# APRENDIZAJE DE MÁQUINAS
## Data Labellin' or Data Annotation

## John Ballesteros

**Profesor Asociado**
**Departamento de Ingeniería Civil**
**jballes@unal.edu.co**

https://github.com/srobles05/3008422-AprendizajeDeMaquinas

Gidia
Grupo de I+D
en Inteligencia Artificial

UNIVERSIDAD
NACIONAL
DE COLOMBIA

# A data – model tandem

a11X1 + a12X2 +…..+a1nXn = w1

a21X1 + a22X2 +…..+a2nXn = w2

.

.

.

an1X1 + an2X2 +…..+annXn = wn

$X + Y = 80$

$X/2 - 3Y = 28$

In a DNN
AX = b, donde: A,X,b son matrices o tensores, A son los parametros o los pesos, b son los bias (interceptos)
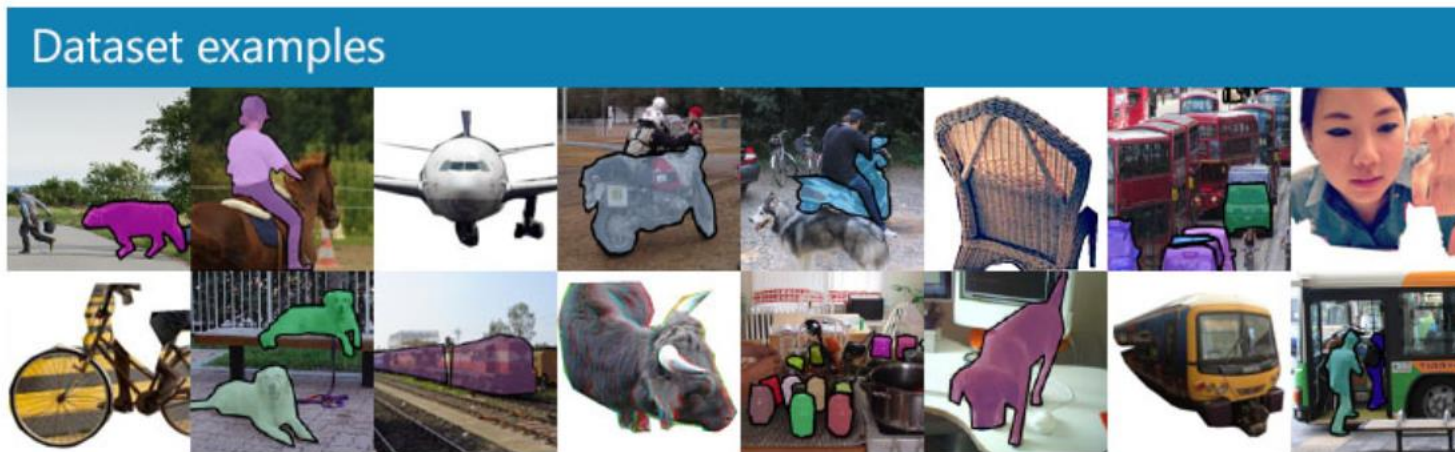
# Training Data

Gidia
Grupo de I+D
en Inteligencia Artificial

UNIVERSIDAD
NACIONAL
DE COLOMBIA

# Famous Datasets



**Mnist: handwritten numbers** (http://yann.lecun.com/exdb/mnist/) Yann LeCun, et all

**Coco: common objects in context for detection & segmentation** (https://cocodataset.org/#home)

**Planet: Multiclassification satellite from amazon rainforest** (https://www.kaggle.com/nikitarom/planets-dataset)

**Camvid: Segmentation Street datasets** (http://mi.eng.cam.ac.uk/research/projects/VideoRec/CamVid/)



Dataset examples

Gidia
Grupo de I+D
en Inteligencia Artificial

UNIVERSIDAD
NACIONAL
DE COLOMBIA

- From the field (primary)


- From the Internet (secondary)

