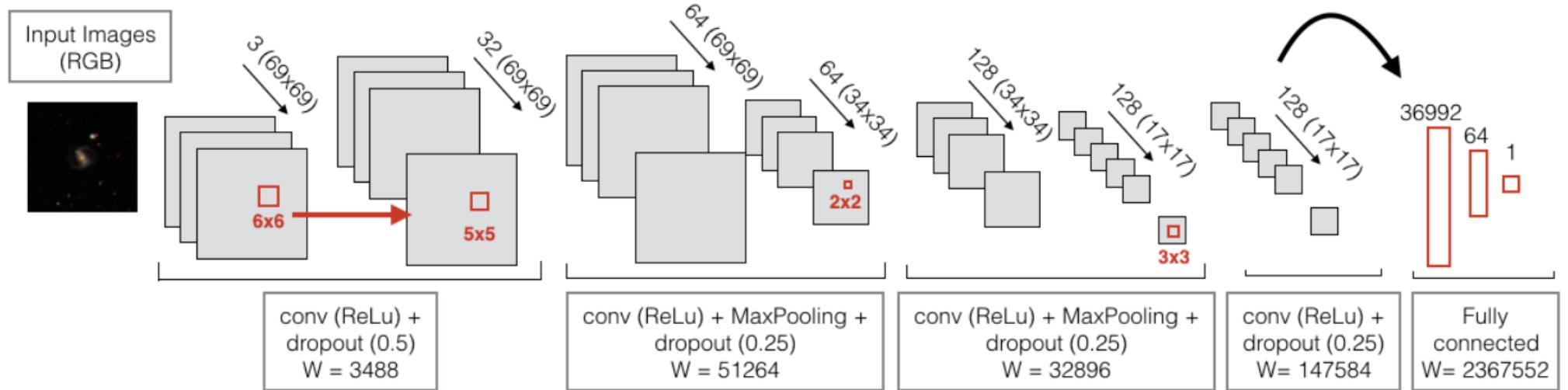


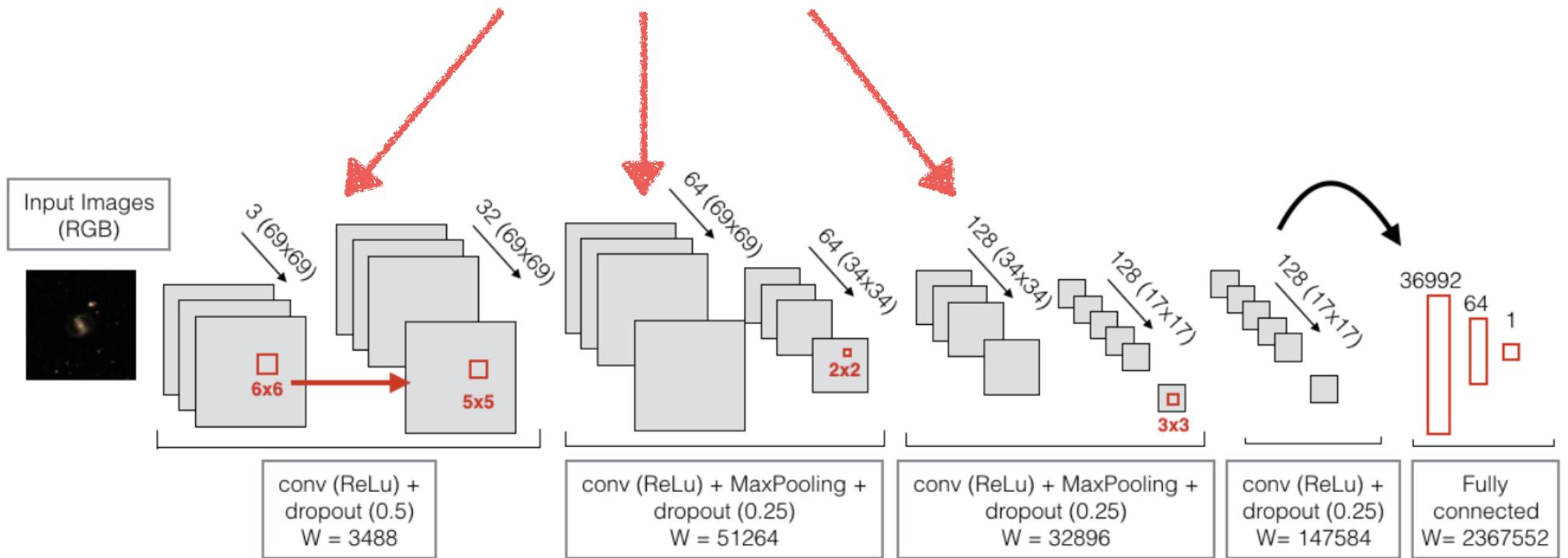
# EXAMPLE OF VERY SIMPLE CNN



Dominguez-Sanchez+18

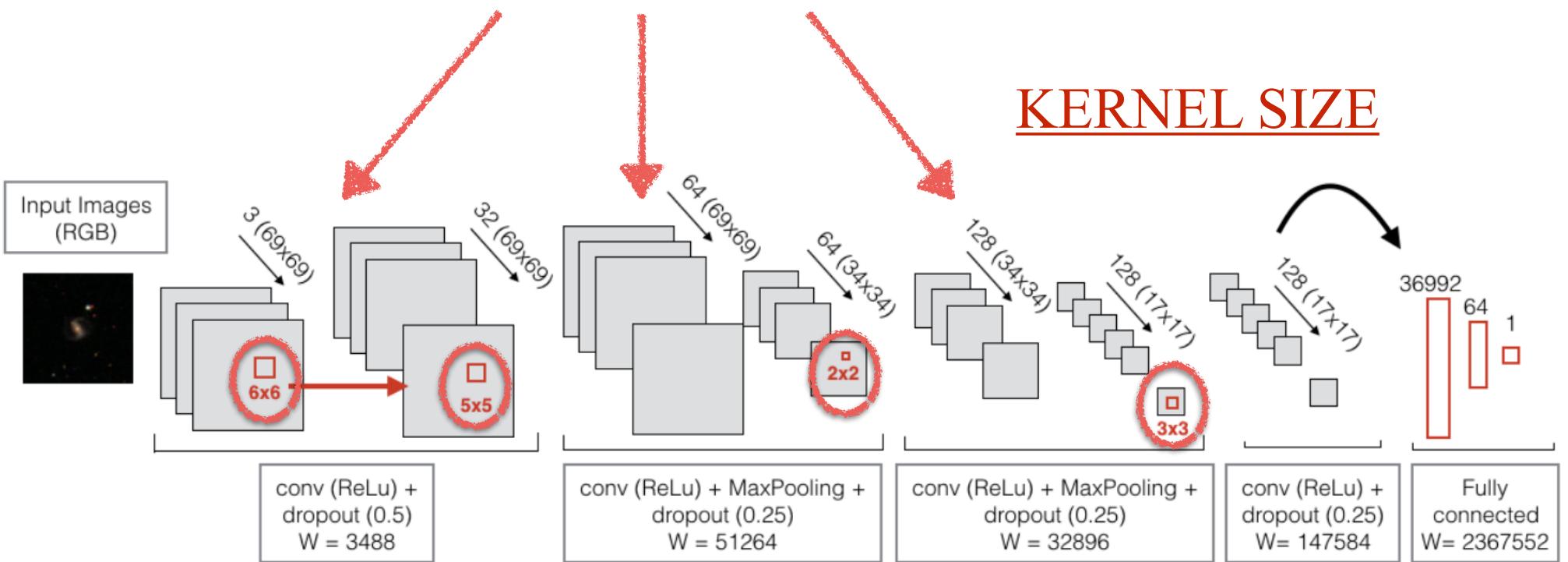
# EXAMPLE OF VERY SIMPLE CNN

## 3 convolutional layers



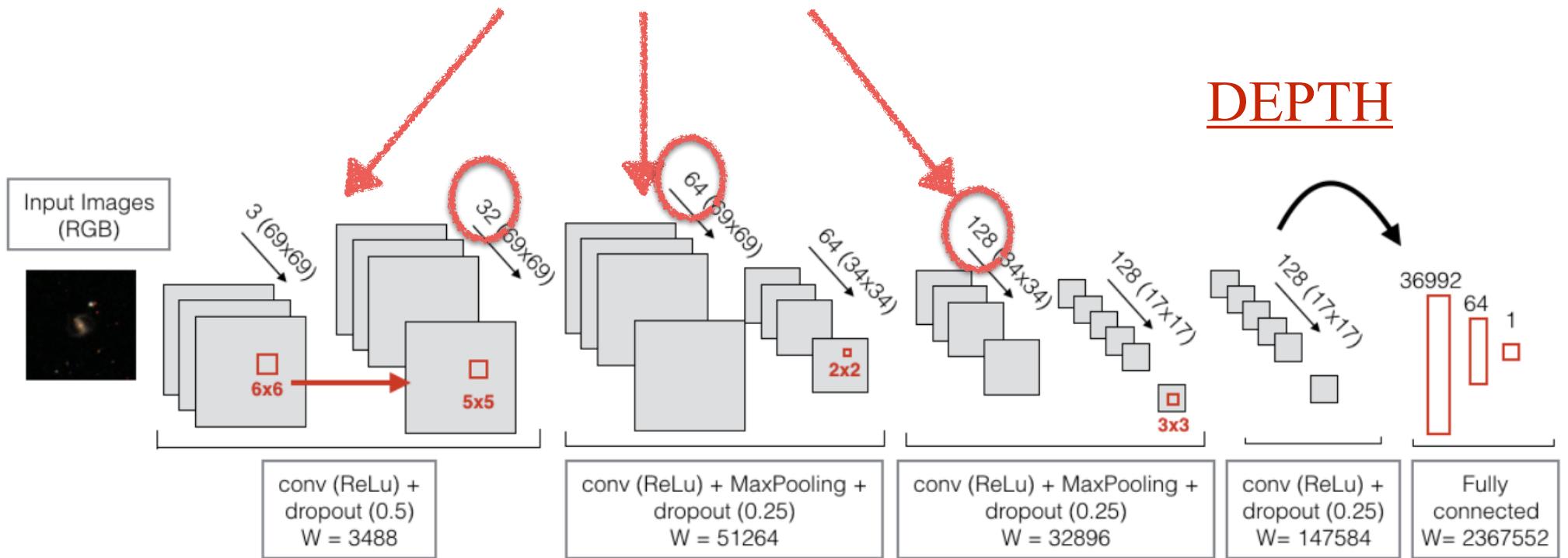
# EXAMPLE OF VERY SIMPLE CNN

3 convolutional layers



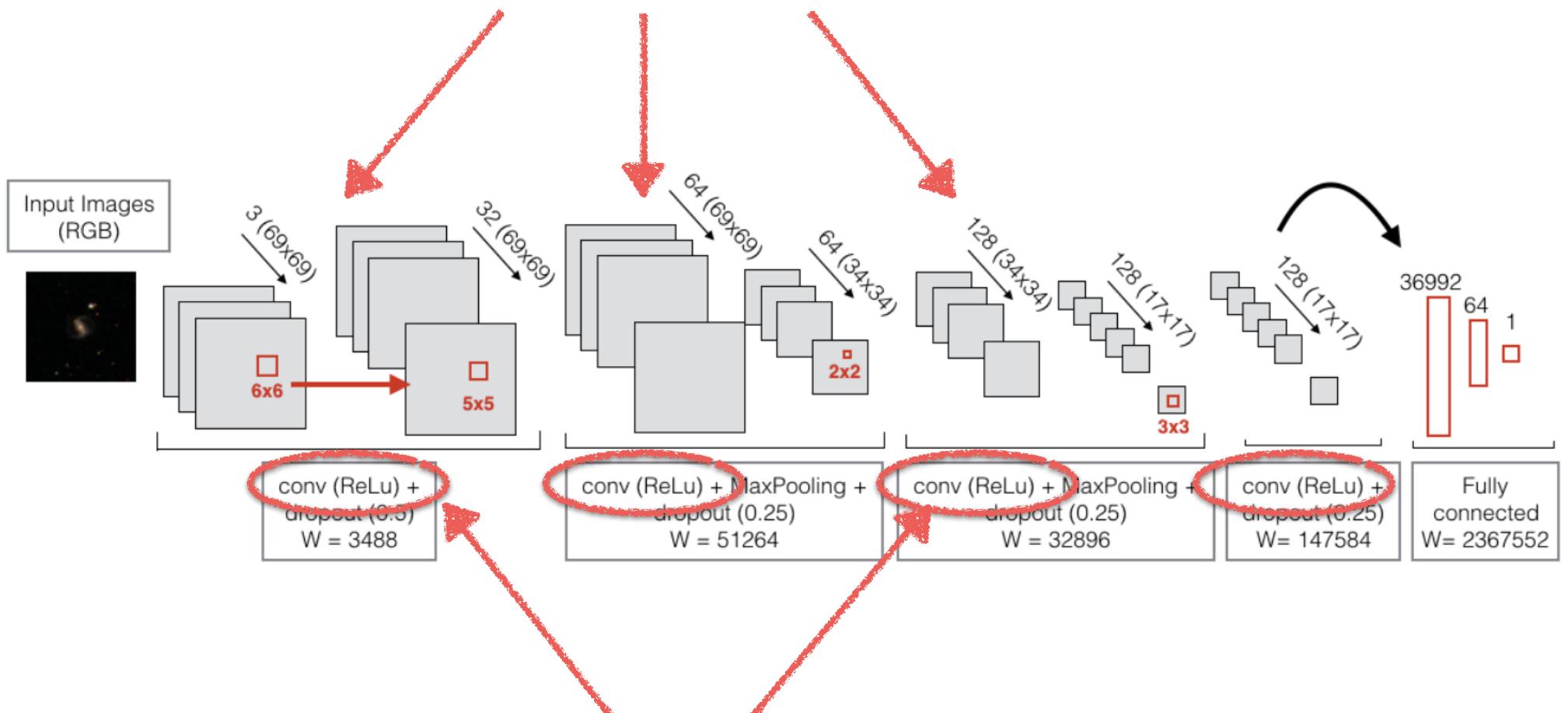
# EXAMPLE OF VERY SIMPLE CNN

3 convolutional layers



# EXAMPLE OF VERY SIMPLE CNN

3 convolutional layers

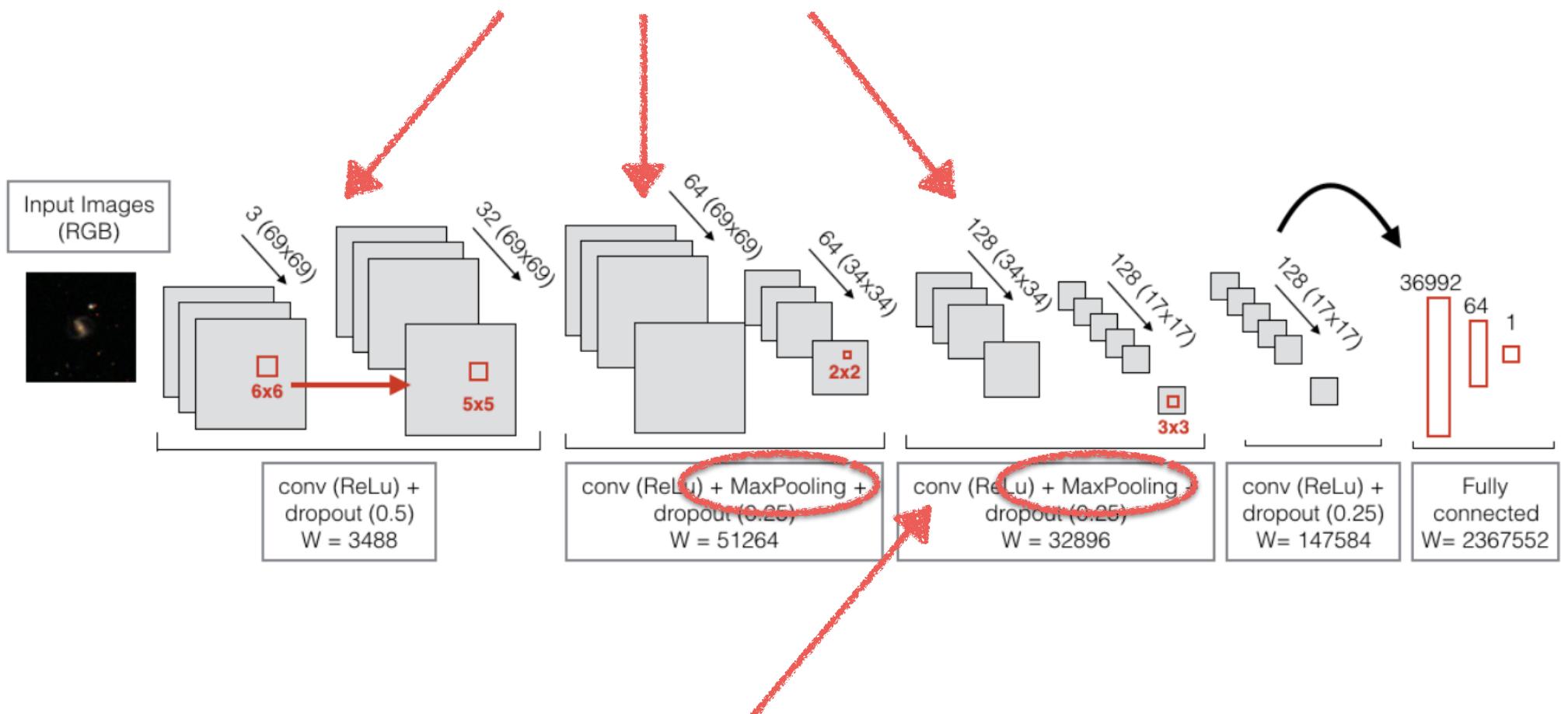


ReLU activation

Dominguez-Sanchez+18

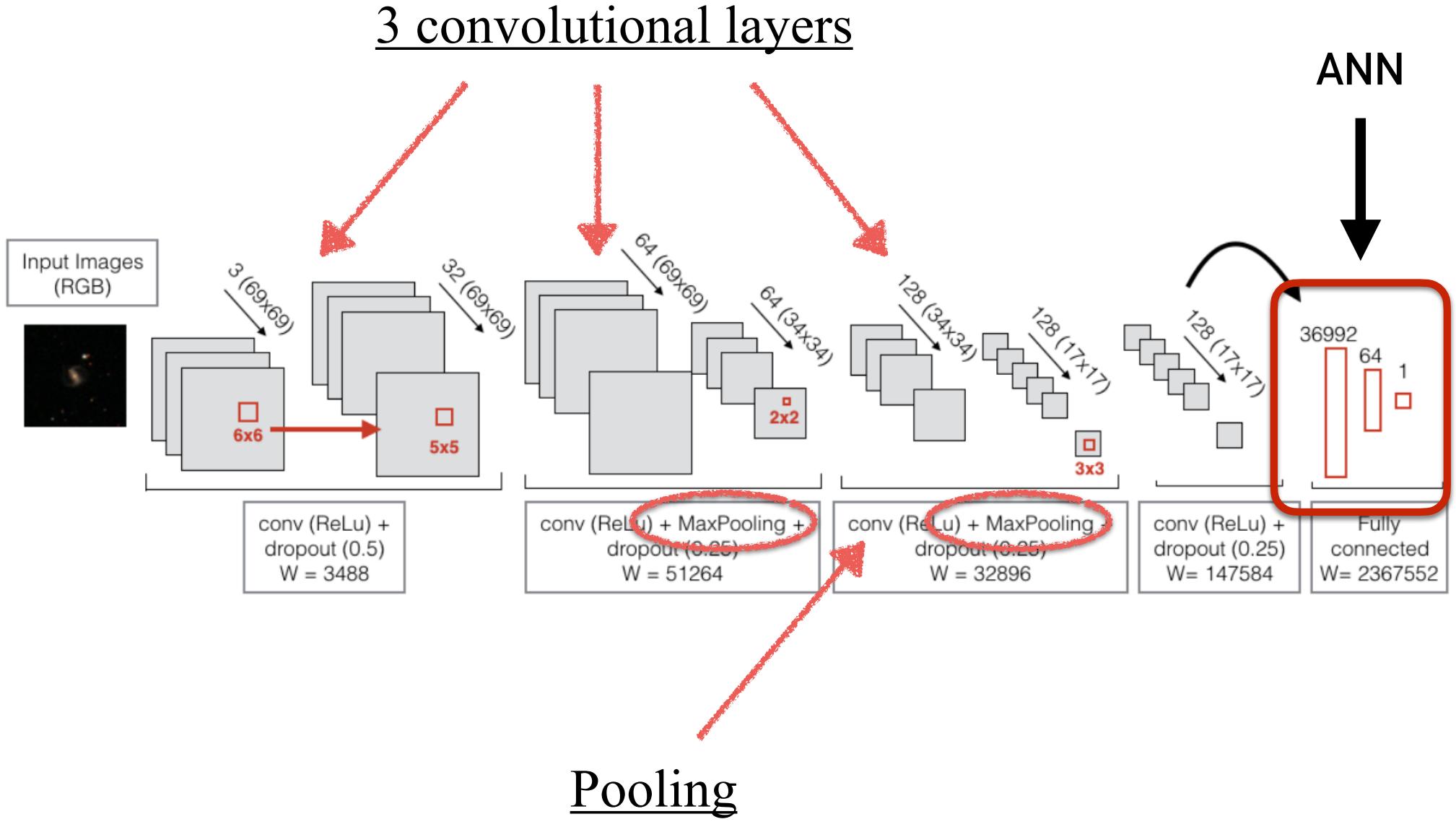
# EXAMPLE OF VERY SIMPLE CNN

3 convolutional layers



Pooling

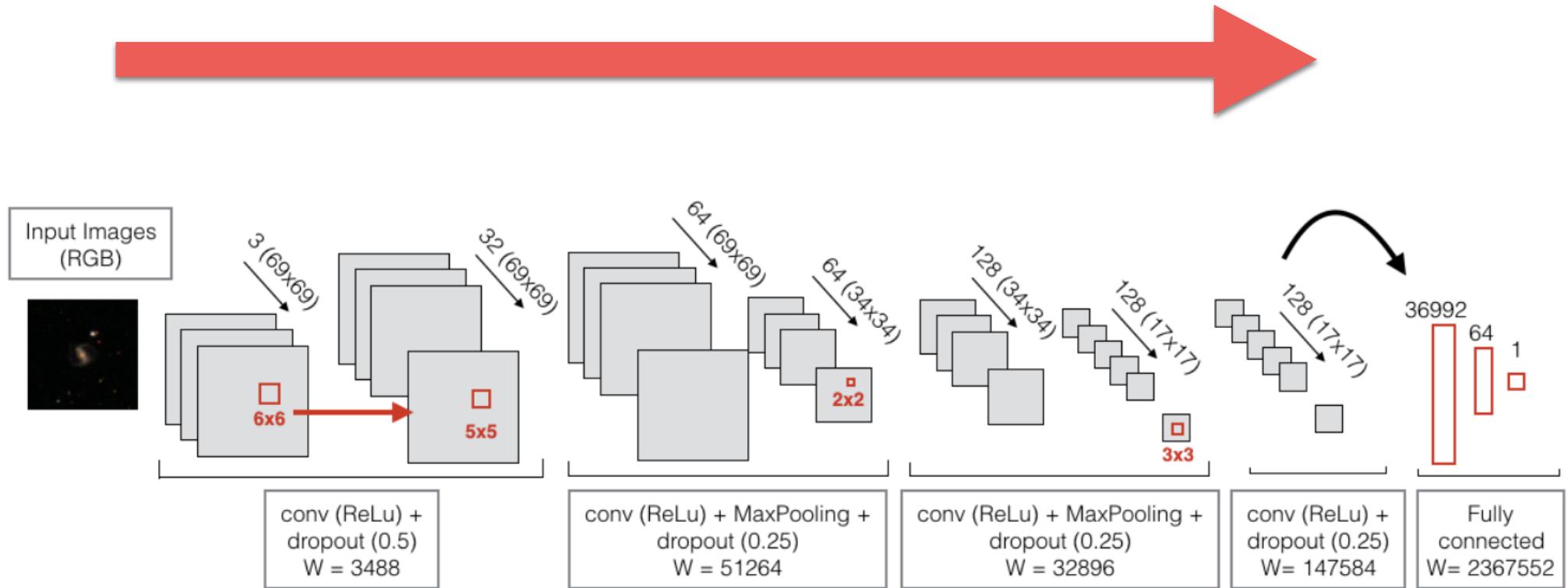
# EXAMPLE OF VERY SIMPLE CNN



# EXAMPLE OF VERY SIMPLE CNN

OVERALL:

- decrease of tensor size
- increase of depth



# IMPLEMENTATION IN KERAS

```
#===== Model definition=====

#Convolutional Layers

model = Sequential()
model.add(Convolution2D(32, 6,6, border_mode='same',
                       input_shape=(img_channels, img_rows, img_cols)))
model.add(Activation('relu'))
model.add(Dropout(0.5))

model.add(Convolution2D(64, 5, 5, border_mode='same'))
