

# Detección e identificación de características en vehículos utilizando algoritmos de Machine Learning.

Instituto Balseiro - Universidad Nacional de Cuyo

Alumno: Enrique Nicanor Mariotti

Noviembre 2020

# Texto



# Texto

## Particularidades

- Permite inferir información sobre el contenido semántico de la escena.
- Particularmente difícil mediante técnicas de OCR, ideadas para la lectura de documentos escritos.
- Gran numero de aplicaciones y amplia disponibilidad.

# Texto

## Problemas asociados:

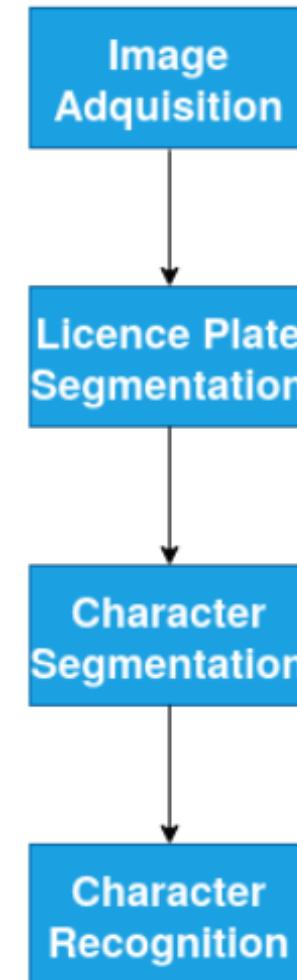
- Variables en lay-out
- Estilos y tipografías distintas
- Iluminación variable
- Oclusiones parciales o totales
- Orientación
- Ruido
- Background FP

# ALPR

## **Automatic Licence Plate Recognition:**

- Control automotor por parte de las fuerzas de la ley.
- Sondeo de información sobre lotes automotores vastos.
- Control de acceso a lugares privados.

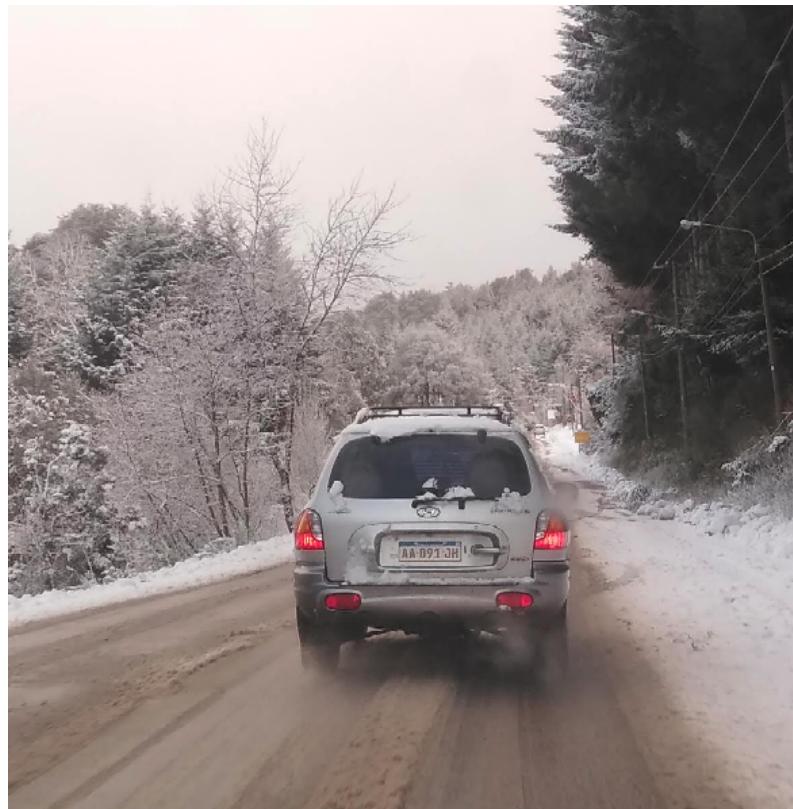
# ALPR



## Problemas:

- Diferentes colores, fuentes y fondos
- Tamaño de caracteres no uniforme
- Factores externos:
  - Velocidad del vehículo
  - Angulo de la cámara
  - Trasfondo de la escena
  - Iluminación
  - Muchos otros...

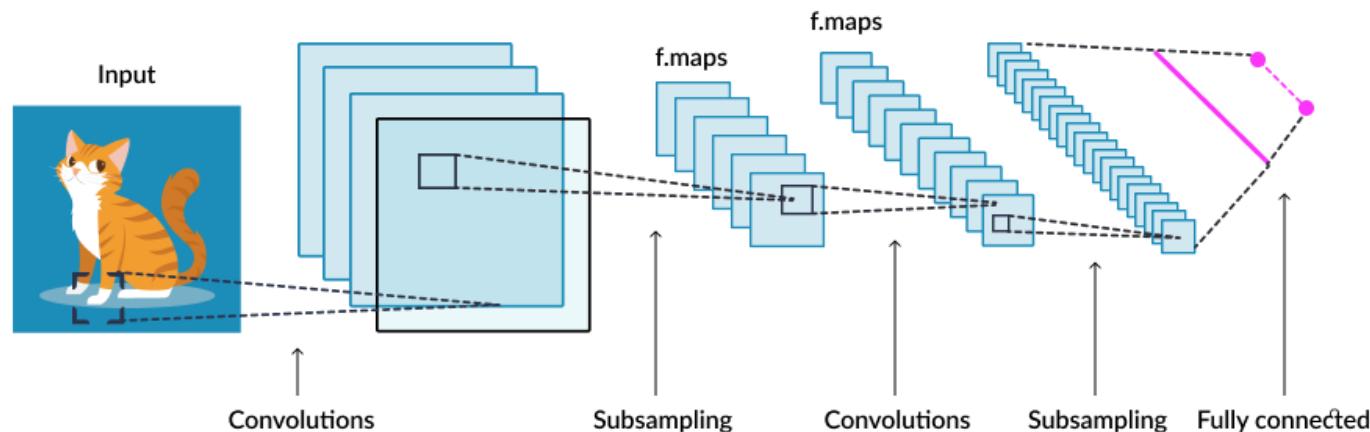
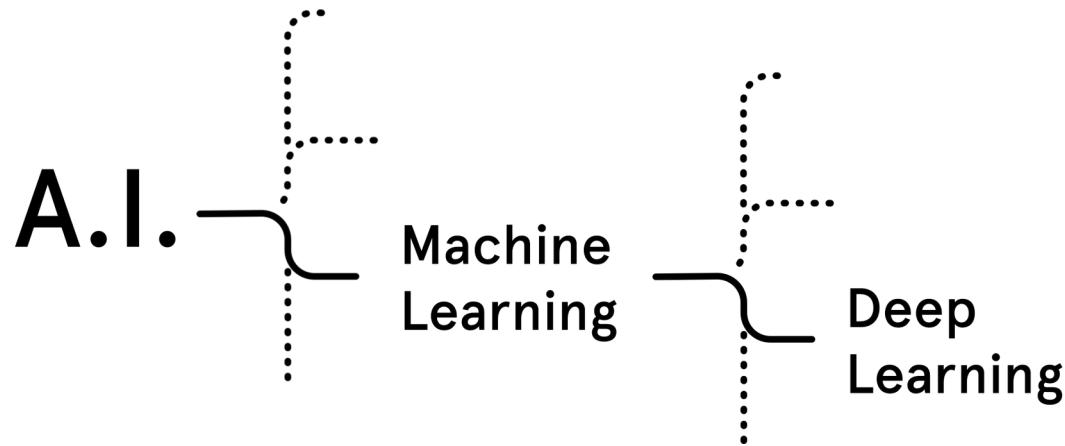
# ALPR



# Objetivos

- Explorar una solución integral al problema **ALPR** usando **Deep Learning (DL)**:
  - Patrones complejos de información.
  - Trabajar con datos no estructurados.
  - La interpretabilidad del resultado no es un factor preponderante.

# Objetivos



# Objetivos

## Problemática:

- Disponibilidad de datos anotados (**ALPR**) y recursos físicos (muy) limitada.

## ¿Solución?:

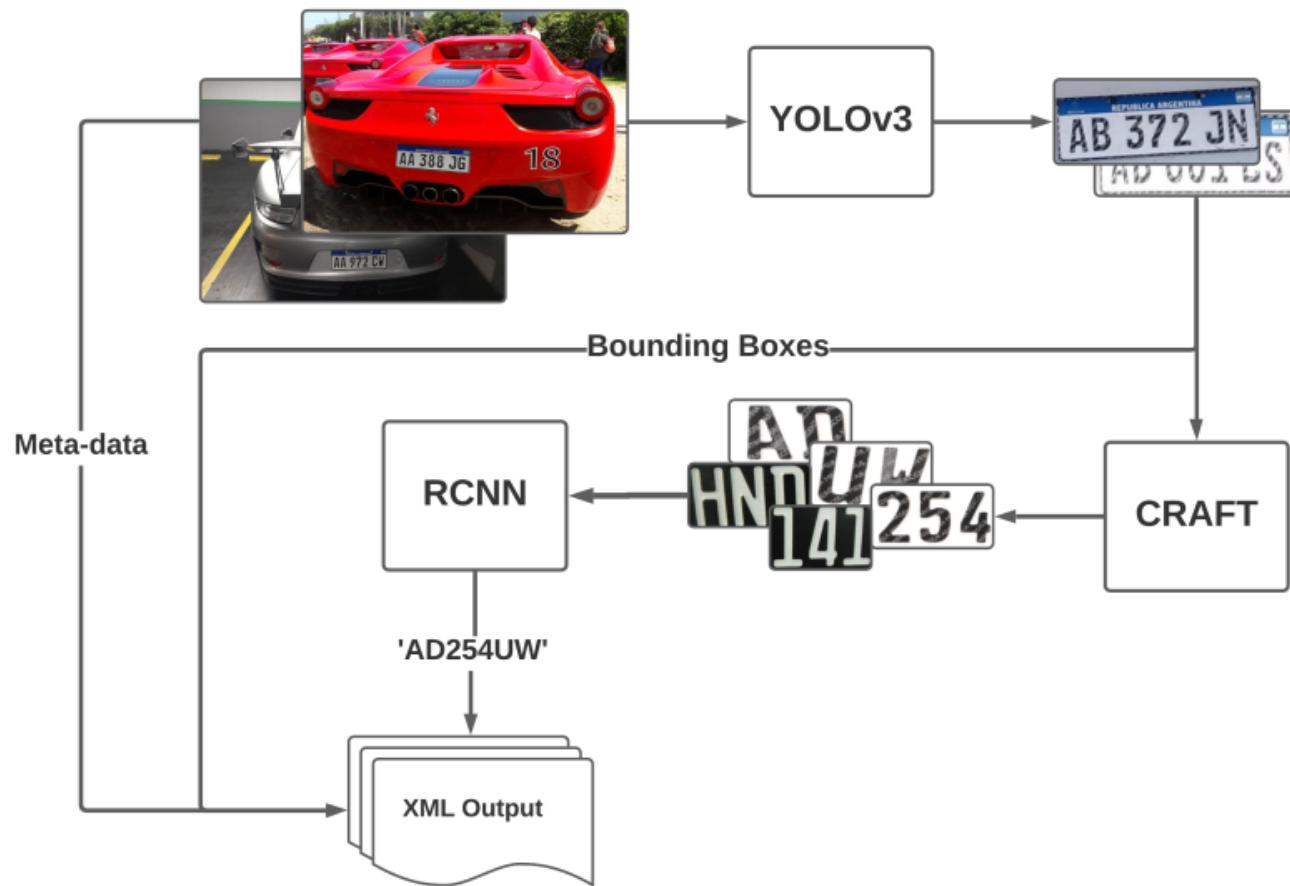
- Priorizar el uso de transfer learning en redes profundas.