- In the file name: it's the most common:  ex. dog.jpg, horse.jpg, cat.jpg

- In a separate file: tipically in a .csv file

-In a field of a table: target

-In the whole image itself:  segmentation
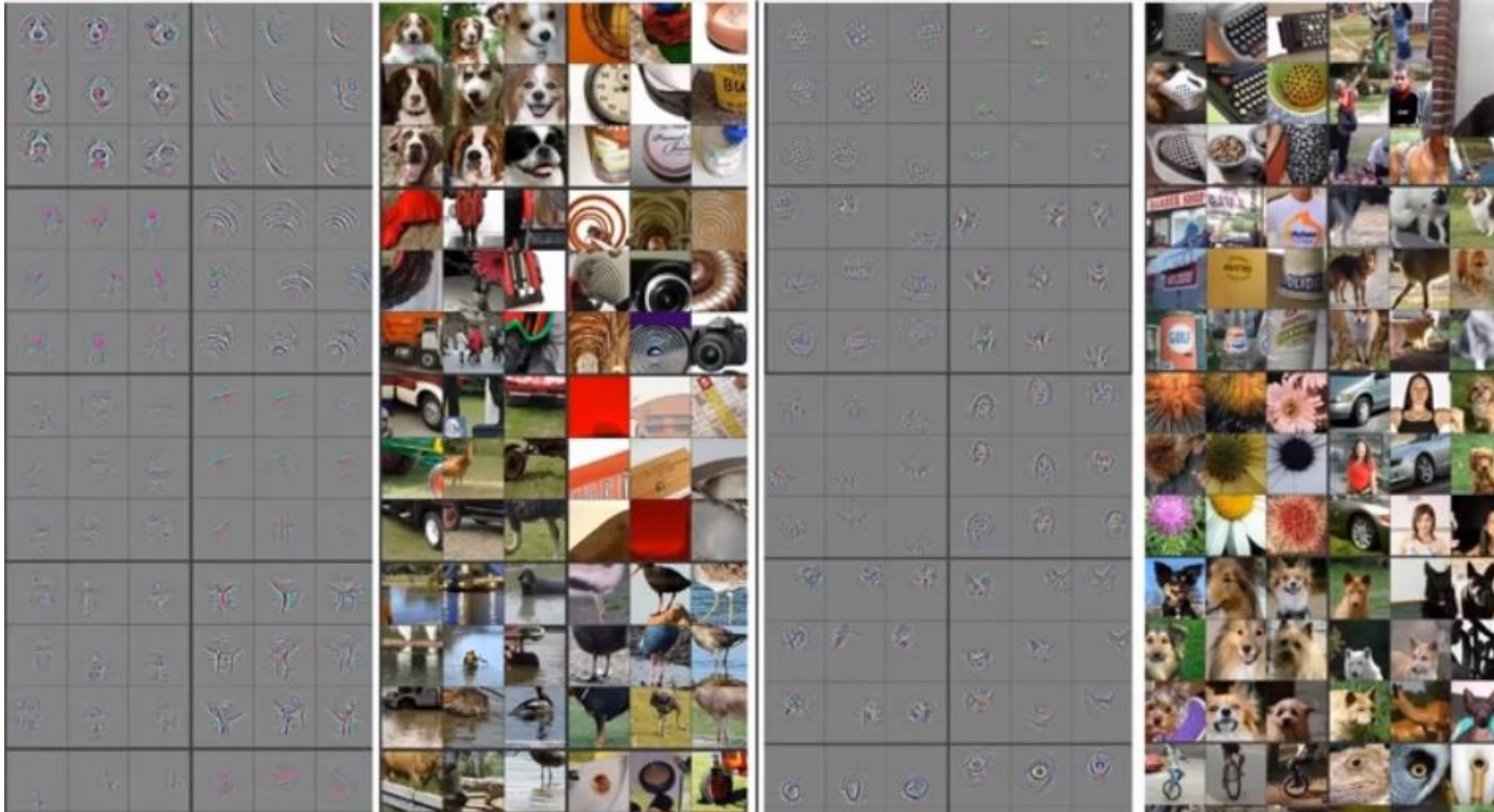
# Search and scroll [Google Images](#)

-Be specific, less manual pruning you will have to do.

-The maximum number of images Google Images shows is 700.

-Put things you want to exclude into the search query using (-)

-Limit your results to show only photos by clicking on Tools and selecting Photos from the Type dropdown.

ctrl+shift+j

# Transfer Learning

```
learn = create_cnn(data, models.resnet34, metrics=error_rate)
```

## What is enough data?

If hyper parameters: architecture, epochs and learning rate are ok but model accuracy is still not good, below 80%

You need more data !!!!

How to get more data?

Gidia
Grupo de I+D
en Inteligencia Artificial

UNIVERSIDAD
NACIONAL
DE COLOMBIA

# Data Augmentation and sintetic data

-Programming is the solution, there is little workaround

-Python is more than ok: PIL, CV2, and imageo modules do the trick!

-Typical transformations are: rotate=flip, scale (affine), warping

This can multiply actual data by 4 !

- Crappify and blurring is easy.

# Data Augmentation = Image Transformations (Geometric, Color, etc)

Original          Rotated          Flip          Warped



```
i = Image.open(os.path.join(mypath, f))
print(f, i.size, i.mode)
i10 = i.rotate(10)
i10.save(os.path.join(mypath,'rot10{}.jpg'.format(fname)))
iflip = i.transpose(Image.FLIP_LEFT_RIGHT)
iflip.save(os.path.join(mypath,'flip{}.jpg'.format(fname)))
ifvert = i.transpose(Image.FLIP_TOP_BOTTOM)
ifvert.save(os.path.join(mypath,'fvert{}.jpg'.format(fname)))
iwarp = i.transform((new_width, height), (Image.AFFINE),(1, m, -xshift if m > 0 else 0, 0, 1, 0),
        (Image.BICUBIC))
```

