

Deep learning methods in population genomics and phylogeography

Manolo Fernandez Perez

Department of Life Sciences, Imperial College London

@ManoloLearning 

manolofperez@gmail.com 

sites.google.com/site/manolofperez 

Goals

- Conceive and simulate genetic data under competing demographic scenarios
- Understand deep learning background and how a CNN works
- How to use deep learning to compare demographic scenarios



Program

Program

- ***Part I: - The building blocks of a CNN script.***
- ***Practical: Comparing demographic scenarios with deep learning.***
- ***Part II: Quick overview of other applications and future perspectives.***
- ***Part III: Wrapup.***

Part I: The building blocks of a CNN script



The diagram illustrates the building blocks of a CNN script. It starts with an 'Input Patch' of size 28x28. This is followed by a 'Conv' layer with a 5x5 kernel, producing 20 feature maps. A 'Max Pool' layer with a 2x2 kernel then reduces these to 10 feature maps. This sequence of convolution and max pooling layers is repeated, with the number of feature maps increasing to 50 in subsequent layers. The final output is passed through a 'Fully Connected' layer and a 'ReLU' activation function.

CNN Script

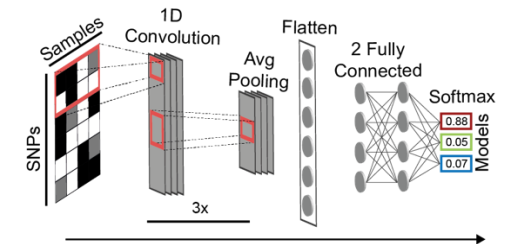


required python modules.



TensorFlow

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

CNN Script

Inputs:

Scenario 1



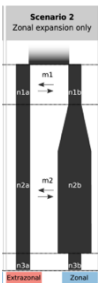
Samples

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

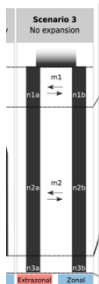
Scenario 2



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

Parameters

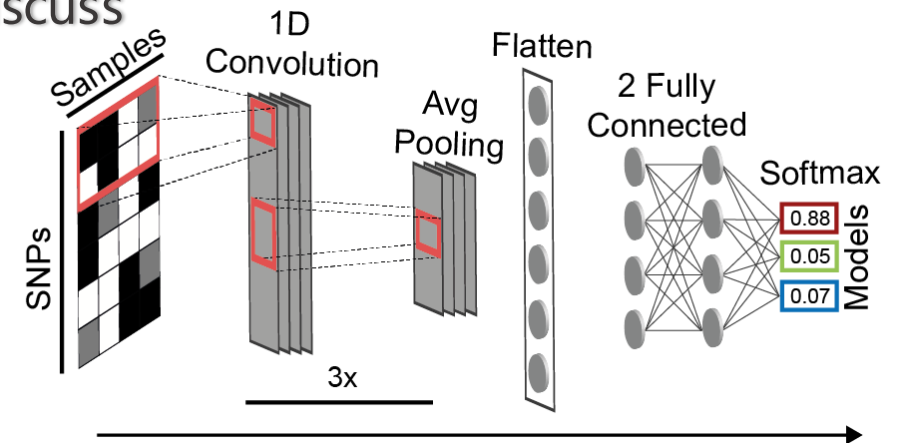
	Theta	T1	T2	T3	Ne
Sim1					
Sim2					
Sim3					
Sim4					

No of simulations

3-D Numpy array

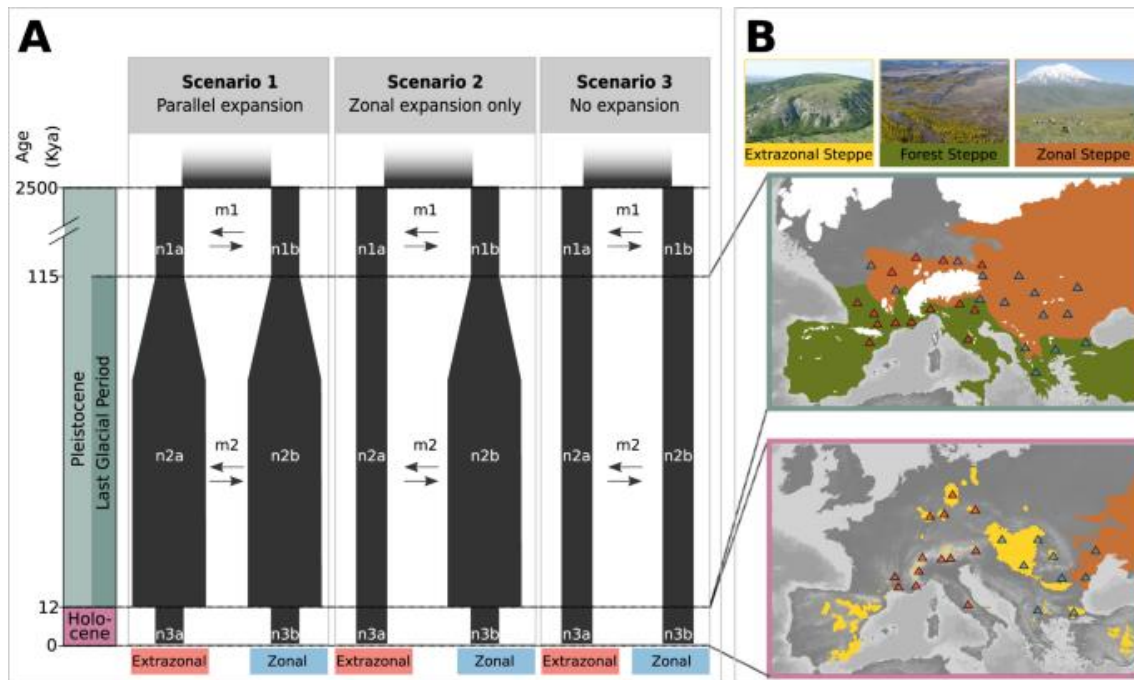
- Practical Exercise 1:

- Go through Section 1 of the notebook:
https://drive.google.com/file/d/1dA2GWzG5I7m0NbInxtE7-LL_Os0Ytmql/view?usp=sharing
- We will 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.

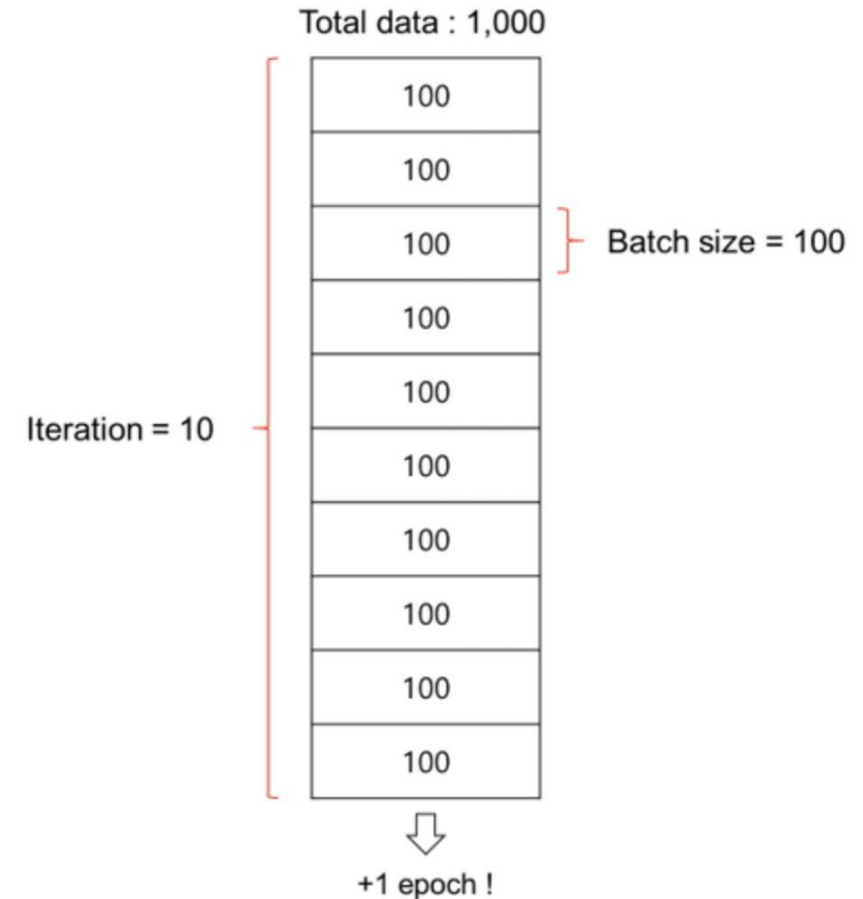



```
# Define parameters for the CNN run.
batch_size = 128
### how much iterations to train the network
epochs = 50
```