# Objetivos

## Requerimientos adicionales:

- Secuencia de etapas modulares (**OOP**). Los módulos u objetos individuales pueden ser usados con otros propósitos.
- **Caso de estudio:** Patente única del Mercosur, para facilitar la comparación de los resultados.

# Lay-Out



# YOLO

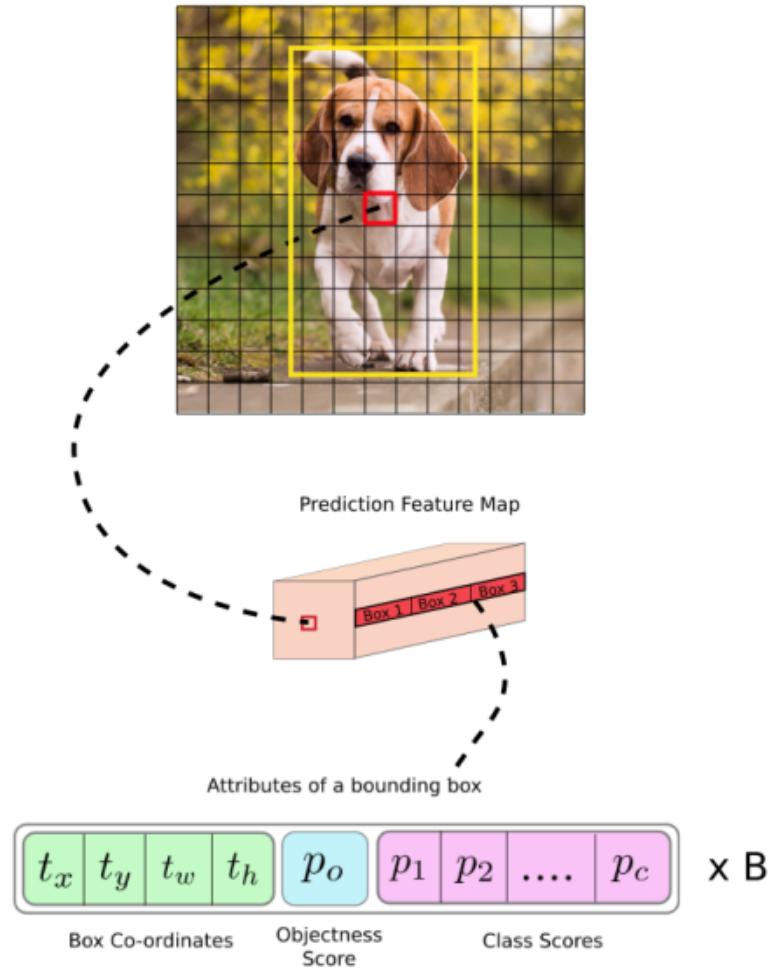
| Type          | Filters       | Size             | Output           |
|---------------|---------------|------------------|------------------|
| Convolutional | 32            | $3 \times 3$     | $256 \times 256$ |
| Convolutional | 64            | $3 \times 3 / 2$ | $128 \times 128$ |
| 1x            | Convolutional | 32               | $1 \times 1$     |
|               | Convolutional | 64               | $3 \times 3$     |
|               | Residual      |                  | $128 \times 128$ |
|               | Convolutional | 128              | $3 \times 3 / 2$ |
| 2x            | Convolutional | 64               | $1 \times 1$     |
|               | Convolutional | 128              | $3 \times 3$     |
|               | Residual      |                  | $64 \times 64$   |
|               | Convolutional | 256              | $3 \times 3 / 2$ |
| 8x            | Convolutional | 128              | $1 \times 1$     |
|               | Convolutional | 256              | $3 \times 3$     |
|               | Residual      |                  | $32 \times 32$   |
|               | Convolutional | 512              | $3 \times 3 / 2$ |
| 8x            | Convolutional | 256              | $1 \times 1$     |
|               | Convolutional | 512              | $3 \times 3$     |
|               | Residual      |                  | $16 \times 16$   |
|               | Convolutional | 1024             | $3 \times 3 / 2$ |
| 4x            | Convolutional | 512              | $1 \times 1$     |
|               | Convolutional | 1024             | $3 \times 3$     |
|               | Residual      |                  | $8 \times 8$     |
|               | Avgpool       |                  | Global           |
|               | Connected     |                  | 1000             |
|               | Softmax       |                  |                  |

Tabla 2.1: CNN detrás del procesamiento de YOLO [1].

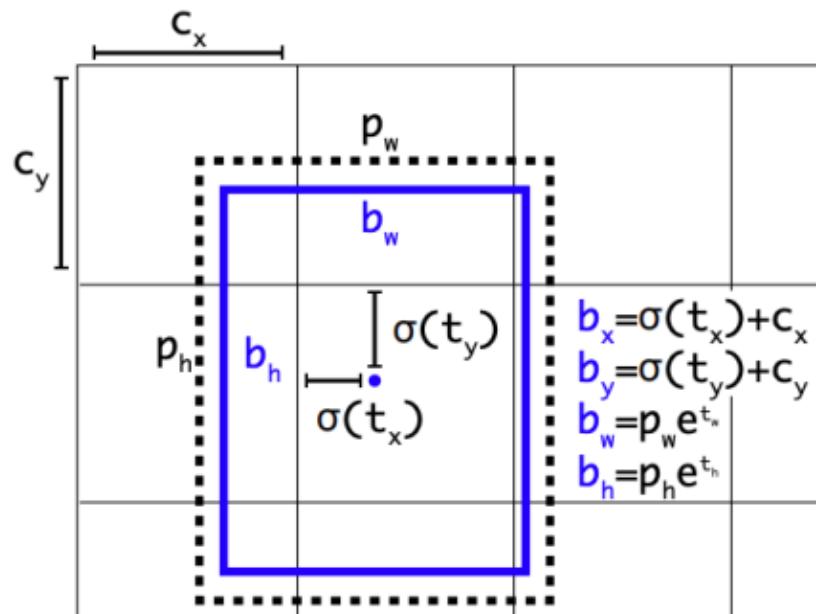
# YOLO

|    | Type          | Filters | Size             | Output           |
|----|---------------|---------|------------------|------------------|
| 1x | Convolutional | 32      | $3 \times 3$     | $256 \times 256$ |
|    | Convolutional | 64      | $3 \times 3 / 2$ | $128 \times 128$ |
|    | Convolutional | 32      | $1 \times 1$     |                  |
|    | Convolutional | 64      | $3 \times 3$     |                  |
| 2x | Residual      |         |                  | $128 \times 128$ |
|    | Convolutional | 128     | $3 \times 3 / 2$ | $64 \times 64$   |
|    | Convolutional | 64      | $1 \times 1$     |                  |
|    | Convolutional | 128     | $3 \times 3$     |                  |
| 8x | Residual      |         |                  | $64 \times 64$   |
|    | Convolutional | 256     | $3 \times 3 / 2$ | $32 \times 32$   |
|    | Convolutional | 128     | $1 \times 1$     |                  |
|    | Convolutional | 256     | $3 \times 3$     |                  |
| 8x | Residual      |         |                  | $32 \times 32$   |
|    | Convolutional | 512     | $3 \times 3 / 2$ | $16 \times 16$   |
|    | Convolutional | 256     | $1 \times 1$     |                  |
|    | Convolutional | 512     | $3 \times 3$     |                  |
| 4x | Residual      |         |                  | $16 \times 16$   |
|    | Convolutional | 1024    | $3 \times 3 / 2$ | $8 \times 8$     |
|    | Convolutional | 512     | $1 \times 1$     |                  |
|    | Convolutional | 1024    | $3 \times 3$     |                  |
|    | Residual      |         |                  | $8 \times 8$     |
|    | Avgpool       |         | Global           |                  |
|    | Connected     |         | 1000             |                  |
|    | Softmax       |         |                  |                  |

Tabla 2.1: CNN detrás del procesamiento de YOLO [1].