Gidia
Grupo de I+D
en Inteligencia Artificial

UNIVERSIDAD
NACIONAL
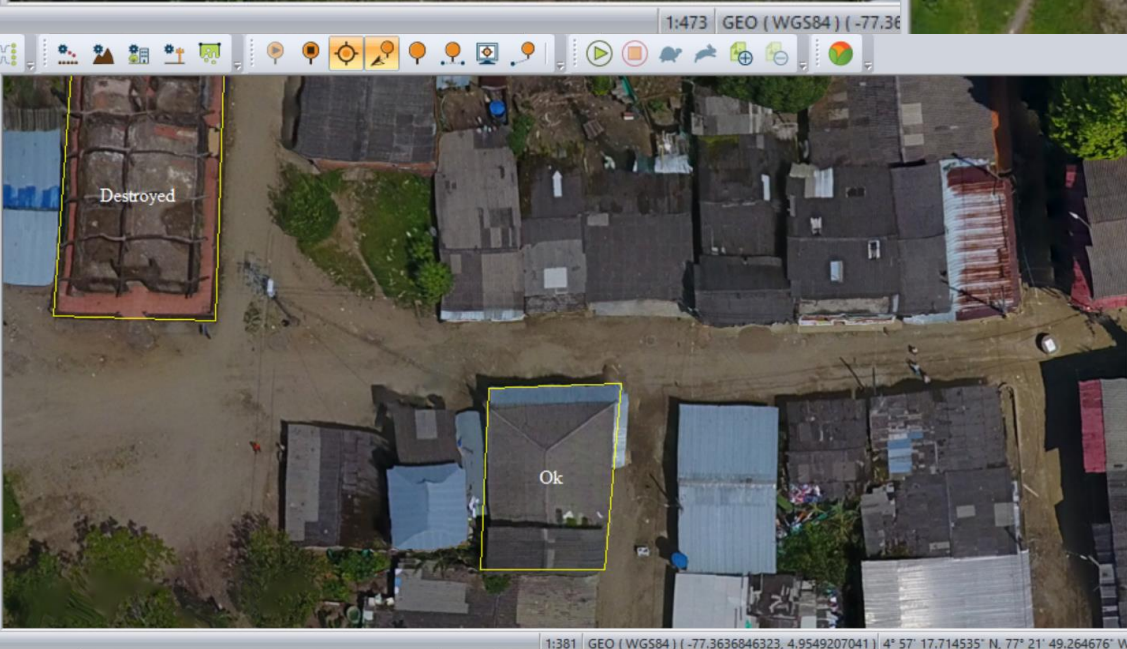