

# Machine learning applied to population genomics

Manolo Fernandez Perez  
Department of Life Sciences, Imperial College London

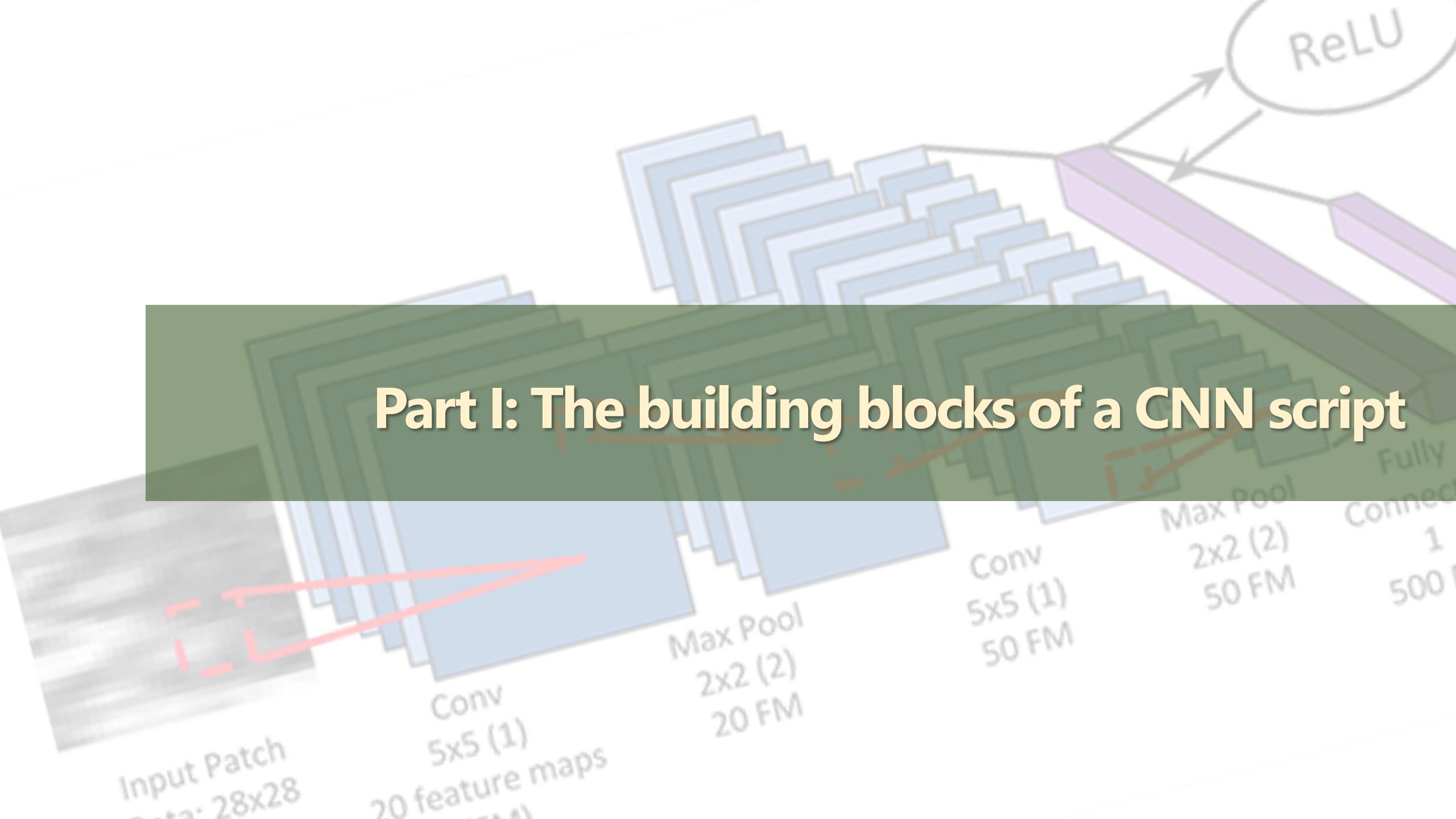
@ManoloLearning   
manolofperez@gmail.com   
[sites.google.com/site/manolofperez](http://sites.google.com/site/manolofperez) 

# Goals

- Conceive and simulate genetic data under competing demographic scenarios 
- Understand deep learning background and how a CNN works 
- Use CNN to detect regions with selective sweeps on real genomes
- How to use deep learning to compare demographic scenarios

# Program

# Part I: The building blocks of a CNN script



# CNN Script



required python modules.

Define the CNN architecture.

Load and process the training data.

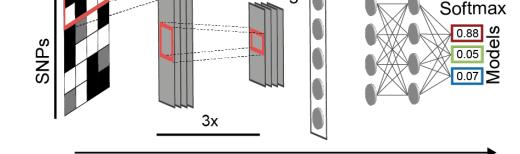
Train the network.

Load the test data and perform cross-validation.

Predict the most likely model for the empirical data



TensorFlow



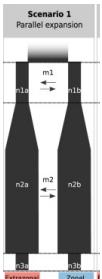
# CNN Script

Inputs:

**Scenario 1**

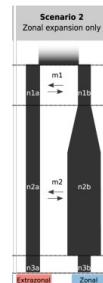
Samples

SNPs	-1 -1 1 0 1 -1 1 1 -1 0 1 -1 -1 1 1 1 1 1 -1 0 0 1 -1 -1	-1 -1 1 0 1 -1 1 1 -1 0 1 -1 -1 1 1 1 1 1 -1 0 0 1 -1 -1
------	---	---



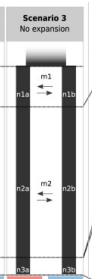
**Scenario 2**

-1 -1 1 0 1 -1 1 1 -1 0 1 -1 -1 1 1 1 1 1 -1 0 0 1 -1 -1
-1 1 1 1 1 1 -1 0 1 -1 1 1 -1 0 1 -1 -1 1 1 1 1 1 -1 0
0 1 -1 -1 -1 0 1 -1 1 1 -1 0 -1 1 1 1 1 -1 1 1 0 1 -1 -1
-1 1 1 1 1 1 -1 0 1 -1 1 1 -1 0 1 -1 -1 1 1 1 1 1 -1 0
0 1 -1 -1 -1 0 1 -1 1 1 -1 0 -1 1 1 1 1 -1 1 1 0 1 -1 -1



**Scenario 3**

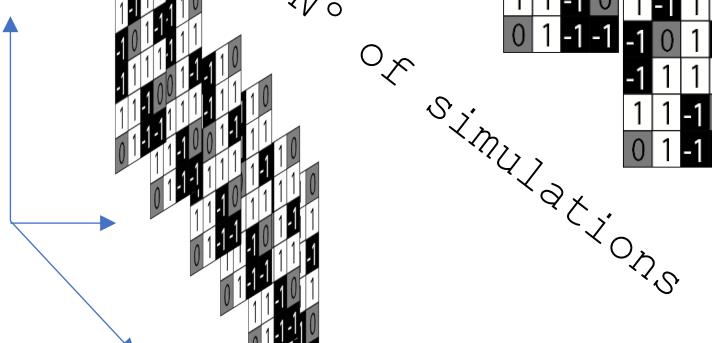
-1 -1 1 0 1 -1 1 1 -1 0 1 -1 -1 1 1 1 1 1 -1 0 0 1 -1 -1
-1 -1 1 0 1 -1 1 1 -1 0 1 -1 -1 1 1 1 1 1 -1 0 0 1 -1 -1
1 -1 1 1 -1 0 1 -1 -1 1 1 1 1 1 -1 0 0 1 -1 -1 -1 1 1 1
1 -1 1 1 -1 0 1 -1 -1 1 1 1 1 1 -1 0 0 1 -1 -1 -1 1 1 1
0 1 -1 -1 -1 0 1 -1 1 1 -1 0 -1 1 1 1 1 -1 1 1 0 1 -1 -1