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



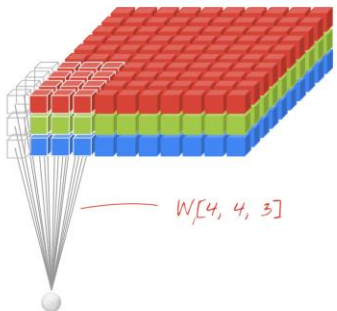
CNN Script



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

CNN Script

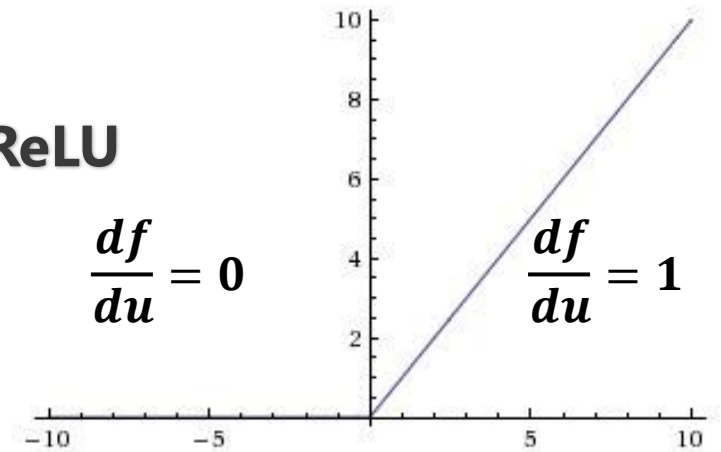
```
# Define the CNN architecture.
def create_cnn(xtest):
    inputShape = (xtest.shape[1], xtest.shape[2])
    ## image size. images need to have EXACTLY the same size
    inputs = Input(shape=inputShape)
    x = inputs
    ## 1D convolution - less computational intensive and is also invariant to the samples order;
    x = Conv1D(256, 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

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

$$\frac{df}{du} = 1$$



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

Kirschner et al. (2022) *Nat Comm*

```
# 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
```

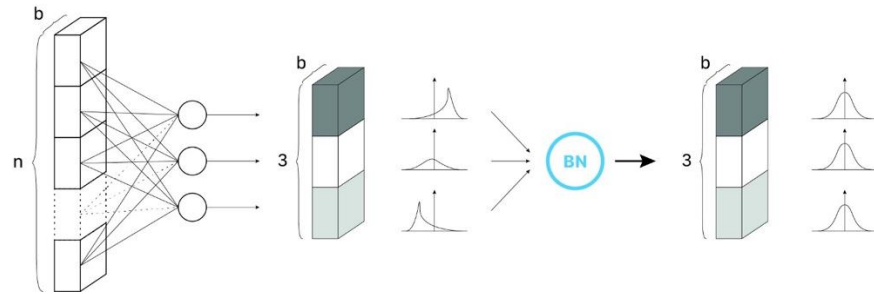
```
    ## 1D convolution - less computational intensive and is also invariant to the samples order;
```

```
    x = Conv1D(256, kernel_size=2, activation='relu', input_shape=(xtest.shape[1], xtest.shape[2]))(x)
```

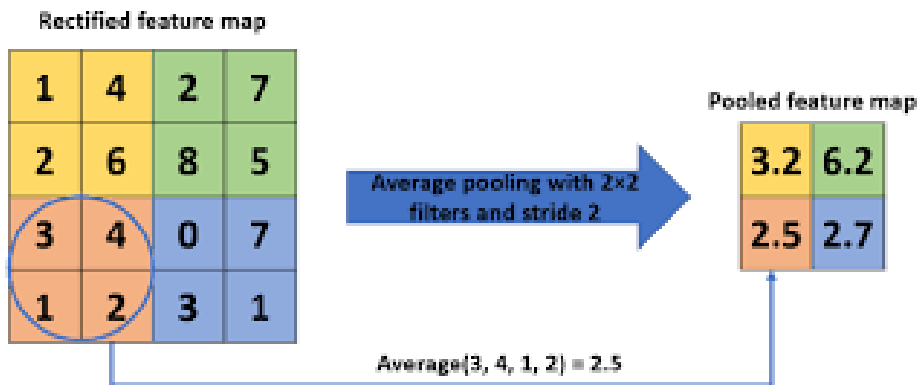
```
    ### Enables the network to learn more complex features / shapes.
```