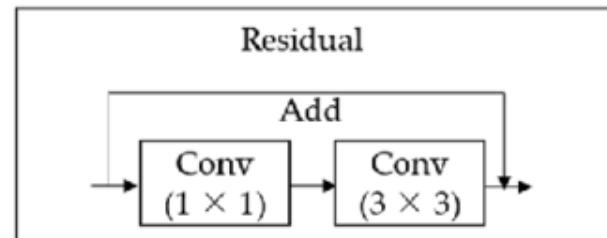
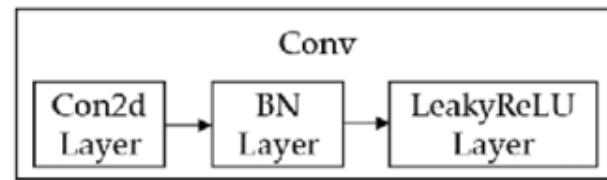
# YOLO



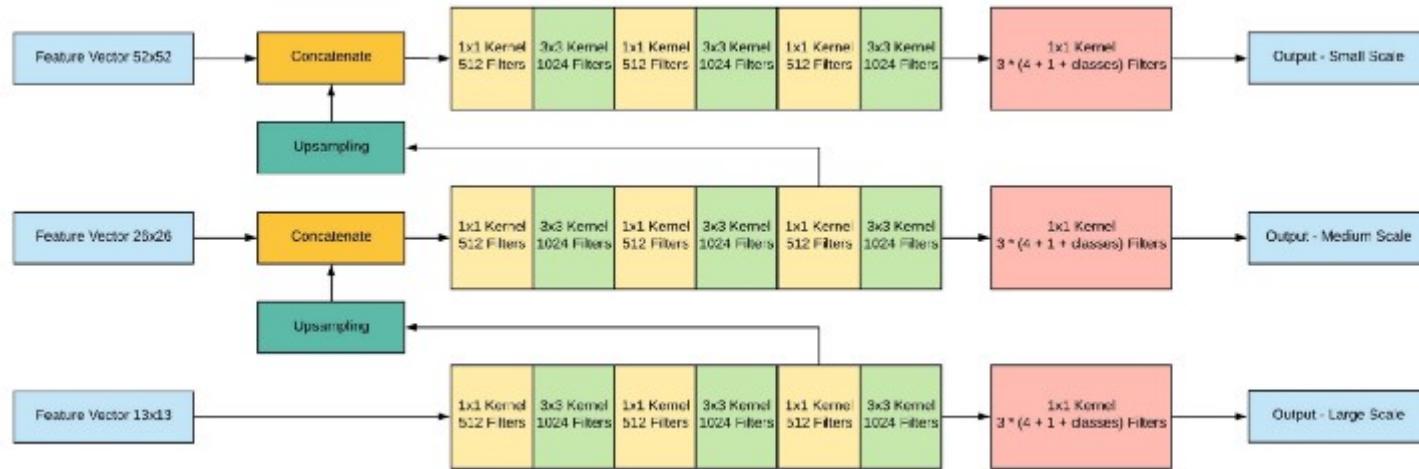
**Figura 2.3:** Predicción de la celda ( $c_x, c_y$ ) en función de las dimensiones predefinidas de los anchors y de una de las predicciones de la red [1].

# YOLO

| Layer    | Filters size        | Repeat   | Output size      |
|----------|---------------------|----------|------------------|
| Image    |                     |          | $416 \times 416$ |
| Conv     | $32 3 \times 3/1$   | 1        | $416 \times 416$ |
| Conv     | $64 3 \times 3/2$   | 1        | $208 \times 208$ |
| Conv     | $32 1 \times 1/1$   | Conv     | $208 \times 208$ |
| Conv     | $64 3 \times 3/1$   | Conv     | $208 \times 208$ |
| Residual |                     | Residual | $208 \times 208$ |
| Conv     | $128 3 \times 3/2$  | 1        | $104 \times 104$ |
| Conv     | $64 1 \times 1/1$   | Conv     | $104 \times 104$ |
| Conv     | $128 3 \times 3/1$  | Conv     | $104 \times 104$ |
| Residual |                     | Residual | $104 \times 104$ |
| Conv     | $256 3 \times 3/2$  | 1        | $52 \times 52$   |
| Conv     | $128 1 \times 1/1$  | Conv     | $52 \times 52$   |
| Conv     | $256 3 \times 3/1$  | Conv     | $52 \times 52$   |
| Residual |                     | Residual | $52 \times 52$   |
| Conv     | $512 3 \times 3/2$  | 1        | $26 \times 26$   |
| Conv     | $256 1 \times 1/1$  | Conv     | $26 \times 26$   |
| Conv     | $512 3 \times 3/1$  | Conv     | $26 \times 26$   |
| Residual |                     | Residual | $26 \times 26$   |
| Conv     | $1024 3 \times 3/2$ | 1        | $13 \times 13$   |
| Conv     | $512 1 \times 1/1$  | Conv     | $13 \times 13$   |
| Conv     | $1024 3 \times 3/1$ | Conv     | $13 \times 13$   |
| Residual |                     | Residual | $13 \times 13$   |



# YOLO



# YOLO

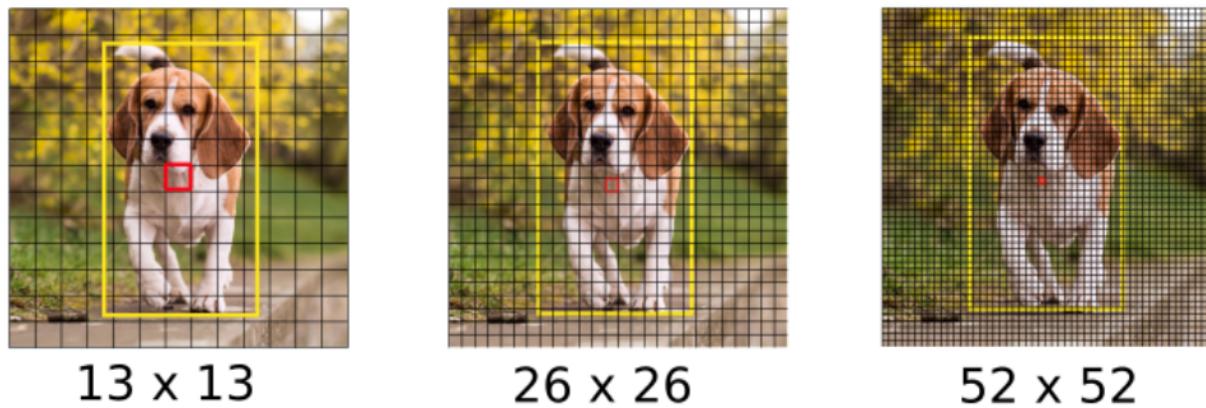
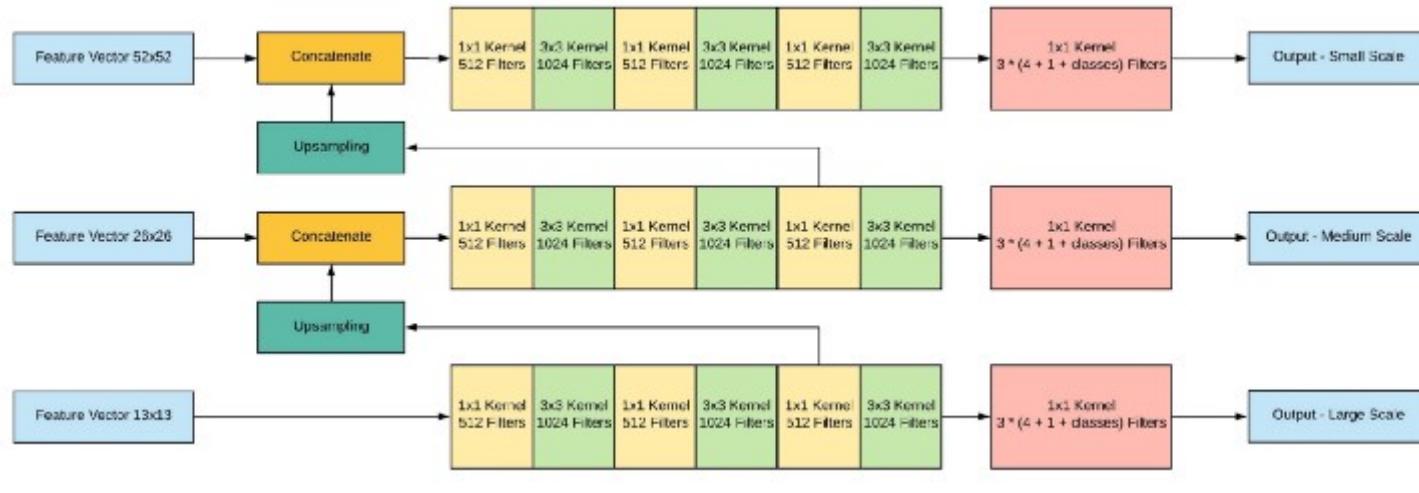
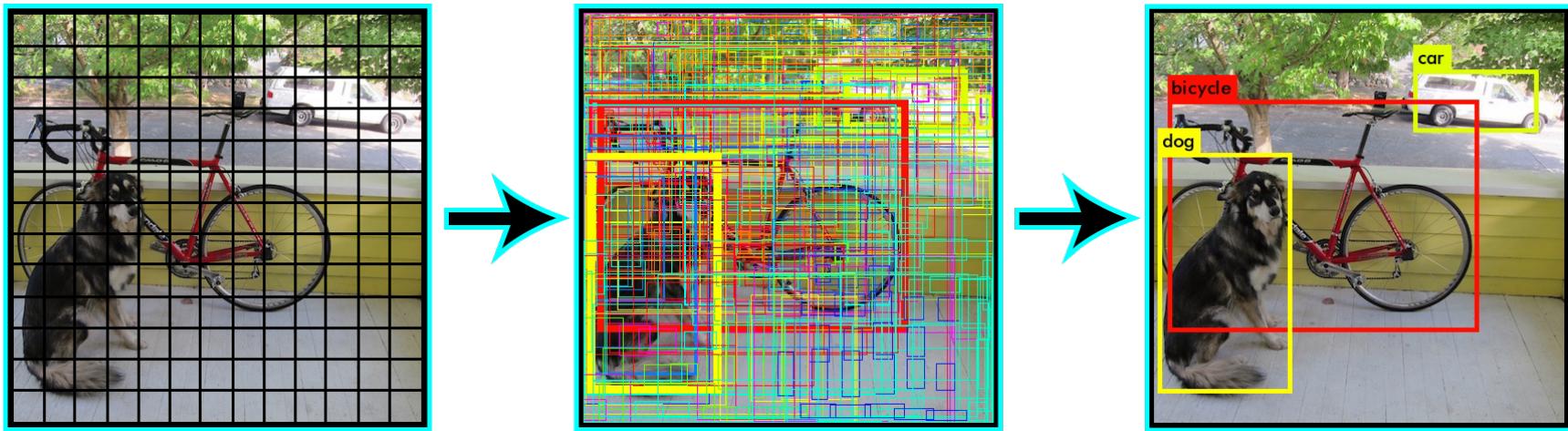


Figura 2.4: Salida multiescala de YOLOv3 [28].

