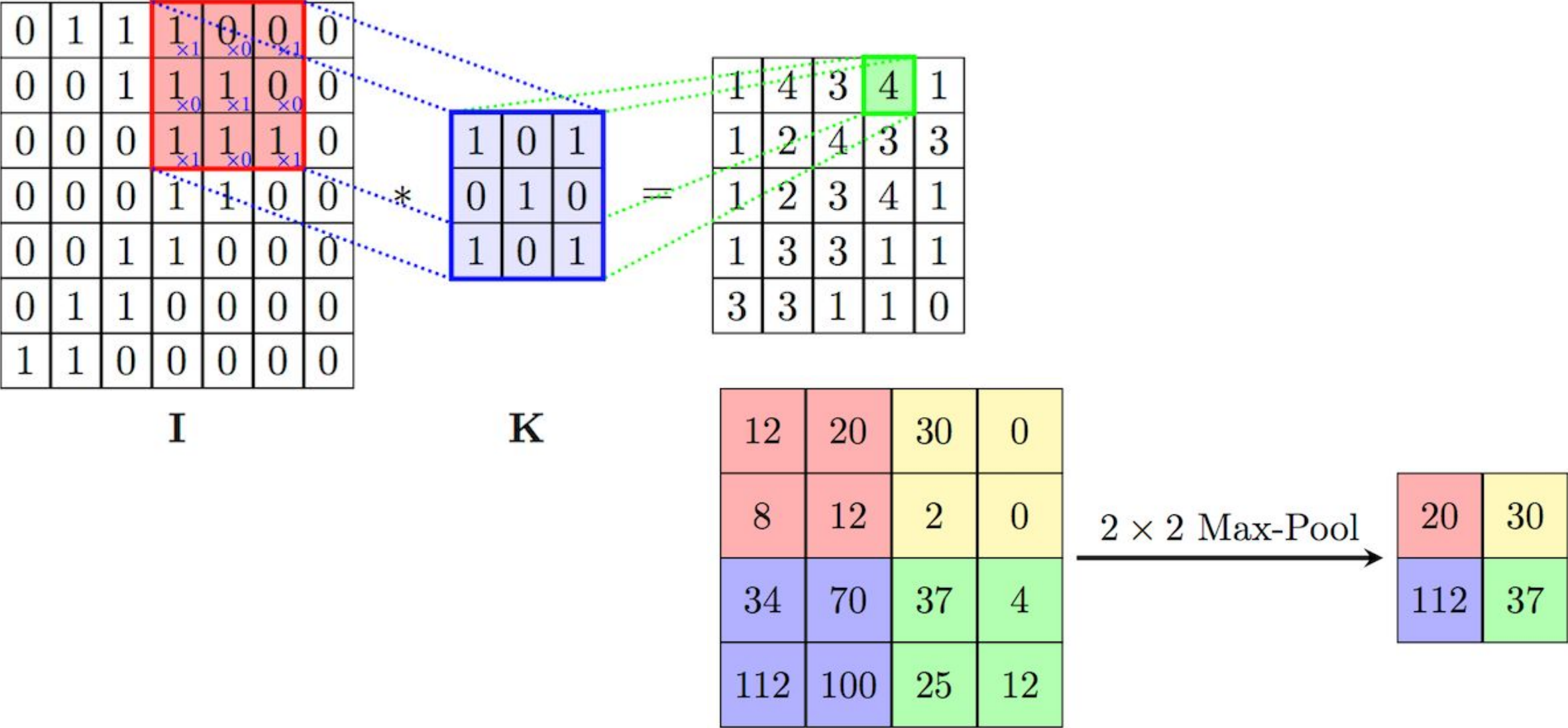
# Intro to Deep Learning



# Convolution Neural Networks in two images - part 1

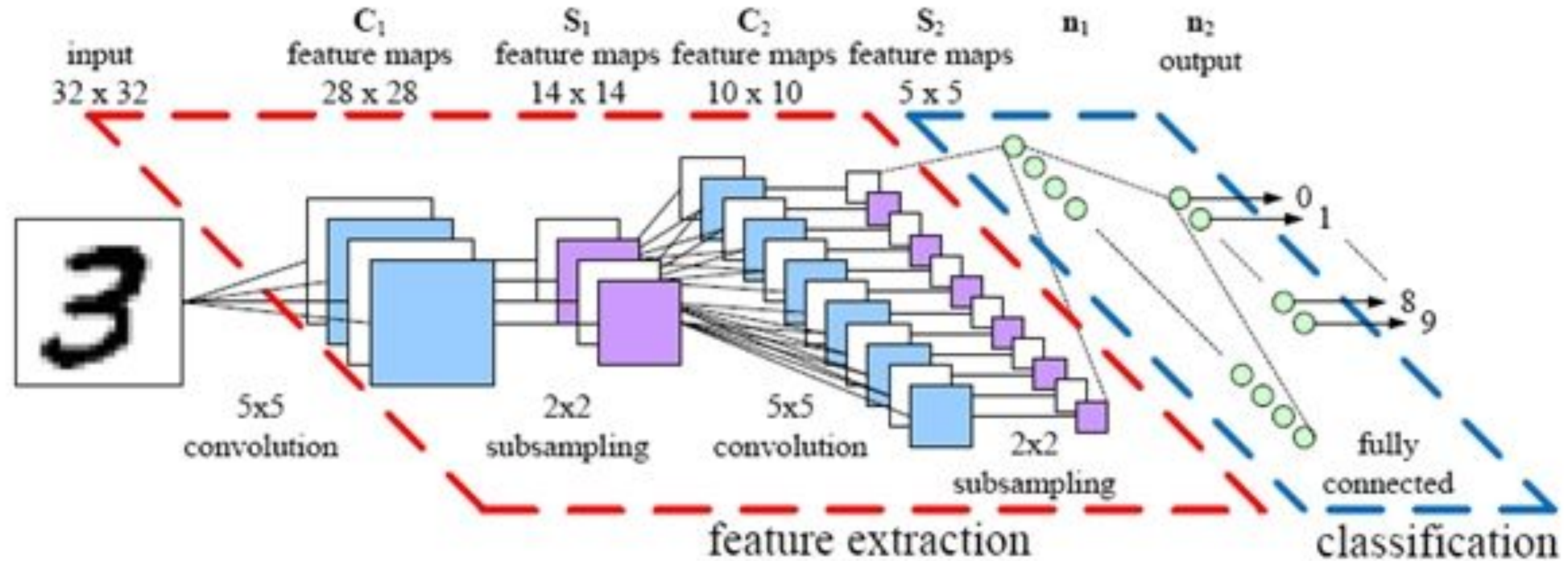