```
    x = AveragePooling1D(pool_size=2)(x)
```

```
    x = BatchNormalization()(x)
```

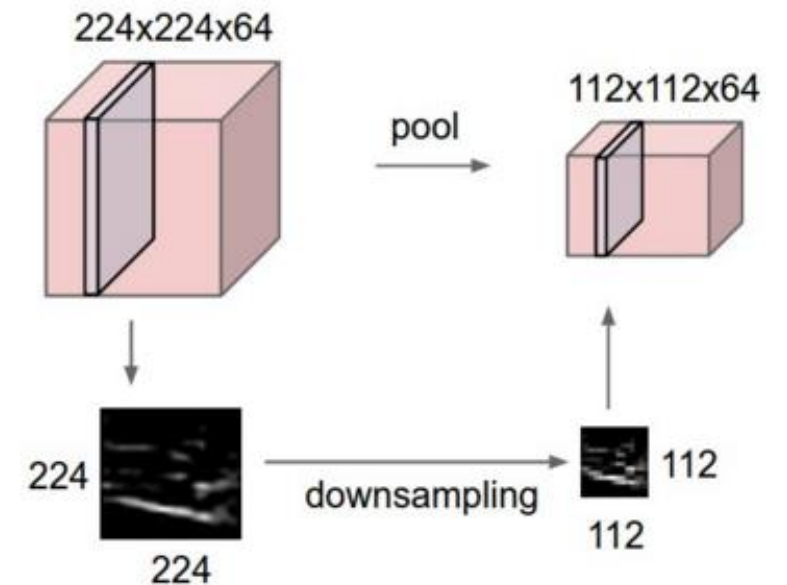


<https://towardsdatascience.com/batch-normalization-in-3-levels-of-understanding-14c2da90a338>



Gholamalinezhad & Khosravi (2020) *arXiv*

CNN Script

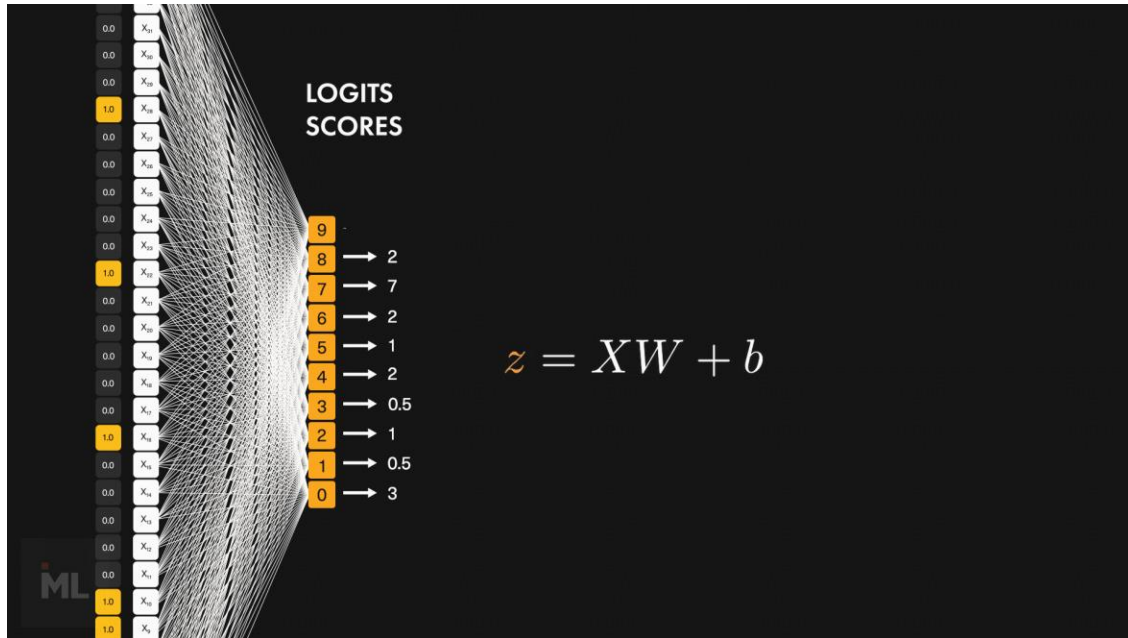


https://leonardoaraujosantos.gitbook.io/artificial-intelligence/machine_learning/deep_learning/pooling_layer

Kirschner et al. (2022) *Nat Comm*

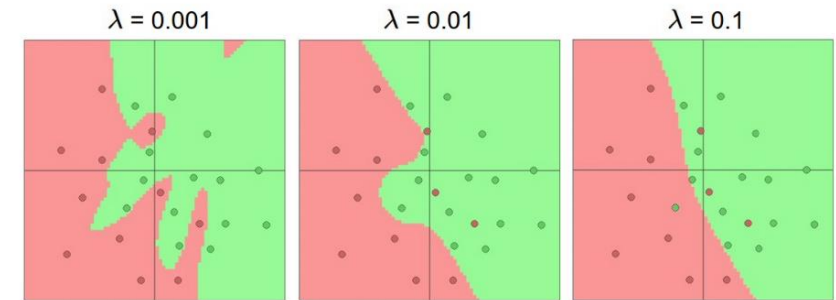
Linearising the image as in the initial step.