model.add(Activation('relu'))

model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Dropout(0.25))

model.add(Convolution2D(128, 2, 2, border_mode='same'))
model.add(Activation('relu'))

model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Dropout(0.25))

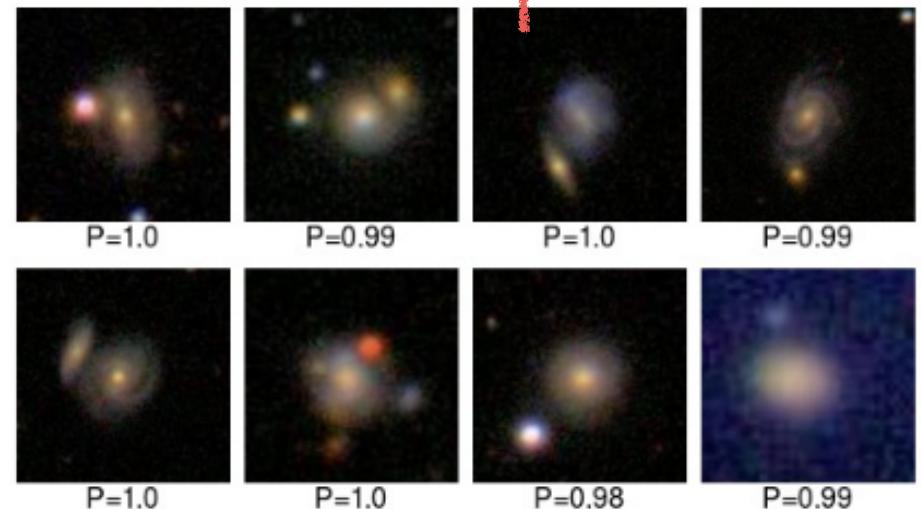
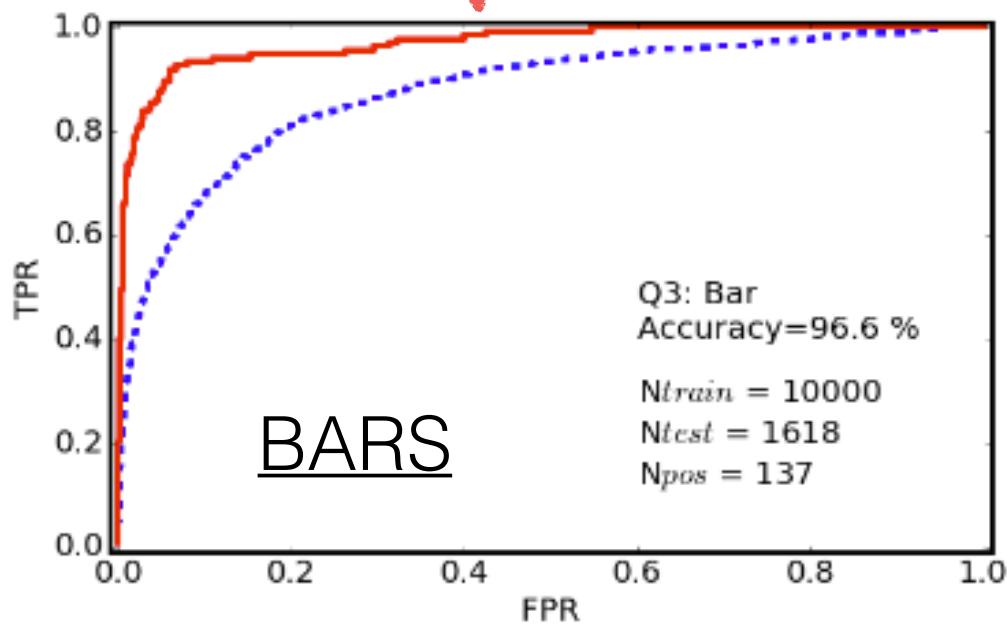
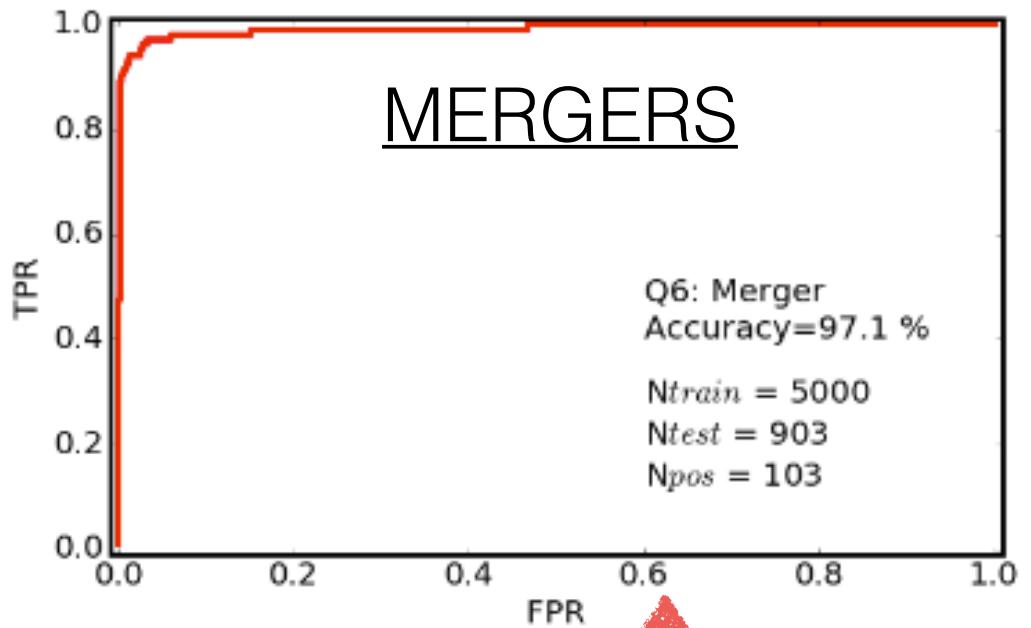
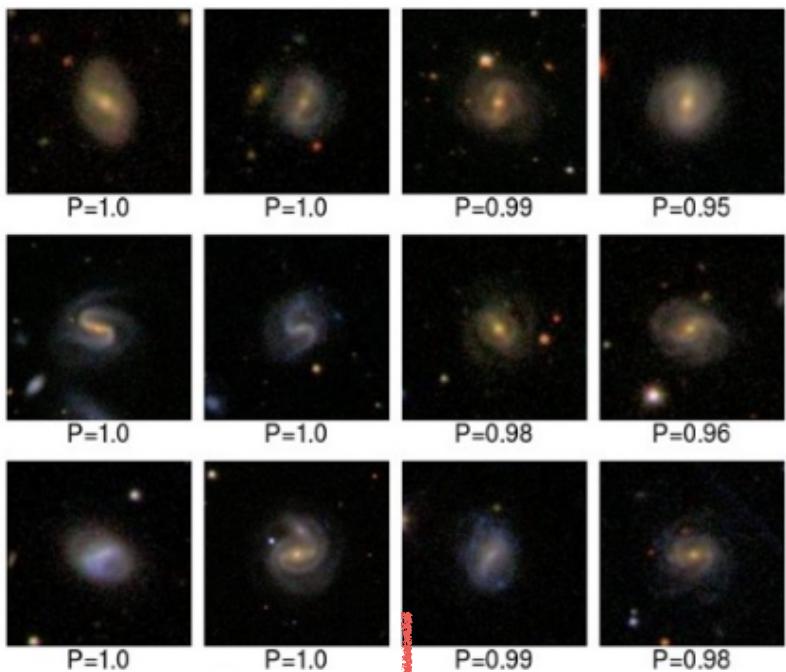
model.add(Convolution2D(128, 3, 3, border_mode='same'))
model.add(Activation('relu'))

model.add(Dropout(0.25))

#Fully Connected start here
#-----#
model.add(Flatten())
model.add(Dense(64, activation='relu'))
model.add(Dropout(.5))
model.add(Dense(1, init='uniform', activation='sigmoid'))

print("Compilation...")

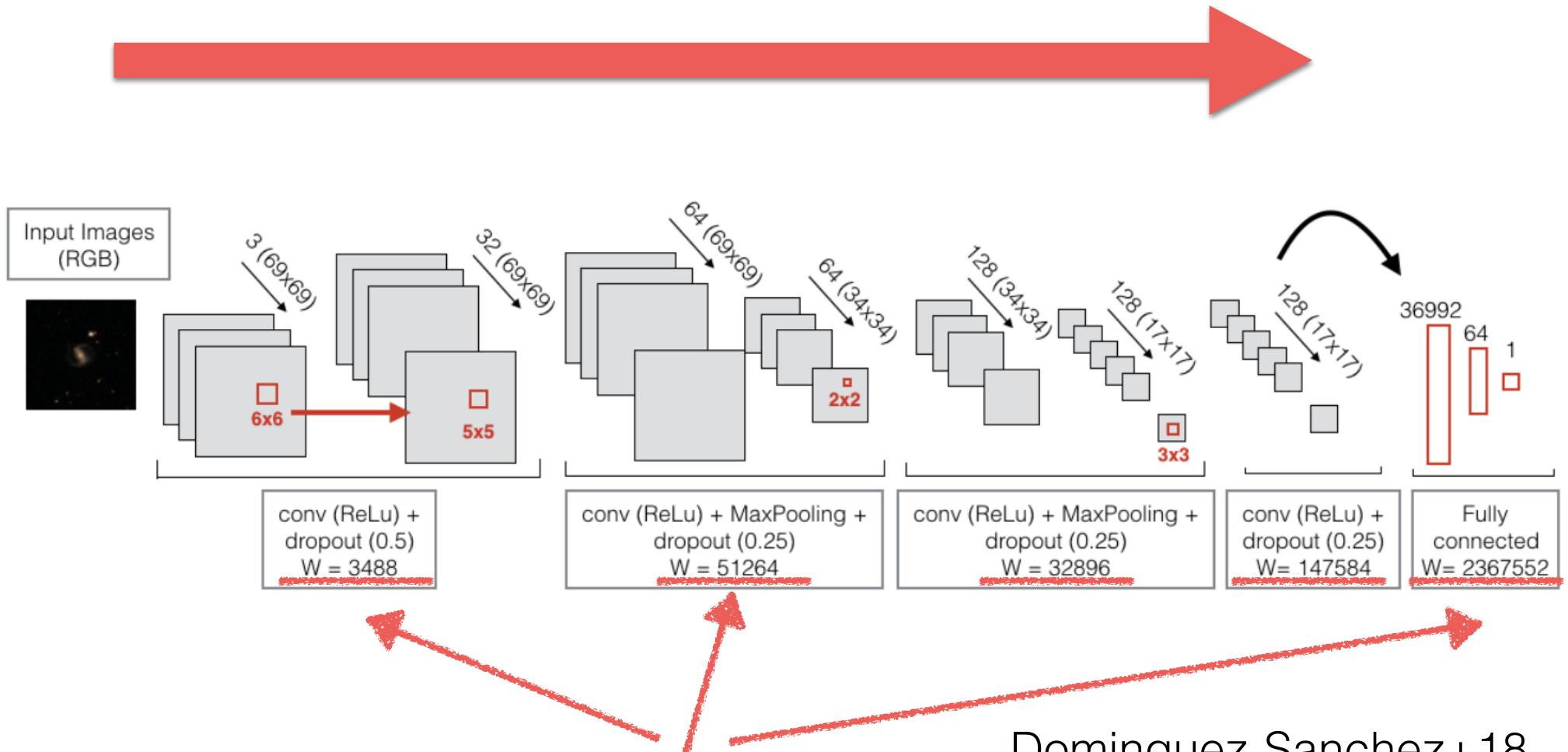
model.compile(loss='binary_crossentropy',optimizer='adam',metrics=['accuracy'])
```



# EXAMPLE OF VERY SIMPLE CNN

OVERALL:

- decrease of tensor size
- increase of depth



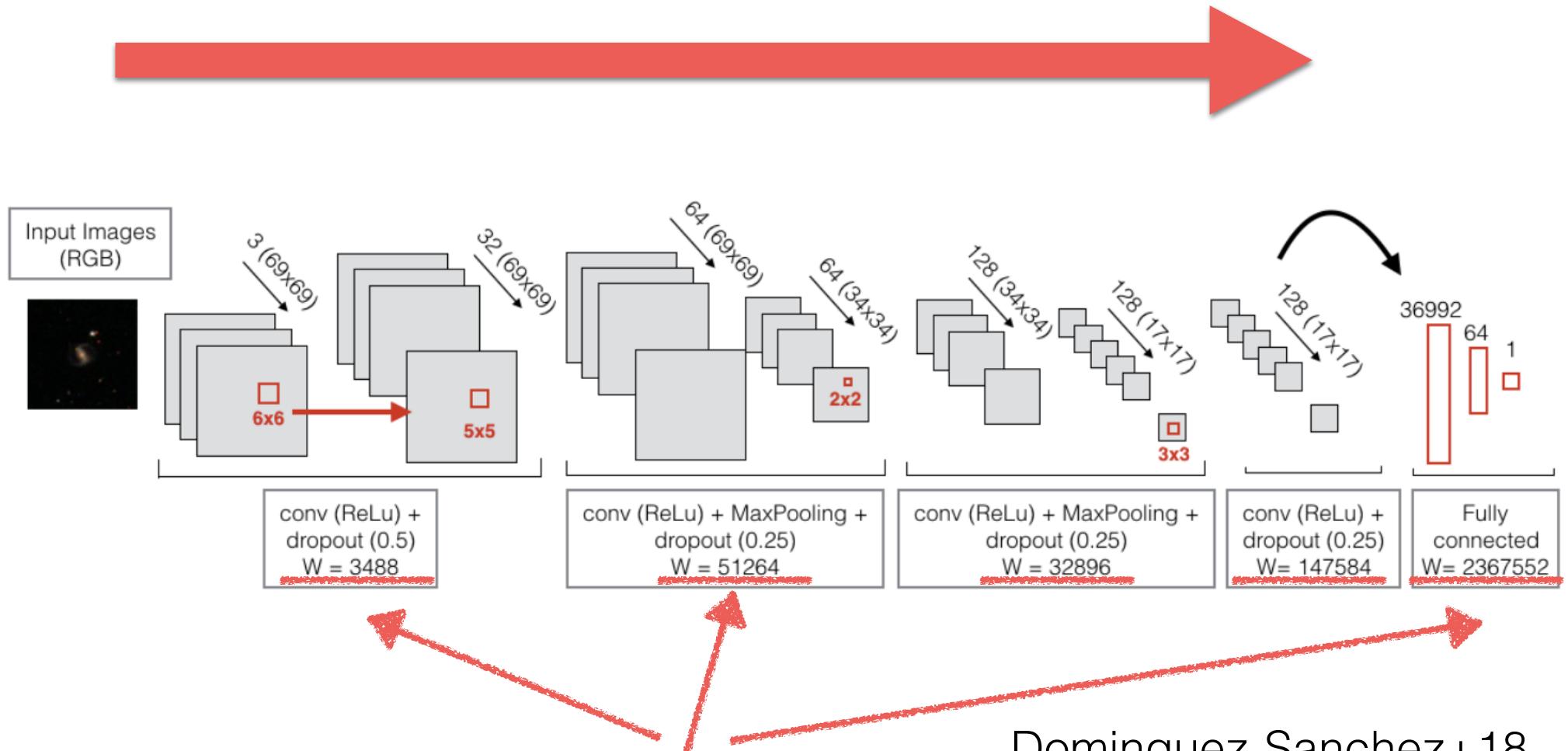
Dominguez-Sanchez+18

Number of parameters

# EXAMPLE OF VERY SIMPLE CNN

OVERALL:

- decrease of tensor size
- increase of depth



2 million of parameters for this very simple network!