# Convolution and Pooling

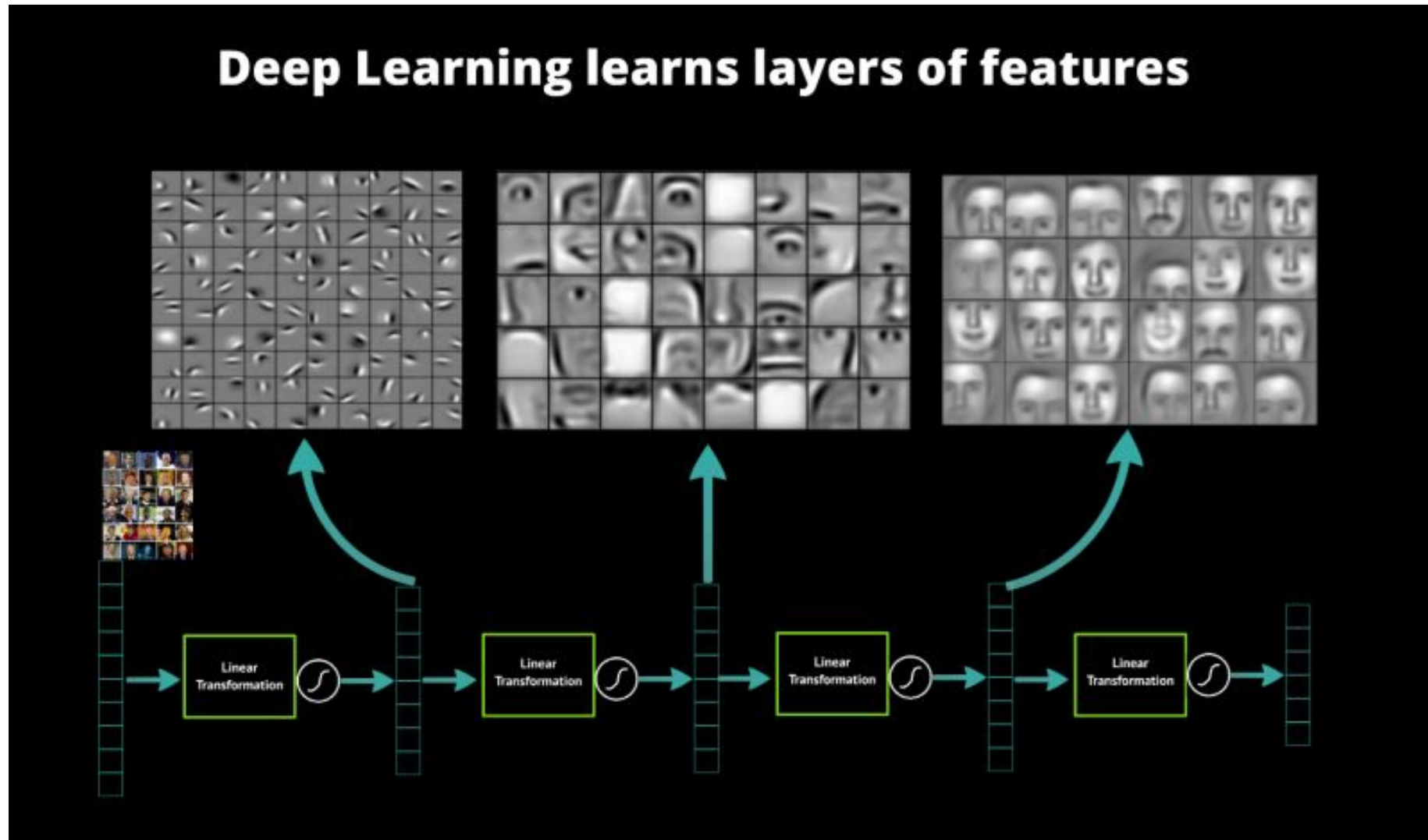




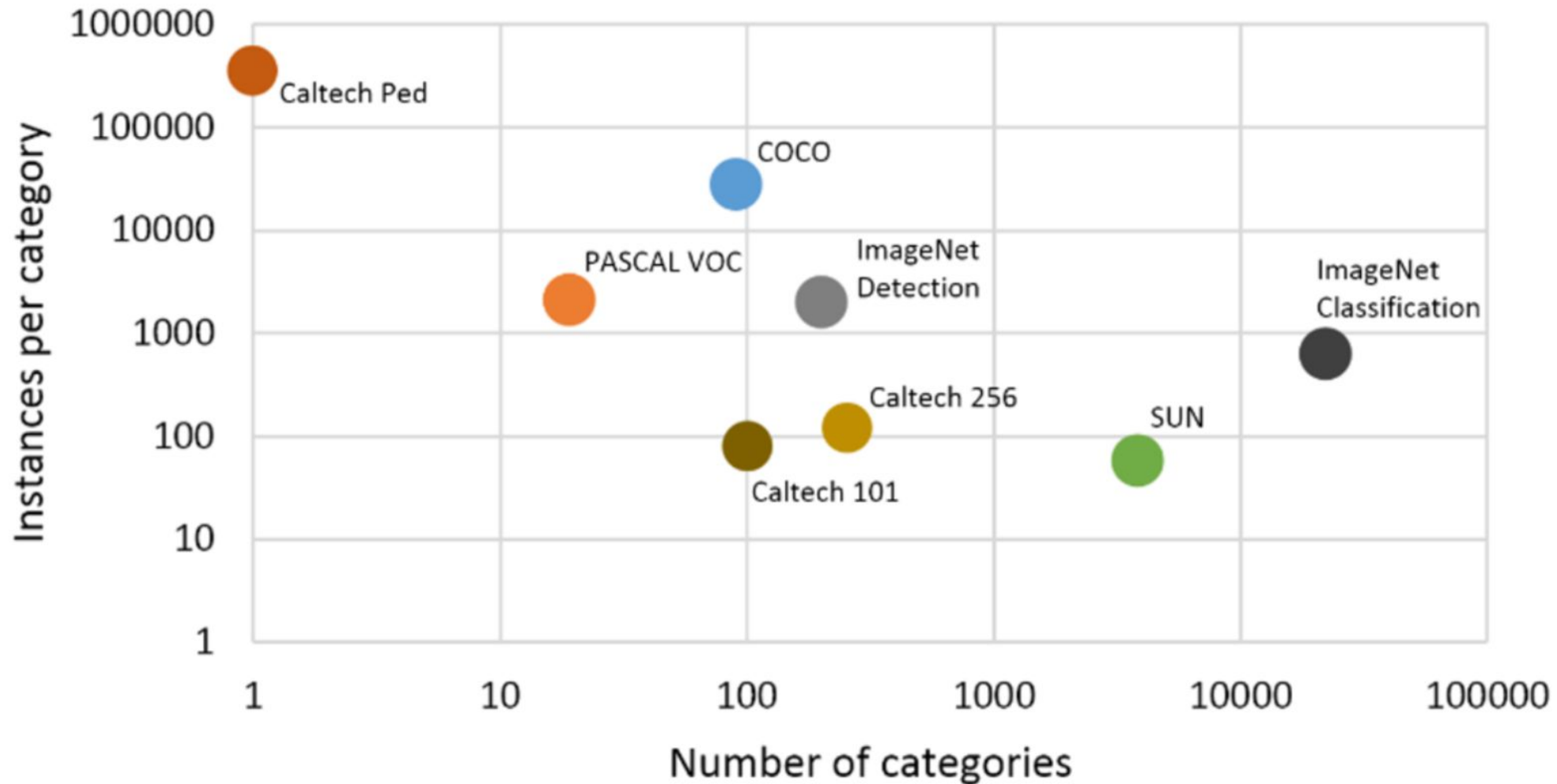
# Convolution Neural Networks in two images - part 2



Why does it work so well?