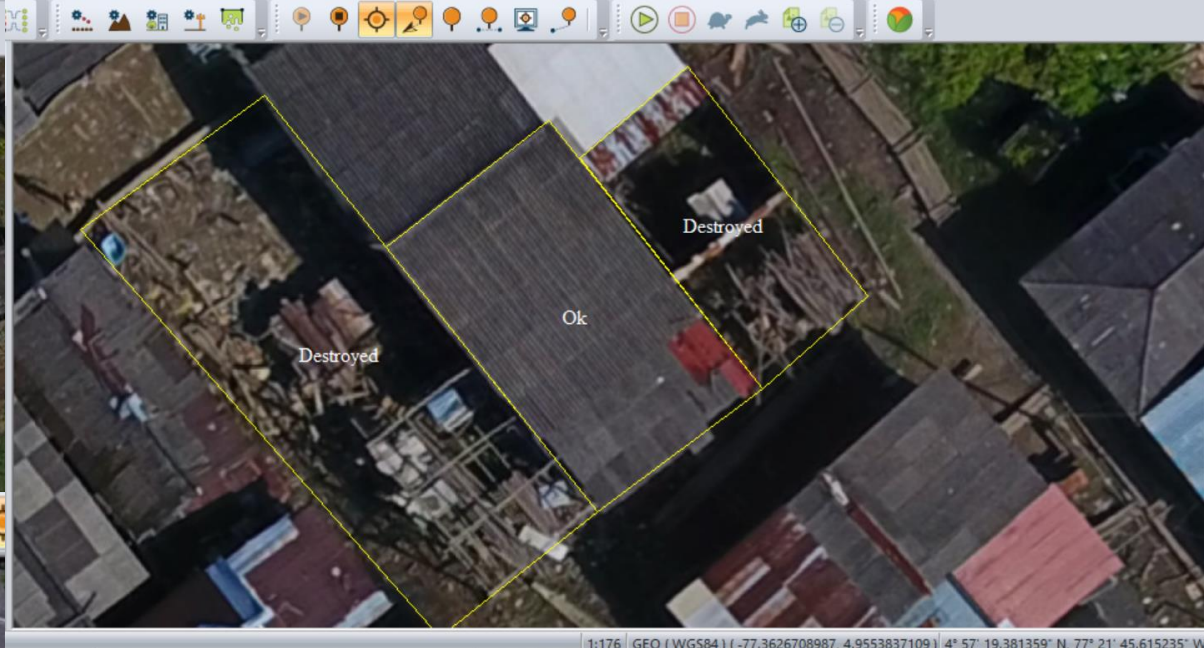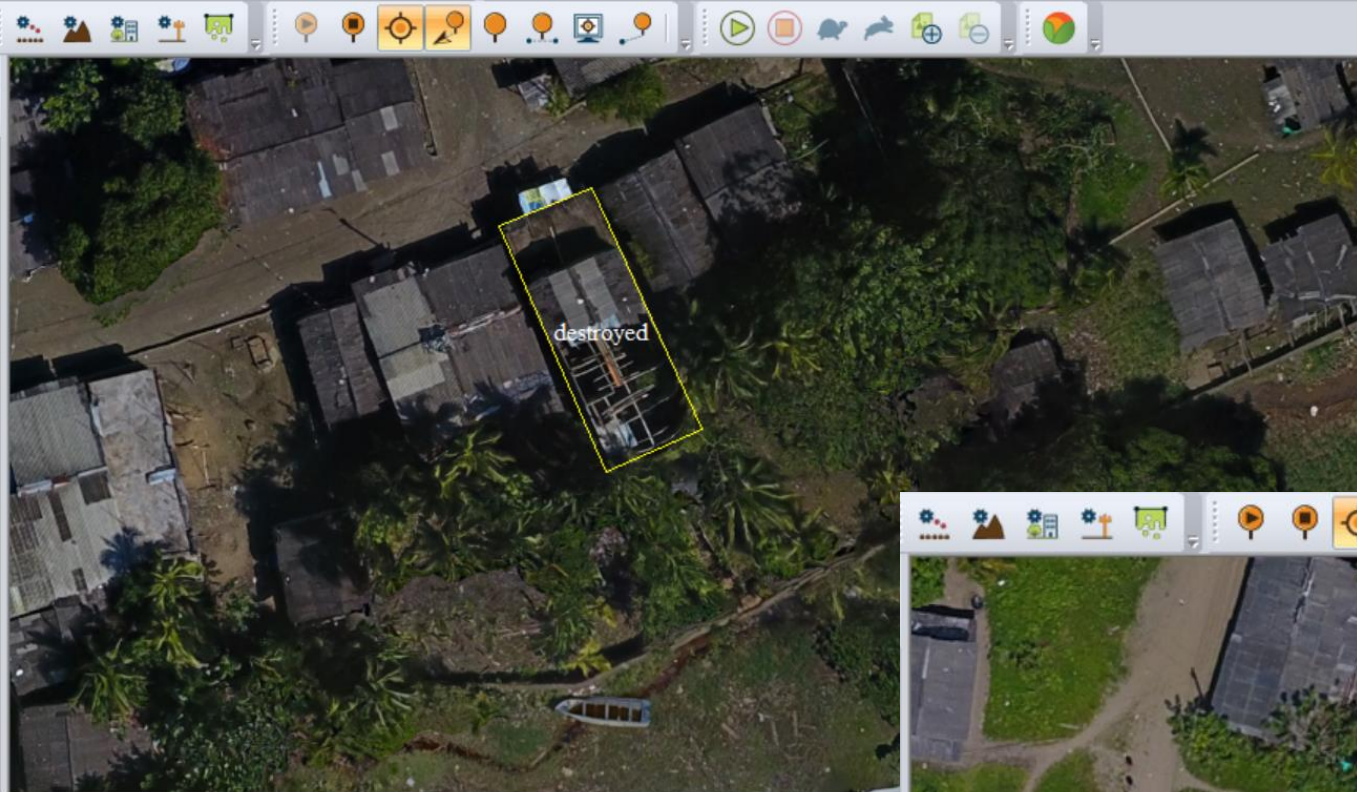DE COLOMBIA