# CHECKING THE NUMBER OF PARAMETERS / LAYERS WITH KERAS

model.summary()

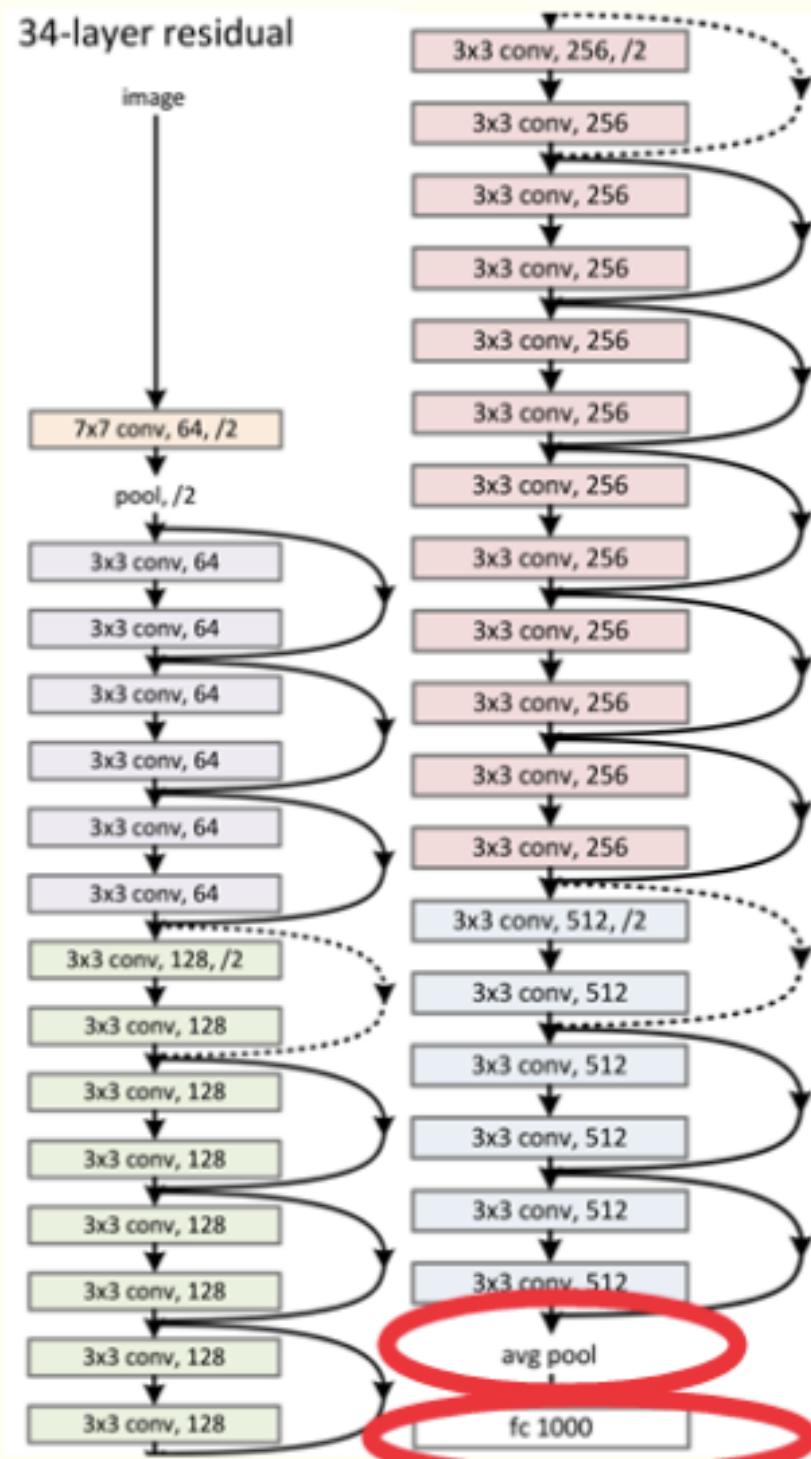


Layer (type)	Output Shape	Param #
=====		
input_1 (InputLayer)	(None, 1, 16, 112, 112)	0
conv3d_1 (Conv3D)	(None, 16, 16, 112, 112)	448
batch_normalization_1 (Batch Normalization)	(None, 16, 16, 112, 112)	448
activation_1 (Activation)	(None, 16, 16, 112, 112)	0
max_pooling3d_1 (MaxPooling3D)	(None, 16, 8, 56, 56)	0
conv3d_2 (Conv3D)	(None, 32, 8, 56, 56)	13856
batch_normalization_2 (Batch Normalization)	(None, 32, 8, 56, 56)	224
activation_2 (Activation)	(None, 32, 8, 56, 56)	0
max_pooling3d_2 (MaxPooling3D)	(None, 32, 4, 28, 28)	0
conv3d_3 (Conv3D)	(None, 64, 4, 28, 28)	55360
batch_normalization_3 (Batch Normalization)	(None, 64, 4, 28, 28)	112
activation_3 (Activation)	(None, 64, 4, 28, 28)	0
max_pooling3d_3 (MaxPooling3D)	(None, 64, 2, 14, 14)	0
activation_12 (Activation)	(None, 64, 2, 14, 14)	0
=====		
Total params: 70,448		
Trainable params: 70,056		
Non-trainable params: 392		

# MODERN CNNs

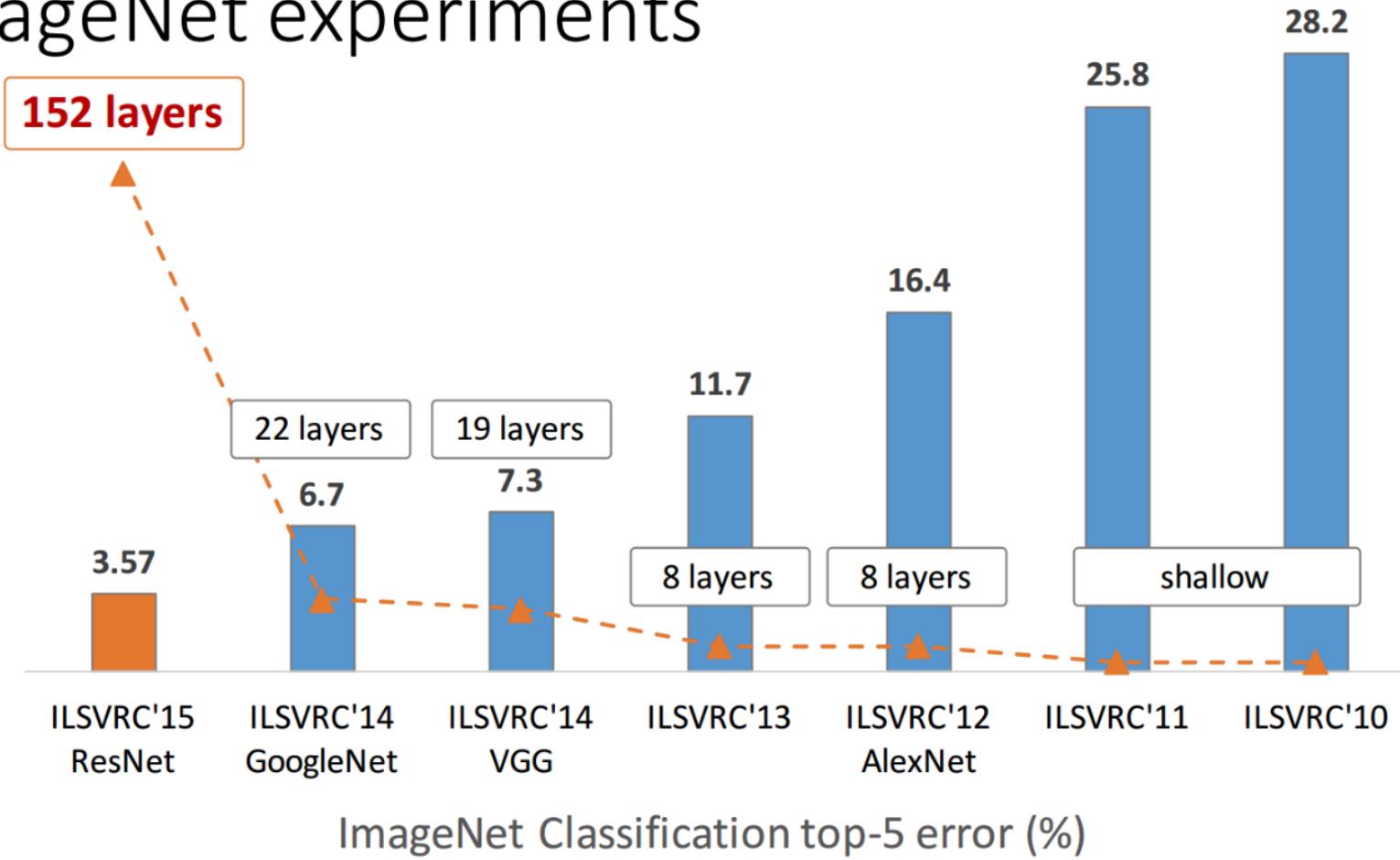
## RESNET

34-layer residual



# DEEPER TENDS TO BE BETTER...

## ImageNet experiments



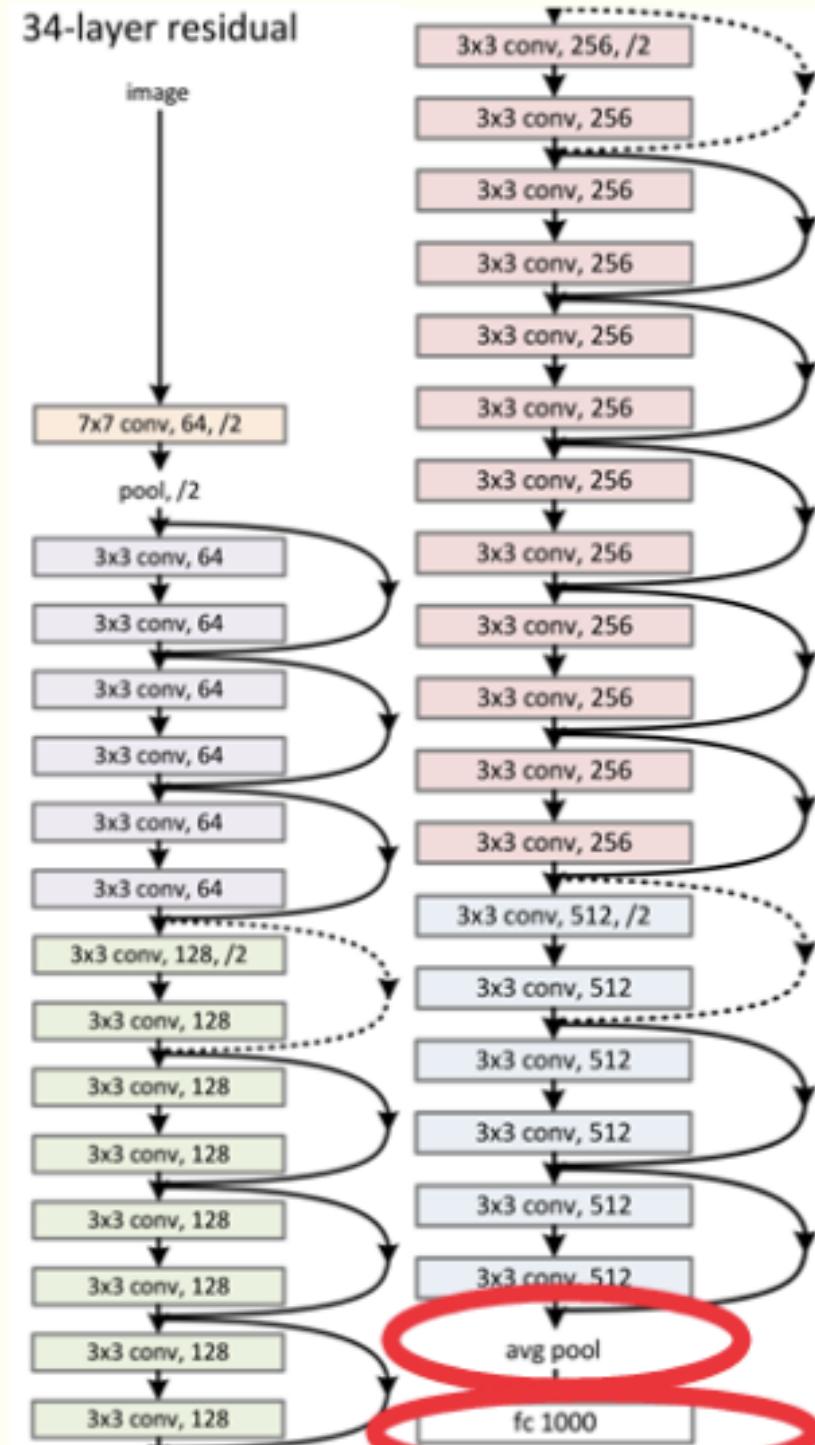
# IN THE REAL LIFE...

## RESNET

DO WE NEED TO GO  
THIS DEEP FOR  
ASTRONOMY  
APPLICATIONS?

[34 layers - authors  
explored up to 1202!]

34-layer residual

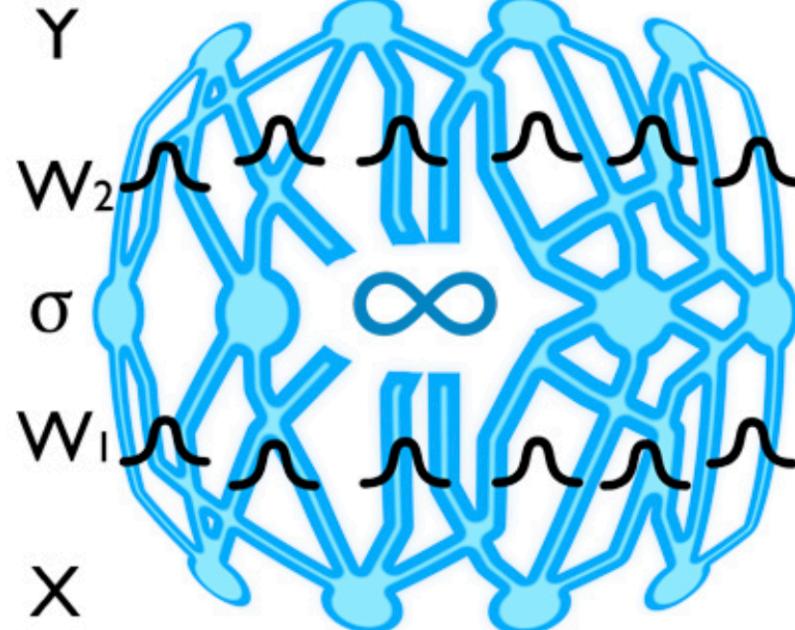


# THE PROBLEMS OF GOING TOO DEEP....

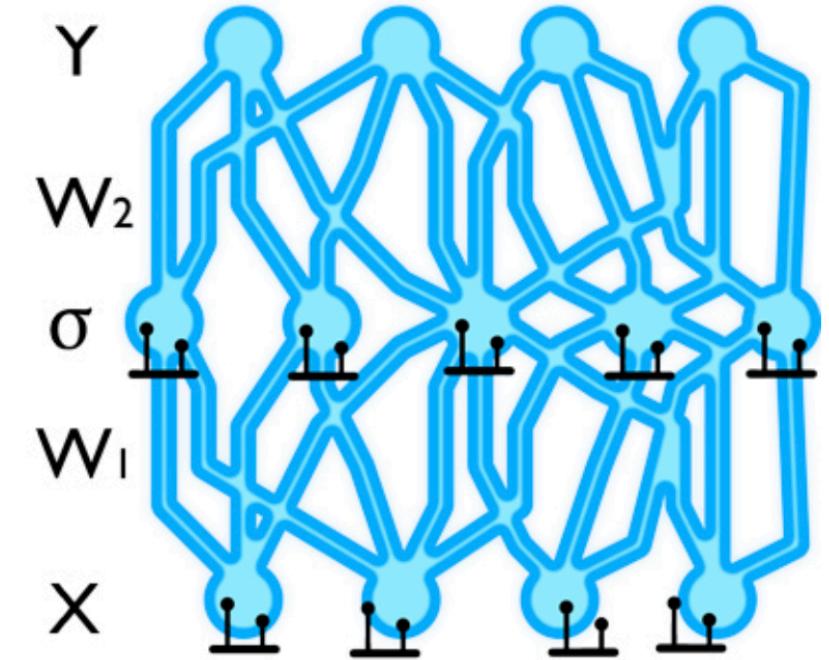
# CAPTURING THE MODEL UNCERTAINTY

## NEURAL NETWORKS AS BAYESIAN MODELS

Denker&Lecun91, Neal+95, Graves+11, Kingma+15, Gal+15...