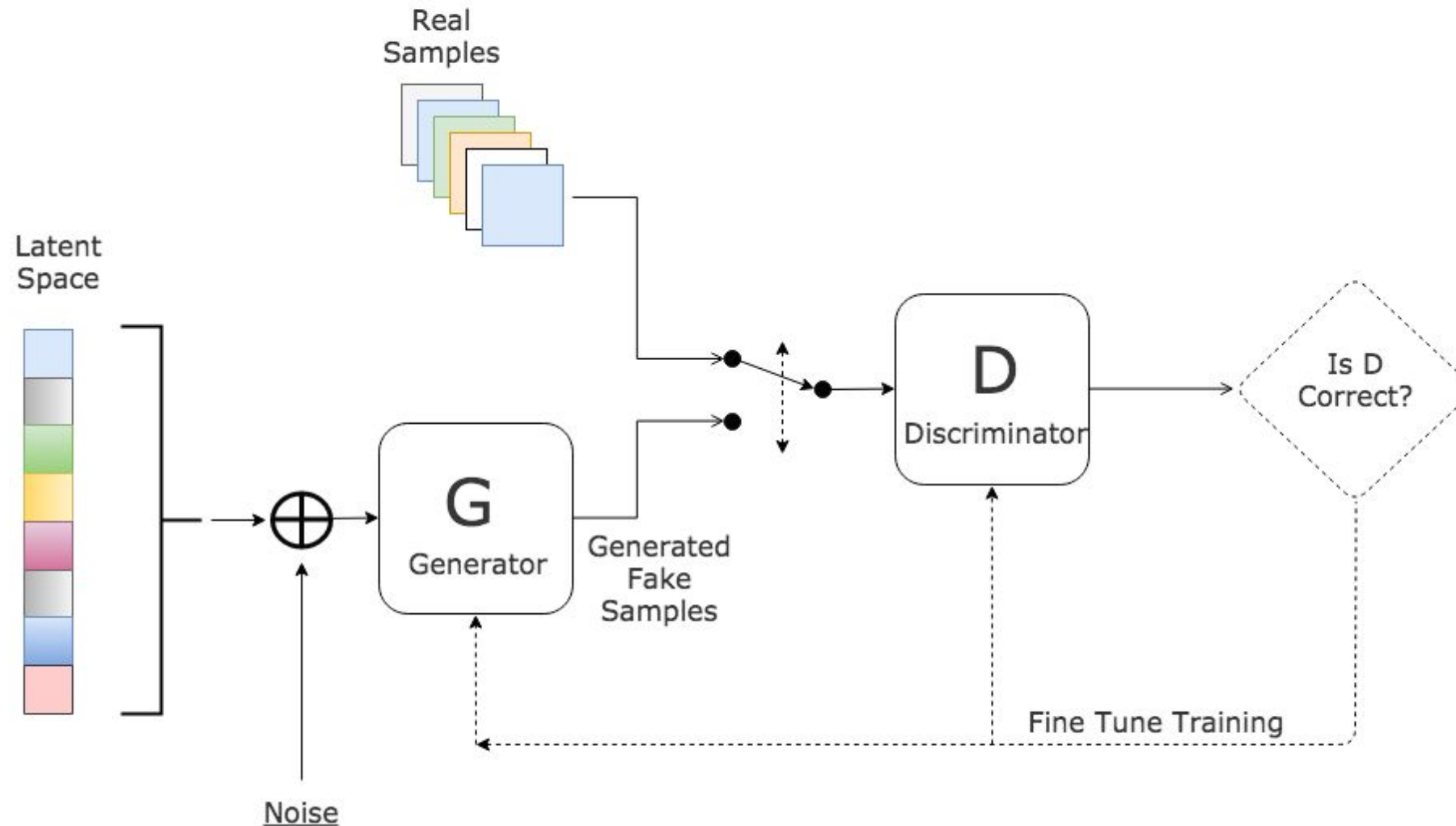


## Number of categories vs. number of instances



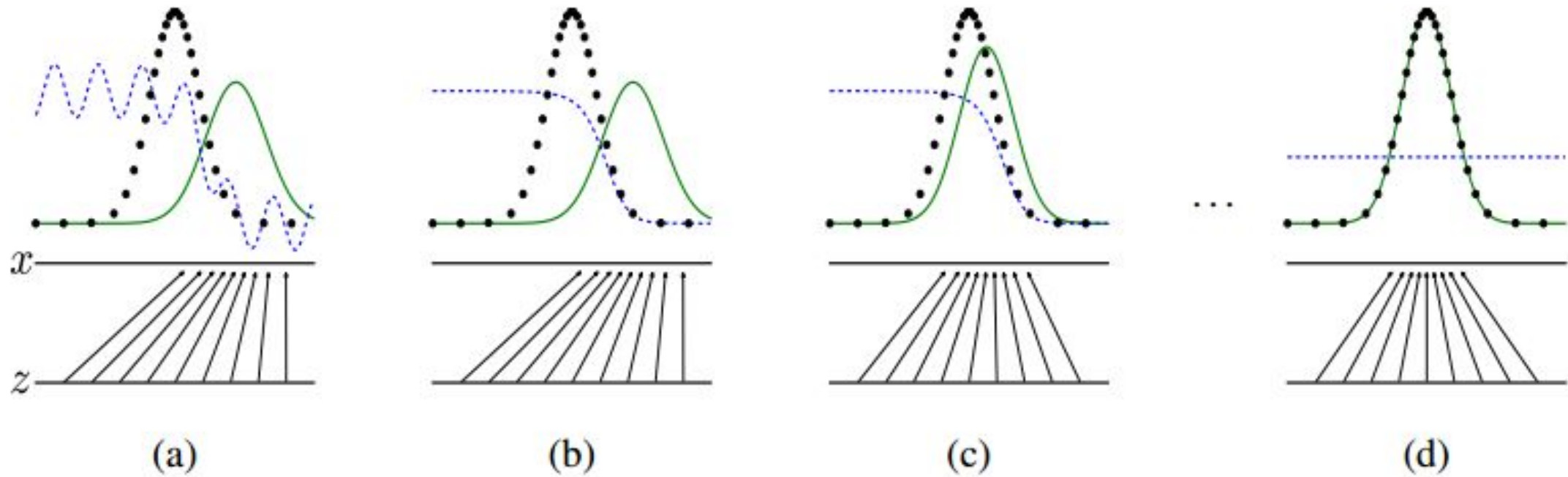
# Generative Adversarial Networks

## Generative Adversarial Network

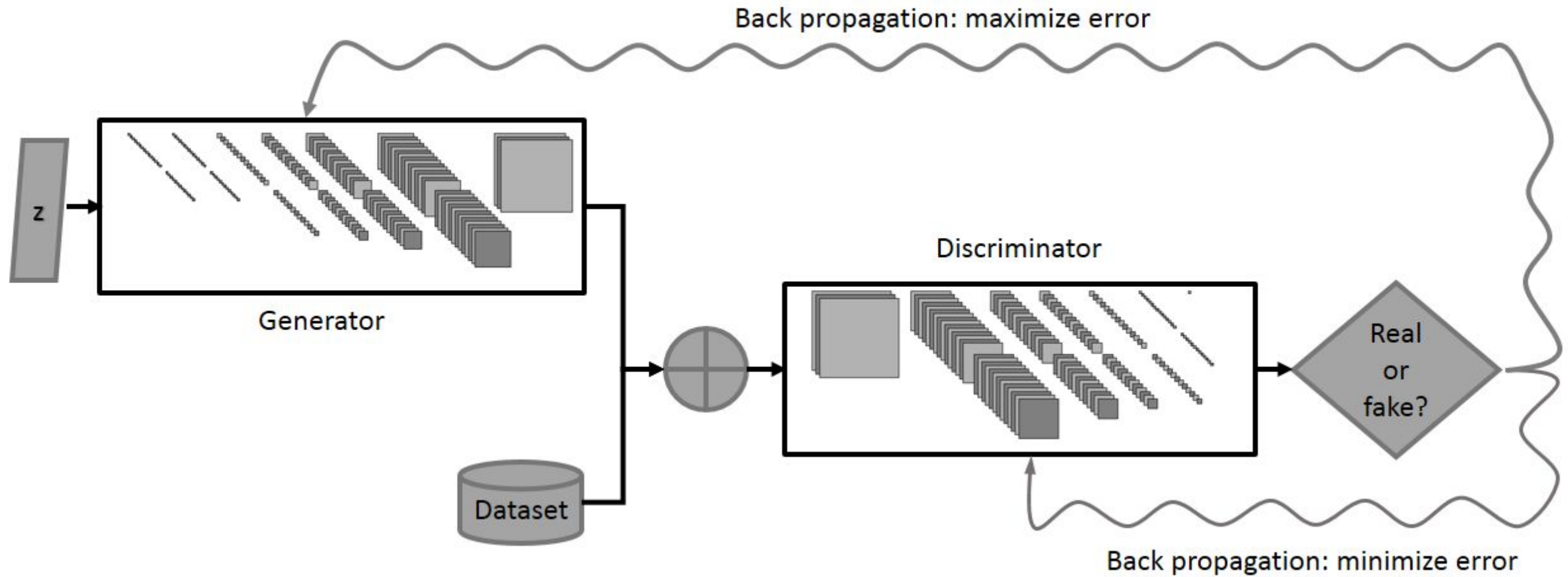




# 1D example of GAN



# GANs in one image



- Mode Collapse: Generator does not generate the full range of plausible samples
- More likely if the generator is optimized while the discriminator is kept constant for many iterations