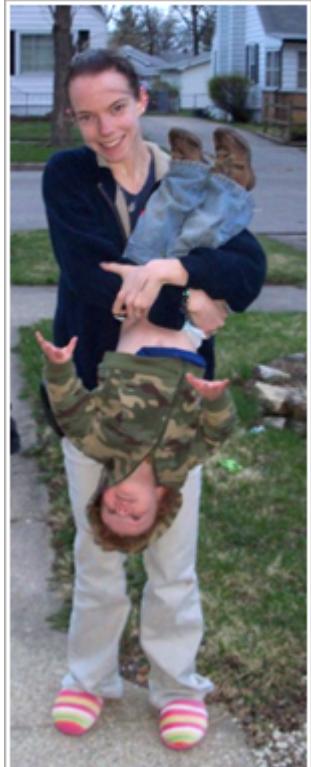
# YOLO



# YOLO

$$\begin{aligned} & \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\ & + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ \left( \sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left( \sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right] \\ & + \lambda_{\text{obj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left( C_i - \hat{C}_i \right)^2 \\ & + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} \left( C_i - \hat{C}_i \right)^2 \\ & + \lambda_{\text{classes}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \end{aligned}$$

# YOLO



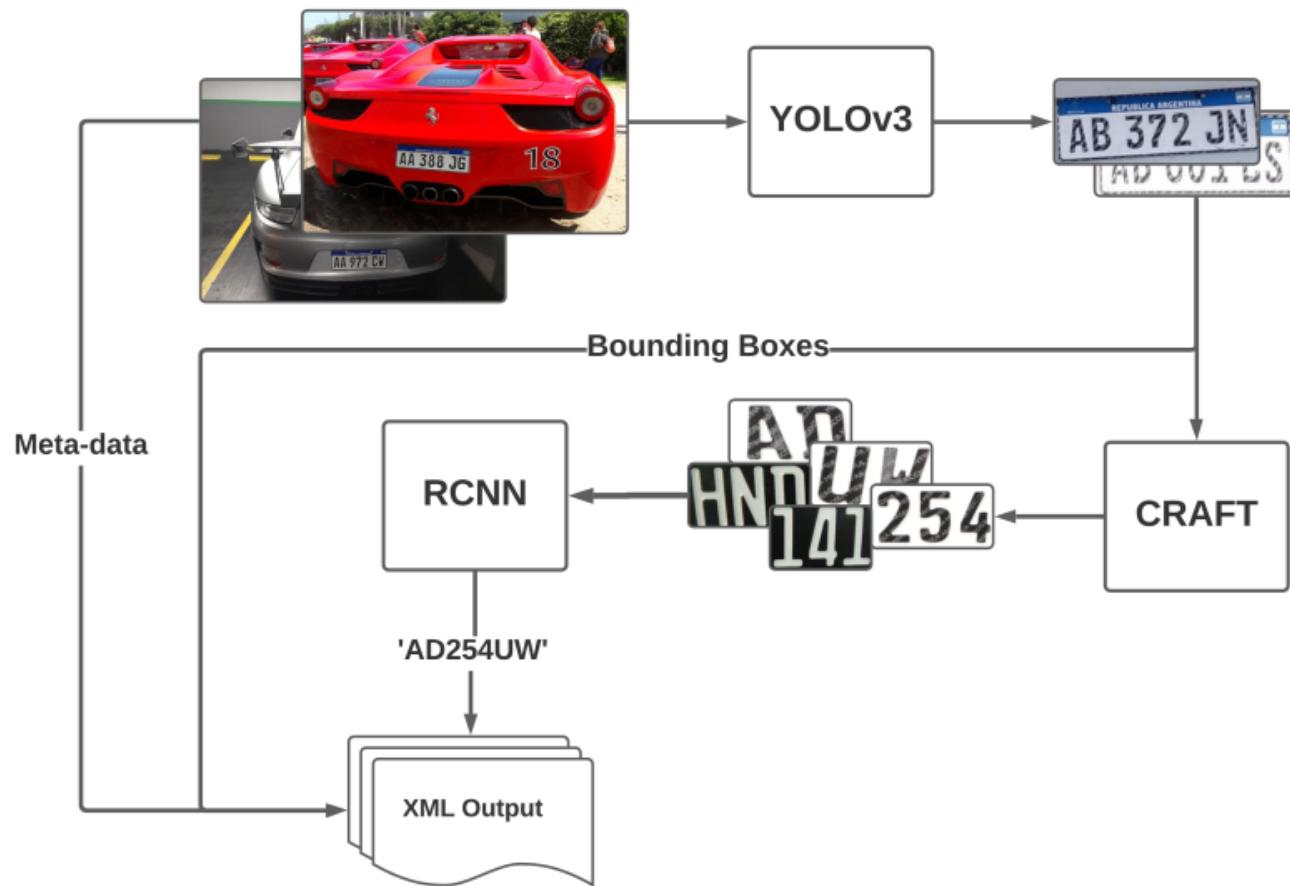
$$\begin{aligned} & \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\ & + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ \left( \sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left( \sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right] \\ & + \lambda_{\text{obj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left( C_i - \hat{C}_i \right)^2 \\ & + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} \left( C_i - \hat{C}_i \right)^2 \\ & + \lambda_{\text{classes}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \end{aligned}$$

# YOLO



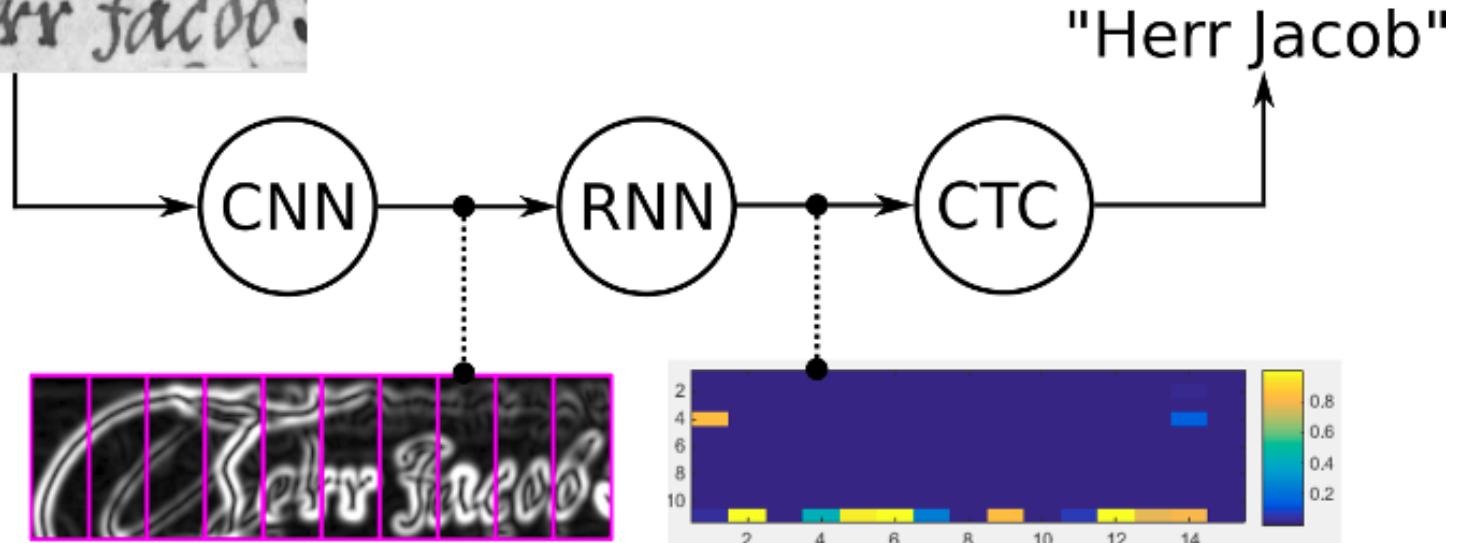
- 100 epochs
- Darknet-53 freezada
- Stochastic Gradient Descent (SGD)
- Learning rate:  $1 \times 10^{-4}$
- Decay:  $5 \times 10^{-4}$
- Momentum:  $9 \times 10^{-1}$
- Batch de entrenamiento: 16 imágenes.
- Aumentación de datos no distorsiva.
- Patente de la Unión Europea (500 imágenes). Anotación con localización

# Lay-Out



# RNN

Herr Jacob.