# Examples of problems that can be solved

**Binary classification:** Damaged/No Damaged--------Prob 1/0 (CNN)

**Multi classification:** Rock classification, Minerals, vipers, trees, outcrops, house assessment--------Prob Treshold above 0.6 (CNN)

**Segmentation:** each pixel is assigned a value into a category (CNN encoding-decoding)

**Image Generation:** Cuttoff, marks & stinks, colorization, improve resolution (more advanced), Neural Style Transfer: GANs

Gidia
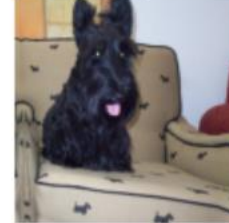Grupo de I+D
en Inteligencia Artificial

UNIVERSIDAD
NACIONAL
DE COLOMBIA

con acabados, garaje, ante jardín, varios pisos, grande, alto valor



con acabados, no garaje, no ante jardín, varios pisos, pequeña, medio valor



sin acabados, no garaje, no ante jardín, varios pisos, mediana, bajo valor



con acabados, no garaje, no ante jardín, 1 piso, pequeña, bajo valor



alta pendiente, no cobertura, alto peso, muy alto riesgo



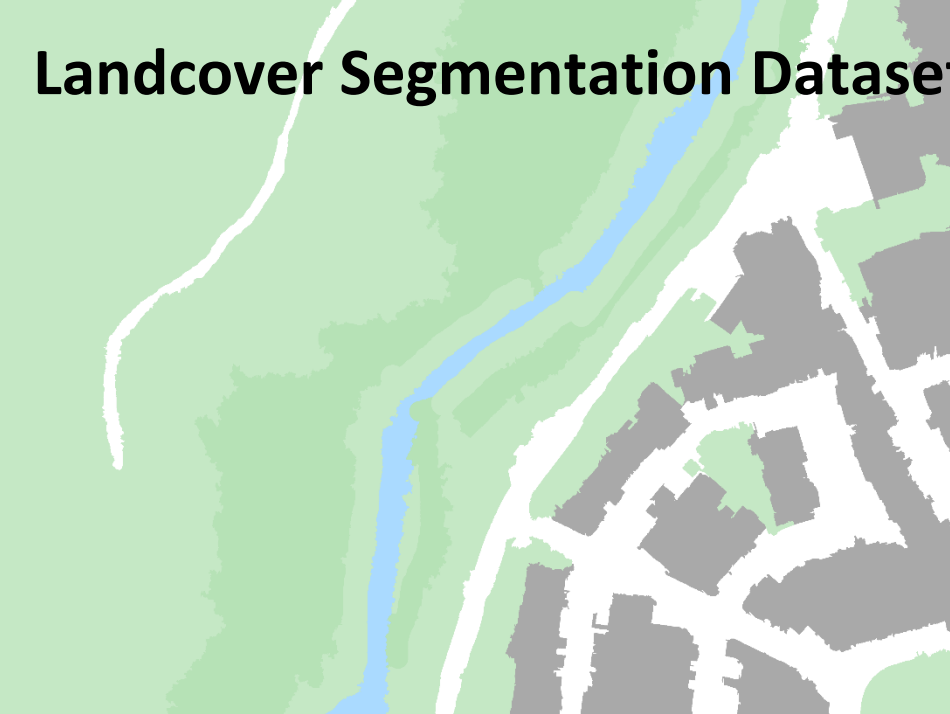alta pendiente, no cobertura, alto peso, muy alto riesgo