- GANs predict the entire sample (e.g. image) at once
- Difficult to predict pixel neighborhoods

- Difficult training process / convergence: Generator and discriminator keep oscillating
- Network does not converge to a stable, (near) optimal solution

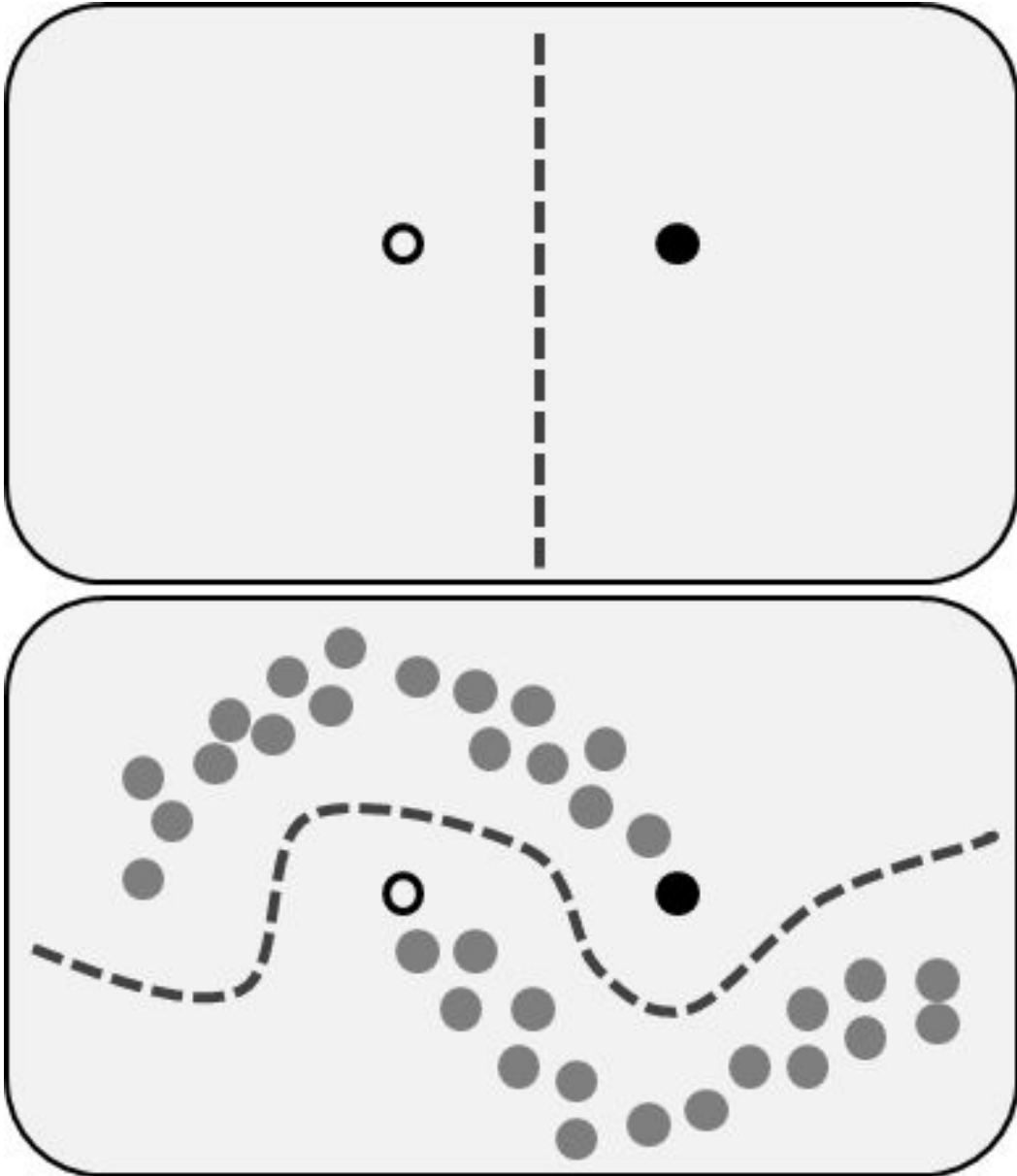
# Semi-supervised learning with GANs



# Why using GANs in a semi-supervised setting?

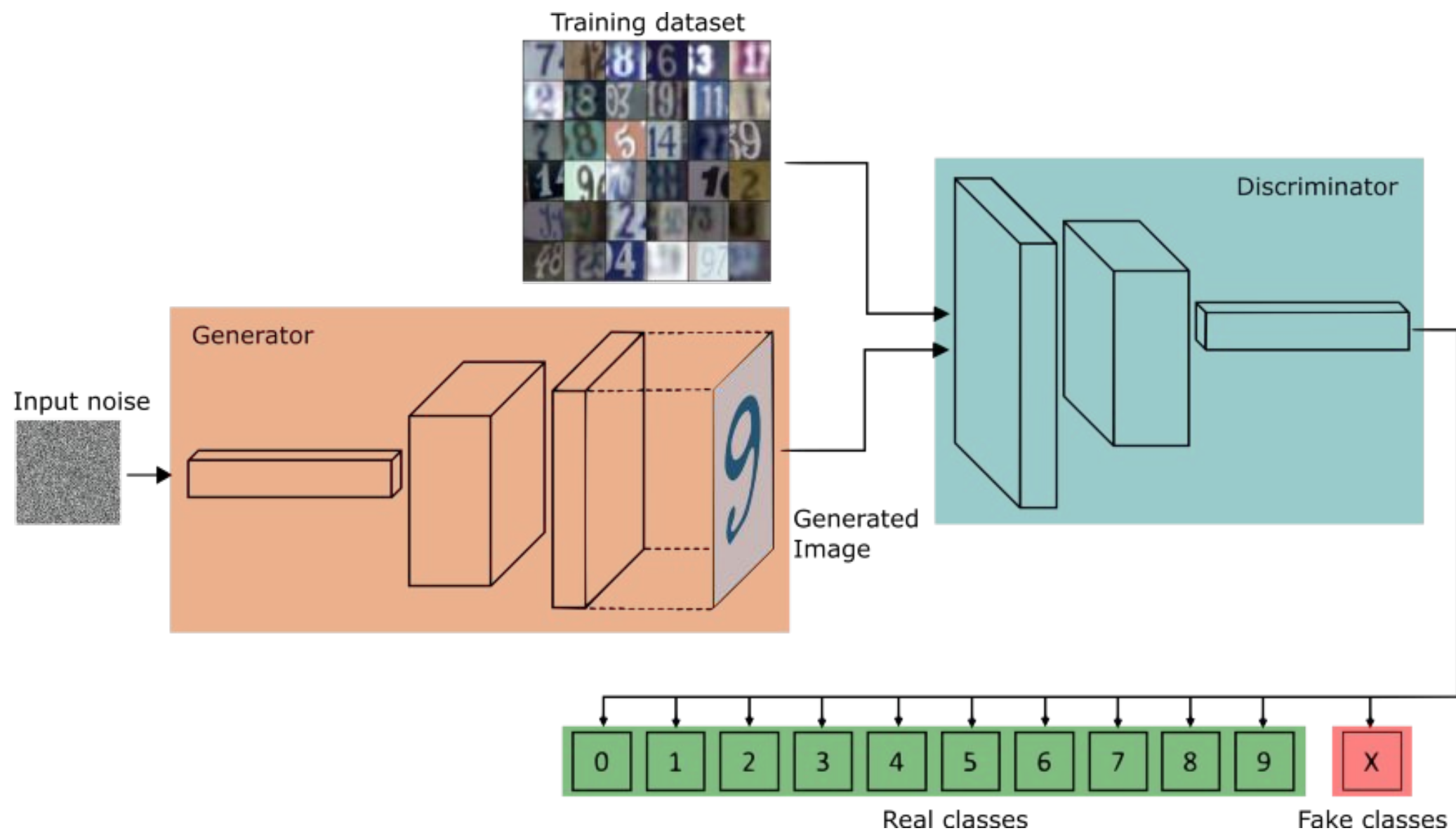
- Create a more diverse set of unlabeled data
- Get a better / more descriptive decision boundary
- Improve generalization when the amount of training samples is small