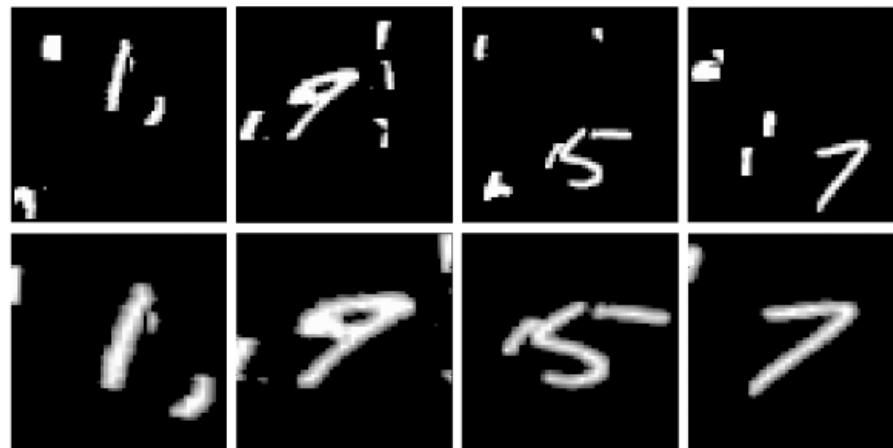


# RNN

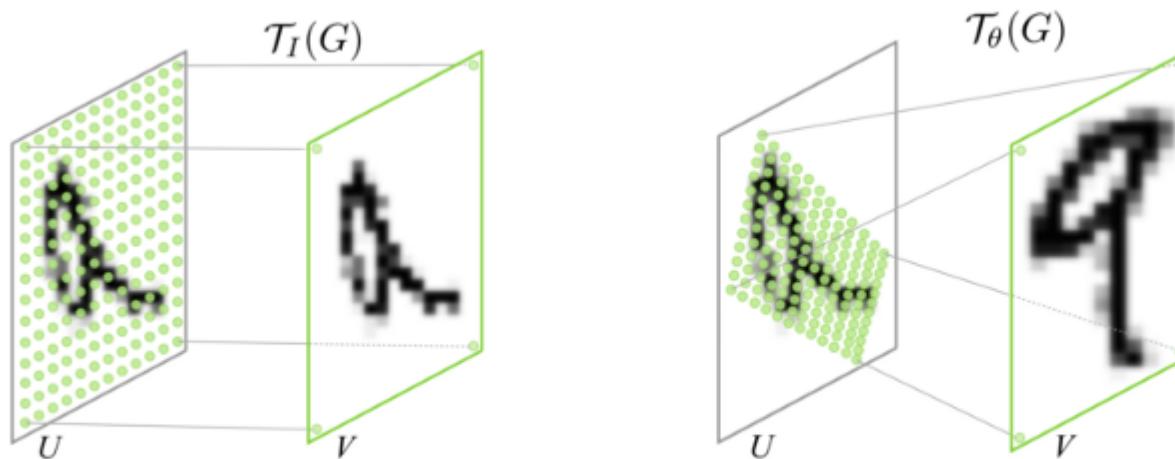
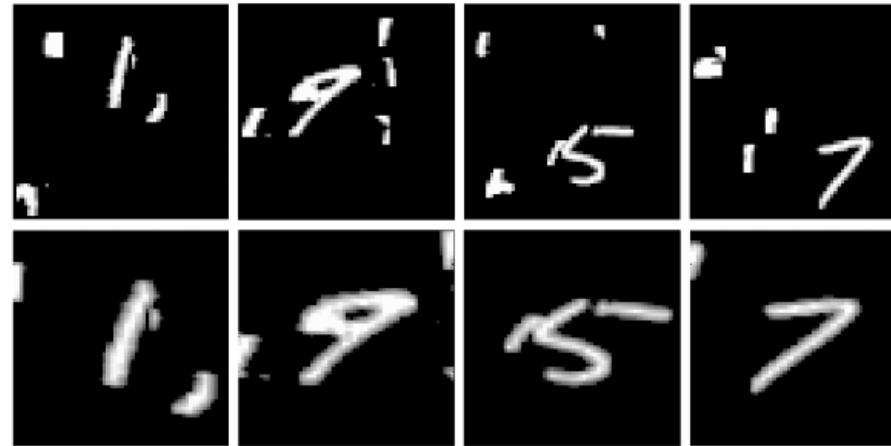
| Bloque | Layer              | Output Size | Features | Kernel Size | Padding | Activation |
|--------|--------------------|-------------|----------|-------------|---------|------------|
| 1      | Input              | 200x31      | -        | -           | -       | -          |
|        | Conv2D             | 200x31x64   | 64       | 3x3         | same    | relu       |
|        | Conv2D             | 200x31x128  | 128      | 3x3         | same    | relu       |
|        | Conv2D             | 200x31x256  | 256      | 3x3         | same    | relu       |
|        | BatchNormalization | 200x31x256  | -        | -           | -       | -          |
|        | MaxPool2D          | 100x15x256  | -        | 2x2         | -       | -          |
| 2      | Conv2D             | 100x15x256  | 256      | 3x3         | same    | relu       |
|        | Conv2D             | 100x15x512  | 512      | 3x3         | same    | relu       |
|        | BatchNormalization | 100x15x512  | -        | -           | -       | -          |
|        | MaxPool2D          | 50x7x512    | -        | 2x2         | -       | -          |
| 3      | Conv2D             | 50x7x512    | 512      | 3x3         | same    | relu       |
|        | Conv2D             | 50x7x512    | 512      | 3x3         | same    | relu       |
|        | BatchNormalization | 50x7x512    | -        | -           | -       | -          |