**Parameters**

Theta T1 T2 T3 Ne

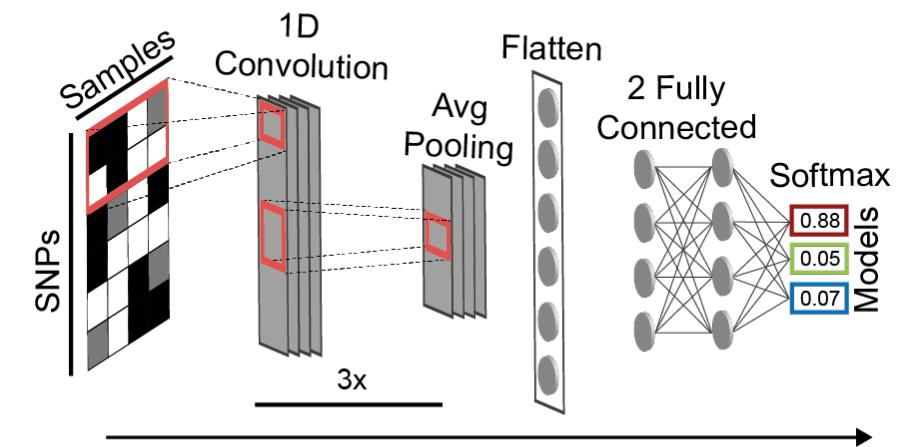
Sim1				
Sim2				
Sim3				
Sim4				



3-D Numpy array

$N_o$   
of simulations

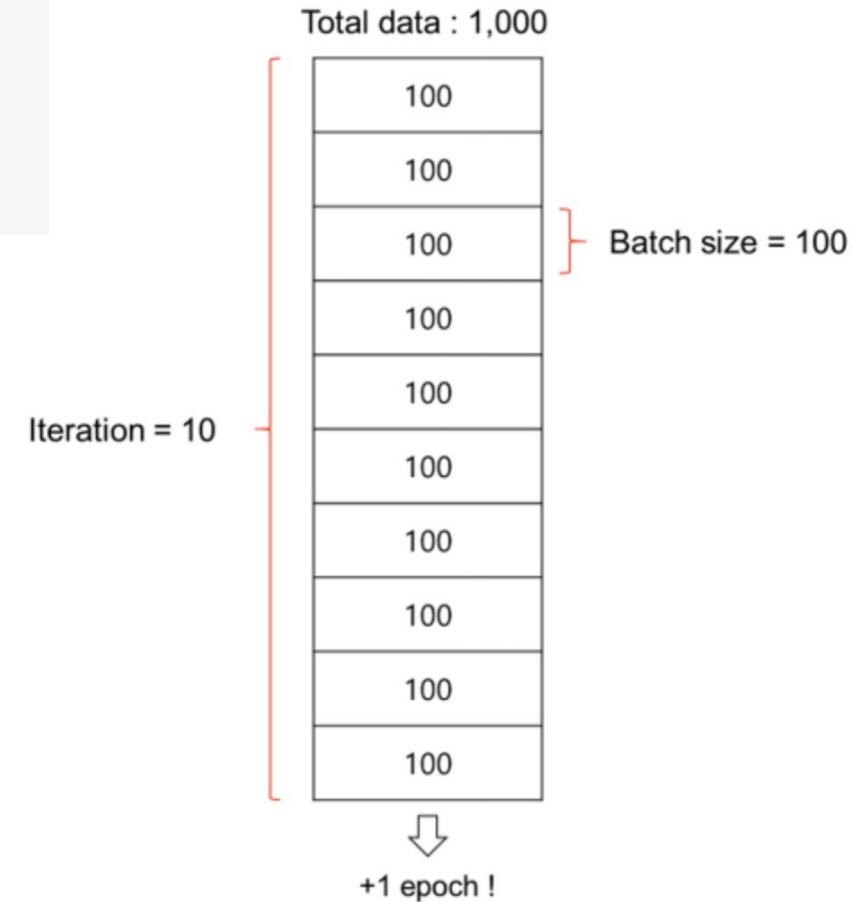
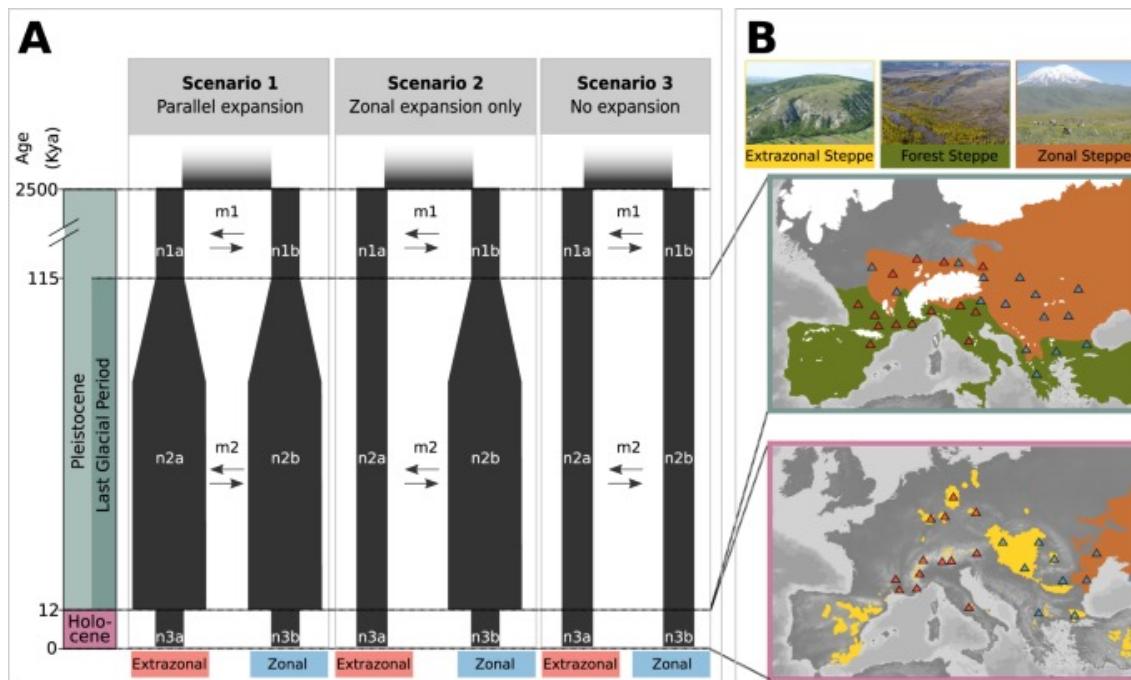
- **Practical Exercise 1:**
- Go through Section 1 of the Part1 script (Demographic models) and try to recognize all the elements of the network. Do you remember the function of each of those elements? Remember that you can add annotations to the code using # and add information that might help you when you get back to the script in the future.
- Now run all the cells until you reach the end of section 2. Your network will be training, so now we will have some time to discuss and do a quick review on the CNN elements.



# CNN Script

```
# Define parameters for the CNN run.
batch_size = 200
### how much interations to train the network
epochs = 100

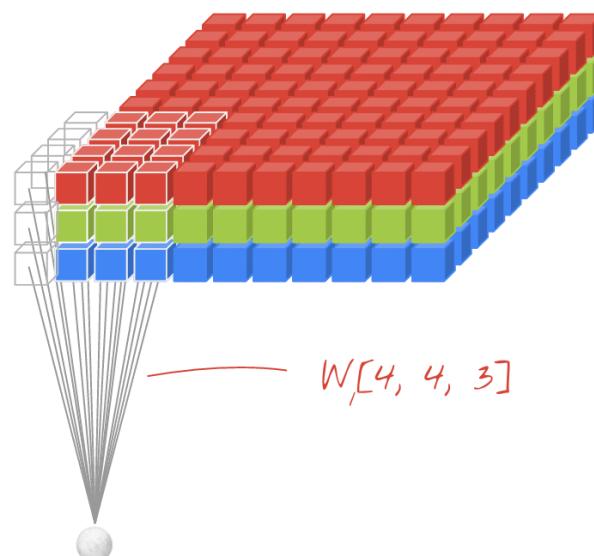
###n of models
num_classes = 3
```