# Influence of unlabeled data in semi-supervised learning



- Continuity assumption
- Cluster assumption
- Manifold assumption

# Example of GAN in a semi-supervised learning role



- Real samples that have labels, similar to normal supervised learning
- Real samples that do not have labels. For those, the discriminator only learns that these images are real
- Samples coming from the generator. For those, the discriminator learns to classify as fake

# What do we have to modify to make it work?

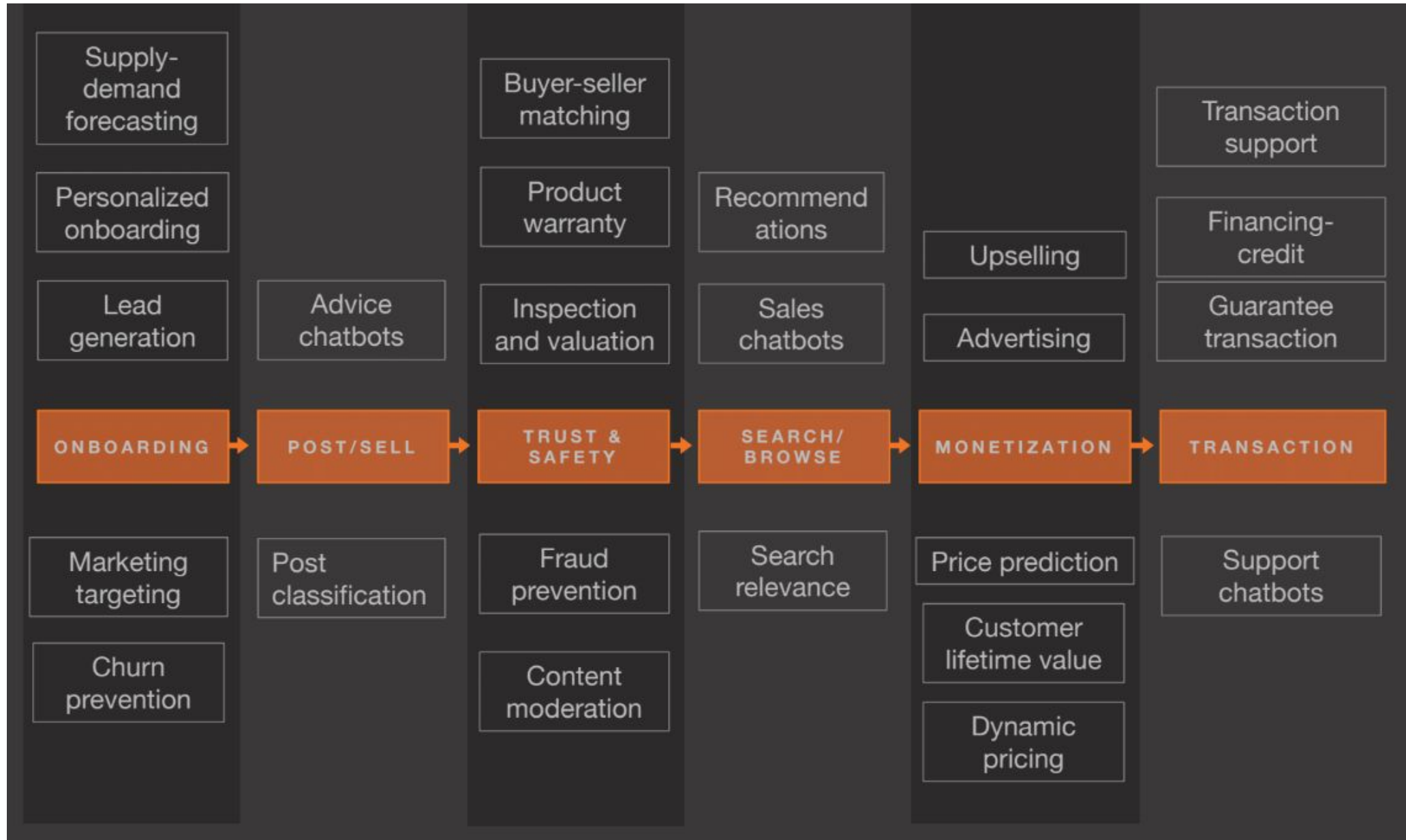
- Adjust the loss functions so that we tackle both problems:
  - as per the original GAN task, use binary cross entropy for the GAN part
  - as per multiclass classification task, use softmax cross entropy (or sparse cross entropy with logits) for the supervised task
  - apply masks to ignore the label not seen in the respective subtask
  - take the average of the GAN and supervised loss
- The discriminator can use the generator's images, as well as the labeled and unlabeled training data, to classify the dataset

- Training GANs is a challenging problem -> feature matching approach
  - a. take the moments for some features for a “real” minibatch
  - b. take the moments for same features for a “fake” minibatch
  - c. minimize mean absolute error between the two



# A use case from Trust and Safety

# Classifieds - Data Science areas of Application



# Classifieds - Data Science areas of Application

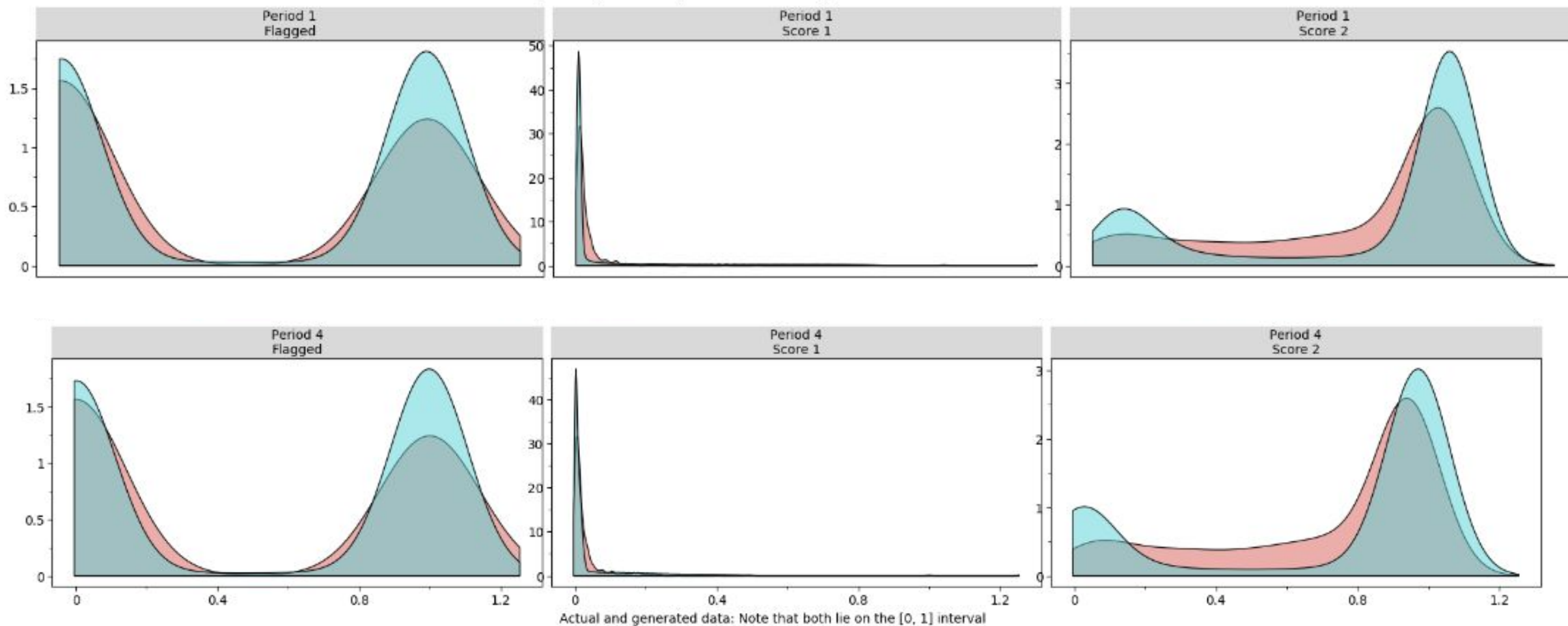
- Spam can be a problem in open platforms like classifieds
- Its annoying and can also be used as a vehicle for fraud
- How can we keep bad users out without annoying good users with stringent policies?