BNNs ADD A PRIOR DISTRIBUTION TO EACH WEIGHT - HARD TO TRAIN



GAL+15 SHOW THAT DROPOUT CAN BE USED TO ESTIMATE UNCERTAINTY

# IMPLEMENTATION IN KERAS / TENSORFLOW

```
#===== Model definition=====

#Convolutional Layers

model = Sequential()
model.add(Convolution2D(32, 6,6, border_mode='same',
                       input_shape=(img_channels, img_rows, img_cols)))
model.add(Activation('relu'))
model.add(Dropout(0.5))

model.add(Convolution2D(64, 5, 5, border_mode='same'))
model.add(Activation('relu'))

model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Dropout(0.25))

model.add(Convolution2D(128, 2, 2, border_mode='same'))
model.add(Activation('relu'))

model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Dropout(0.25))

model.add(Convolution2D(128, 3, 3, border_mode='same'))
model.add(Activation('relu'))

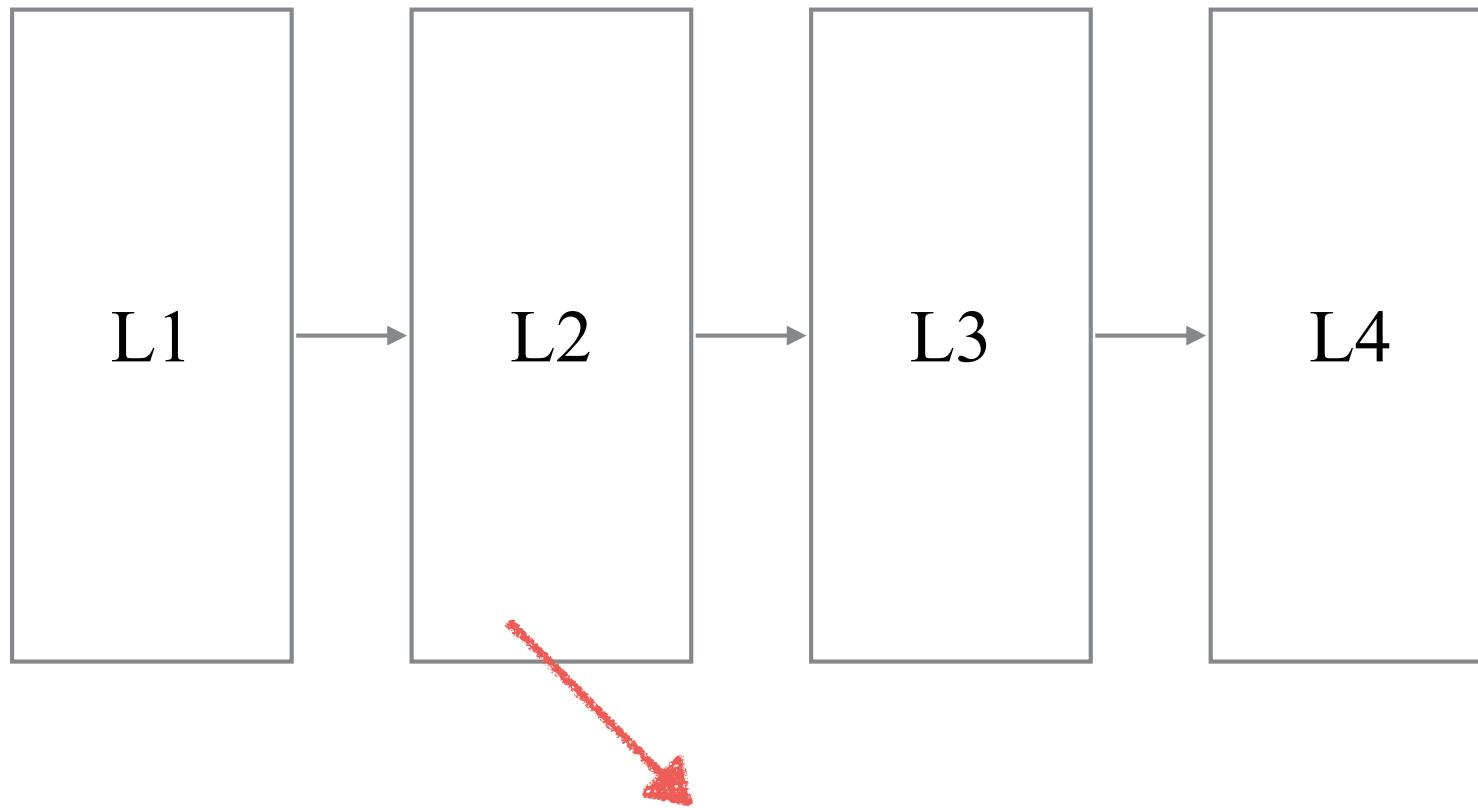
model.add(Dropout(0.25))

#Fully Connected start here
#-----#
model.add(Flatten())
model.add(Dense(64, activation='relu'))
model.add(Dropout(.5))
model.add(Dense(1, init='uniform', activation='sigmoid'))

print("Compilation...")

model.compile(loss='binary_crossentropy',optimizer='adam',metrics=['accuracy'])
```

# CONVNET OR CNN



EACH BLOCK TYPICALLY MADE OF:

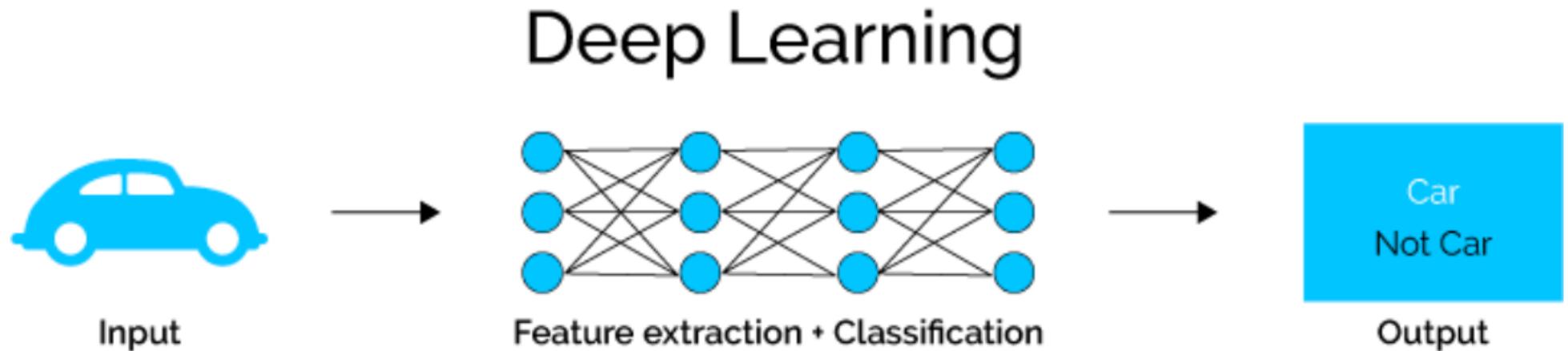
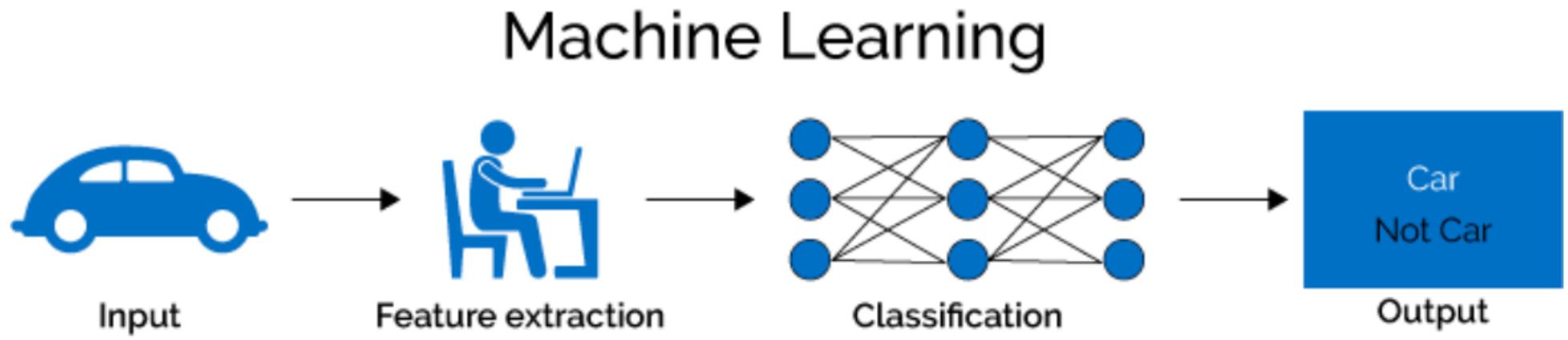
CONV

ACTIVATION

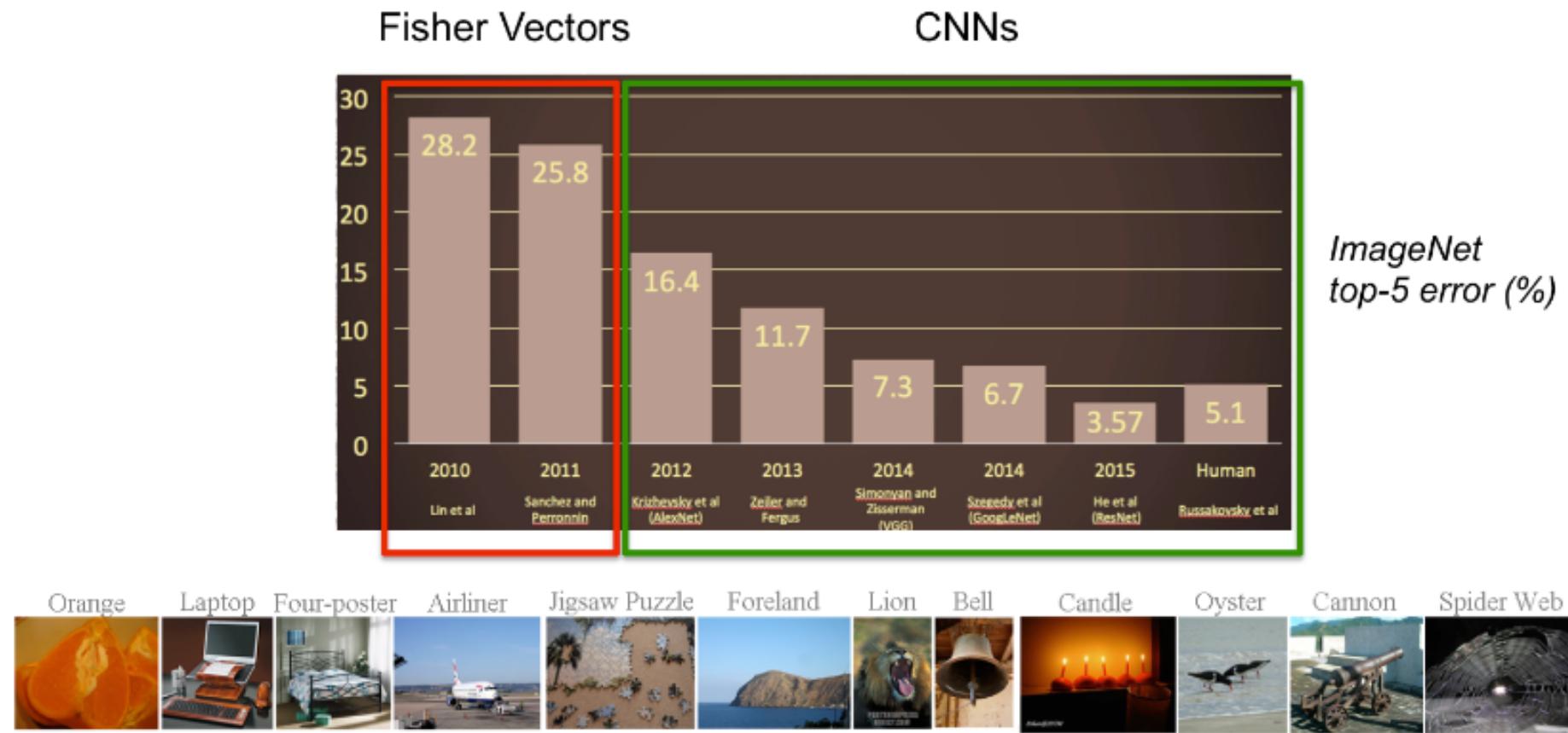
POOLING

(+dropout  
for training)

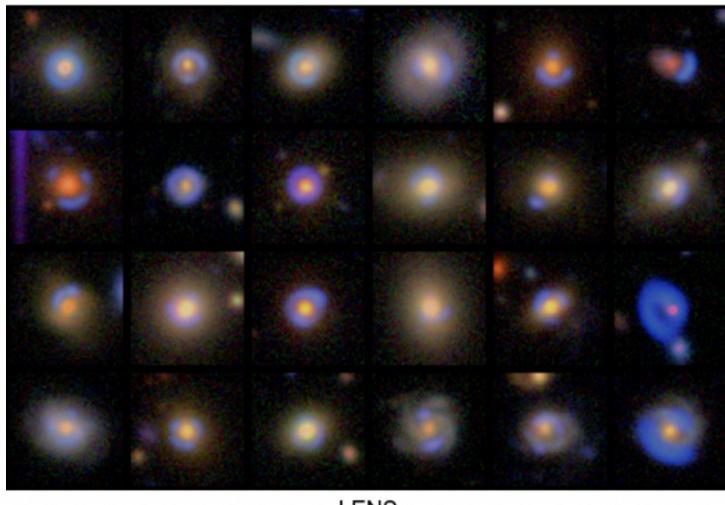
# THIS IS A CHANGE OF PARADIGM!



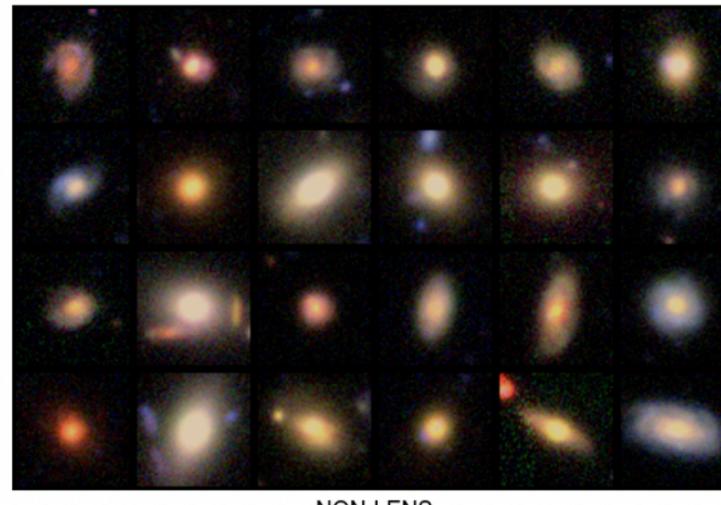
# THIS IS A CHANGE OF PARADIGM!



# 1. Classification



LENS

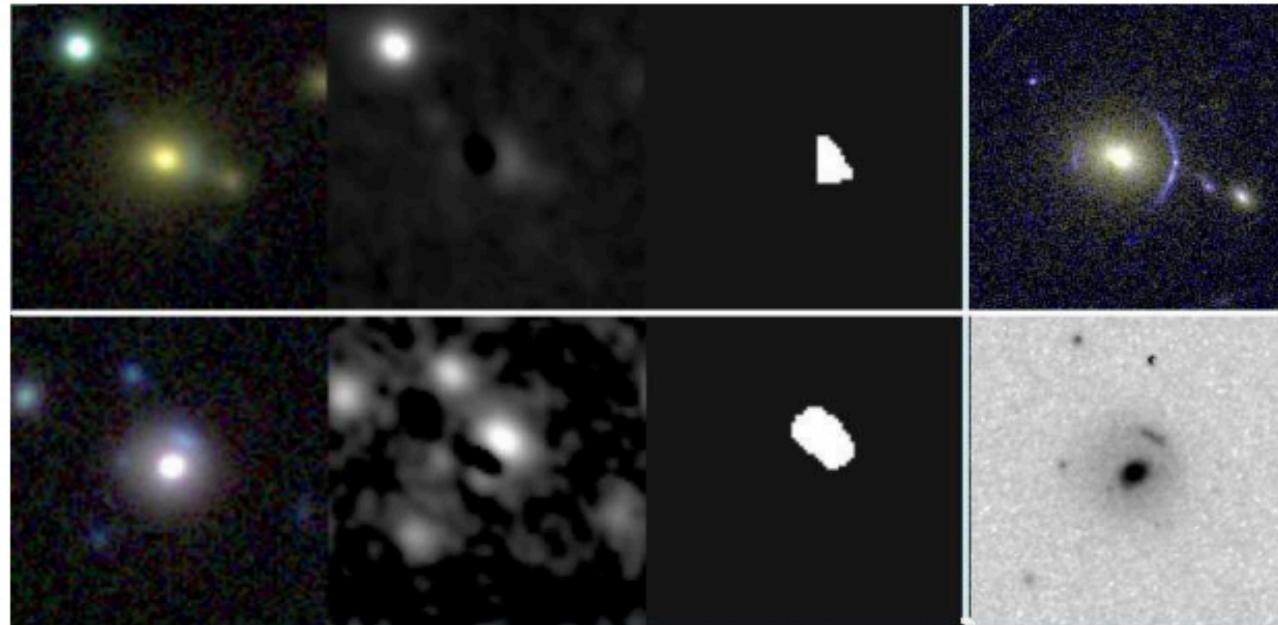


NON-LENS

**Detection of Strong Lenses  
Valuable information of  
Dark Matter properties**

**Future surveys will  
increase the samples by  
orders of magnitude.**

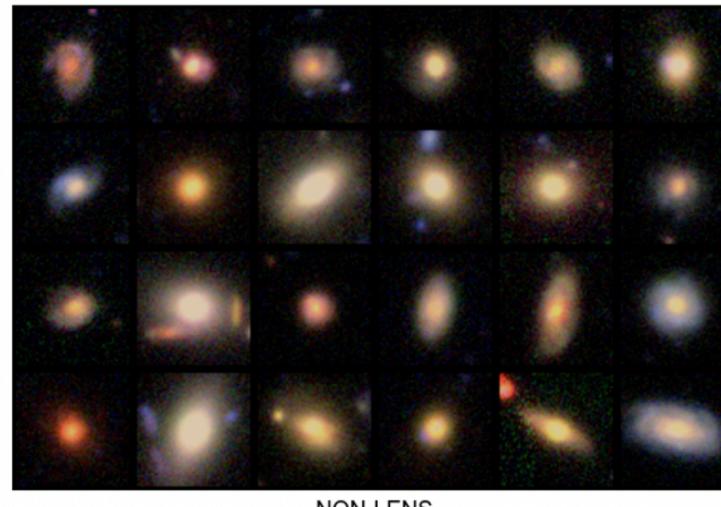
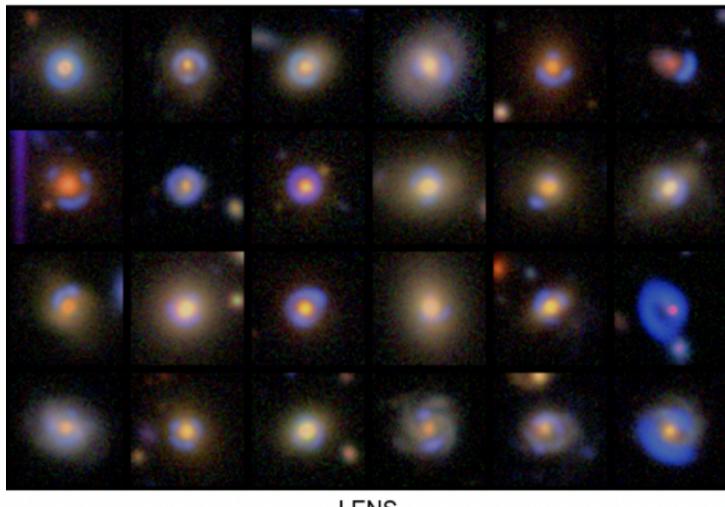
**Jacobs+17**