<https://jerryan.medium.com/batch-size-a15958708a6>

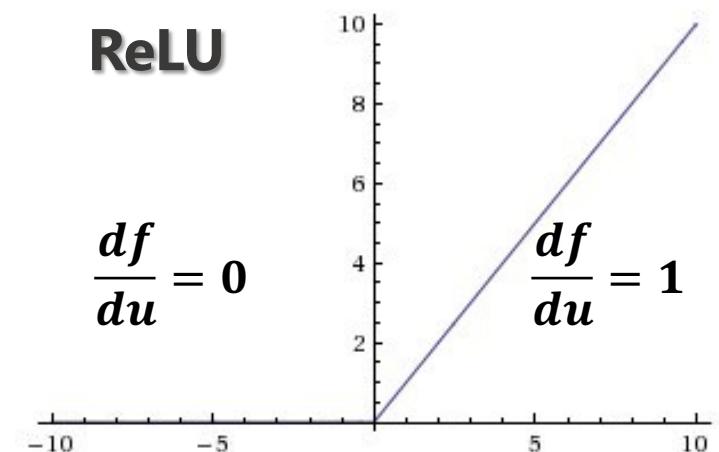
# CNN Script

```
# Define the CNN architecture.
def create_cnn(xtest):
    inputShape = (xtest.shape[1], xtest.shape[2])
    ## image size. images need to have EXACTLY the same size
    inputs = Input(shape=inputShape)
    x = inputs
    ## 1D convolution - less computational intensive and also treats snps as independent;
    x = Conv1D(250, kernel_size=2, activation='relu',input_shape=(xtest.shape[1], xtest.shape[2]))(x)
    ### Enables the network to learn more complex features / shapes.
    x = AveragePooling1D(pool_size=2)(x)
    x = BatchNormalization()(x)
```



ReLU

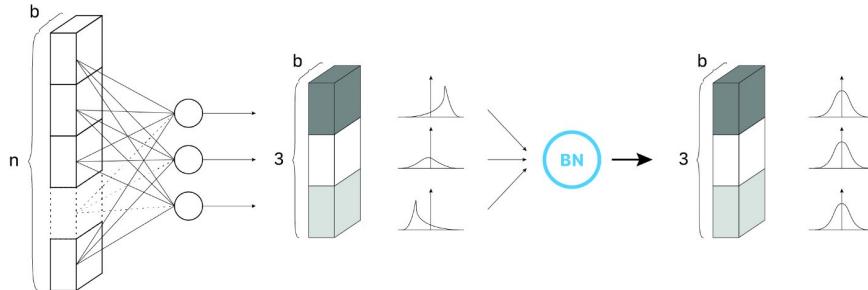
$$\frac{df}{du} = 0$$



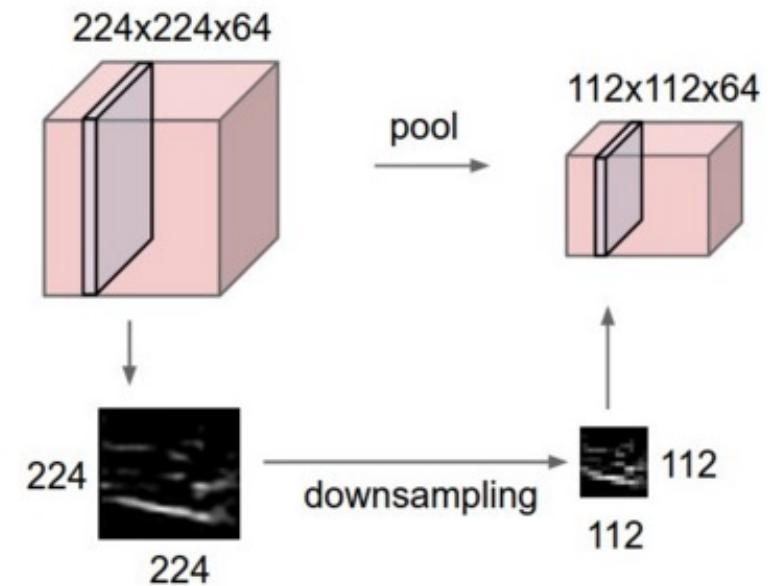
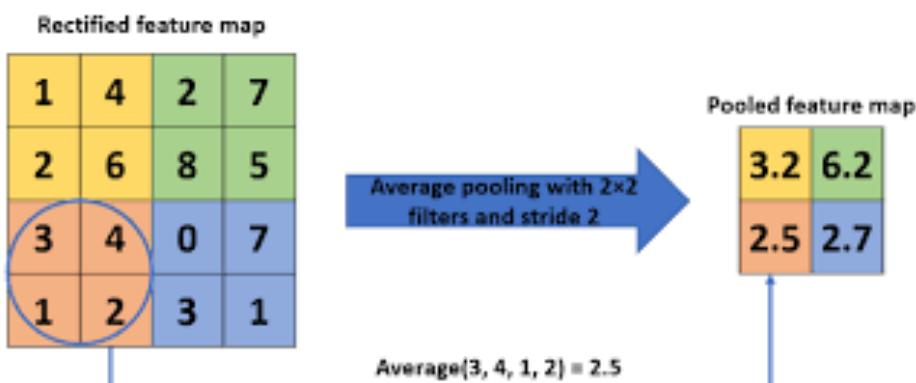
<https://www.kaggle.com/code/dansbecker/rectified-linear-units-relu-in-deep-learning/notebook>

# CNN Script

```
# Define the CNN architecture.
def create_cnn(xtest):
    inputShape = (xtest.shape[1], xtest.shape[2])
    ## image size. images need to have EXACTLY the same size
    inputs = Input(shape=inputShape)
    x = inputs
    ## 1D convolution - less computational intensive and also treats snps as independent;
    x = Conv1D(250, kernel_size=2, activation='relu',input_shape=(xtest.shape[1], xtest.shape[2]))(x)
    ### Enables the network to learn more complex features / shapes.
    x = AveragePooling1D(pool_size=2)(x)
    x = BatchNormalization()(x)
```



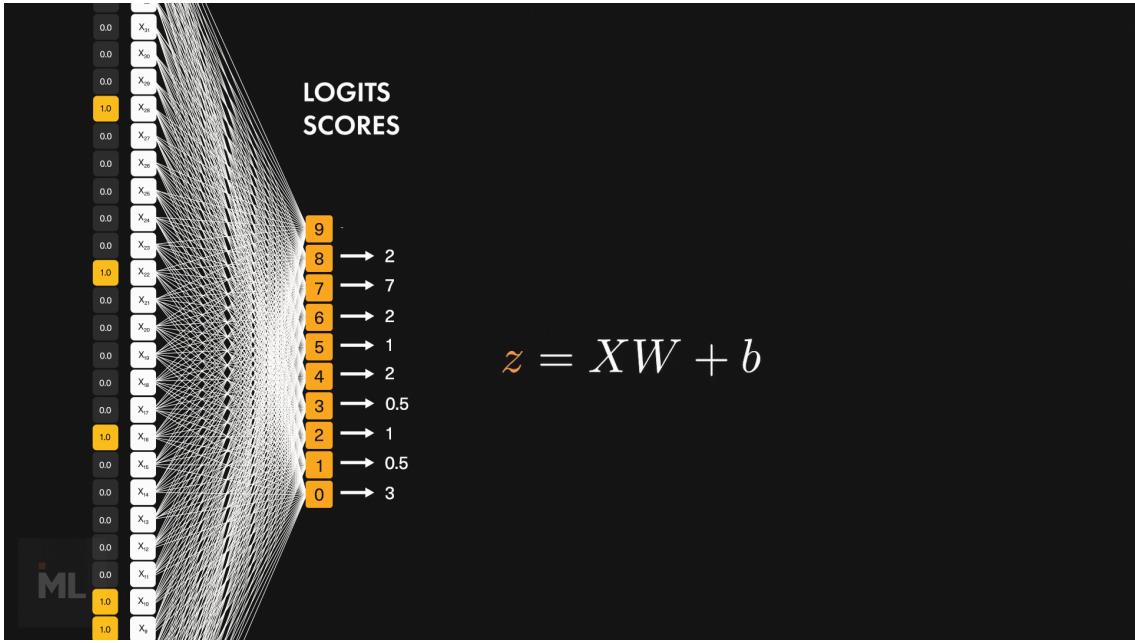
<https://towardsdatascience.com/batch-normalization-in-3-levels-of-understanding-14c2da90a338>