```
x = Flatten()(x)
x = Dense(128, activation='relu')(x)
x = Dropout(0.5)(x)
x = Dense(128, activation='relu')(x)
x = Dropout(0.5)(x)
x = Dense(num_classes, activation="softmax")(x)
```

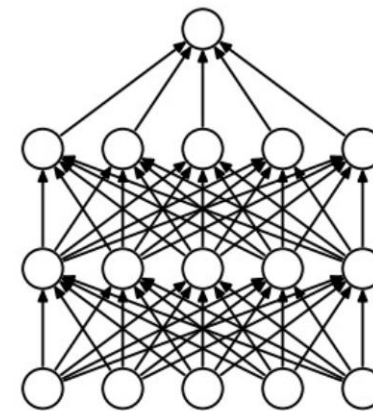


CNN Script

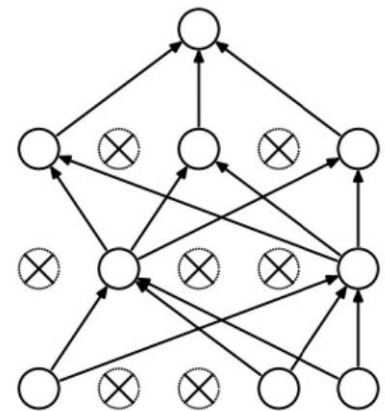
Regularization



Options:



(a) Standard Neural Net



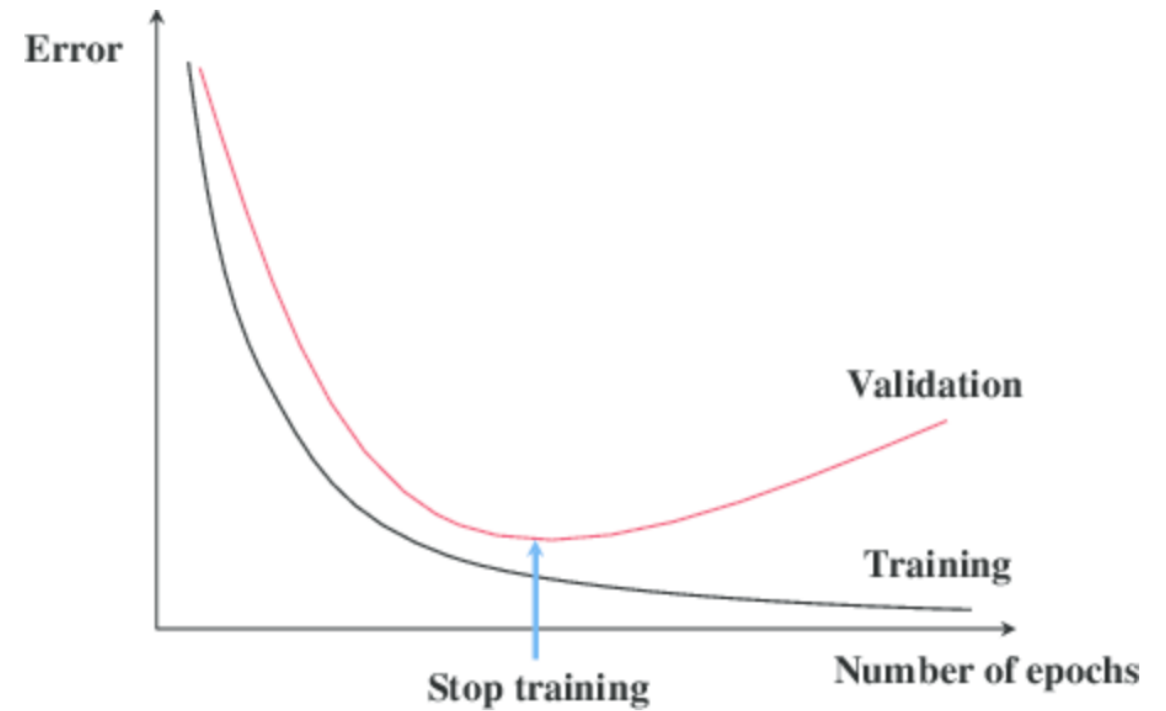
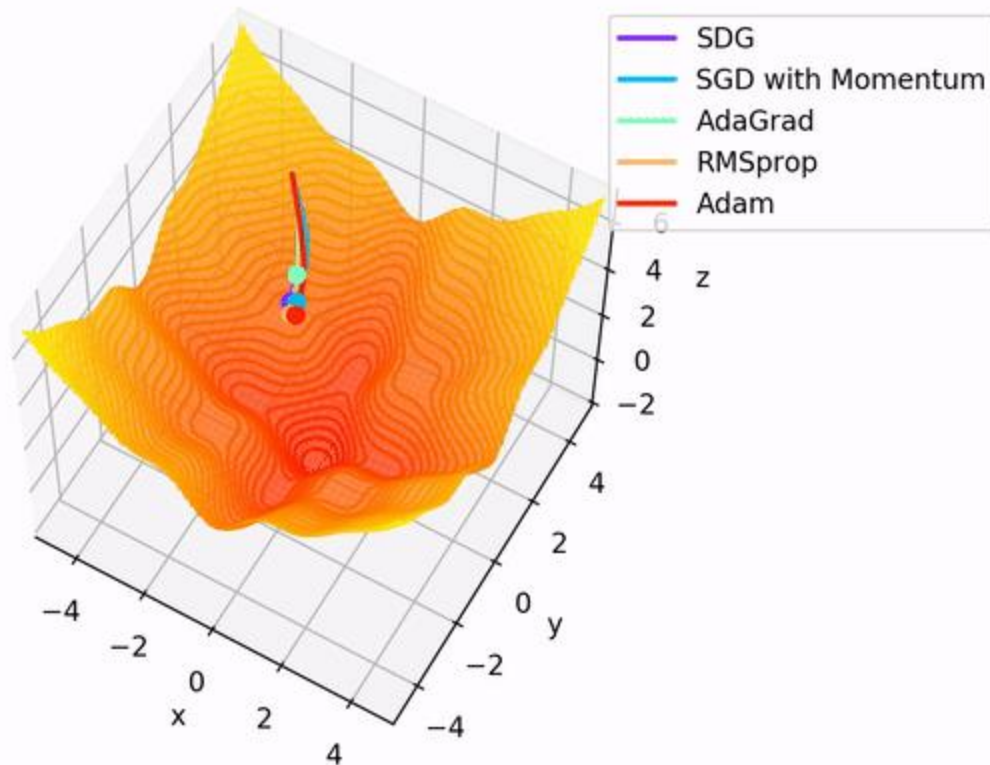
(b) After applying dropout.

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=25, verbose=0, mode='max', restore_best_weights=True)
### stop training when validation error increases (wait 25 epochs to see if there is any improvement).
```