# Why is it a challenge?

- It's a recurring battle, with spammers becoming more inventive to circumvent the algorithms
- We may have only a few cases of cases of annotated examples for certain types of bad content

# Generating Synthetic Samples with GANs

A simple sanity check comparison of actual and generated transactions distributions



# Results

- In the presented example we only have a few tens of a certain type of spam
- We are able to quickly explore the input space and map to a latent space combining the different types of signal we have available from text, image, user behaviour, etc.
- Semi supervised learning with GANs does better (AUC = 99.9) vs Random Forest (AUC = 99.7)

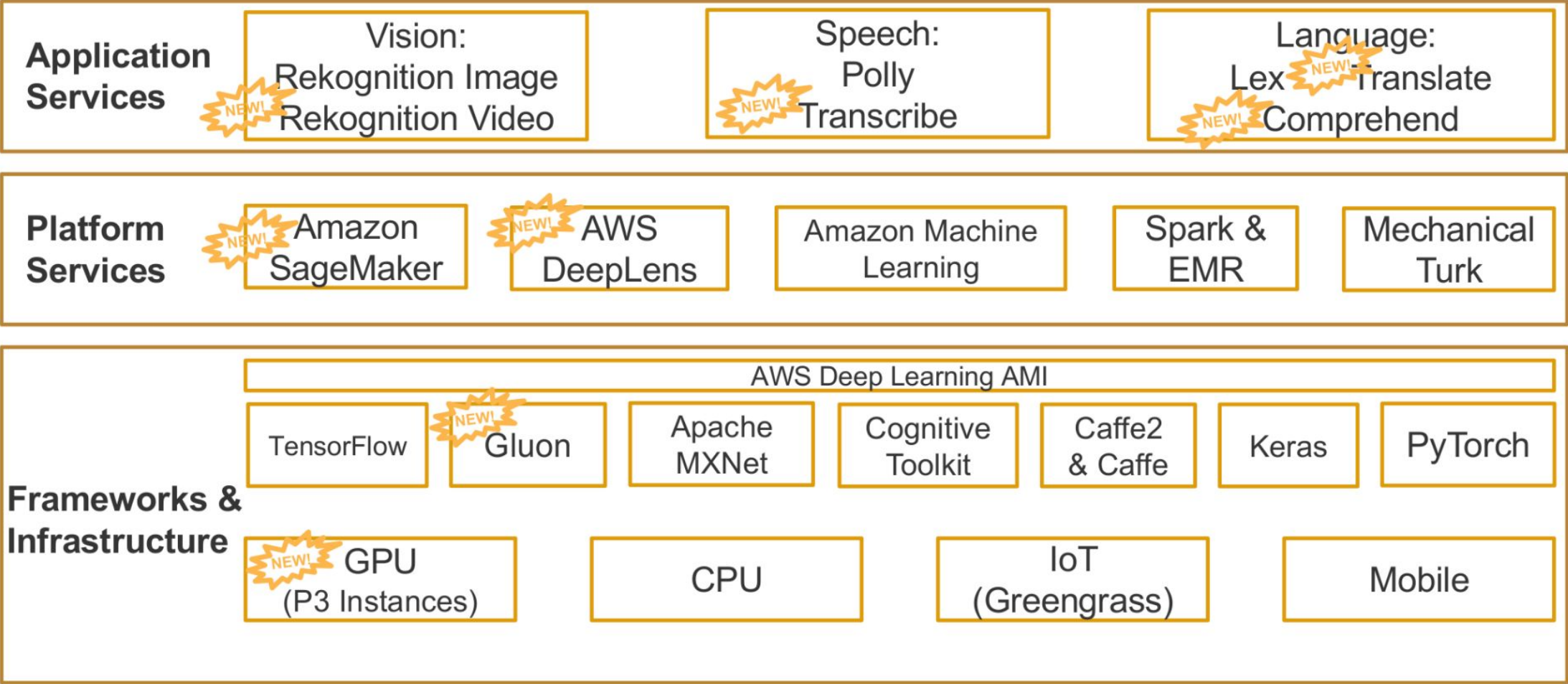




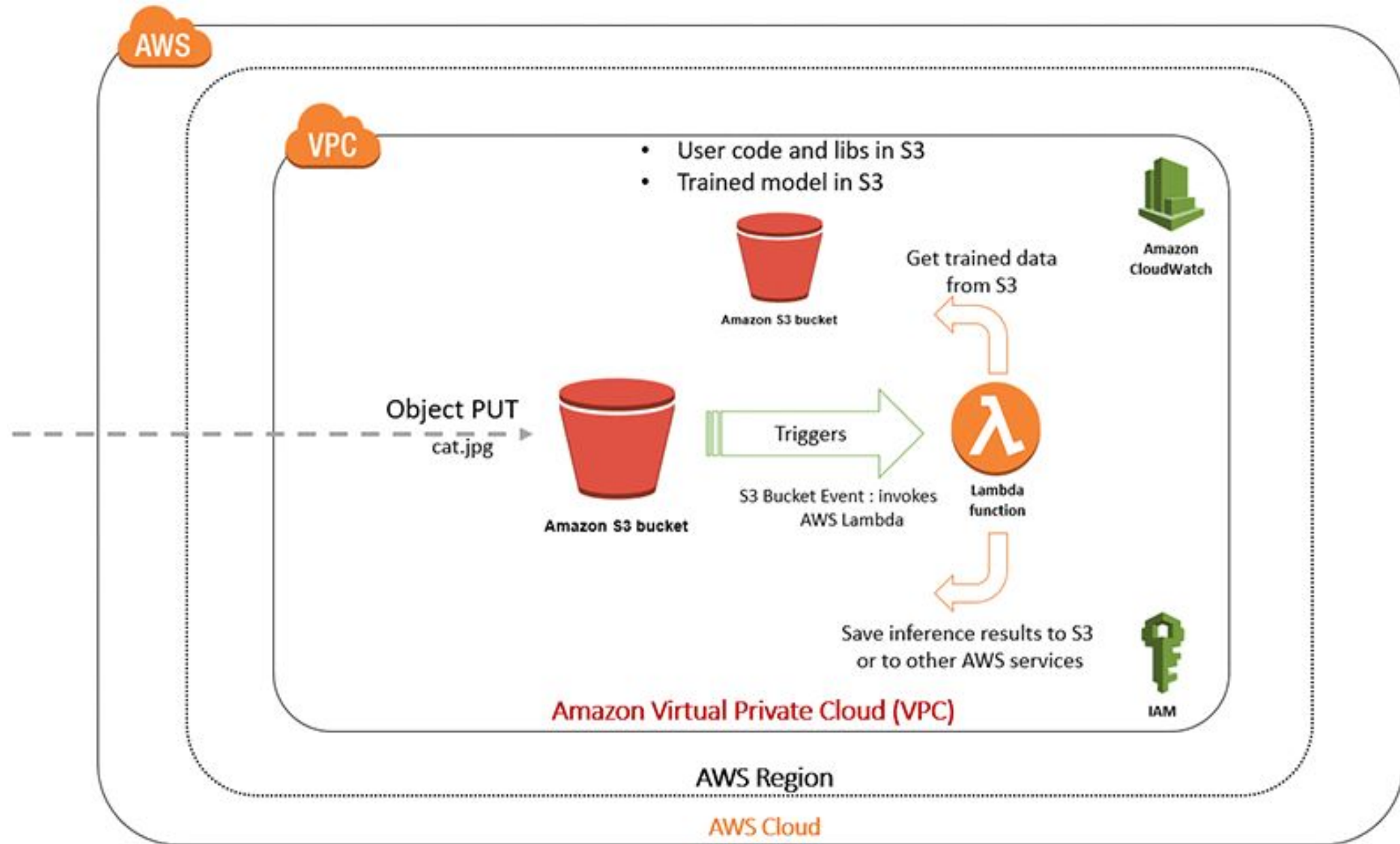
**Questions and Discussion**

# Model Deployment Options - Making endpoints

# AWS Sagemaker in the ML Stack



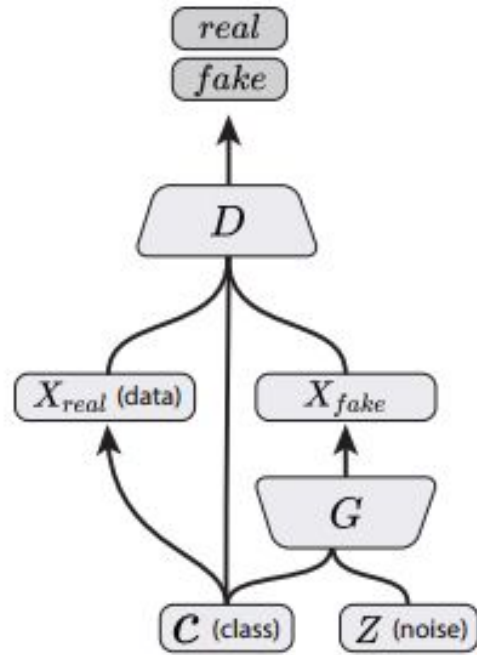
# AWS Lambda with Tensorflow



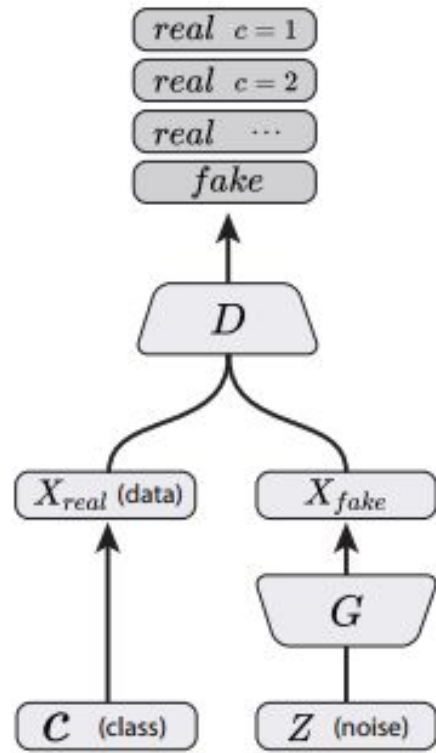
- **Using Chalice to serve SageMaker predictions**
  - <https://medium.com/@julsimon/using-chalice-to-serve-sagemaker-predictions-a2015c02b033>
- **SageMaker examples**
  - <https://github.com/aws-labs/amazon-sagemaker-examples>