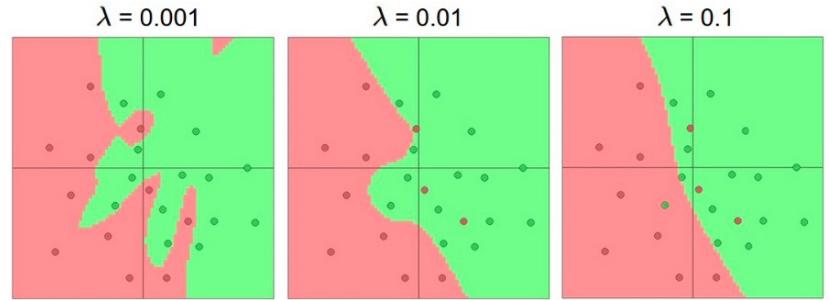
[https://leonardoaraujosantos.gitbook.io/artificial-intelligence/machine\\_learning/deep\\_learning/pooling\\_layer](https://leonardoaraujosantos.gitbook.io/artificial-intelligence/machine_learning/deep_learning/pooling_layer)

# CNN Script

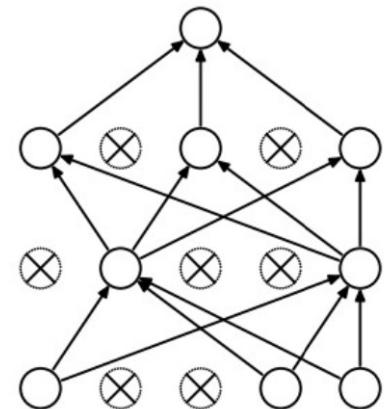
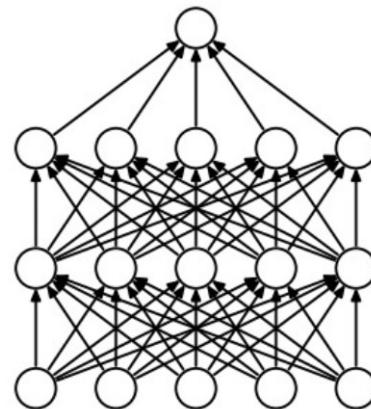
```
### Linearising the image as in the initial step. From this point on the network behaves as a Multi-Layer Perceptron
x = Flatten()(x)
x = Dense(125, activation='relu',kernel_regularizer=l2(1e-3), bias_regularizer=l2(1e-3))(x)
x = Dropout(0.5)(x)
x = Dense(125, activation='relu',kernel_regularizer=l2(1e-3), bias_regularizer=l2(1e-3))(x)
x = Dropout(0.5)(x)
x = Dense(num_classes, activation="softmax")(x)
```



Regularization



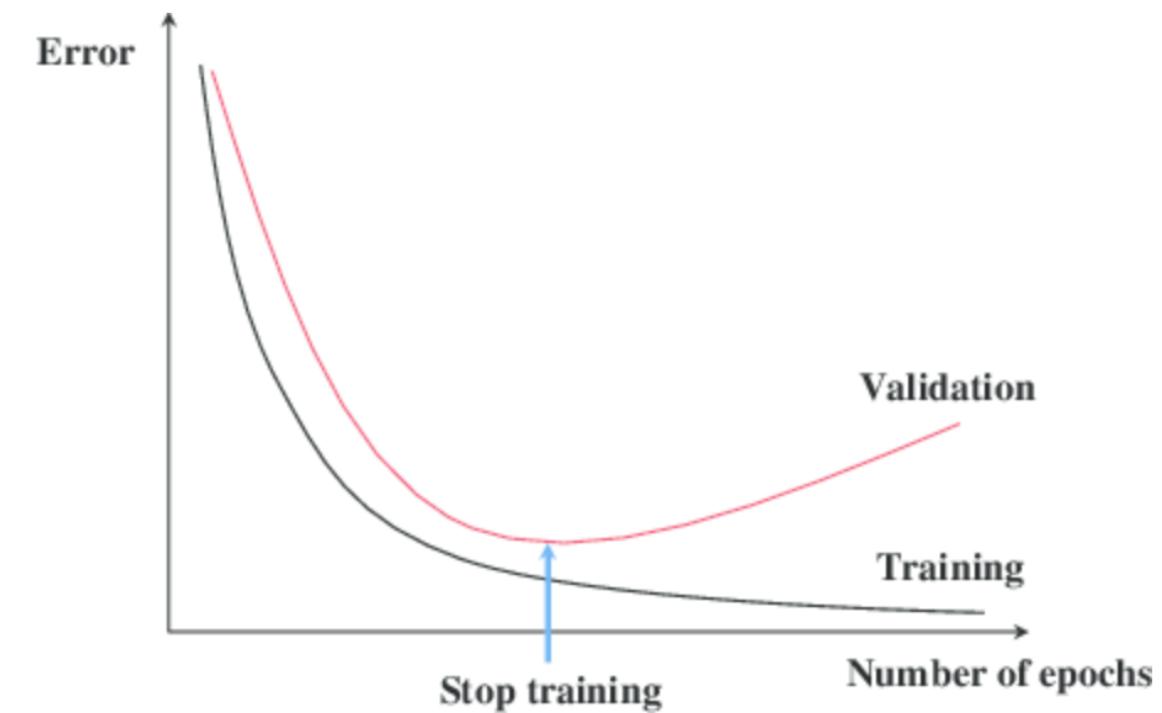
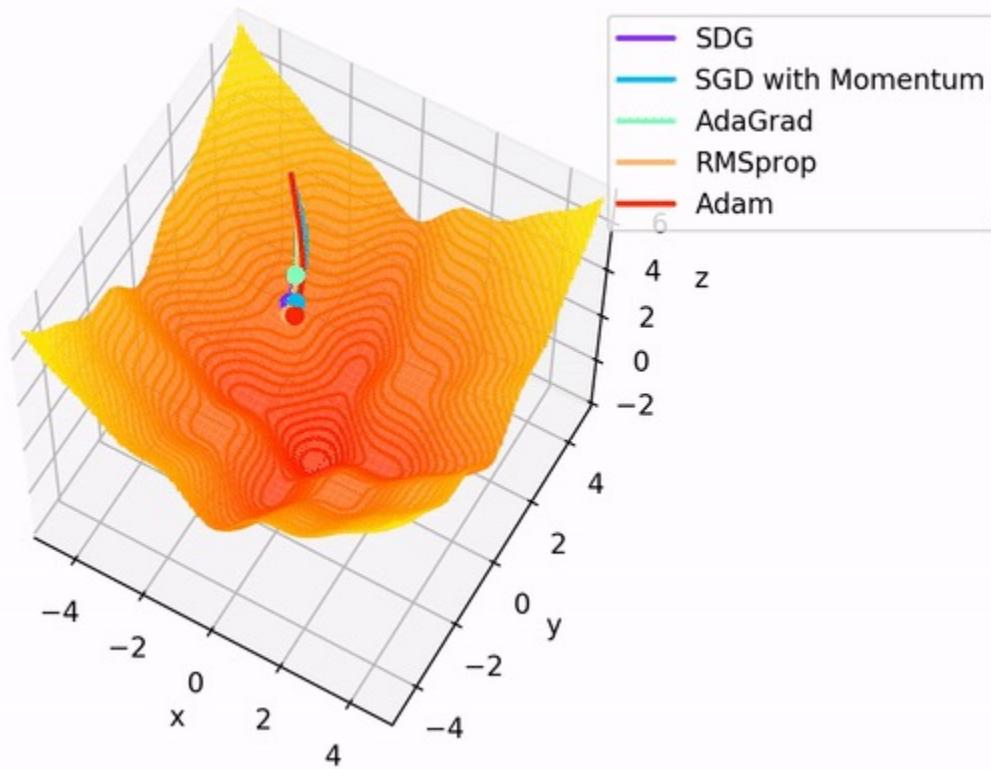
Options: L2, L1, maxnorm and dropout.