Optimizer Comparison



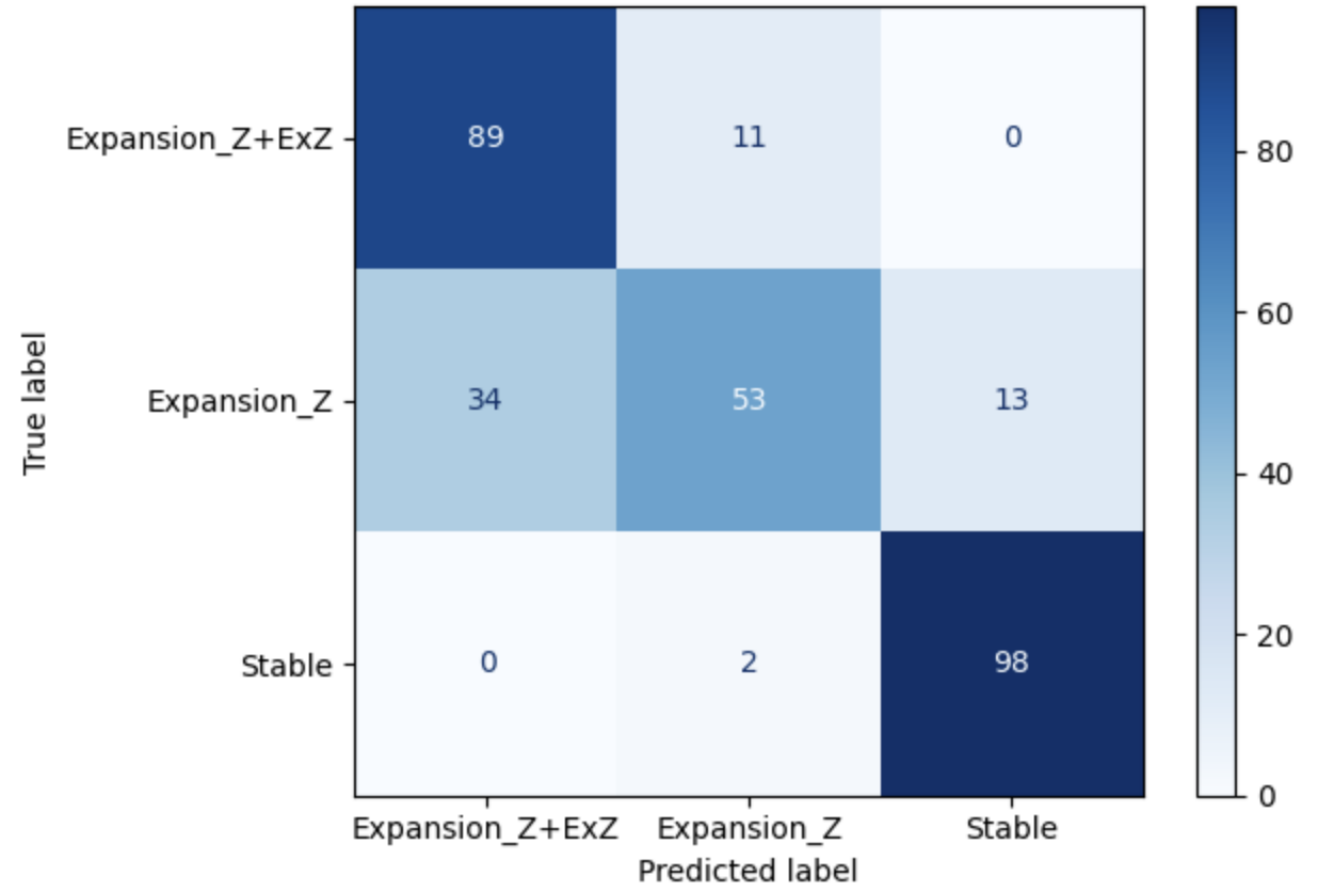
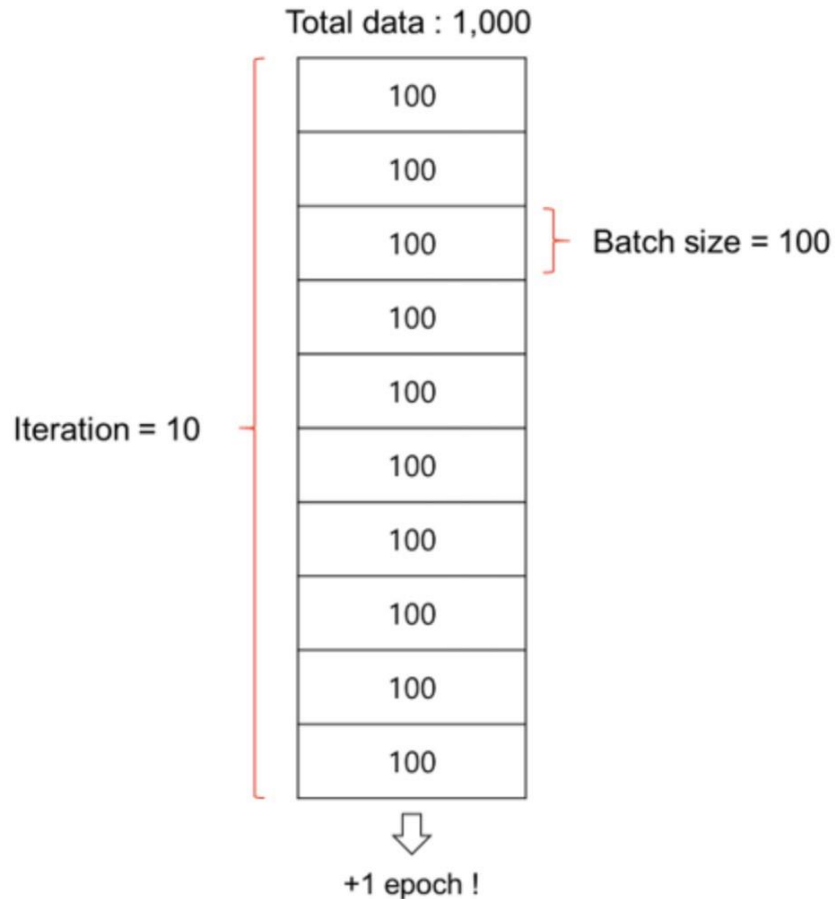
<https://towardsdatascience.com/complete-guide-to-adam-optimization-1e5f29532c3d>