**Tabla 3.1:** Arquitectura empleada.

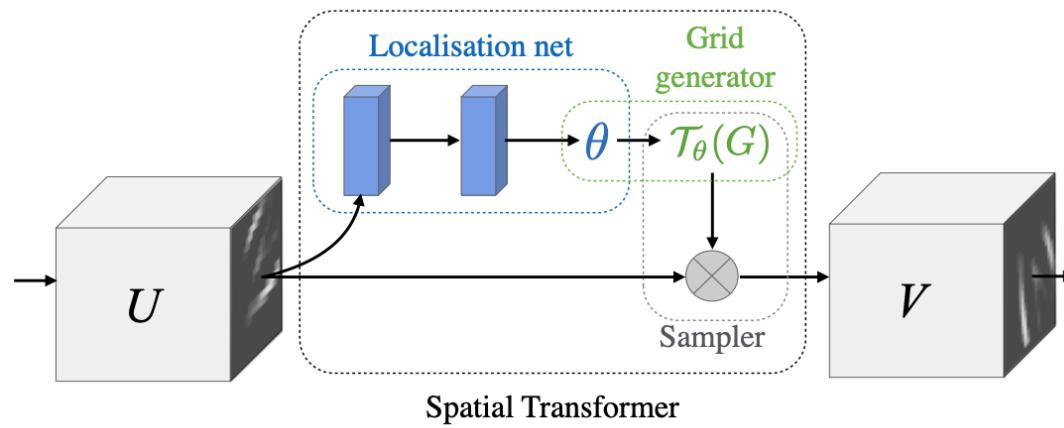
# RNN



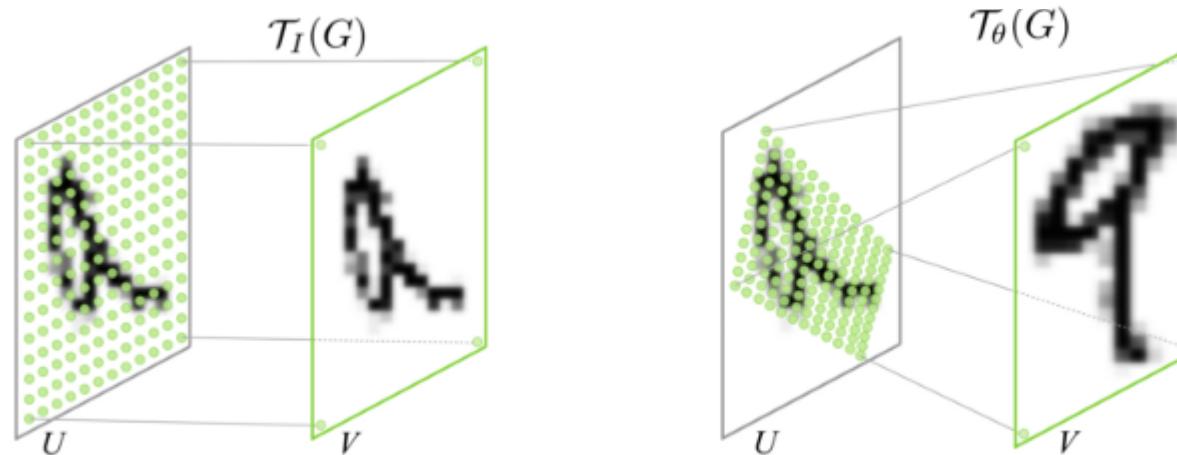
# RNN



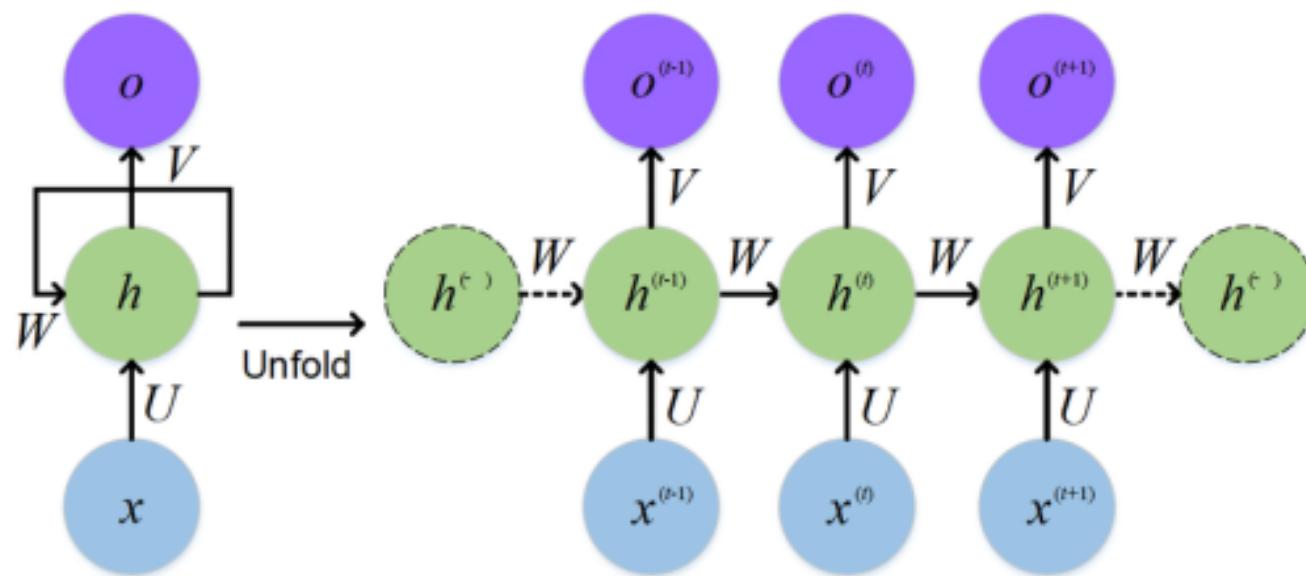
# RNN



Spatial Transformer

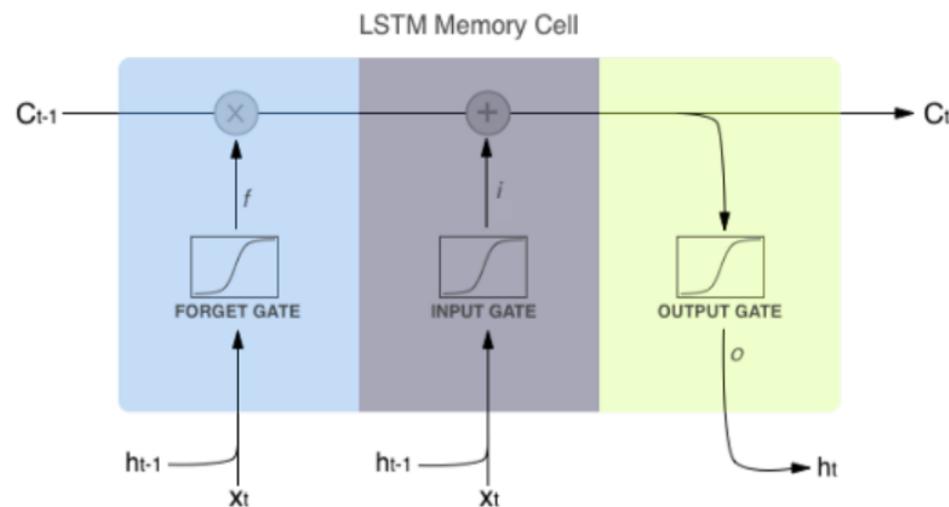


# RNN

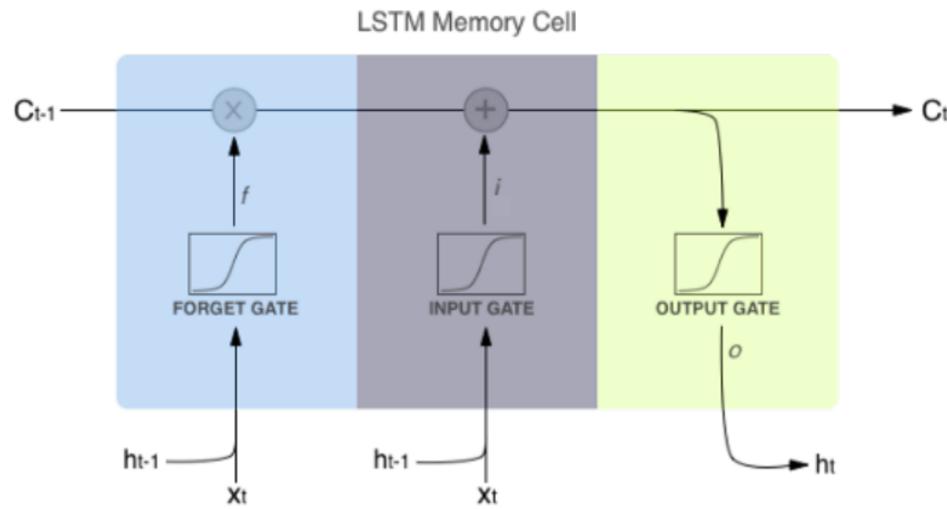


**Figura 3.4:** Desglose del funcionamiento de una RNN [7].

# RNN



# RNN

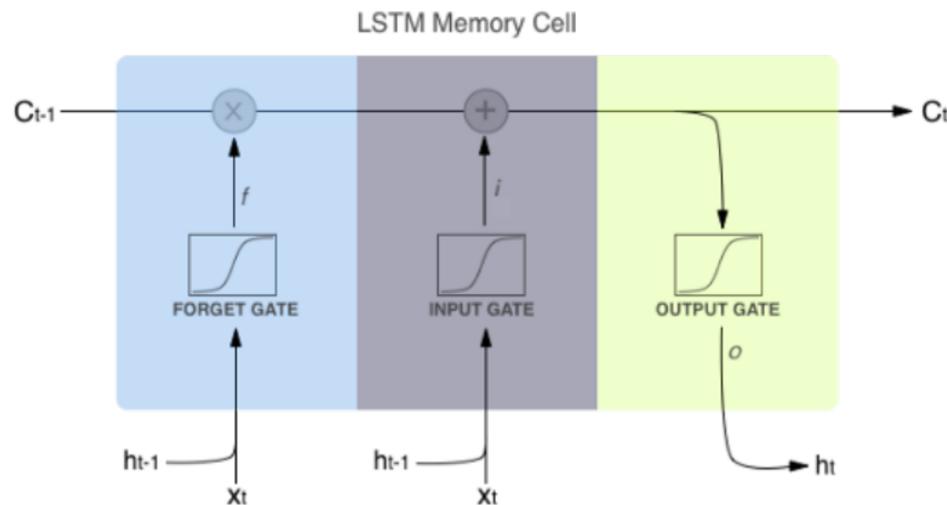


$$i_t = \sigma (w_i [h_{t-1}, x_t] + b_i)$$

$$f_t = \sigma (w_f [h_{t-1}, x_t] + b_f)$$

$$o_t = \sigma (w_o [h_{t-1}, x_t] + b_o)$$

# RNN



$$i_t = \sigma (w_i [h_{t-1}, x_t] + b_i)$$

$$f_t = \sigma (w_f [h_{t-1}, x_t] + b_f)$$

$$o_t = \sigma (w_o [h_{t-1}, x_t] + b_o)$$

$$\tilde{c}_t = \tanh (w_c [h_{t-1}, x_t] + b_c)$$

$$c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$$

$$h_t = o_t * \tanh (c_t)$$

# RNN

h h e ε ε l l l ε l l o

h e ε l ε l l o

h e l l l o

h e l l o

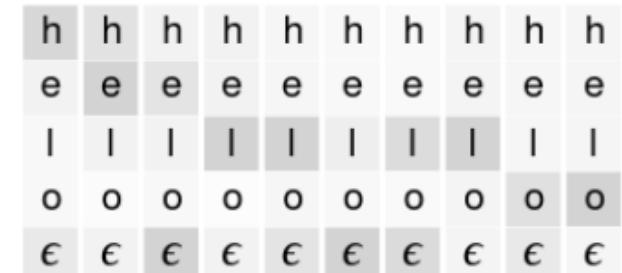
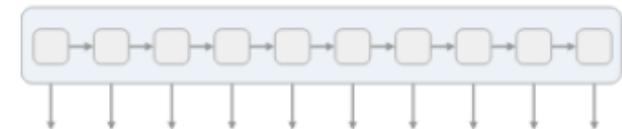
# RNN

h h e ε ε l l l ε l l o

h e ε l ε l o

h e l l o

h e l l o



h e ε l l ε l l o o  
h h e l l l ε ε l ε o  
ε e ε l l l ε ε ε l o o

h e l l o  
e l l o  
h e l o

# RNN

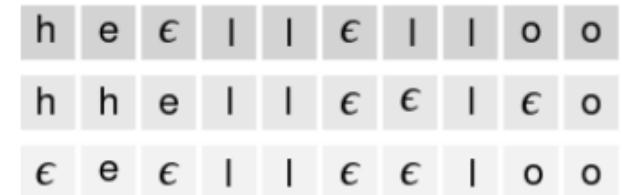
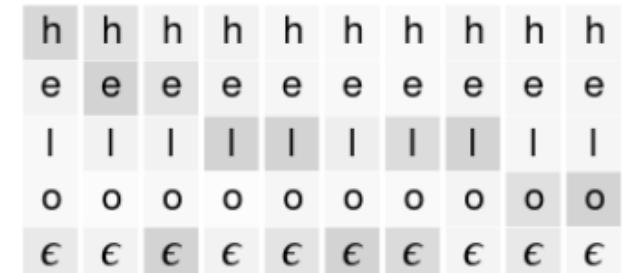
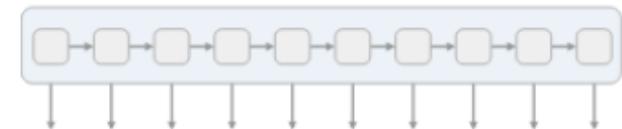
h h e ε ε l l l ε l l o

h e ε l ε l o

h e l o

h e l l o

$$p(\mathbf{Y}|\mathbf{X}) = \sum_{A \in \mathcal{A}_{X,Y}} \prod_{t=1}^T p(\alpha_t|\mathbf{X})$$