# CNN Script

```
# Compile the CNN.  
model.compile(loss=keras.losses.categorical_crossentropy,  
              optimizer='Adam',  
              metrics=['accuracy'])  
  
# We will use early stopping and save the model with the best val_accuracy.  
earlyStopping = EarlyStopping(monitor='val_accuracy', patience=150, verbose=0, mode='max', restore_best_weights  
### stop training when validation increases error
```

Optimizer Comparison



<https://towardsdatascience.com/complete-guide-to-adam-optimization-1e5f29532c3d>

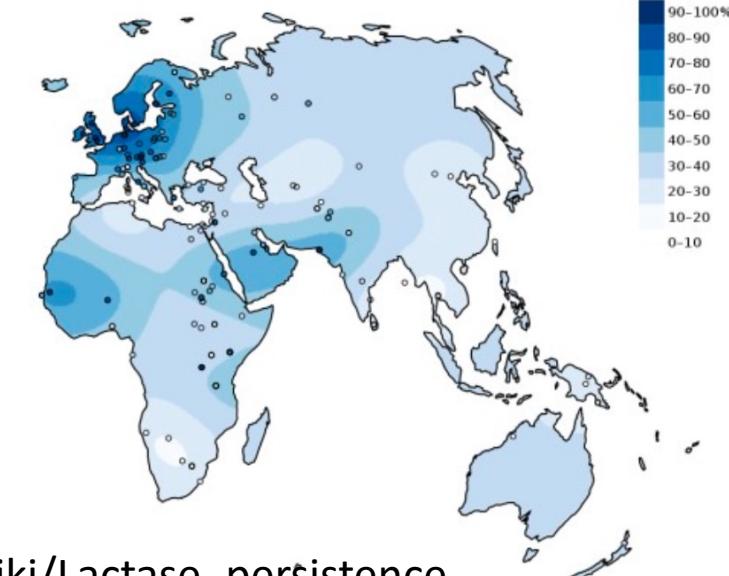
<https://towardsdatascience.com/a-practical-introduction-to-early-stopping-in-machine-learning-550ac88bc8fd>

- **Practical Exercise 2:**

Now open the Part 2 script. It uses Matteo's software (ImaGene) that is a CNN approach to infer selection at the LCT locus.

It uses a "simulation-on-the-fly"-like approach: training data is continuously generated by simulations to avoid the network to see the same data twice and therefore to reduce overfitting. This is a valuable consideration since, when reliable simulators are available, we have access to theoretically infinite training data, the latter being constrained by computing time only.

There are also functions to automate some of the steps we did in the previous example. You can compare the strategies and the architectures adopted.



## Part II: Quick overview of other applications and future perspectives.

## Deep Learning in Population Genetics

Kevin Korfmann<sup>1</sup>, Oscar E. Gaggiotti<sup>2</sup>, and Matteo Fumagalli  <sup>3,\*</sup>

<sup>1</sup>Professorship for Population Genetics, Department of Life Science Systems, Technical University of Munich, Germany

<sup>2</sup>Centre for Biological Diversity, Sir Harold Mitchell Building, University of St Andrews, Fife KY16 9TF, UK

<sup>3</sup>Department of Biological and Behavioural Sciences, Queen Mary University of London, UK

\*Corresponding author: E-mail: m.fumagalli@qmul.ac.uk.

Accepted: 16 January 2023

Review Article | Published: 04 September 2023

## Harnessing deep learning for population genetic inference

[Xin Huang](#) , [Aigerim Rymbekova](#), [Olga Dolgova](#), [Oscar Lao](#)  & [Martin Kuhlwilm](#) 

[Nature Reviews Genetics](#) **25**, 61–78 (2024) | [Cite this article](#)

8148 Accesses | 4 Citations | 41 Altmetric | [Metrics](#)

## The Unreasonable Effectiveness of Convolutional Neural Networks in Population Genetic Inference

Lex Flagel,<sup>1,2</sup> Yaniv Brandvain,<sup>2</sup> and Daniel R. Schrider<sup>\*3</sup>

<sup>1</sup>Monsanto Company, Chesterfield, MO

<sup>2</sup>Department of Plant and Microbial Biology, University of Minnesota, St. Paul, MN

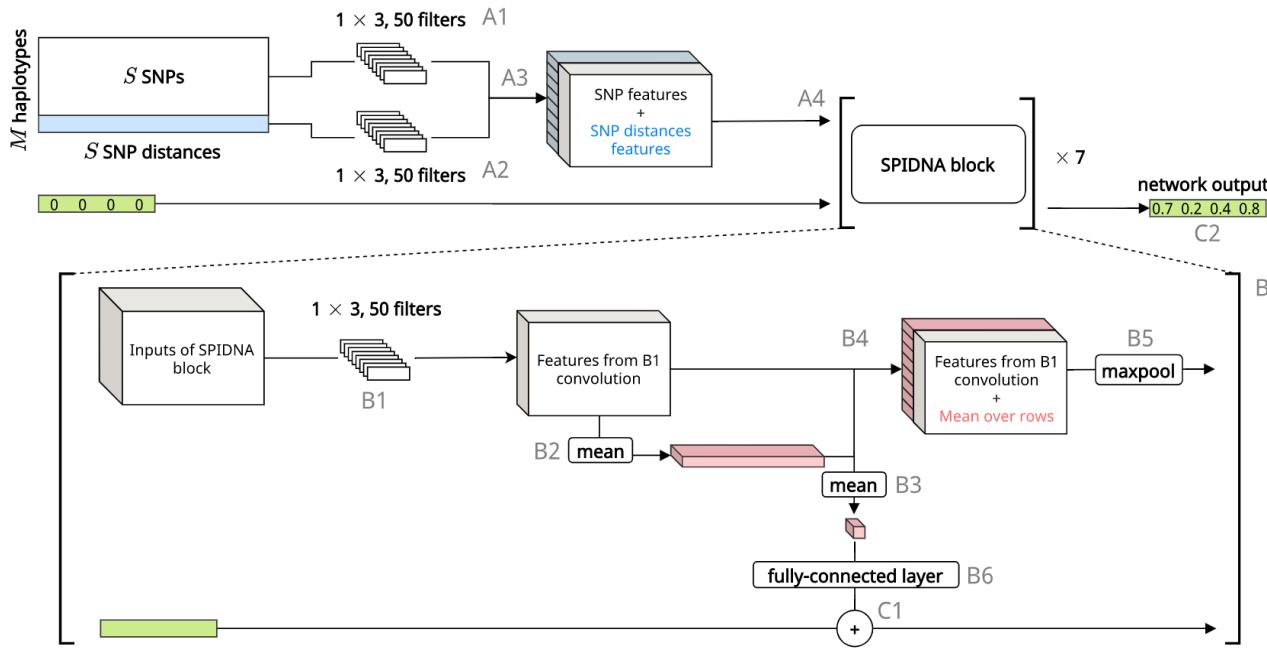
<sup>3</sup>Department of Genetics, University of North Carolina, Chapel Hill, NC

\*Corresponding author: E-mail: drs@unc.edu.