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

#Run the CNN

```
history = model.fit(xtrain, ytrain, batch_size=batch_size,  
                    epochs=epochs,  
                    verbose=1,  
                    validation_data=(xval, yval), callbacks=[earlyStopping])
```

CNN Script



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

The background is a dark, textured surface with various colorful, pixelated lines and shapes. A prominent yellow line with black and white segments runs diagonally across the upper half. Below it, a blue line with white segments runs diagonally. In the lower half, there are several rectangular blocks with colorful patterns, including one with a blue and white grid and another with a red and white grid. The overall aesthetic is digital and abstract.

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

Deep Learning in Population Genetics

Kevin Korfmann¹, Oscar E. Gaggiotti², and Matteo Fumagalli ^{3,*}

¹Professorship for Population Genetics, Department of Life Science Systems, Technical University of Munich, Germany

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

³Department of Biological and Behavioural Sciences, Queen Mary University of London, UK

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

Accepted: 16 January 2023

The Unreasonable Effectiveness of Convolutional Neural Networks in Population Genetic Inference

Lex Flagel,^{1,2} Yaniv Brandvain,² and Daniel R. Schrider^{*,3}

¹Monsanto Company, Chesterfield, MO

²Department of Plant and Microbial Biology, University of Minnesota, St. Paul, MN