- **How to Deploy Deep Learning Models with AWS Lambda and Tensorflow**
  - <https://aws.amazon.com/blogs/machine-learning/how-to-deploy-deep-learning-models-with-aws-lambda-and-tensorflow/>

- **Getting started with ML (Google)**
  - <https://drive.google.com/file/d/0B0phzgXZajPFZFRYNFRicjlZTVk/view>
- **Semi-supervised learning with Generative Adversarial Networks (GANs)**
  - <https://towardsdatascience.com/semi-supervised-learning-with-gans-9f3cb128c5e>
- **CNNs, three things you need to know**
  - <https://www.mathworks.com/solutions/deep-learning/convolutional-neural-network.html>
- **Semi-supervised learning**
  - [https://en.wikipedia.org/wiki/Semi-supervised\\_learning](https://en.wikipedia.org/wiki/Semi-supervised_learning)
- **Variants of GANs**
  - <https://www.slideshare.net/thinkingfactory/variants-of-gans-jaejun-yoo>
- **From GAN to WGAN**
  - <https://lilianweng.github.io/lil-log/2017/08/20/from-GAN-to-WGAN.html>

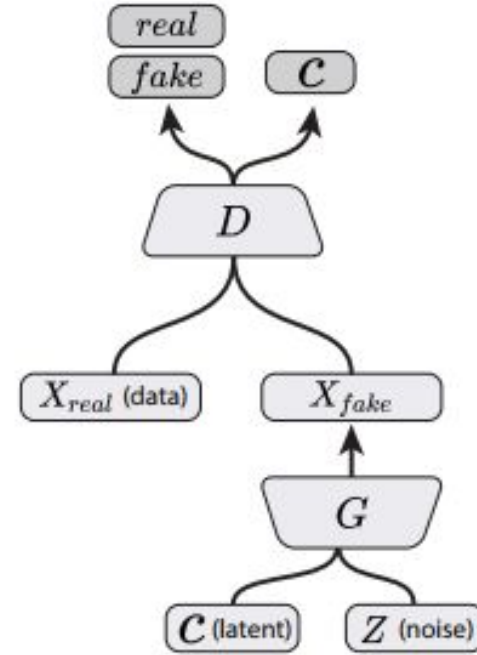
# Types of GANs



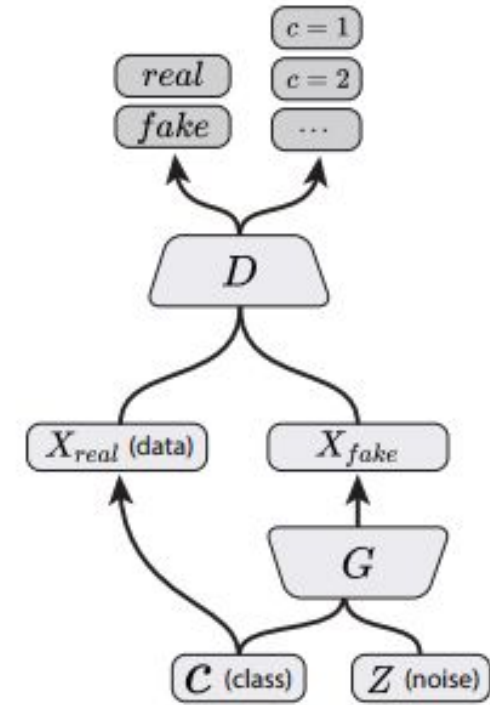
Conditional GAN  
(Mirza & Osindero, 2014)



Semi-Supervised GAN  
(Odena, 2016; Salimans, et al., 2016)



InfoGAN  
(Chen, et al., 2016)