Associate editor: Yuseob Kim

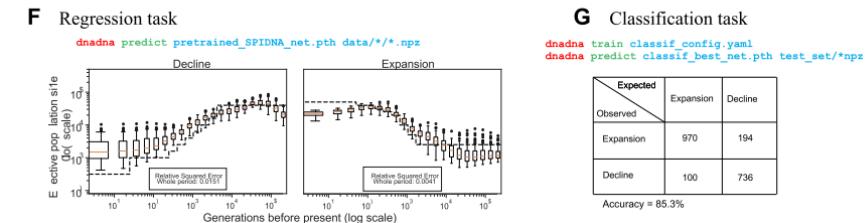
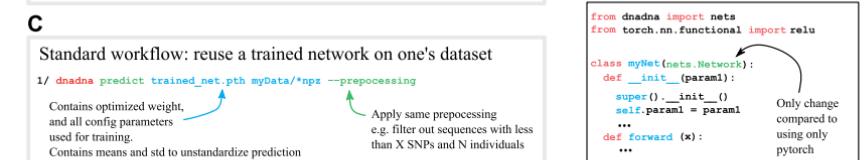
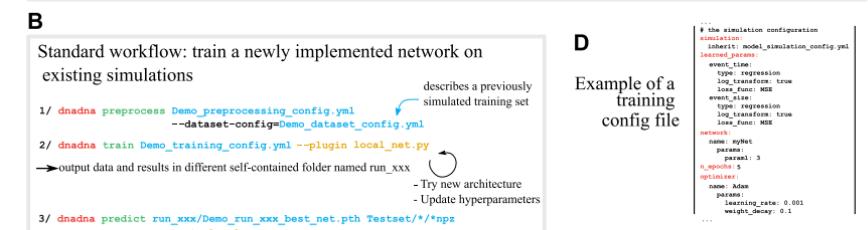
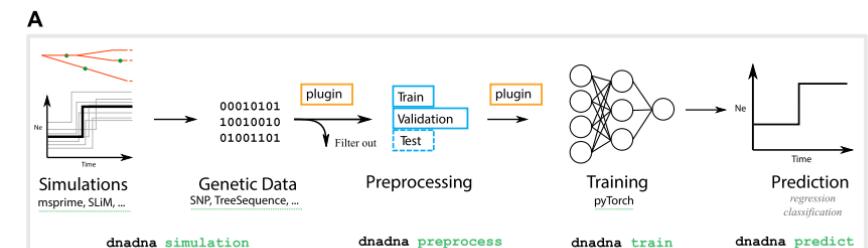
# Deep learning for population size history inference: Design, comparison and combination with approximate Bayesian computation

Théophile Sanchez | Jean Cury | Guillaume Charpiat | Flora Jay



## Genetics and population analysis dnadna: a deep learning framework for population genetics inference

Théophile Sanchez<sup>1†</sup>, Erik Madison Bray<sup>1†</sup>, Pierre Jobic<sup>1,2</sup>, Jérémie Guez<sup>1,3</sup>, Anne-Catherine Letournel<sup>1</sup>, Guillaume Charpiat<sup>1</sup>, Jean Cury <sup>1,4\*</sup>‡ and Flora Jay <sup>1,\*†</sup>



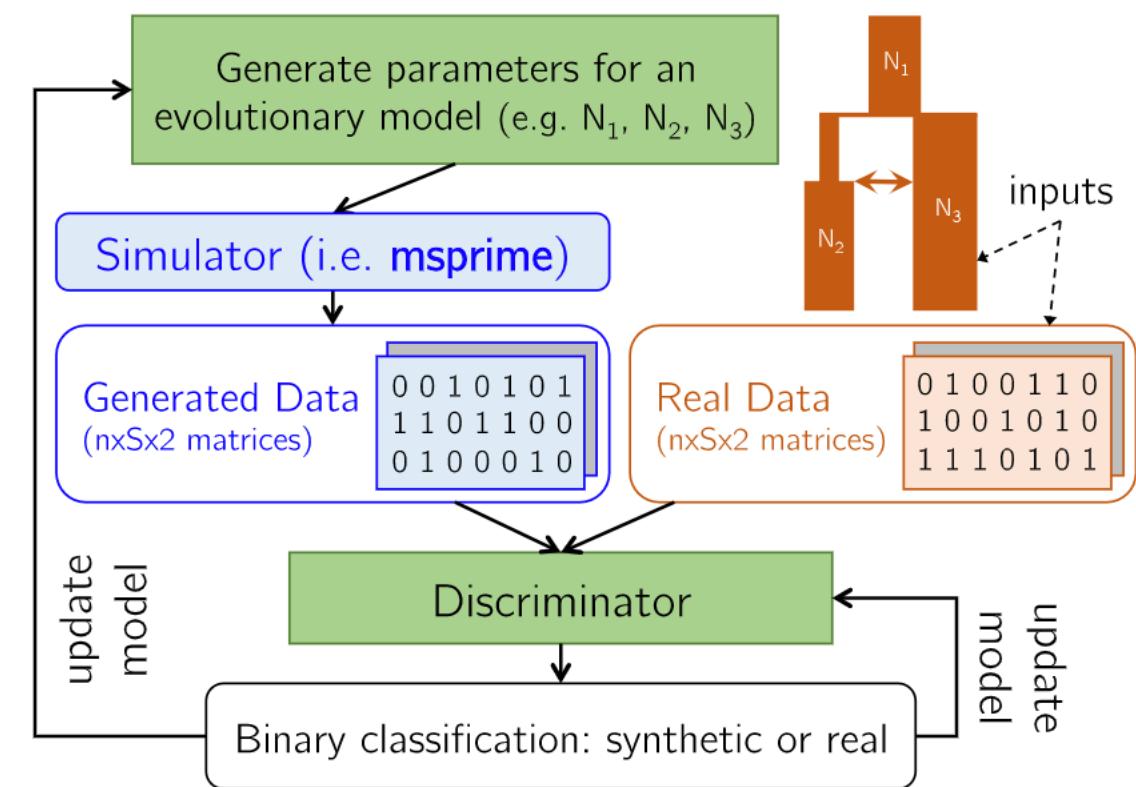
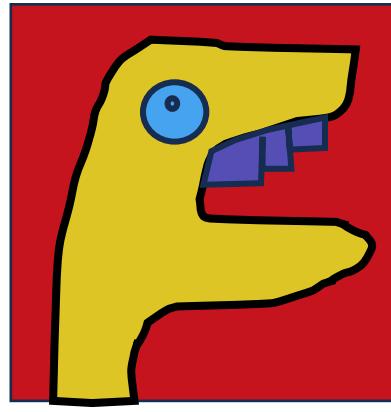
```
# the simulation configuration
# dnadna simulation
# dnadna preprocess
# dnadna train
# dnadna predict
```

```
# the training configuration
# dnadna train
# dnadna predict
# dnadna preprocess
```

```
# the network configuration
# dnadna predict
```

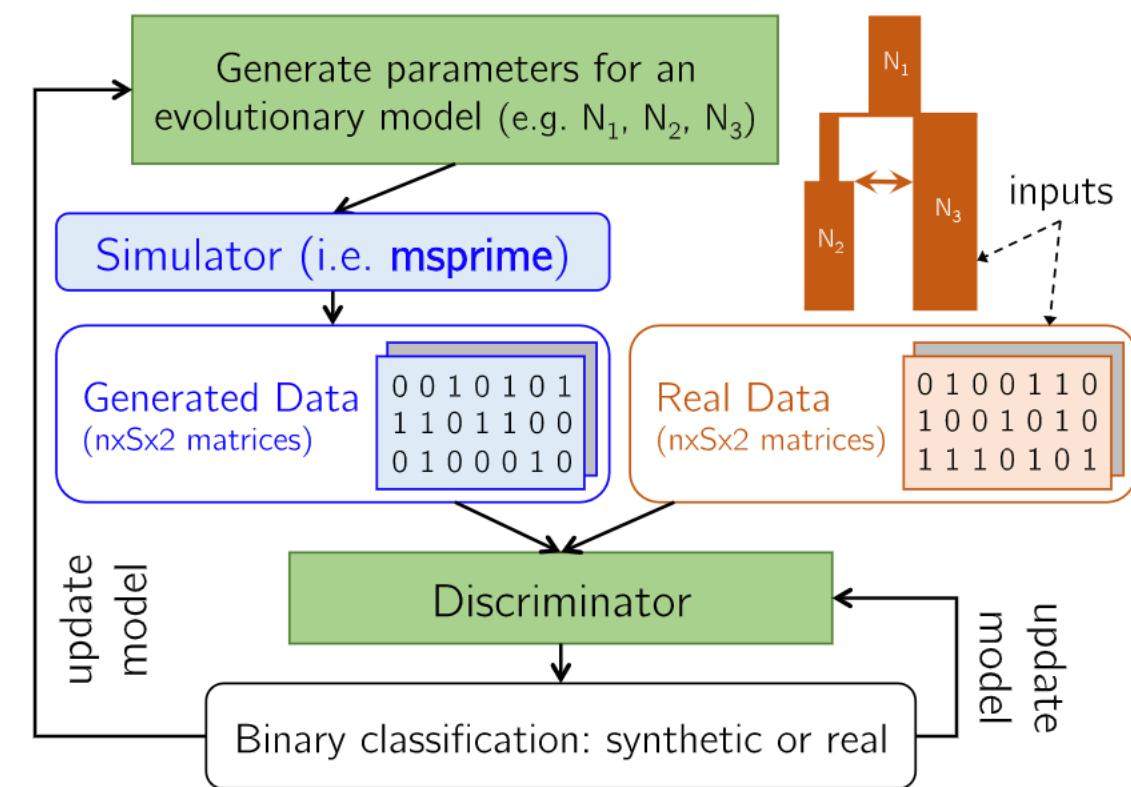
# Automatic inference of demographic parameters using generative adversarial networks