alta pendiente, no cobertura, bajo peso, muy bajo riesgo

Gidia
Grupo de I+D
en Inteligencia Artificial

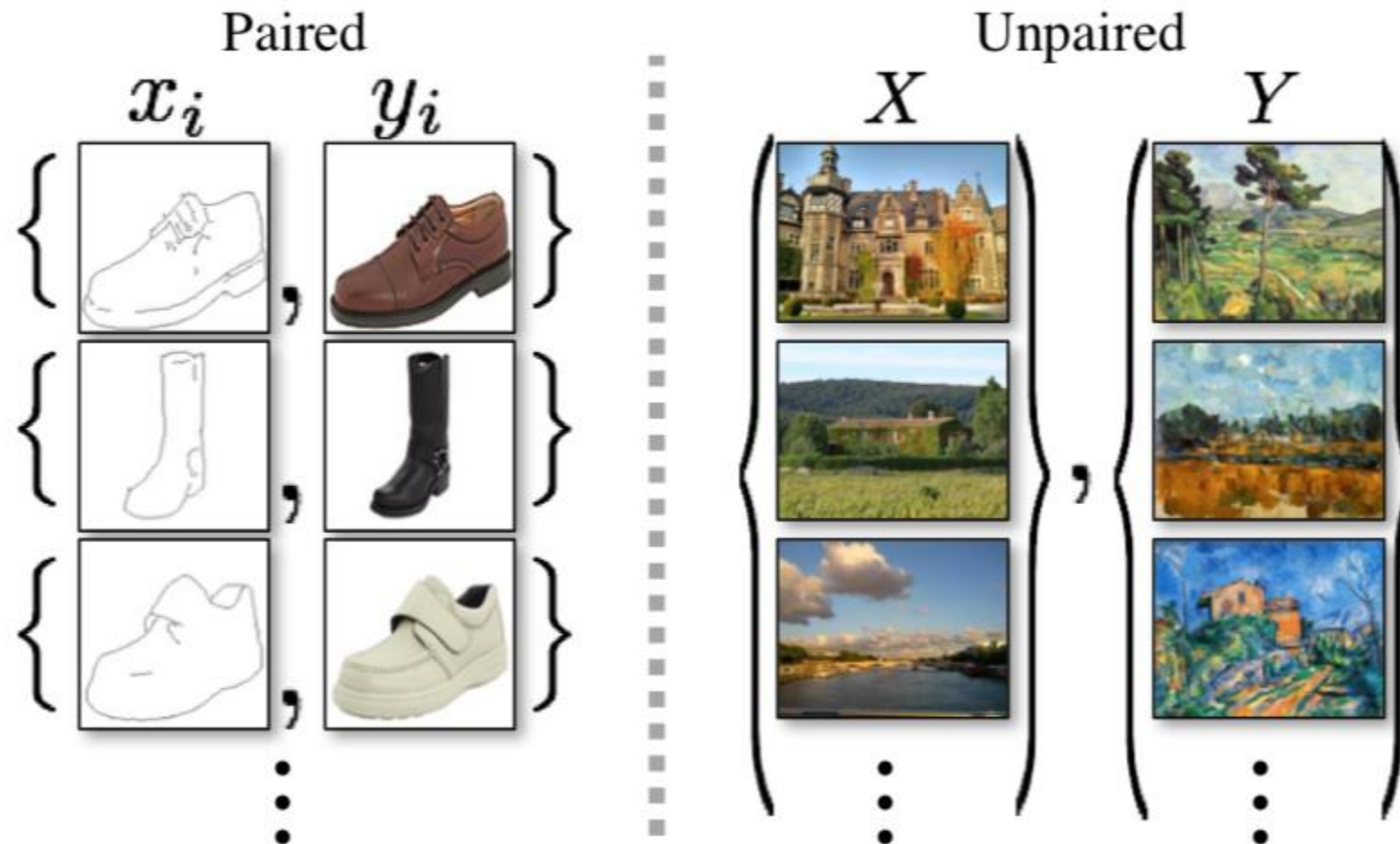UNIVERSIDAD NACIONAL DE COLOMBIA

# Camvid: Streets Segmentation Dataset

Landcover Segmentation Dataset

# Types of training datasets for image generation

Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks
Jun-Yan Zhu*, Taesung Park*, Phillip Isola, Alexei A. Efros
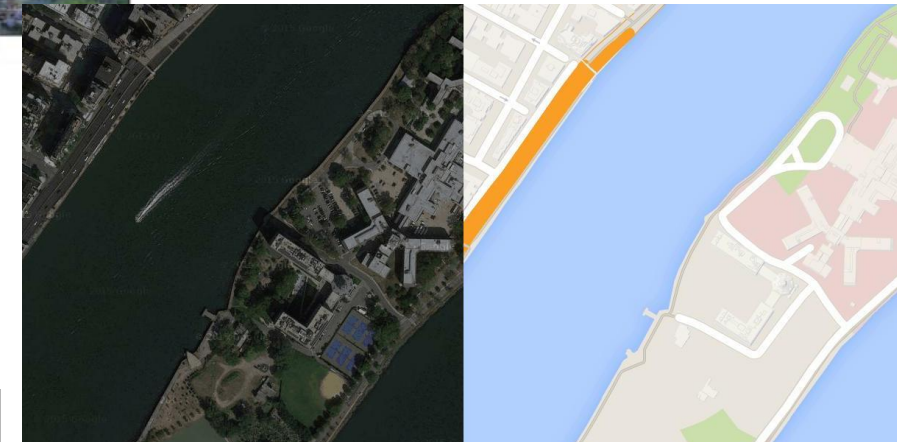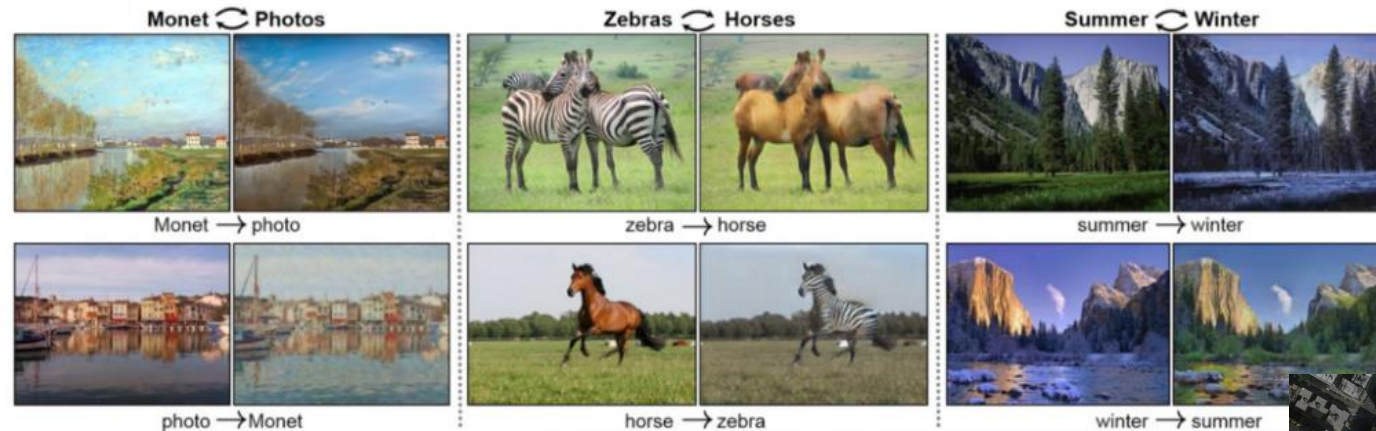Berkeley AI Research Lab, UC Berkeley
In ICCV 2017. (* equal contributions)

# Image Translation: Cuttoff-marks & stinks, colorization, improve resolution, NST (neural style transfer)

*Pix to pix*

*Supervised Learning !*

Easy to learn and train, difficult to get data



Pix to pix image translation, Phillip Isola et al
In ICCV 2017.