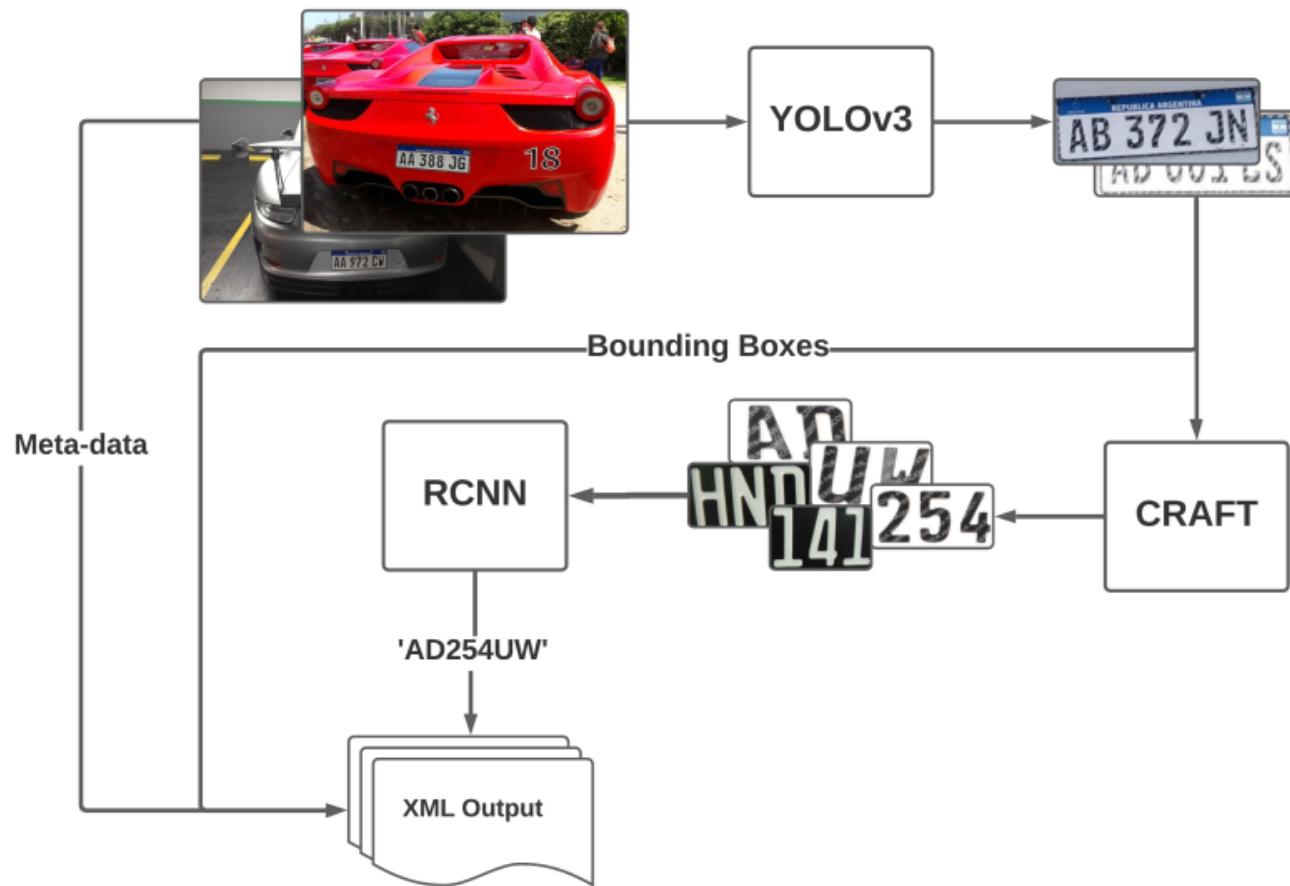
h e l l o  
e l l o  
h e l o

# RNN

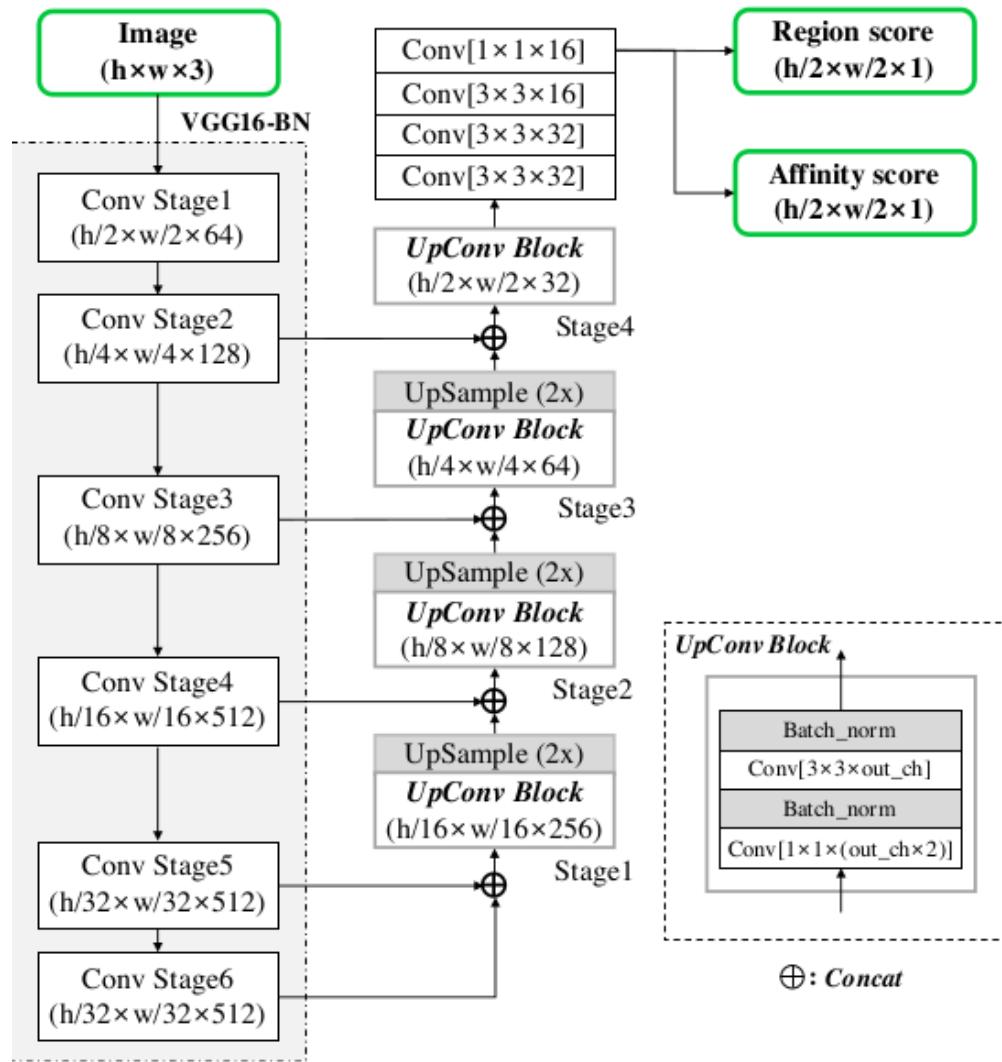


- 50 epochs
- Adam
- Learning rate:  $1 \times 10^{-3}$
- $\beta_1 = \beta_2 = 9 \times 10^{-1}$
- $\epsilon = 1 \times 10^{-8}$
- Decay:  $5 \times 10^{-3}$
- Batch de entrenamiento: 128 imágenes.
- Sin aumentación de datos.
- Imágenes sintéticas ( $1,5 \times 10^5$ ). **Padding**.

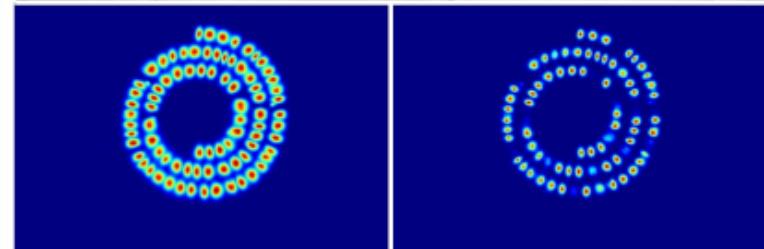
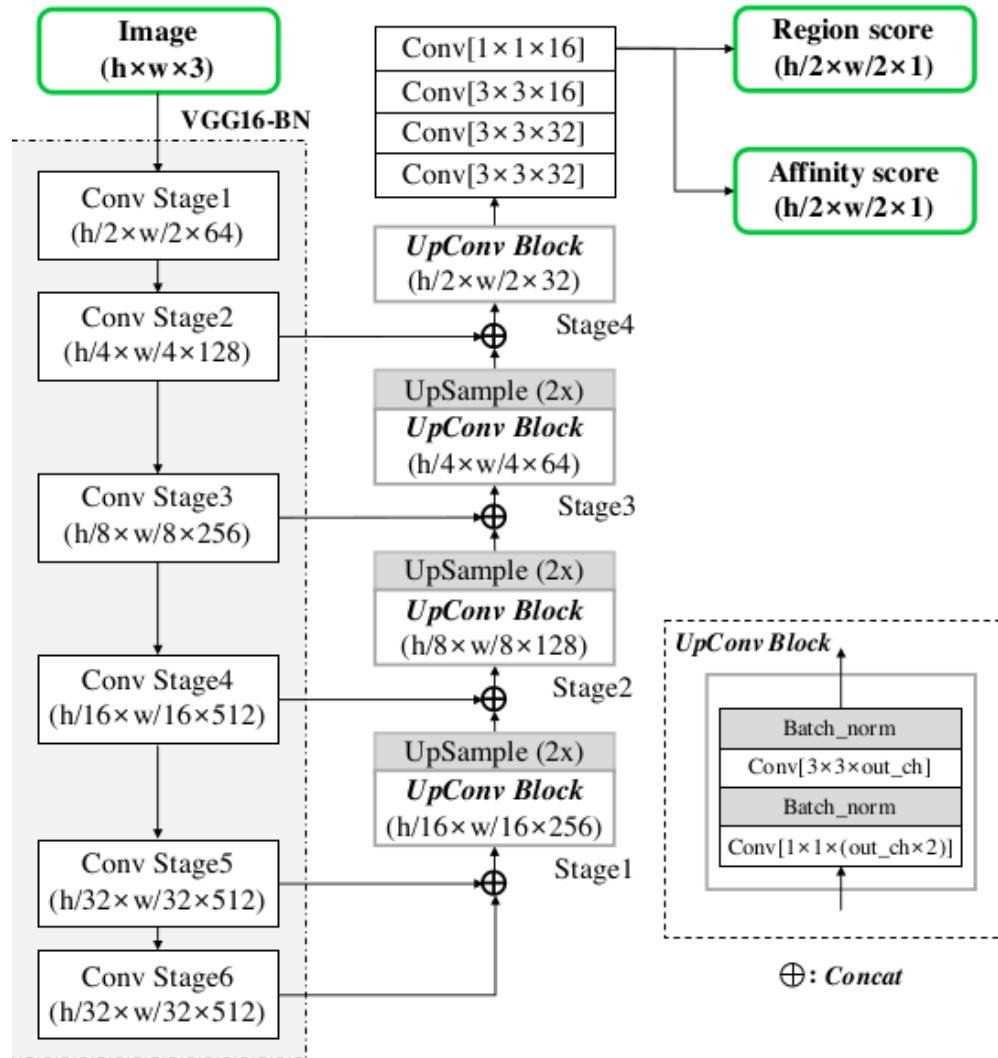
# Lay-Out



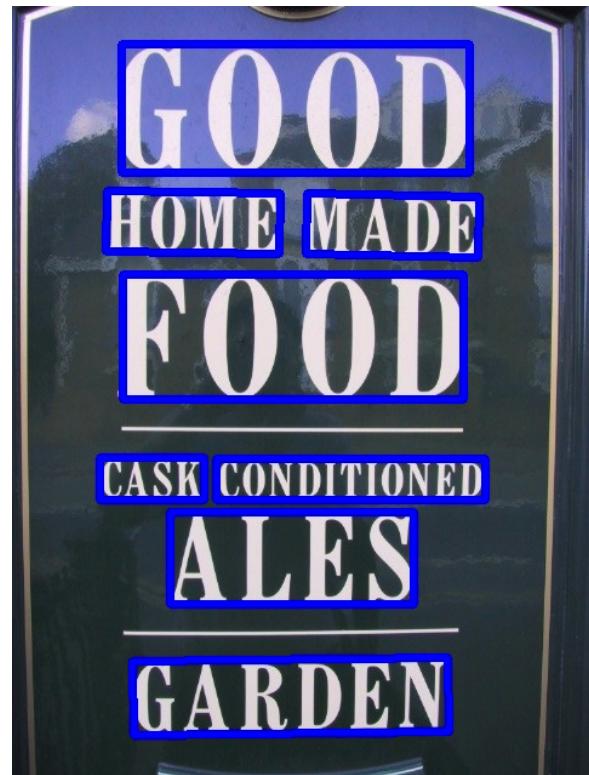
# CRAFT



# CRAFT



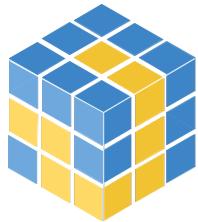
# CRAFT



# CRAFT

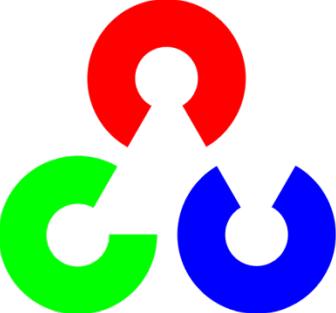


# Tools



*NumPy*

**matplotlib** 

  
**OpenCV**

 **Keras**



# Tools

```
enmariotti@raxamd:~/CV$ tree . -d
```

```
.
```

- |- binder
- |- data
  - |- footage
    - |- aerial\_footage
    - |- street\_footage
    - |- text\_footage
  - |- images
  - |- test
  - |- train
    - |- 00negative
    - |- cars
    - |- text
- |- dev
  - |- cam\_examples
    - |- tests
  - |- yolo\_examples
- |- doc
- |- models
  - |- \_\_pycache\_\_
- |- models-checkpoints
  - |- cam\_model
    - |- resnet50\_car-text\_non-freeze
    - |- resnet50\_text\_non-freeze
- |- src
  - |- \_\_pycache\_\_

25 directories



# Hardware

**CPU:** 2 x Intel(R) Xeon(R) CPU E5-2620 v4 @ 2.10GHz 8 cores por procesador, 16 subprocesos.

**MEM:** 32 GB DDR2

**GPU:** NVIDIA Tesla K-40 12 GB GDRR5

**GPU:** NVIDIA Quadro M4000 8GB GDDR5



# Dataset



Figura 4.1: Imagen original sin modificar.