Zhanpeng Wang<sup>1</sup> | Jiaping Wang<sup>1</sup> | Michael Kourakos<sup>2</sup> | Nhunghoang<sup>2</sup> |  
Hyong Hark Lee<sup>2</sup> | Iain Mathieson<sup>3</sup> | Sara Mathieson<sup>1</sup> 



# Automatic inference of demographic parameters using generative adversarial networks

Zhanpeng Wang<sup>1</sup> | Jiaping Wang<sup>1</sup> | Michael Kourakos<sup>2</sup> | Nhung Hoang<sup>2</sup> |  
Hyong Hark Lee<sup>2</sup> | Iain Mathieson<sup>3</sup> | Sara Mathieson<sup>1</sup> 



# Peer Community Journal

Section: Evolutionary Biology

Research article

Published  
2024-03-18

Cite as

Kevin Korfmann, Thibaut Paul Patrick Sellinger, Fabian Freund, Matteo Fumagalli and Aurélien Tellier (2024) Simultaneous Inference of Past Demography and Selection from the Ancestral Recombination Graph under the Beta Coalescent, Peer Community Journal, 4: e33.

Kevin Korfmann, Thibaut Paul Patrick Sellinger, Fabian Freund, Matteo Fumagalli and Aurélien Tellier (2024) Simultaneous Inference of Past Demography and Selection from the Ancestral Recombination Graph under the Beta Coalescent, Peer Community Journal, 4: e33.

## Simultaneous Inference of Past Demography and Selection from the Ancestral Recombination Graph under the Beta Coalescent

Kevin Korfmann , Thibaut Paul Patrick Sellinger , Fabian Freund , Matteo Fumagalli , and Aurélien Tellier 

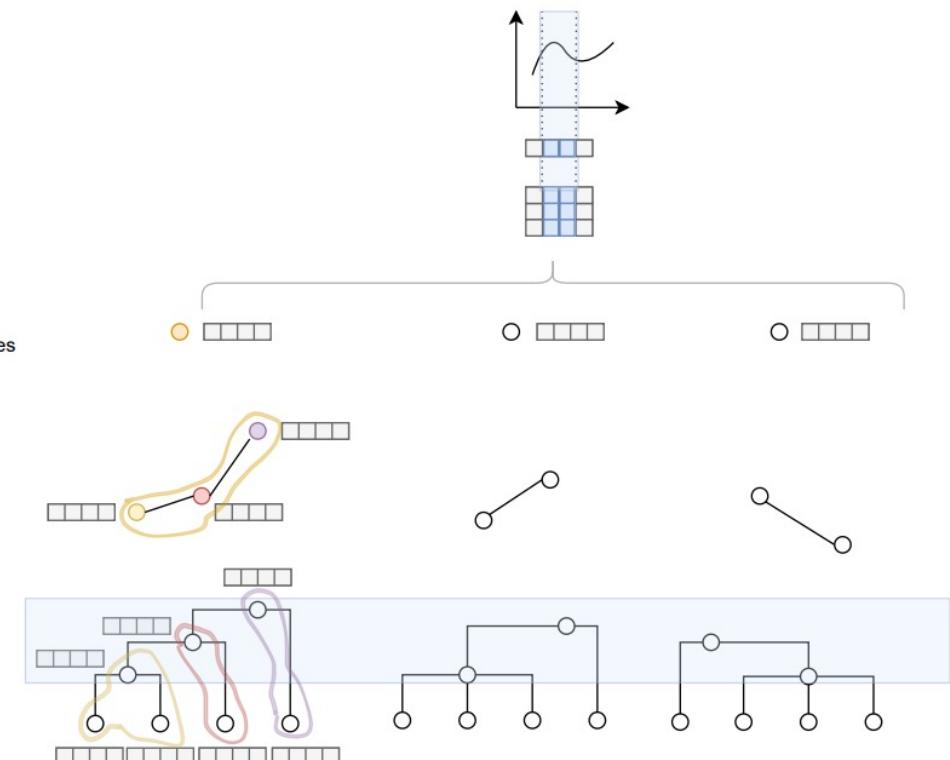
5. Visualization of inferred variables

4. Masking of time-relevant regions and column-wise mean

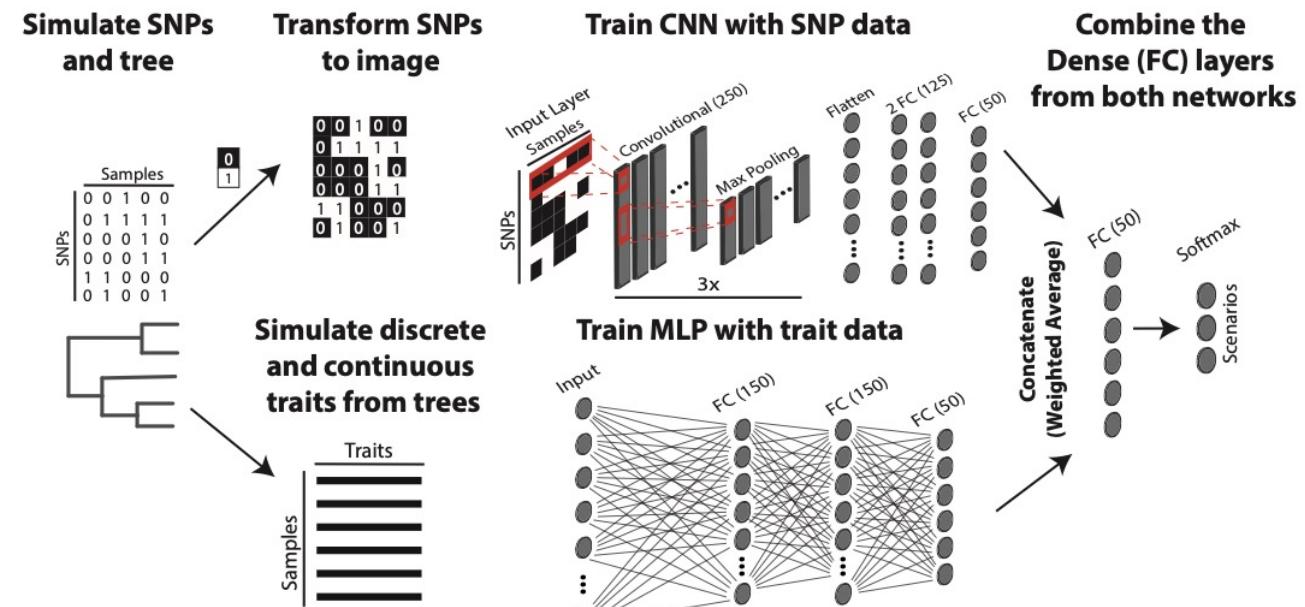
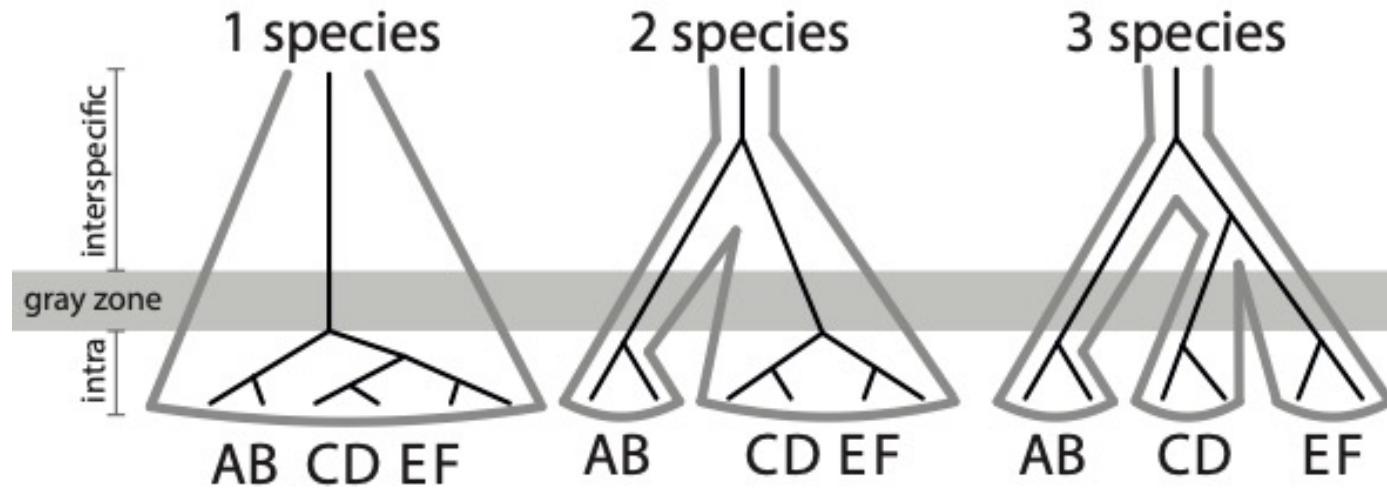
3. Last pooling step with feature vector containing inferred variables