Source training set     Target training set

Input     CycleGAN

*Object transfiguration*

*Unsupervised Learning !*

Difficult to train and learn but easy to get training data

Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks
Jun-Yan Zhu*, Taesung Park*, Phillip Isola, Alexei A. Efros
Berkeley AI Research Lab, UC Berkeley
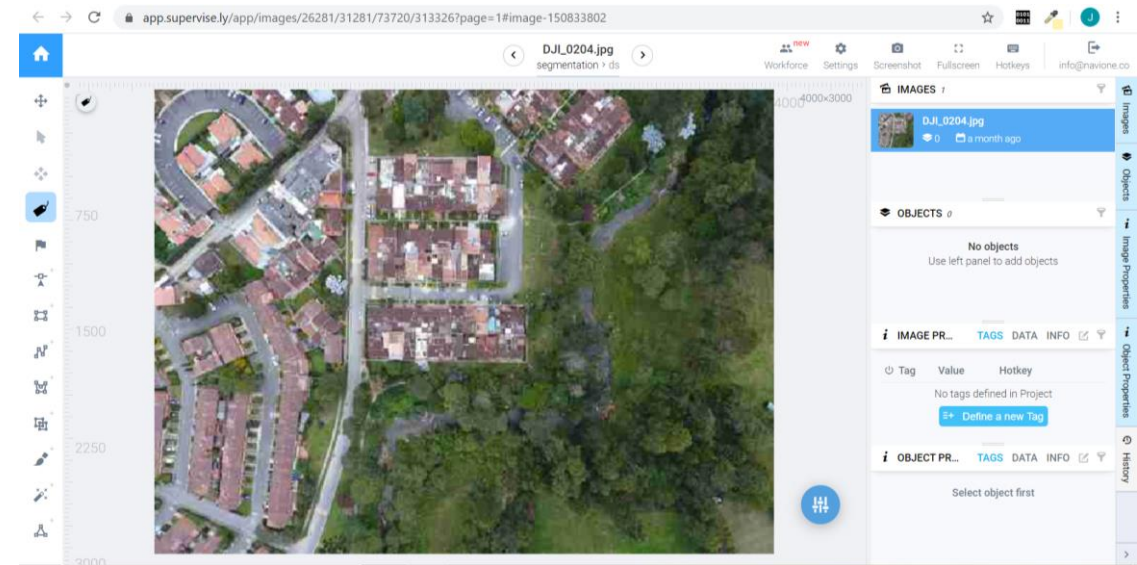In ICCV 2017. (* equal contributions)

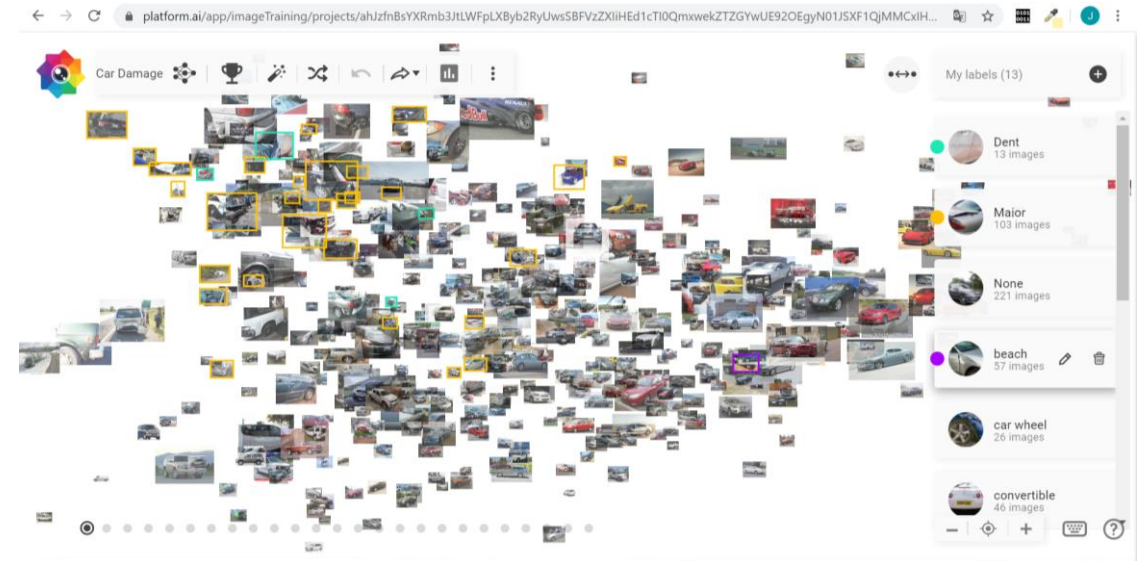# Labelling Time (research topic)