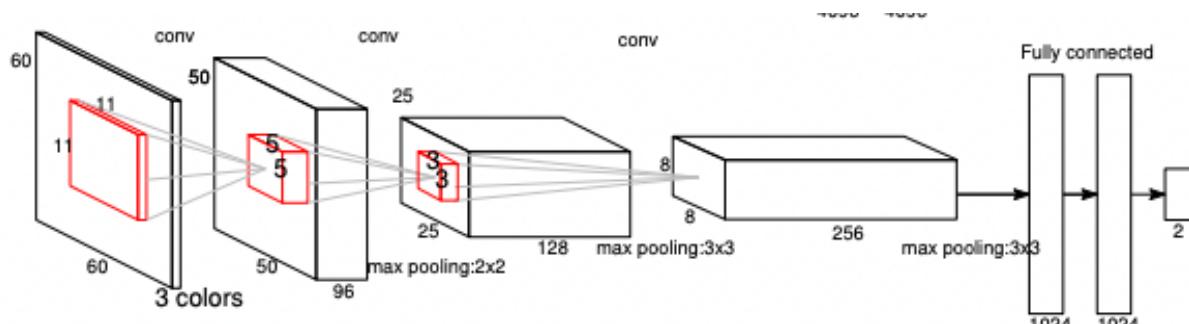
"Pre Deep Learning"  
Approach

**Gavazzi+17**

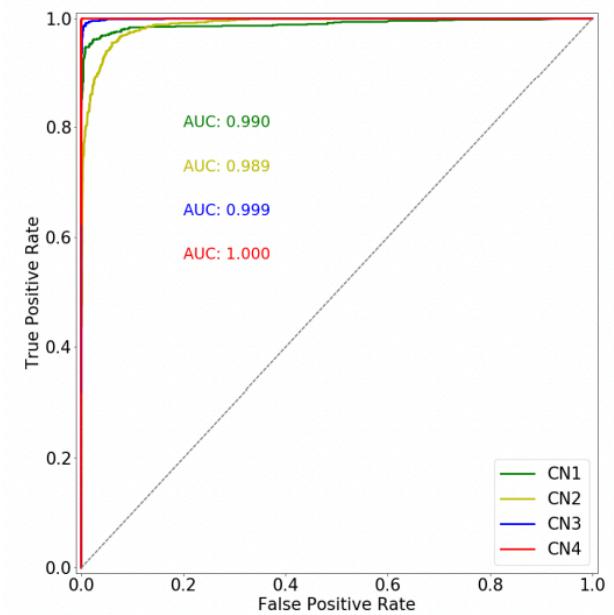
# 1. Classification



**Jacobs+17**



**It illustrates the change of paradigm from an algorithmic centric focus to a purely data driven approach to data**



# 1. Classification

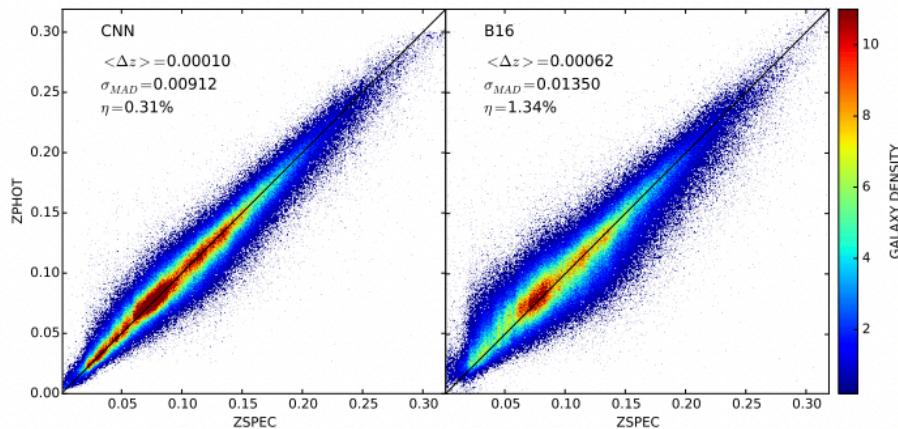
**CNN based classifications reach unprecedented accuracy**

Name	type	AUROC	TPR <sub>0</sub>	TPR <sub>10</sub>	short description
CMU-DeepLens-Resnet-ground3	Ground-Based	0.98	0.09	0.45	CNN
CMU-DeepLens-Resnet-Voting	Ground-Based	0.98	0.02	0.10	CNN
LASTRO EPFL	Ground-Based	0.97	0.07	0.11	CNN
CAS Swinburne Melb	Ground-Based	0.96	0.02	0.08	CNN
AstrOmatic	Ground-Based	0.96	0.00	0.01	CNN
Manchester SVM	Ground-Based	0.93	0.22	0.35	SVM / Gabor
Manchester2	Ground-Based	0.89	0.00	0.01	Human Inspection
ALL-star	Ground-Based	0.84	0.01	0.02	edges/gradiants and Logistic Reg.
CAST	Ground-Based	0.83	0.00	0.00	CNN / SVM
YattaLensLite	Ground-Based	0.82	0.00	0.00	SExtractor
LASTRO EPFL	Space-Based	0.93	0.00	0.08	CNN
CMU-DeepLens-Resnet	Space-Based	0.92	0.22	0.29	CNN
GAMOCLASS	Space-Based	0.92	0.07	0.36	CNN
CMU-DeepLens-Resnet-Voting	Space-Based	0.91	0.00	0.01	CNN
AstrOmatic	Space-Based	0.91	0.00	0.01	CNN
CMU-DeepLens-Resnet-aug	Space-Based	0.91	0.00	0.00	CNN
Kapteyn Resnet	Space-Based	0.82	0.00	0.00	CNN
CAST	Space-Based	0.81	0.07	0.12	CNN
Manchester1	Space-Based	0.81	0.01	0.17	Human Inspection
Manchester SVM	Space-Based	0.81	0.03	0.08	SVM / Gabor
NeuralNet2	Space-Based	0.76	0.00	0.00	CNN / wavelets
YattaLensLite	Space-Based	0.76	0.00	0.00	Arcs / SExtractor
All-now	Space-Based	0.73	0.05	0.07	edges/gradiants and Logistic Reg.
GAHEC IRAP	Space-Based	0.66	0.00	0.01	arc finder

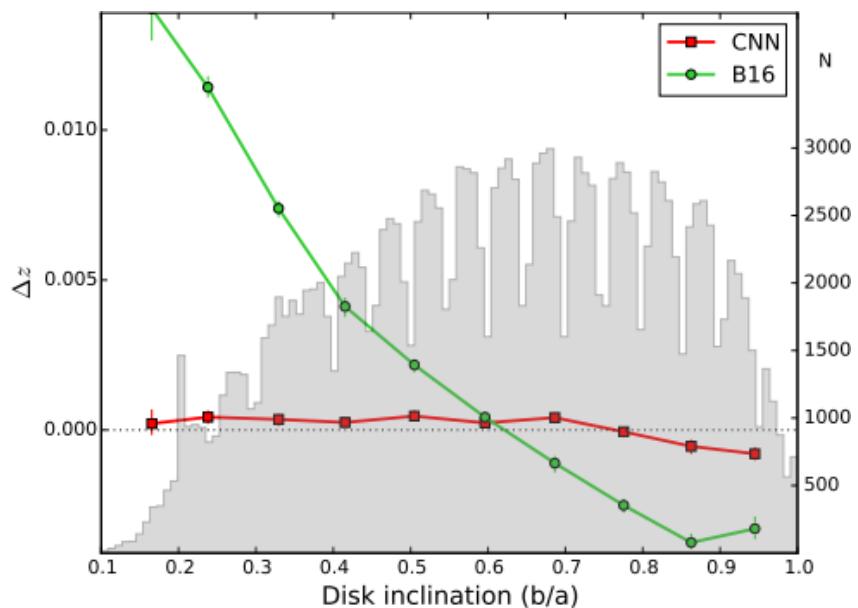
**Metcalf+19**

# Photometric Redshifts

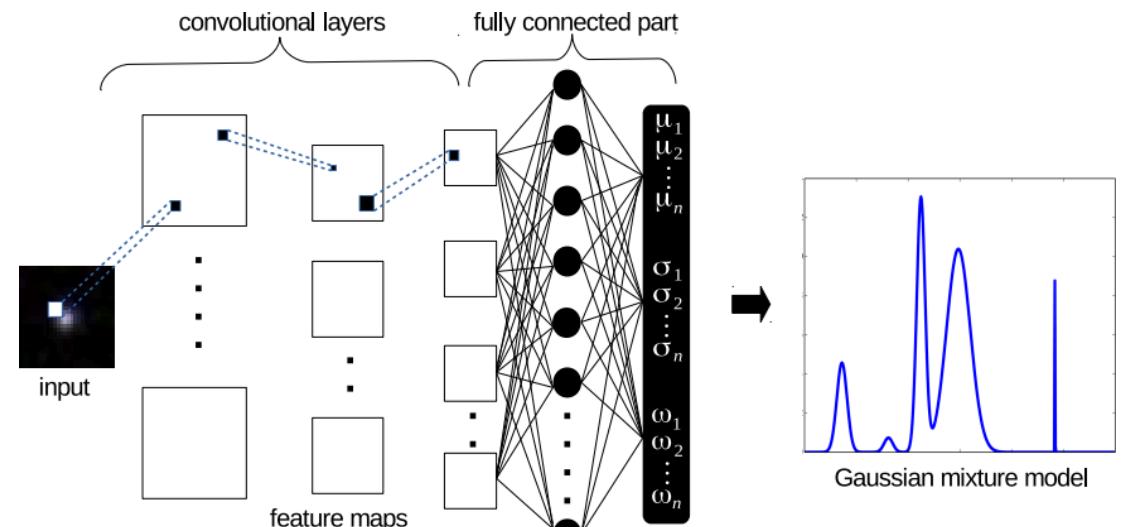
## Deep Learning    Classical approach



Geometric Effects are automatically considered  
(beyond photometry)