2. Learned subgraph with updated feature vectors

1. Coalescent trees with feature vectors



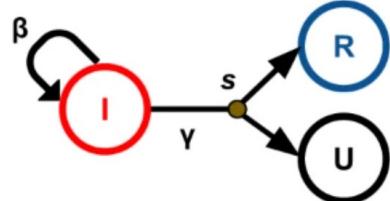
# Integrative Deep Learning species delimitation



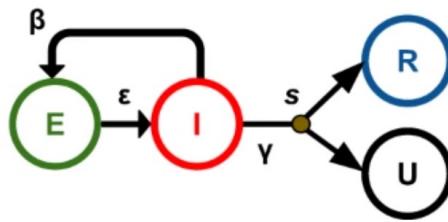
# Deep Learning for phylodynamics and macroevolution

Perez & Gascuel (in prep. for *Syst. Biol.*)

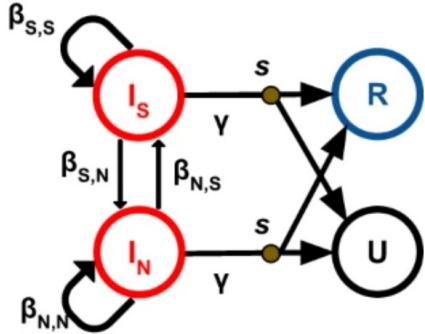
**BD**



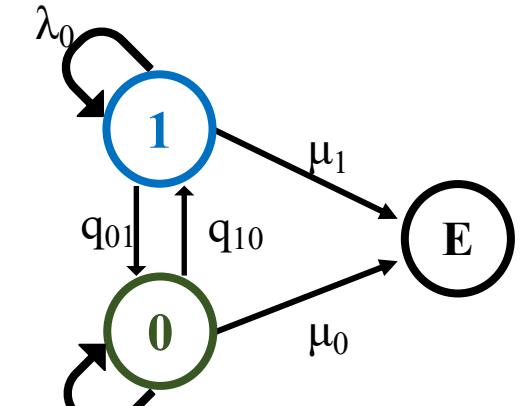
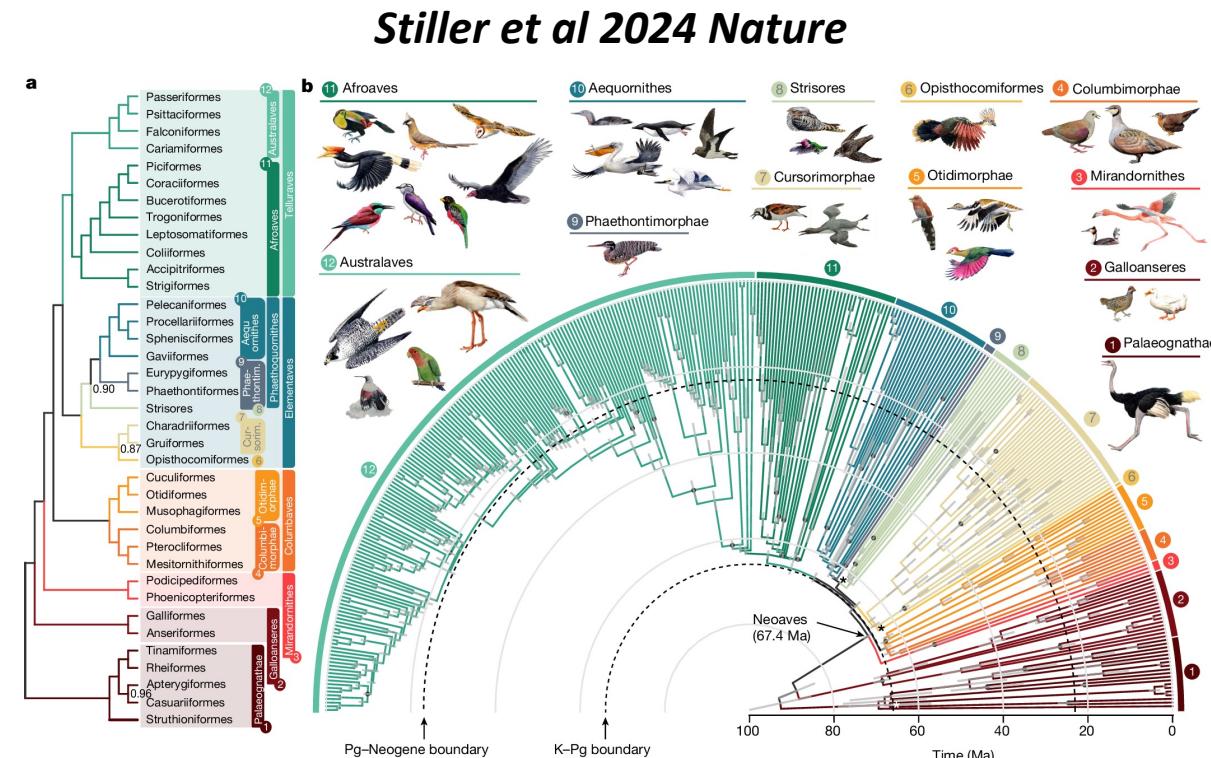
**BDEI**



**BDSS**



Voznica et al. (2022) Nat Comm



$$q = q_{01} = q_{10}$$

$$\epsilon = \mu_0 / \lambda_0 = \mu_1 / \lambda_1$$

What's next?

Your  
Project  
Here



# Population Genomics: background and tools

18 – 26 May 2017 | Napoli, Italy

**Part III: Wrapup.**

# Goals

- Conceive and simulate genetic data under competing demographic scenarios 
- Understand deep learning background and how a CNN works 
- Use CNN to detect regions with selective sweeps on real genomes 
- How to use deep learning to compare demographic scenarios 