# What about data in text and tables?

| Dir. Inicial | Dirección | Complemento | Barrio / Vereda / I | POI | Pto. geográf | Ciudad | Depto | País |
|---|---|---|---|---|---|---|---|---|
| CRR 14 #7A 34 SUR BARRIO SAN CARLOS O RETEN SUR | Carrera 14 7A Sur Sur 34 | | Barrio San Carlos | | | | | |
| CARRERA 26 A # 1A-43 BARRIO SANTA ISABEL | Carrera 26 A 1A 43 | | Barrio Santa Isabel | | | | | |
| KR 19B1 9 - 25 LOS CORTIJOS | Carrera 19B1 9 25 | | Los Cortijos | | | | | |
| CALLE 7 NO 11-24 | Calle 7 11 24 | | | | | | | |
| APTO 304 T1 EDIF SIEMPRE VERDE CRA 17 #118-32 | Carrera 17 118 32 | Apartamento 304 | | Edificio Siempre Verde | | | | |
| CRA 42 NO 85 A 95 ITAGUI-ITAGUI-ITAGUI-ANTIOQUIA | Carrera 42 85 A 95 | | | | | Itagui | Antioquia | |
| CR 7, 10-87, | Carrera 7 10 87 | | | | | | | |
| CALLE 59 # 56-63, NRO. TORRE 6, APTO. 721 | Calle 59 56 . 63 | Apartamento. 721 | | | | | | |
| CALLE69SUR 46A 64 CALLE DEL BANCO  TORRE ASÍS II APT 701, S, | Calle69Sur 46A 64 | Apartamento 701 | | | | Sabaneta | Antioquia | Colombia |
| MANZANA J5 CASA 11 B. VILLA CATALINA | Manzana J5 | Casa 11 | Villa Catalina | | | | | |
| CRA91#44A39,AMÉRICA NIZA 2PISO MEDELLÍN   BARRIO: AMÉRICA | Cra91 44A39 | | America Niza | | | Medellin | | |
| CR 34 # 10 581 ACOPI YIMBO | Carrera 34 10 581 | | | Acopi | | | Yimbo | |
| CALLE 22 #1-140 AV BOLIVAR FERRETERIA ARGENTINA | Calle 22 1 140 | | | Ferreteria Argentina | | | | |
| TR 58 BIS # 2 C 60 B CAMELIA | Transversal 58 Bis 2 C 60 | | B Camelia | | | | | |
| KILOMETRO 3981 ANIKLO VIAL RIO FRIO ZOBA FRANCA SANTANDE | Kilometro 3981 Anillo Vial Rio Frio | | | Edificio Baiachara | | Santander | | |
| CARRERA 4 #2-03 BARRIO CHAPINERO | Carrera 4 2 03 | | Barrio Chapinero | | | | | |
| CRA 14F 76B 57 SUR  INT 1 MARICHUELA USME BOGOTÁ, D.C. BOG | Carrera 14F 76B Sur 57 | Interior 1 | | | | Bogota D.C. | Bogota | Col |
| CR 10 A # 11 75 LOCL 103 PASAGE GOMEZ | Carrera 10 A 11 75 | Local 103 | | Pasage Gomez | | | | |
| CL 24 # 6 67 CENTRO APTO 301 PEREIRA BARRIO: CENTRO | Calle 24 6 67 | Apartamento 301 | Pereira Barrio Centro | | Centro | | | |
| CLL 82 # 67 A -51 | Calle 82 67 A 51 | | | | | | | |
| CL 21 # 6 36 CENTRO MONTERÍA CÓRDOBA COL | Calle 21 6 36 | | | | Centro | Monteria | Cordoba | Col |
| CL 71 # 6-21 OF 301HIDROCARBUROS DEL CASANARE S.A.S | Calle 71 6 21 | | | Hidrocarburos Del Casanare S.A.S | | | | |
| CL 20 A # 12 70 METRO SECCION ELECTRO | Calle 20 A 12 70 | | | Metro Seccion Electro | | | | |
| CARRERA 22B NUMERO 13A 47 BARRIO GUAYAQUIL   BARRIO GUA | Carrera 22B 13A 47 | | Barrio Guayaquil | Barrio Guayaquil | | | | |
| CR 1 # 12 -118 CC PLAZA BOCAGRANDE LC 105 | Carrera 1 12 118 | Local 105 | | Cc Plaza Bocagrande | | | | |

# Software for Labelling

Supervisely: [www.app.supervise.ly](www.app.supervise.ly)



Platform.ai: [www.platform.ai](www.platform.ai)



*PixelAnnotationTool (github)*

*Gidia*
Grupo de I+D
en Inteligencia Artificial

UNIVERSIDAD
NACIONAL
DE COLOMBIA

multilabel.herokuapp.com

Gidia
Grupo de I+D
en Inteligencia Artificial

UNIVERSIDAD
NACIONAL
DE COLOMBIA