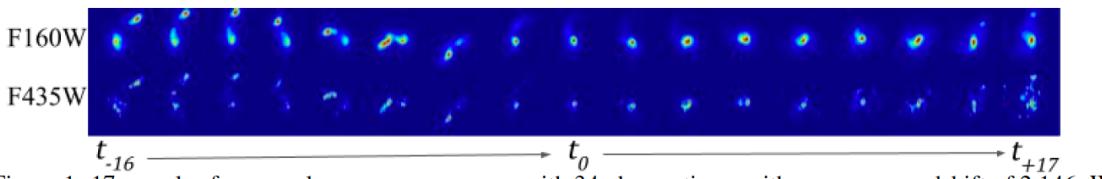
Pasquet+18



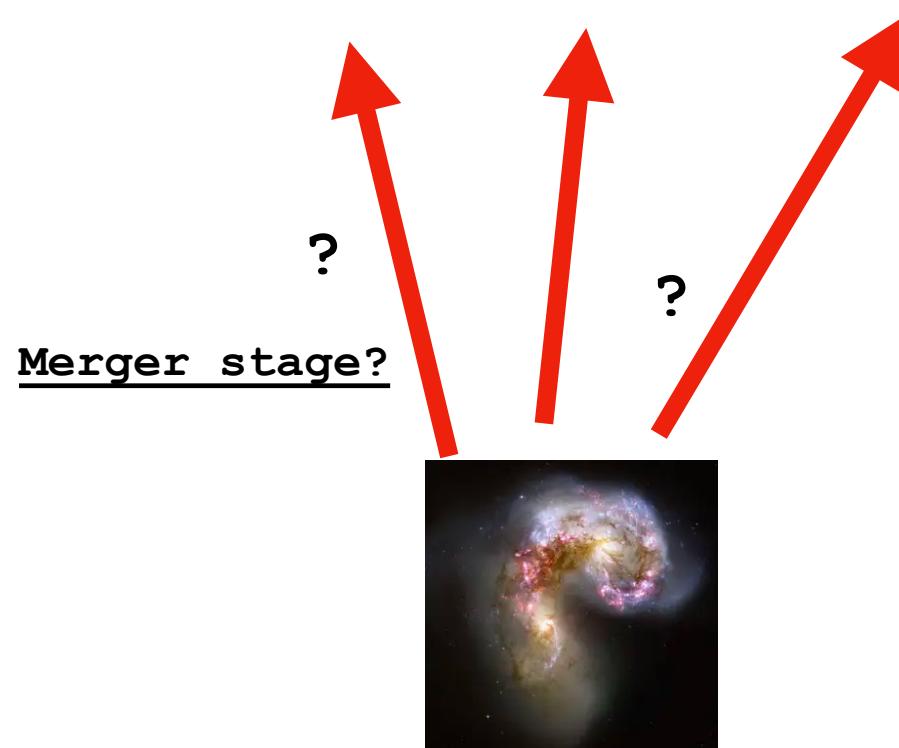
Disanto+18

Uncertainty quantification through Mixture Density Networks

# Mergers of Galaxies

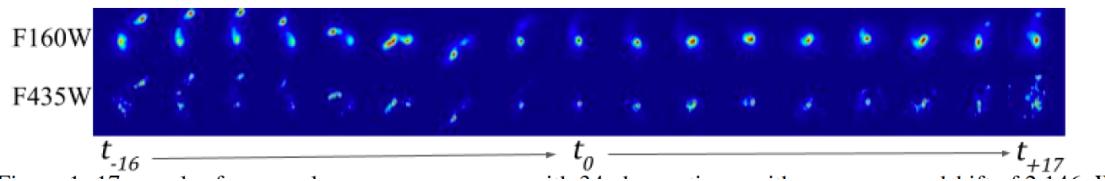


Merger of galaxies sequence  
from cosmological simulations

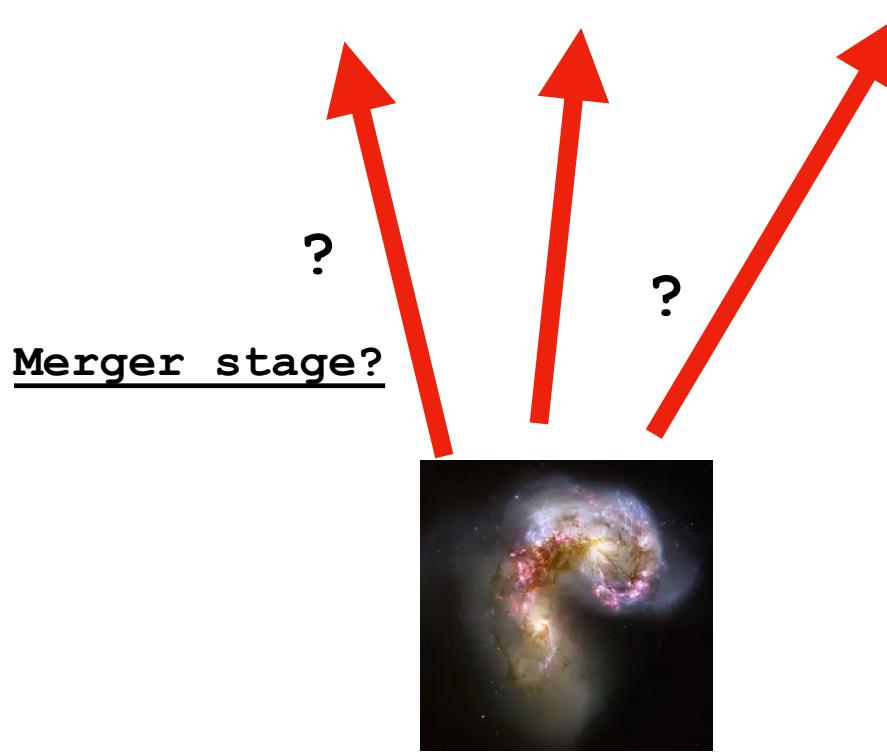


**Neural Networks to find  
relations between observables  
and physical processes**

# Mergers of Galaxies

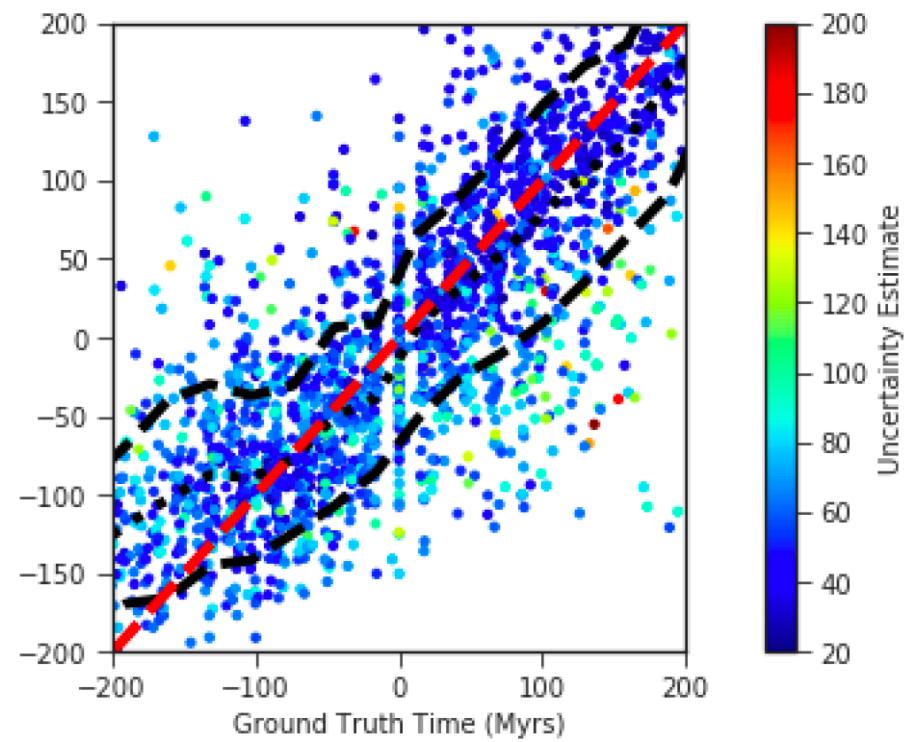


Merger of galaxies sequence  
from cosmological simulations



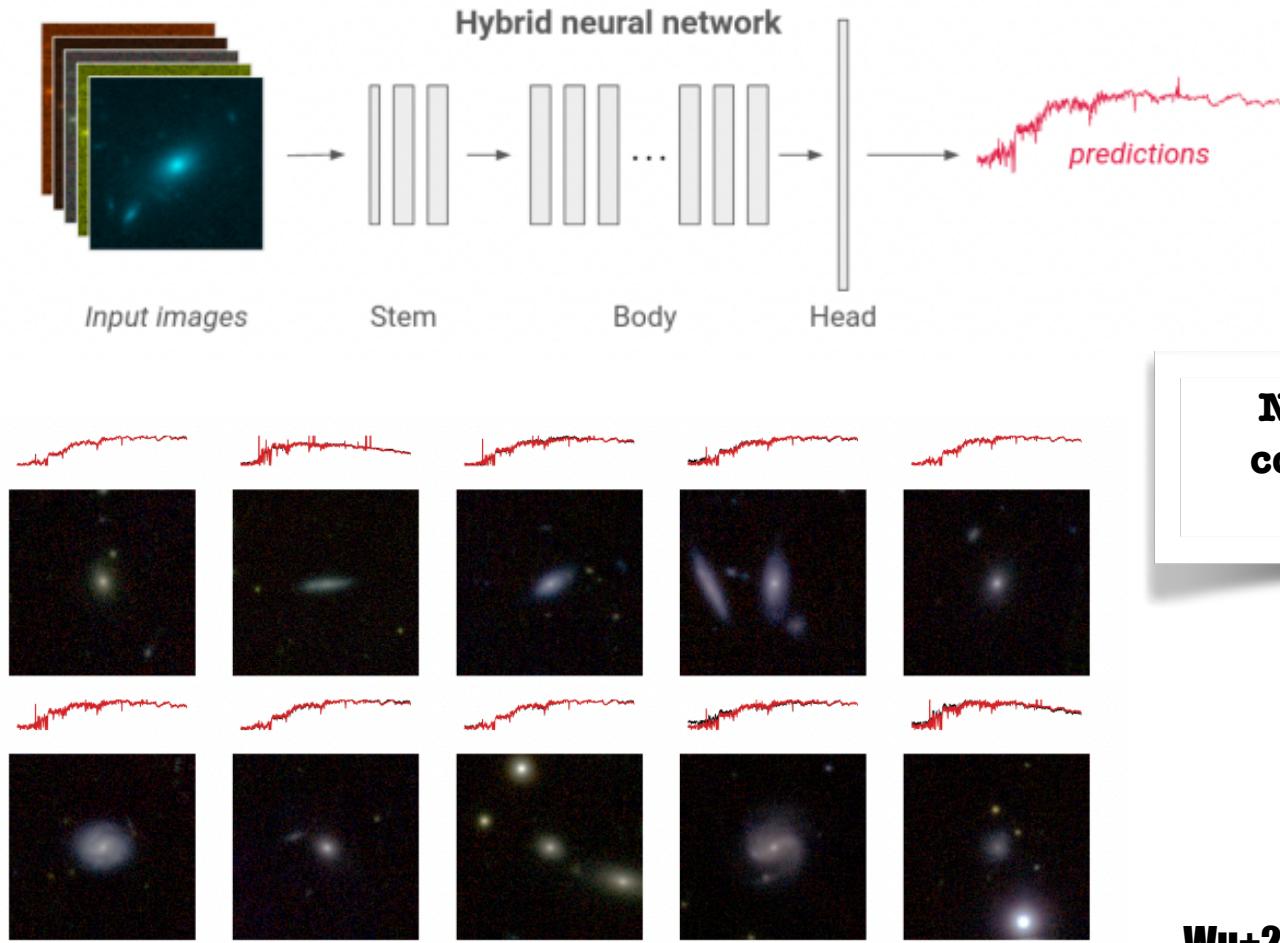
Merger stage?

Neural Networks to find  
relations between observables  
and physical processes



Koppula+21

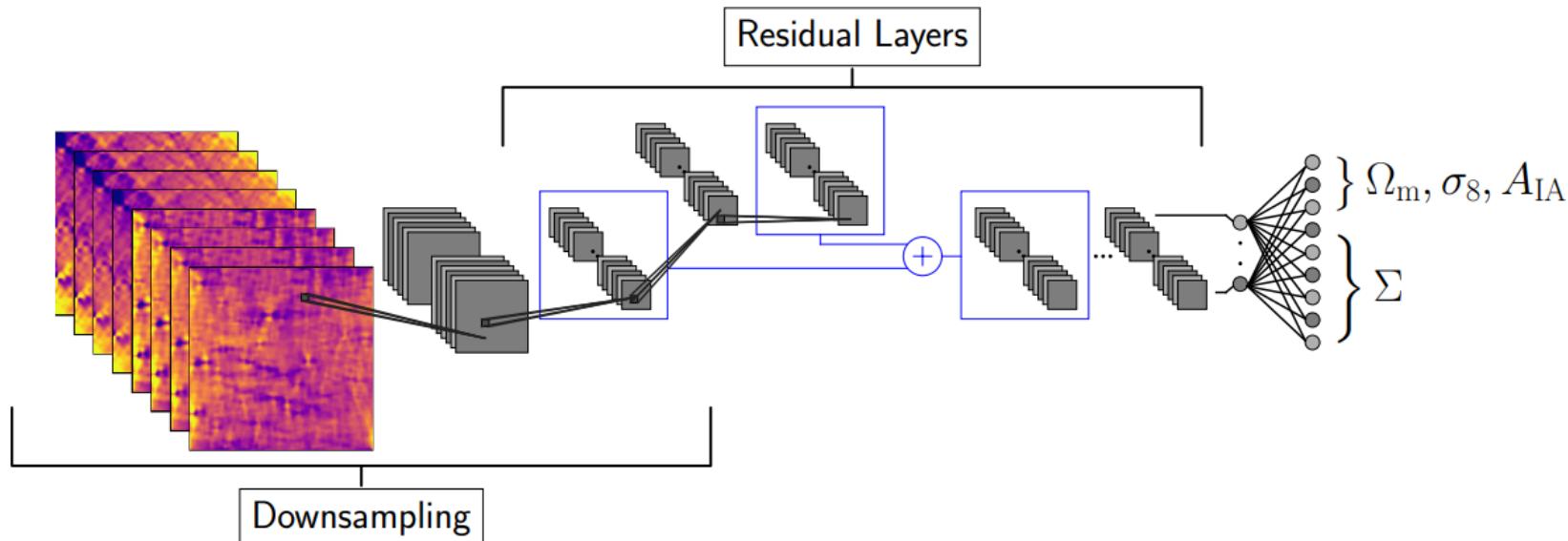
# Virtual Observatory



**Neural Networks to learn complex mapping between observables**

**Wu+21**

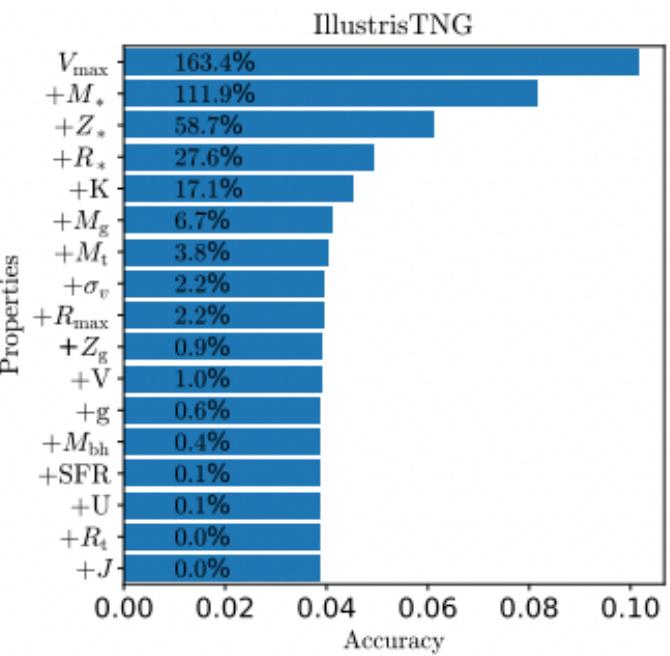
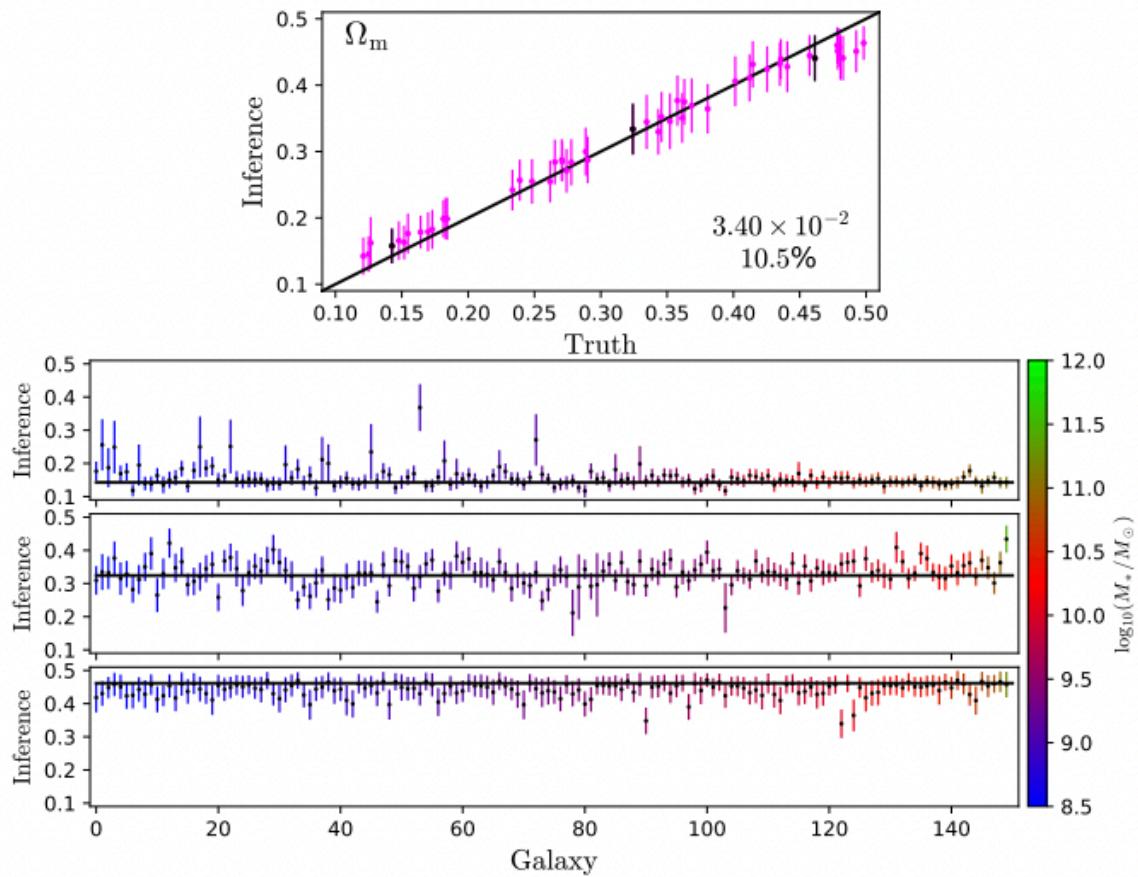
# Deep Learning for Cosmological Inference



Motivation: Generalize comparison of observations with theory, beyond basic summary statistics

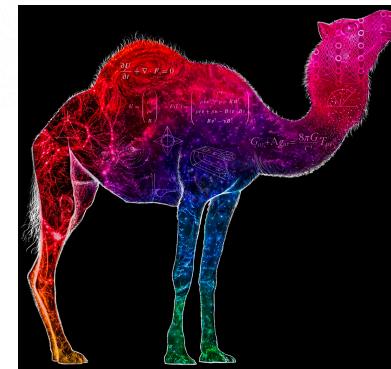
Neural Networks are used as efficient feature extractors

# Cosmological Inference



Villaescusa-Navarro+22

CAMELS



**Neural Networks to learn new relations between models and observables**

Ravanbakhsh+17, Brehmer+19, Ribli+19, Pan+19, Ntampaka+19, Alexander+20, Arjona+20, Coogan+20, Escamilla-Rivera+20, Hortua+20, Vama+20, Vernardos+20, Wang+20, Mao+20, Arico+20, Villaescusa\\_navarro+20, Singh+20, Park+21, Modi+21, Villaescusa-Navarro+21ab, Moriwaki+21, DeRose+21, Makinen+21, Villaescusa-Navarro+22

# ALSO FOR GALAXY MORPHOLOGY

SVMs

CNNs

[HUERTAS-COMPANY+14]

AUTOMATIC

Late-Type

13

Early-Type

87

75

25

Early-Type



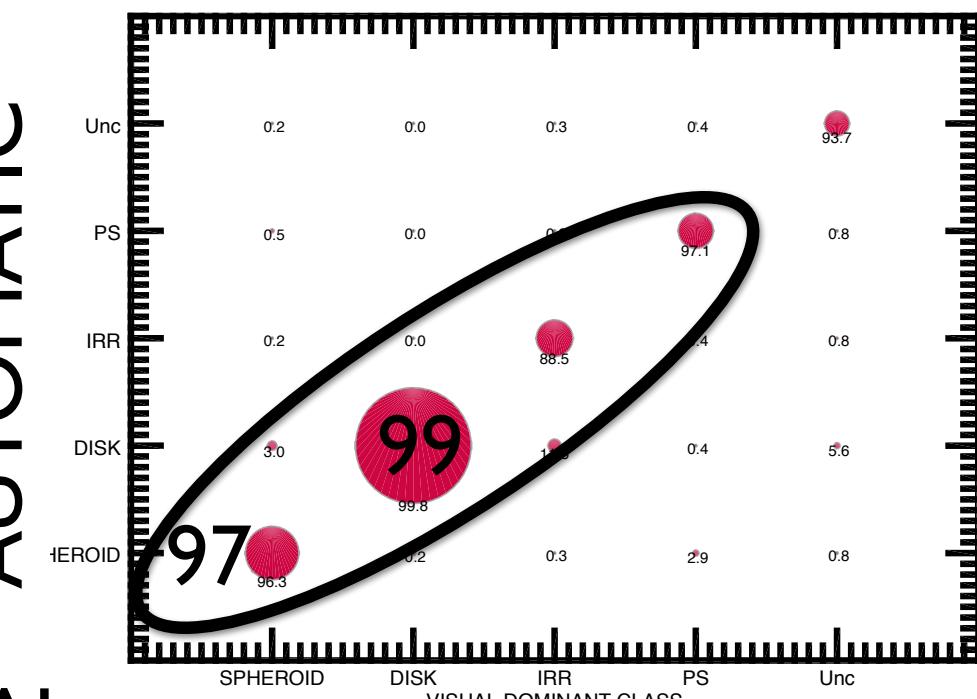
Late-Type



AUTOMATIC

VISUAL

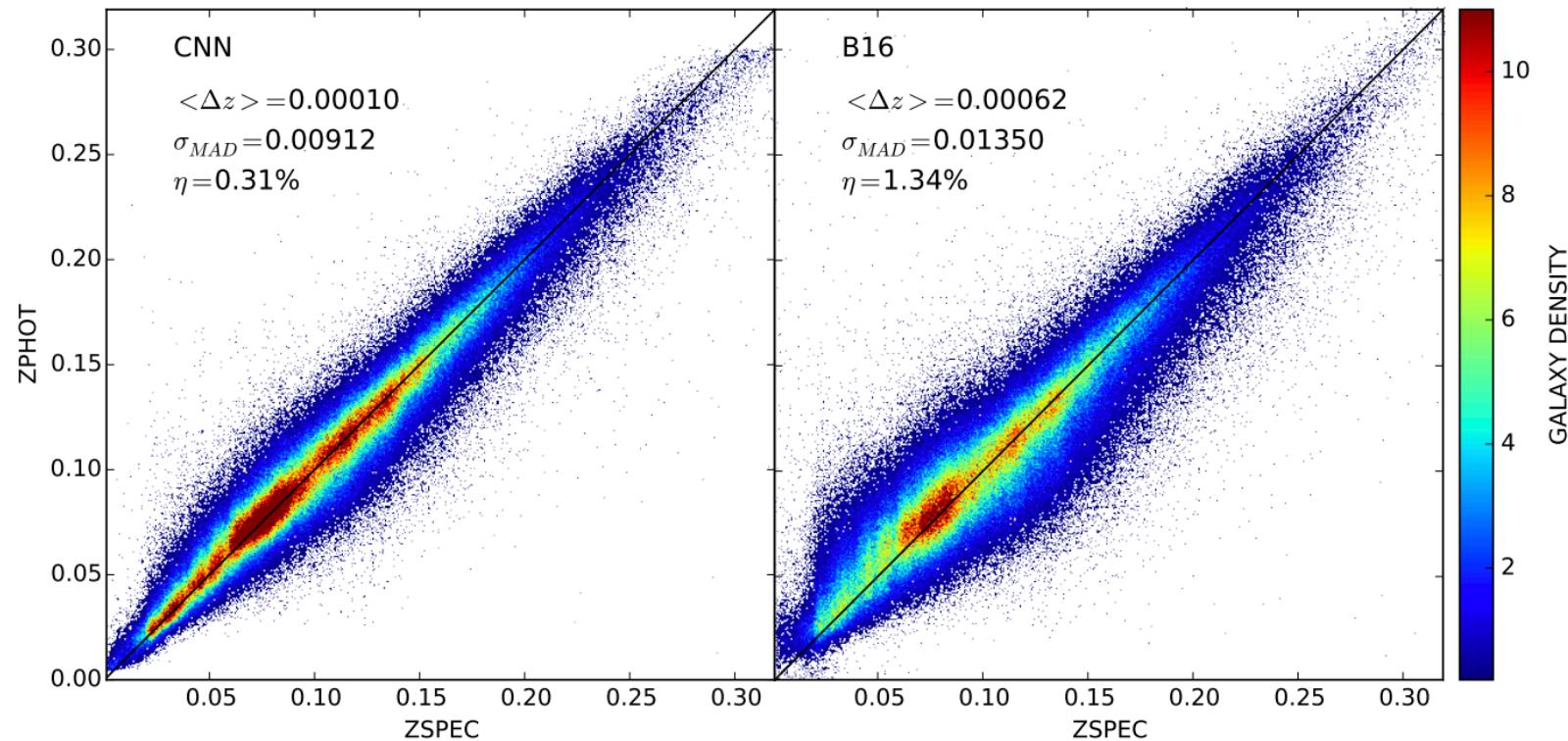
[HUERTAS-COMPANY+15b]