# Dataset



Figura 4.2: Imagen modificada con Albumentations [12].



Figura 4.3: Imagen modificada con Albumentations [12].

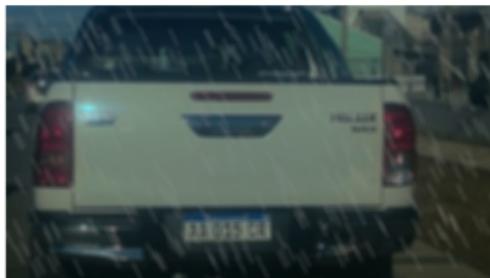


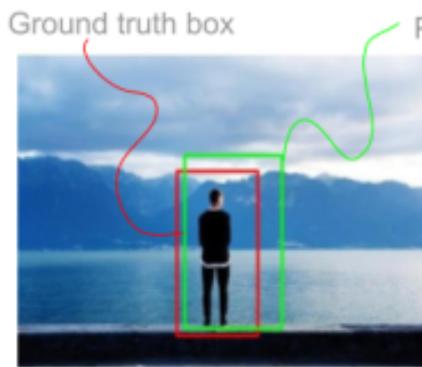
Figura 4.4: Imagen modificada con Albumentations [12].



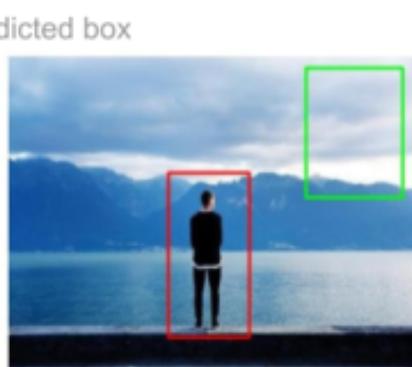
Figura 4.1: Imagen original sin modificar.

# Metricas

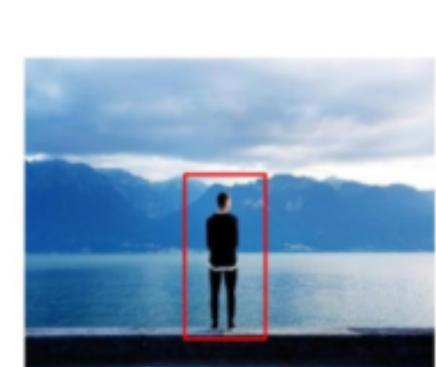
True Positive - TP



False Positive - FP



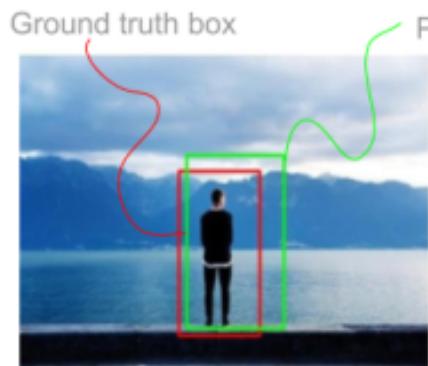
False Negative - FN



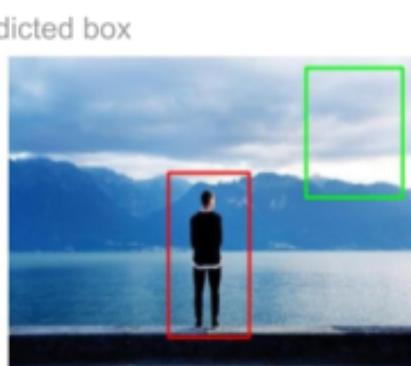
**Figura 4.5:** Ejemplo de clasificación en categorías no triviales [13].

# Metricas

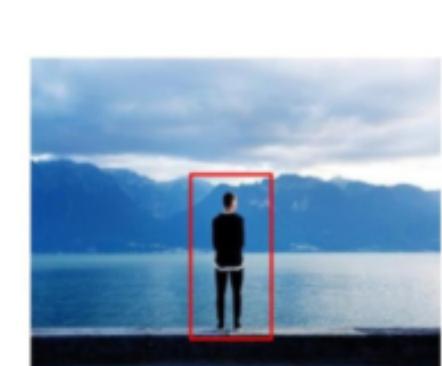
True Positive - TP



False Positive - FP



False Negative - FN



**Figura 4.5:** Ejemplo de clasificación en categorías no triviales [13].

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN}$$

$$Precision = \frac{TP}{TP + FP} = \frac{\text{true object detection}}{\text{all detected boxes}}$$

$$Recall = \frac{TP}{TP + FN} = \frac{\text{true object detection}}{\text{all ground truth boxes}}$$

# Metricas

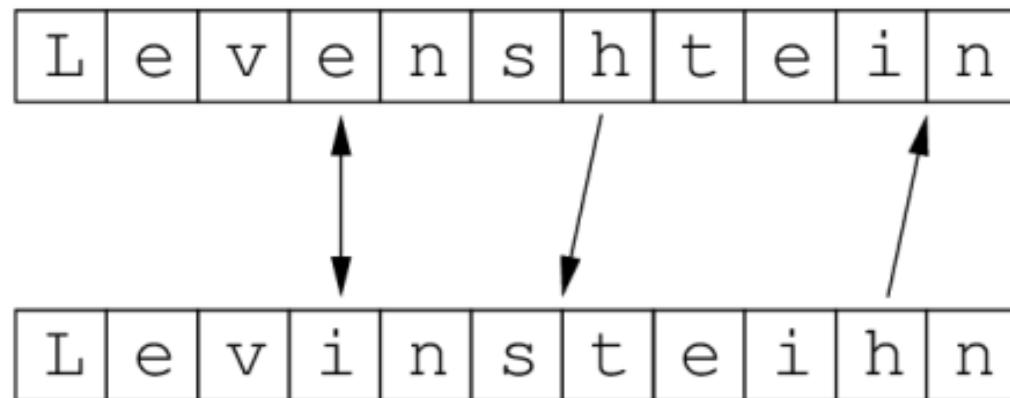


Predictions: [['aa', '205', 'nx'], ['gdt', '465']]

# Metricas

## Text Distance Infographics

| CATEGORY                | APPROACH   | ⊕ PROS   | ⊖ CONS  | ALGORITHM  | AREAS   |
|-------------------------|--|--|---|--|---|
| Edit based Similarities | compare two strings by counting the minimum number of operations required to transform one string into the other | <ul style="list-style-type: none"> <li>+ Simplicity</li> <li>+ Good for short strings</li> </ul> | <ul style="list-style-type: none"> <li>- Inefficient for longer strings</li> <li>- Computationally expensive</li> <li>- Doesn't take into account the semantic meaning</li> </ul> | Hamming<br>Levenshtein<br>Damerau-Levenshtein<br>Jaro-Winkler<br>Strcmp95<br>Needleman-Wunsch<br>Gotoh<br>Smith-Waterman<br>MLIPNS | Spelling correction<br>Duplication search<br>DNA analysis<br>Correction systems for OCR<br>Measuring the linguistic distance of the languages |



# Resultados

| Dataset  | Pipeline |           |        |      | OpenALPR |           |        |      |
|----------|----------|-----------|--------|------|----------|-----------|--------|------|
|          | Accuracy | Precision | Recall | LS*  | Accuracy | Precision | Recall | LS*  |
| Normal   | 0.38     | 0.6       | 0.52   | 0.46 | 0.40     | 0.83      | 0.43   | 0.17 |
| Soleado  | 0.55     | 0.78      | 0.65   | 0.34 | 0.44     | 0.79      | 0.50   | 0.21 |
| Lluvioso | 0.16     | 0.35      | 0.23   | 0.36 | 0.34     | 0.75      | 0.39   | 0.18 |
| Ruidoso  | 0.42     | 0.48      | 0.76   | 0.49 | 0.61     | 0.90      | 0.65   | 0.16 |

**Tabla 4.1:** Comparación de resultados. (\*) Promedio de la similitud normalizada de *Levenshtein*.

# Resultados



Predictions: [['castle', 'stirling', 'lodgings', 'argylls']]



Predictions: [['mea', '959']]

# Resultados



# Resultados

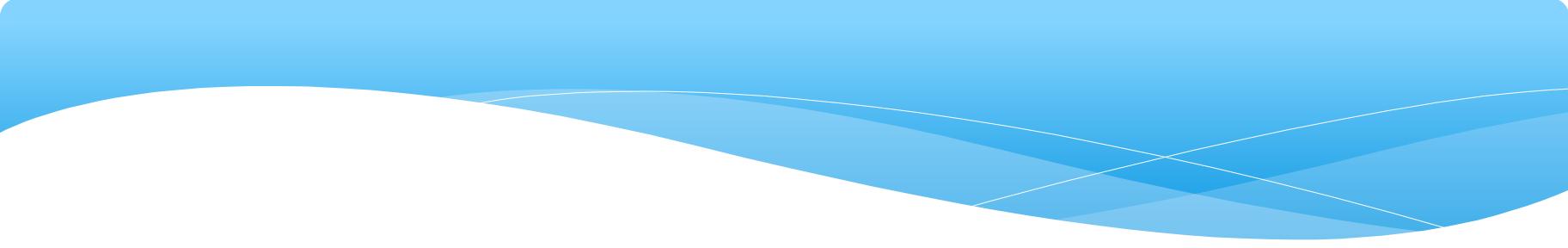
**(VIDEO & OUTPUT)**

# Conclusiones

- Resolvimos el problema **ALPR** con un sistema modular del tipo *detección-detección-reconocimiento* usando deep learning.
- Evaluamos nuestro sistema frente al estandar comercial, **OpenALPR**
- Dependiendo de la iluminación de la escena, uno u otro método entrega mejores resultados.
- Algunos resultados cualitativos atisban una transferencia buena a otras aplicaciones.

# Mejoras

- Re-entrenar YOLOv3 con un dataset de patentes locales para mejorar su desempeño.
- Re-entrenar la RCNN con un pre-filtro convolucional a la entrada, sobre imágenes de sentencias de matriculas de origen natural anotadas.
- Implementar otros mecanismos de atención y entrenar.



**¡GRACIAS!**

....

**¿PREGUNTAS?**