AC-GAN  
(Present Work)

- **Amazon Sagemaker**
- **AWS Lambda and Tensorflow**

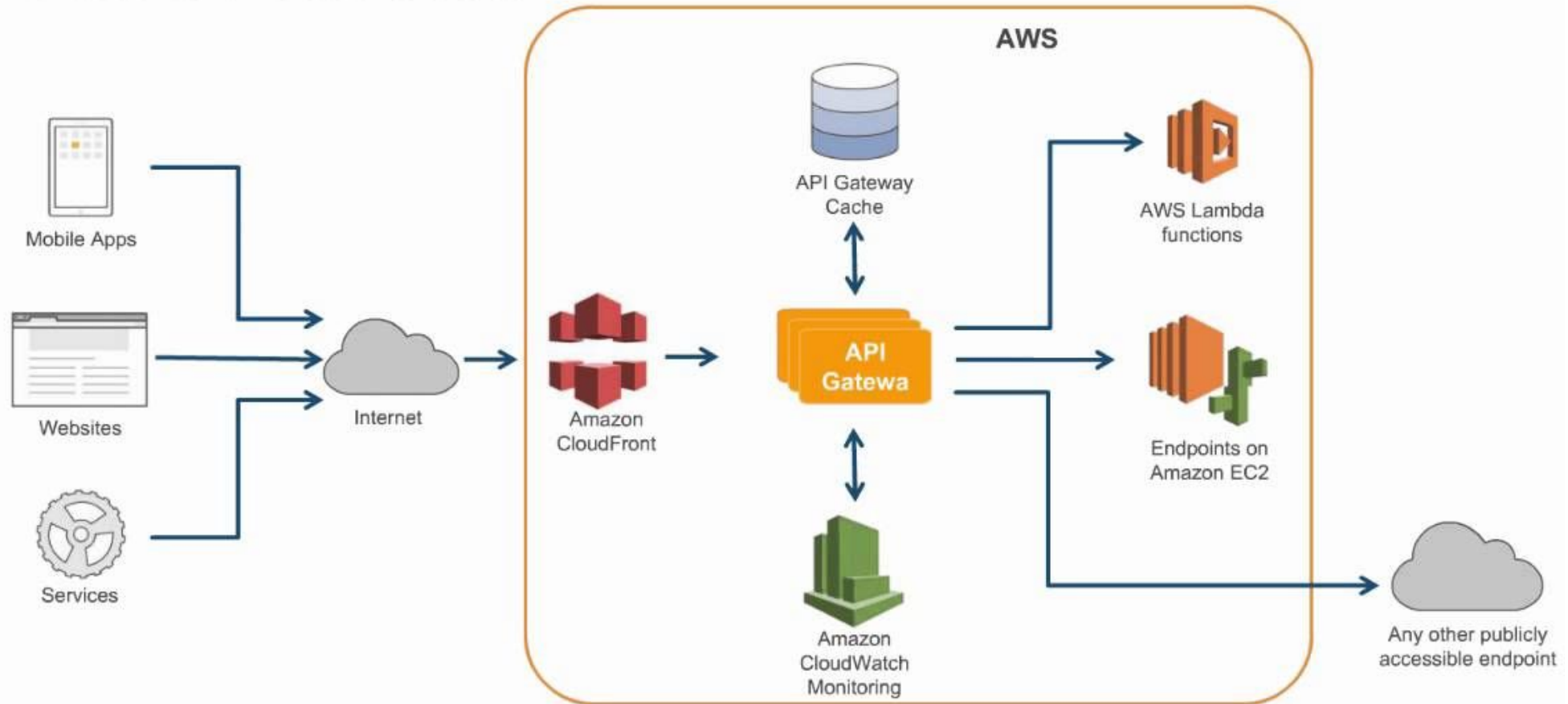


# Amazon SageMaker



A **fully managed service** that enables **data scientists** and **developers** to quickly and easily **build** machine-learning based models **into production** smart applications.

## An API Call Flow




# Google Colaboratory

# Google Colaboratory: Tensorflow + GPU + GDrive + Notebook + Free


A Medium Corporation [US] | <https://medium.com/deep-learning-turkey/google-colab-free-gpu-tutorial-e113627b9f5d>

**M** Deep Learning Turkey [Follow](#) [Twitter](#) [Facebook](#)

 **fuat** [Follow](#)  
Love to Train Deep Neural Networks  
Jan 26 · 7 min read

## Google Colab Free GPU Tutorial

Now you can develop **deep learning** applications with [Google Colaboratory](#) - on the **free Tesla K80 GPU**- using [Keras](#), [Tensorflow](#) and [PyTorch](#).



5K [47](#) [Twitter](#) [Facebook](#) [Bookmark](#)

### Notebook settings

Runtime type

Python 3

Hardware accelerator

GPU

☐ Omit code cell output when saving this notebook

CANCEL

SAVE



- **What is Google Colaboratory?**
- **Tensorflow + GPU + GDrive + Notebook + Free**
- **Free (limited) use of GPUs for training**

# Interface - Show me the code

The screenshot displays the Google Colaboratory web interface. At the top, there's a header with the Colab logo, a 'Hello, Colaboratory' message, and a menu (File, Edit, View, Insert, Runtime, Tools, Help). On the right, there are 'SHARE' and 'CONNECT' buttons. Below the header, a toolbar shows 'CODE', 'TEXT', 'CELL', and 'COPY TO DRIVE' options. A left sidebar contains a 'Table of contents' and a 'Code snippets' section. The main content area features a 'Welcome to Colaboratory!' message, followed by sections for 'Local runtime support', 'Python 3', and 'TensorFlow execution'. The 'Python 3' section includes a code snippet that prints a versioned greeting, which has been executed, showing the output 'Hello, Colaboratory from Python 3!'. The 'TensorFlow execution' section includes a mathematical matrix addition example and a corresponding code snippet.

**Welcome to Colaboratory!**

Colaboratory is a Google research project created to help disseminate machine learning education and research. It's a Jupyter notebook environment that requires no setup to use and runs entirely in the cloud.

Colaboratory notebooks are stored in [Google Drive](#) and can be shared just as you would with Google Docs or Sheets. Colaboratory is free to use.

For more information, see our [FAQ](#).

**Local runtime support**

Colab also supports connecting to a Jupyter runtime on your local machine. For more information, see our [documentation](#).

**Python 3**

Colaboratory supports both Python2 and Python3 for code execution.

- When creating a new notebook, you'll have the choice between Python 2 and Python 3.
- You can also change the language associated with a notebook; this information will be written into the `.ipynb` file itself, and thus will be preserved for future sessions.

```
[ ] 1 import sys
    2 print('Hello, Colaboratory from Python {}'.format(sys.version_info[0]))
```

Hello, Colaboratory from Python 3!

**TensorFlow execution**

Colaboratory allows you to execute TensorFlow code in your browser with a single click. The example below adds two matrices.

$$\begin{bmatrix} 1. & 1. & 1. \\ 1. & 1. & 1. \end{bmatrix} + \begin{bmatrix} 1. & 2. & 3. \\ 4. & 5. & 6. \end{bmatrix} = \begin{bmatrix} 2. & 3. & 4. \\ 5. & 6. & 7. \end{bmatrix}$$

```
[ ] 1 import tensorflow as tf
    2 import numpy as np
    3
    4 with tf.Session():
    5     input1 = tf.constant(1.0, shape=[2, 3])
    6     input2 = tf.constant(np.reshape(np.arange(1.0, 7.0, dtype=np.float32), (2, 3)))
    7     output = tf.add(input1, input2)
    8     result = output.eval()
    9
   10 result
```

<https://colab.research.google.com/drive/1CqIRvhhR3bT-NUgTmgH5S9OSHn2ABVYt#scrollTo=RA4eBhnzBZO3>