VISUAL

VISUAL

# PHOTOMETRIC REDSHIFTS



Pasquet+18

AUTOMATICALLY COMBINING MORPHOLOGY AND COLOR  
FOR PHOTOZ ESTIMATION

Making deep neural networks more robust

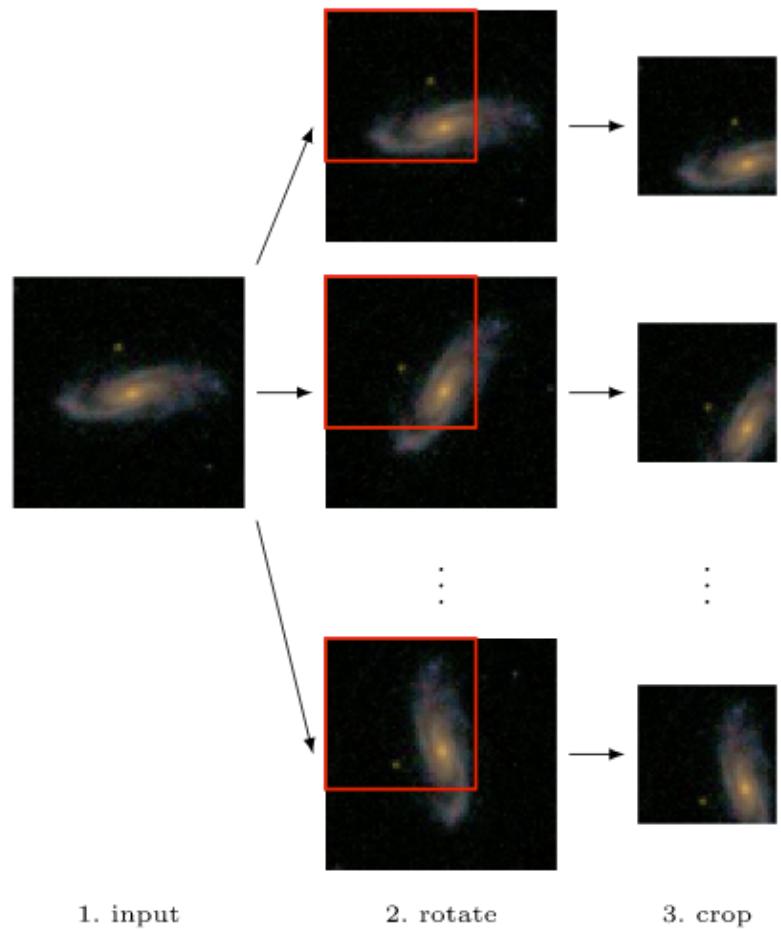
# DATA AUGMENTATION

**CNNs are invariants to translations**

FOR MANY APPLICATIONS THE CLASSIFICATION SHOULD  
BE INDEPENDENT TO:

- TRANSLATIONS
- ROTATIONS
- SCALINGS
- ETC...

# DATA AUGMENTATION



Dieleman+15

FOR MANY APPLICATIONS THE CLASSIFICATION SHOULD  
BE INDEPENDENT TO:  
- TRANSLATIONS  
- ROTATIONS  
- SCALINGS  
- ETC...

# DATA AUGMENTATION



FOR MANY APPLICATIONS THE CLASSIFICATION SHOULD  
BE INDEPENDENT TO:  
- TRANSALTIONS  
- ROTATIONS  
- SCALINGS  
- ETC...

# DATA AUGMENTATION

CAN BE DONE “ON THE FLY” DURING THE TRAINING PHASE

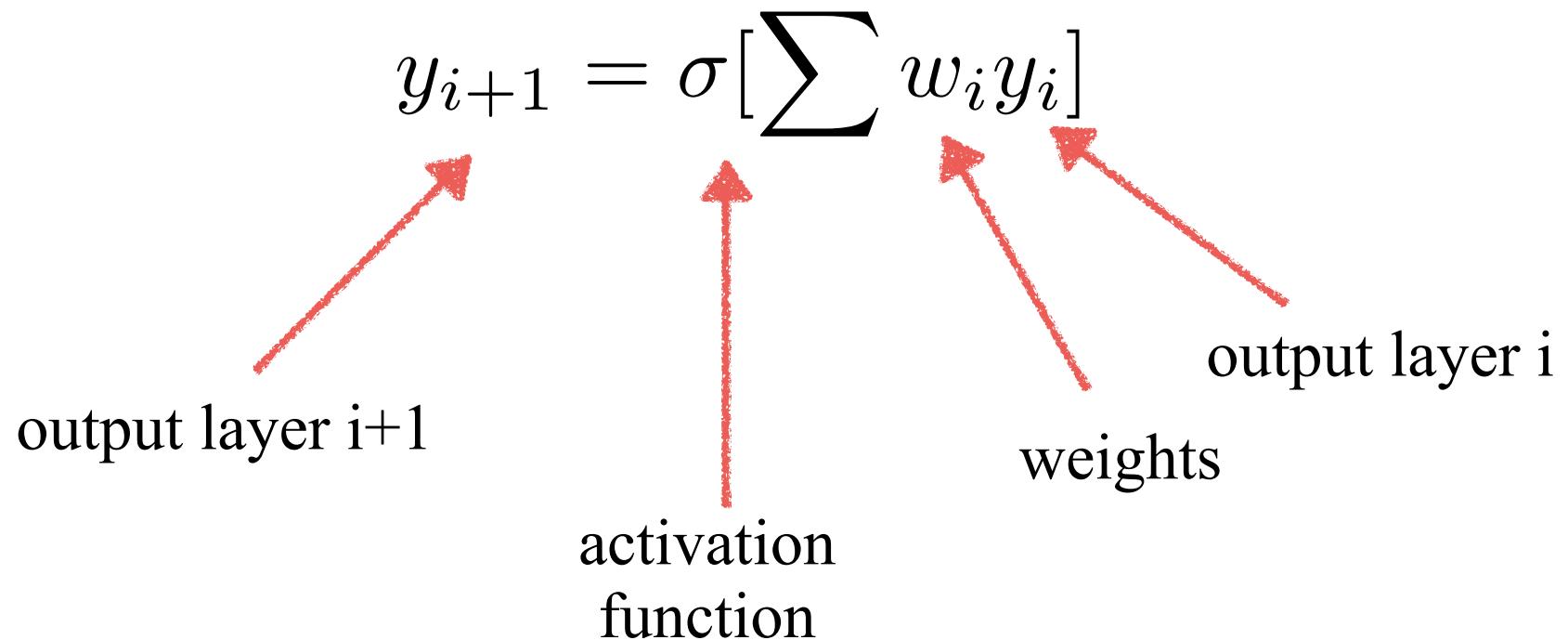
## KERAS IMPLEMENTATION:

DEFINES  
THE RANGE OF  
PERTURBATIONS  
APPLIED TO  
IMAGES DURING  
TRAINING

```
datagen = ImageDataGenerator(  
    featurewise_center=False,  
    samplewise_center=False,  
    featurewise_std_normalization=False,  
    samplewise_std_normalization=False,  
    zca_whitening=False,  
    rotation_range=45,  
    width_shift_range=0.05,  
    height_shift_range=0.05,  
    horizontal_flip=True,  
    vertical_flip=True,  
    zoom_range=[0.75,1.3])
```

# VANISHING / EXPLODING GRADIENT PROBLEM

REMEMBER THAT:



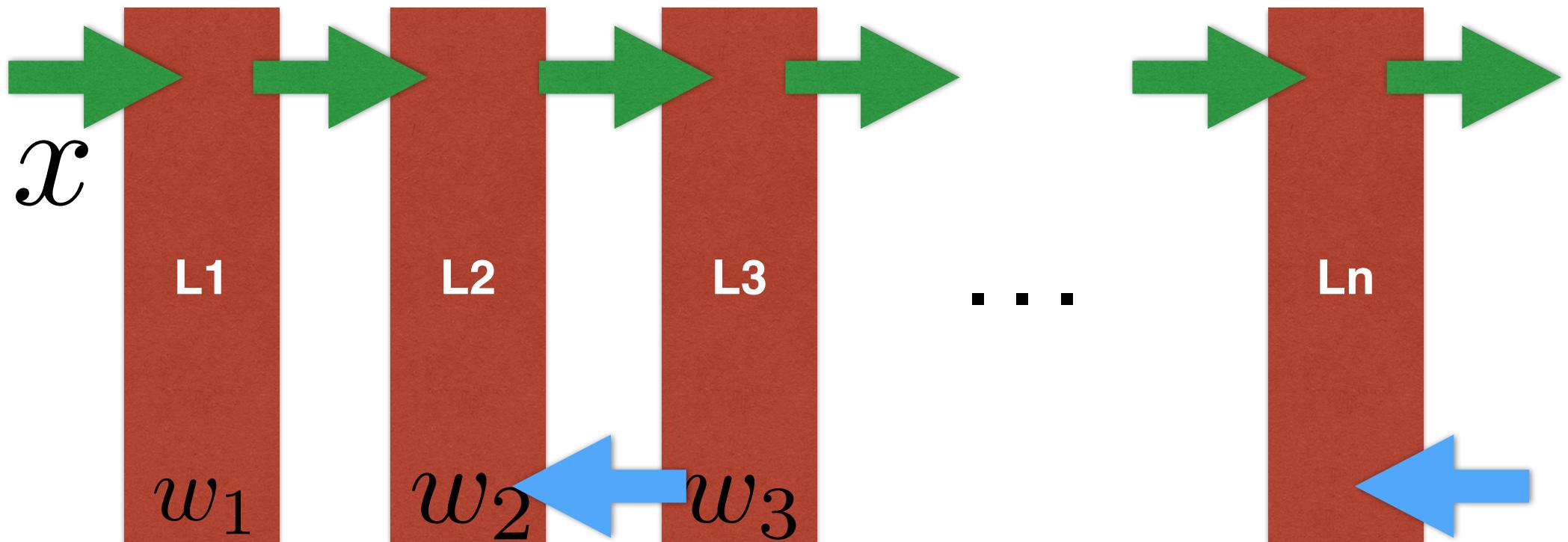
# VANISHING / EXPLODING GRADIENT PROBLEM

WITH MANY LAYERS:

$$y_n = \sigma \left( \dots \sigma \left( \dots \sigma \left( \sum w_0 x \right) \right) \right)$$

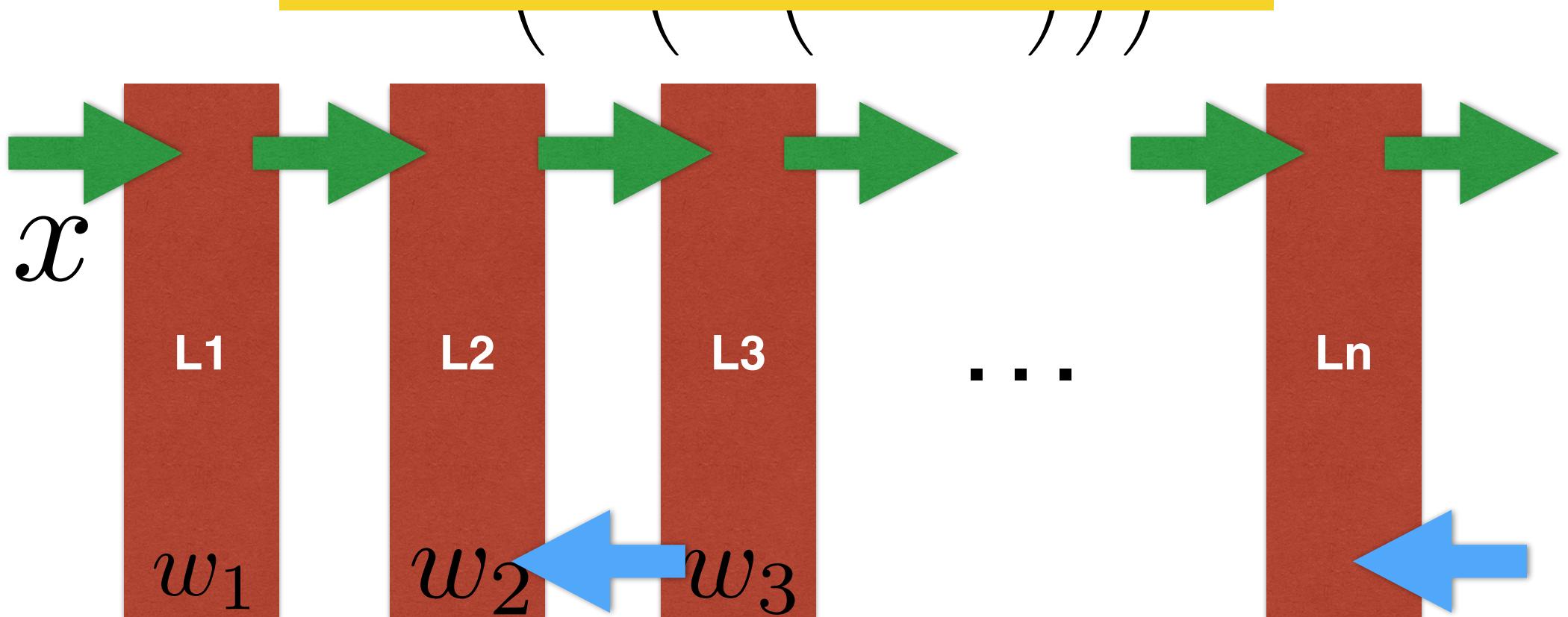
# VANISHING/EXPLODING GRADIENT PROBLEM

$$y_n = \sigma\left(\dots\sigma\left(\dots\sigma\left(\sum w_0x\right)\right)\right)$$

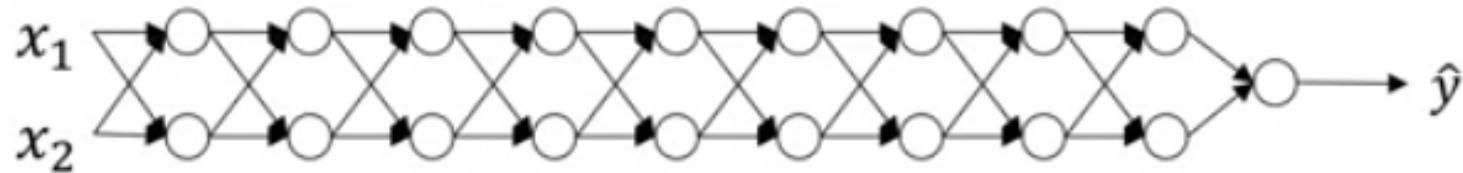


# VANISHING/EXPLODING GRADIENT PROBLEM

TRAINING BECOMES UNSTABLE  
VERY SLOW OR NO CONVERGENCE



# VANISHING/EXPLODING GRADIENT PROBLEM



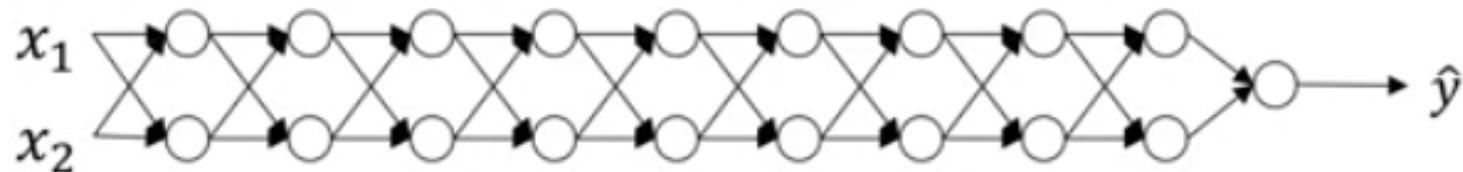
IF WE ASSUME AN IDENTITY ACTIVATION FUNCTION:

$$\hat{y} = x \prod_n w_i$$

with:

$$w_i = \begin{pmatrix} w_i^0 & 0 \\ 0 & w_i^1 \end{pmatrix} \quad x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

# VANISHING/EXPLODING GRADIENT PROBLEM



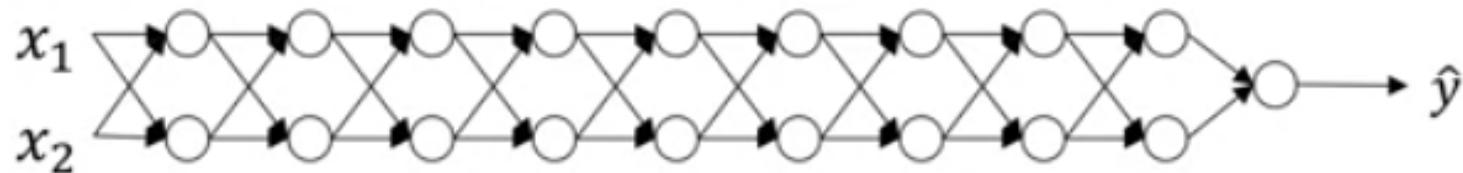
$$w_i = \begin{pmatrix} w_i^0 & 0 \\ 0 & w_i^1 \end{pmatrix} \quad \hat{y} = x \prod_n w_i$$

$$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \quad \text{IF WEIGHTS ARE ALL INITIALIZED TO VALUES } \ll 1:$$

$$w_i^L \rightarrow 0$$

VANISHING GRADIENT

# VANISHING/EXPLODING GRADIENT PROBLEM



$$w_i = \begin{pmatrix} w_i^0 & 0 \\ 0 & w_i^1 \end{pmatrix} \quad \hat{y} = x \prod_n w_i$$

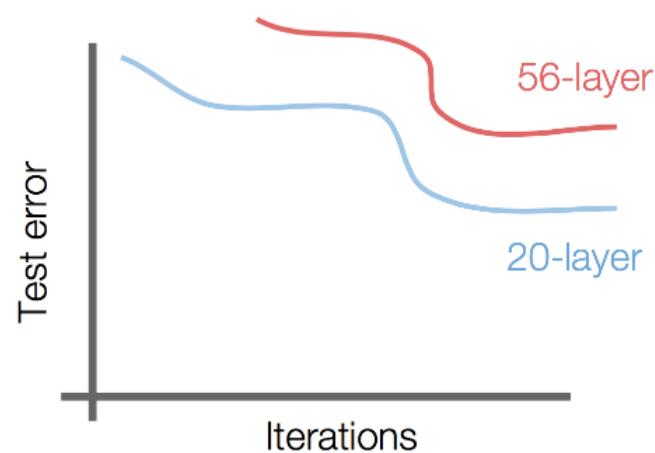
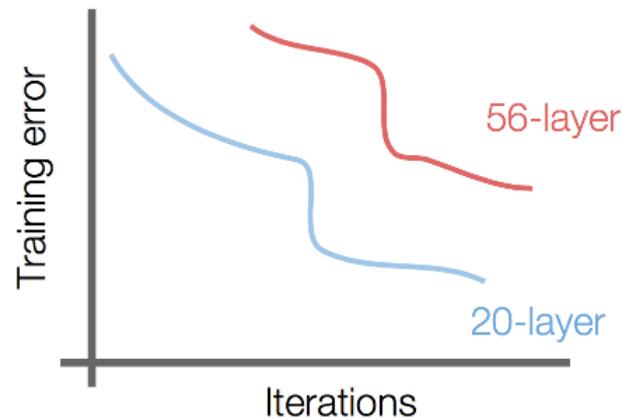
$$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \quad \text{IF WEIGHTS ARE ALL INITIALIZED TO VALUES } > 1:$$

$$w_i^L \rightarrow \infty$$

EXPLODING GRADIENT

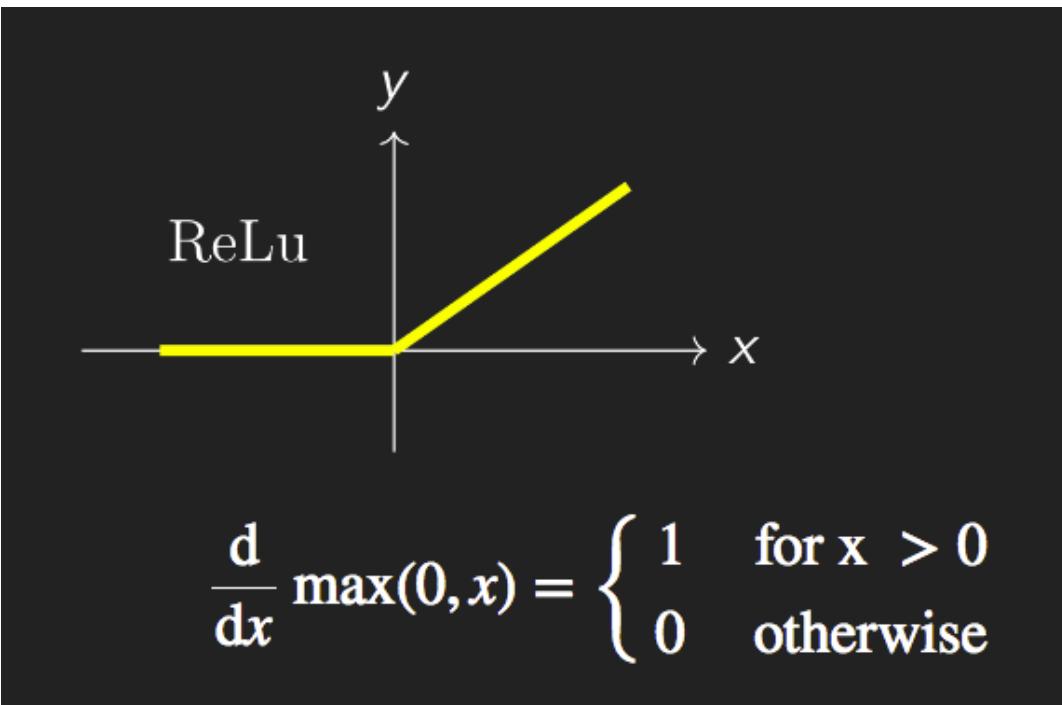
# VANISHING/EXPLODING GRADIENT PROBLEM

**TRAINING BECOMES UNSTABLE  
VERY SLOW OR NO CONVERGENCE**



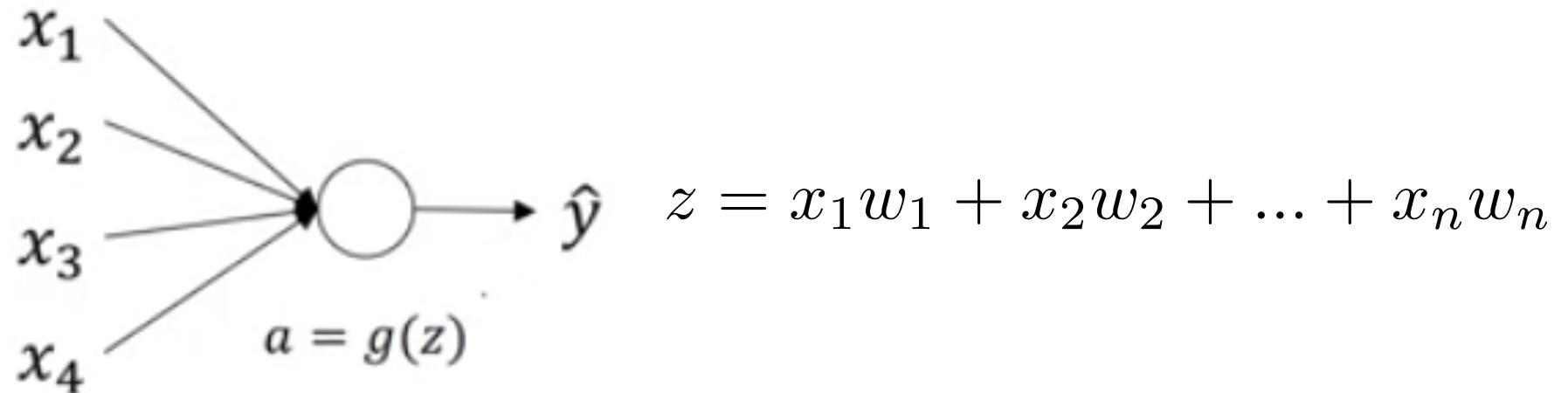
# VANISHING GRADIENT PROBLEM

THE ReLU ACTIVATION FUNCTION HELPS PREVENTING VANISHING GRADIENTS

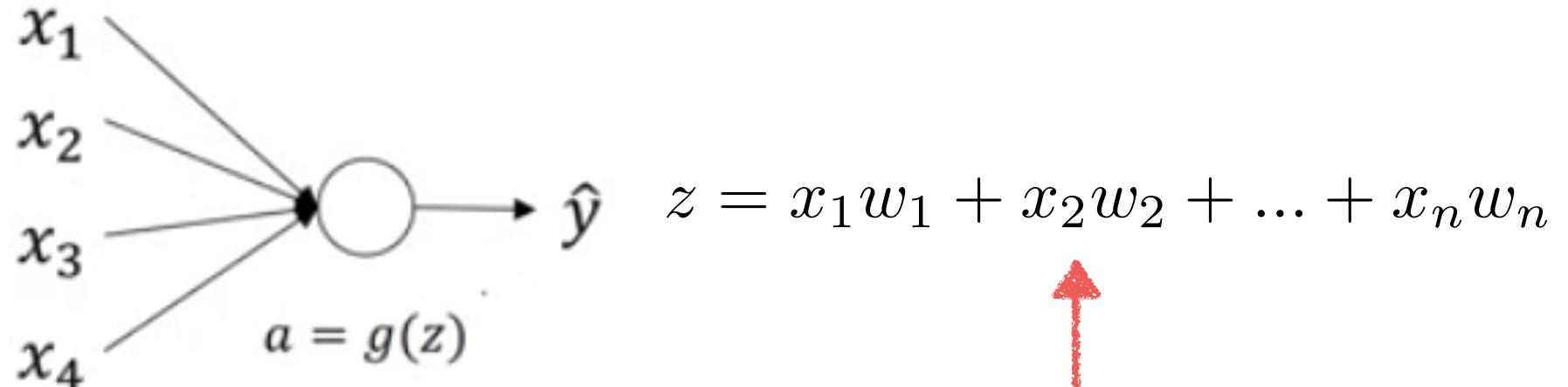


MOST WIDELY USED  
IN DEEP NETWORKS

# WEIGHT INITIALIZATION IS A KEY POINT...

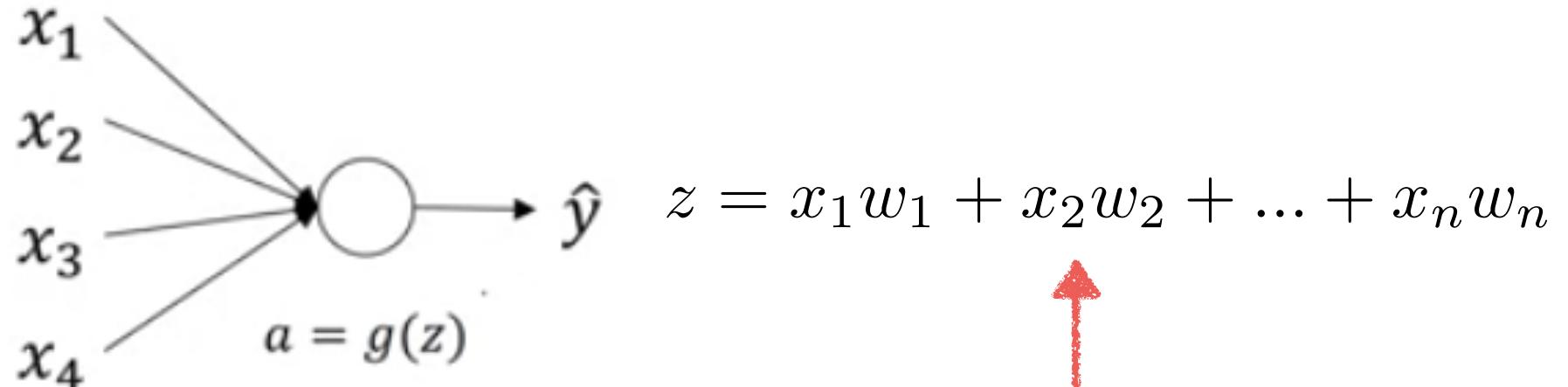


# WEIGHT INITIALIZATION IS A KEY POINT...



THE LARGER  $n$ , THE SMALLER  
WEIGHTS SHOULD BE...

# WEIGHT INITIALIZATION IS A KEY POINT...



THE LARGER  $n$ , THE SMALLER  
WEIGHTS SHOULD BE...

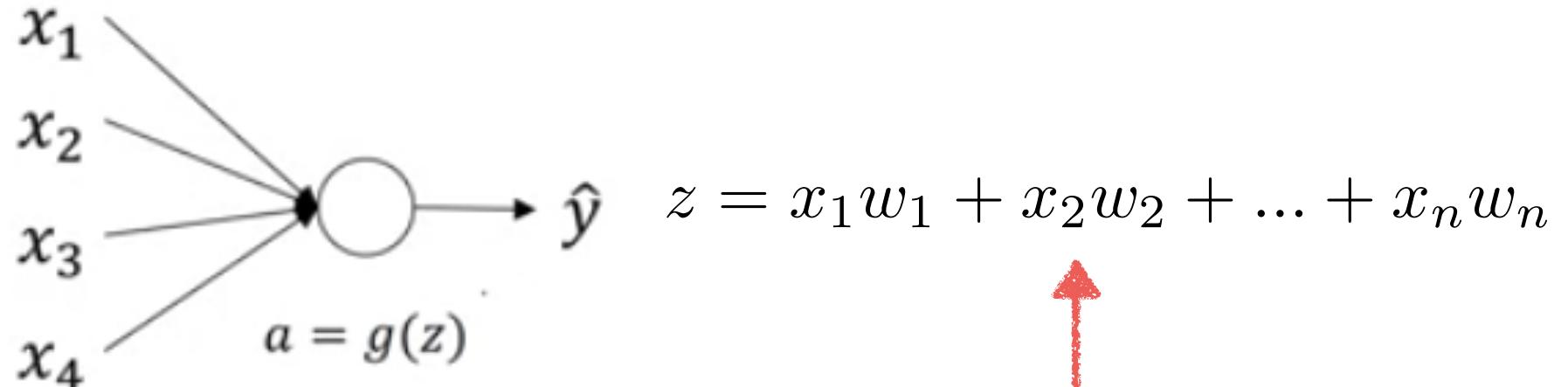
ONE SIMPLE SOLUTION:

$$\sigma^2(w_i) = \frac{1}{n}$$

number  
of  
inputs

A mathematical equation for weight initialization is shown:  $\sigma^2(w_i) = \frac{1}{n}$ . A red arrow points from the word "variance" above the equation to the fraction  $\frac{1}{n}$ . Another red arrow points from the text "number of inputs" below the equation to the denominator  $n$ .

# WEIGHT INITIALIZATION IS A KEY POINT...



THE LARGER  $n$ , THE SMALLER  
WEIGHTS SHOULD BE...

ONE SIMPLE SOLUTION:

$$\sigma^2(w_i) = \frac{1}{n}$$

← number of inputs

# WEIGHT INITIALIZATION IS A KEY POINT...

For ReLU activation functions we typically use:

$$\sigma^2(w_i) = \frac{2}{n}$$

[He initialization, He+15]

# WEIGHT INITIALIZATION IS A KEY POINT...

## IMPLEMENTATION IN KERAS:

```
initialization = 'he_normal'  
act = 'relu'  
  
model = Sequential()  
model.add(Convolution2D(depth, conv_size, conv_size, activation=act, border_mode='same',  
name = "conv%i"%(layer_n), init=initialization, W_constraint=constraint))
```

# WEIGHT INITIALIZATION IS A KEY POINT...

## IMPLEMENTATION IN KERAS:

```
initialization = 'he_normal'  
act = 'relu'  
  
model = Sequential()  
model.add(Convolution2D(depth, conv_size, conv_size, activation=act, border_mode='same',  
name = "conv%i"%(layer_n), init=initialization, W_constraint=constraint))
```

MANY OTHER INITIALIZATIONS AVAILABLE:

keras.initializers



<https://keras.io/initializers/>

# BATCH NORMALIZATION

[SZEGEDY+15]

ANOTHER SOLUTION TO KEEP REASONABLE VALUES OF  
THE ACTIVATIONS IN DEEP NETWORKS

**BATCH NORMALIZATION** PREVENTS LOW OR LARGE  
VALUES BY RE-NORMALIZING THE VALUES BEFORE  
ACTIVATION FOR EVERY BATCH

$$\hat{y}_i = \gamma \frac{y_i - E(y_i)}{\sigma(y_i)} + \beta$$

INPUT      NORMALIZED INPUT      SCATTER

The diagram illustrates the components of the Batch Normalization formula. A red arrow labeled "INPUT" points to the term  $y_i$ . Another red arrow labeled "NORMALIZED INPUT" points to the term  $\hat{y}_i$ . A third red arrow labeled "SCATTER" points to the term  $\beta$ .

# BATCH NORMALIZATION

[SZEGEDY+15]

**BATCH NORMALIZATION** SPEEDS UP AND STABILIZES TRAINING

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;  
Parameters to be learned:  $\gamma, \beta$   
**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

# BATCH NORMALIZATION

[SZEGEDY+15]

IN KERAS, IT IS IMPLEMENTED AS AN ADDITIONAL LAYER

## BatchNormalization

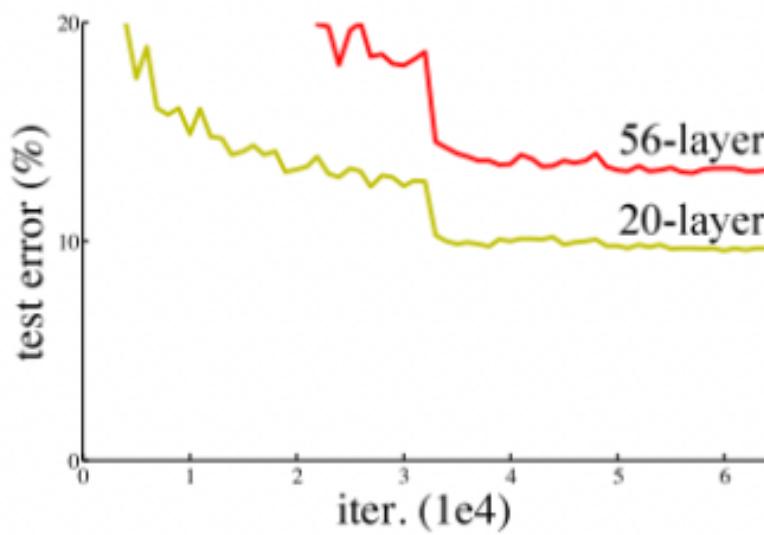
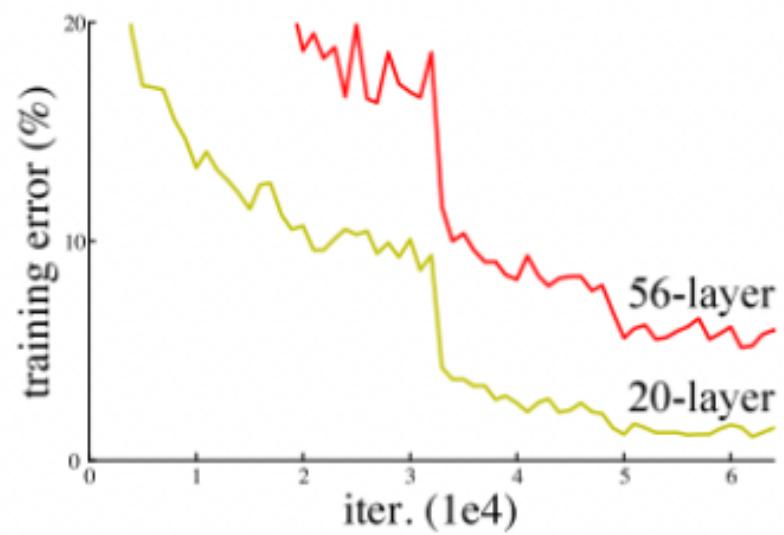
[\[source\]](#)

```
keras.layers.BatchNormalization(axis=-1, momentum=0.99, epsilon=0.001, center=True, scale=True, beta_initializer
```

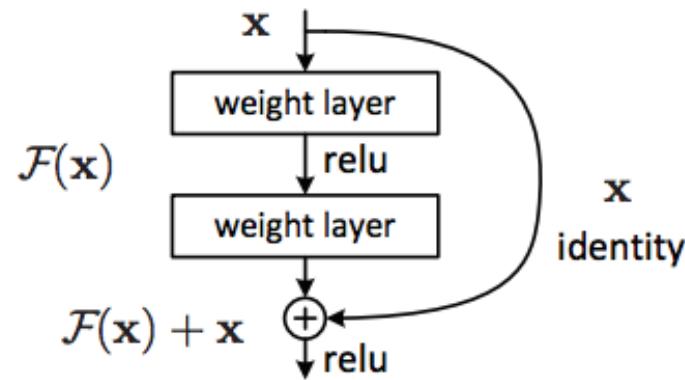
Batch normalization layer (Ioffe and Szegedy, 2014).

Normalize the activations of the previous layer at each batch, i.e. applies a transformation that maintains the mean activation close to 0 and the activation standard deviation close to 1.

## Arguments



# Residual Network



Enables an uninterrupted gradient flow

Learns the residual mapping instead of direct mapping