³Department of Genetics, University of North Carolina, Chapel Hill, NC

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

Associate editor: Yuseob Kim

Review Article | Published: 04 September 2023

Harnessing deep learning for population genetic inference

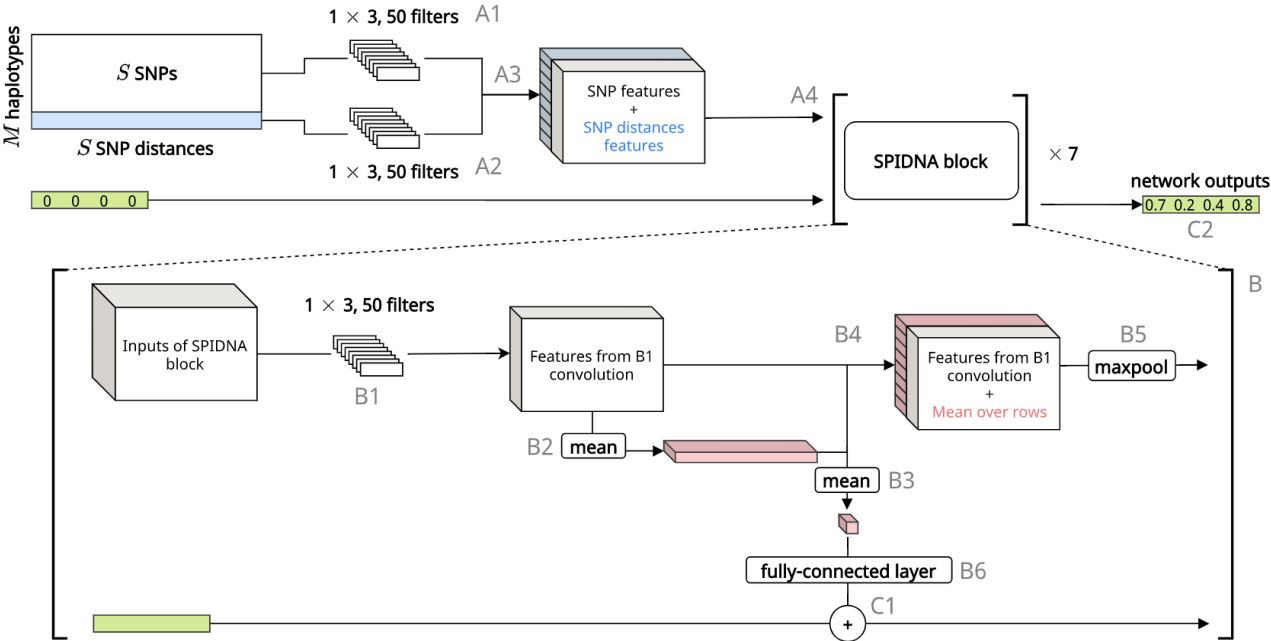
[Xin Huang](#) , [Aigerim Rymbekova](#), [Olga Dolgova](#), [Oscar Lao](#)  & [Martin Kuhlwilm](#) 

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

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

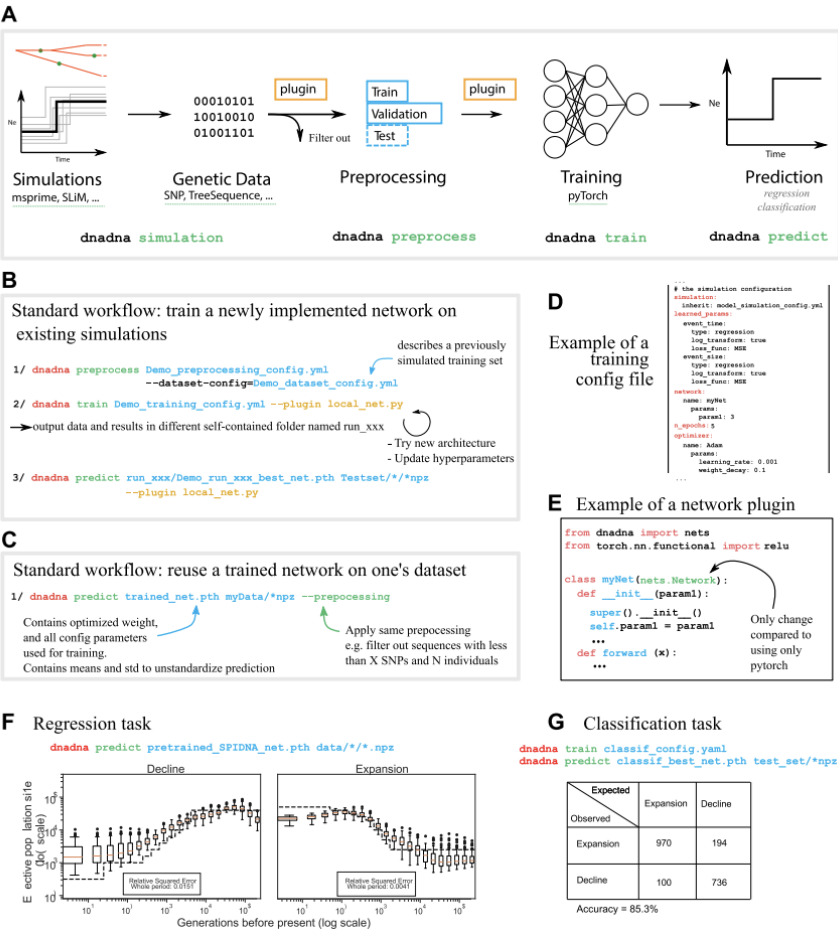
Deep learning for population size history inference: Design, comparison and combination with approximate Bayesian computation

Théophile Sanchez  | Jean Cury  | Guillaume Charpiat | Flora Jay 



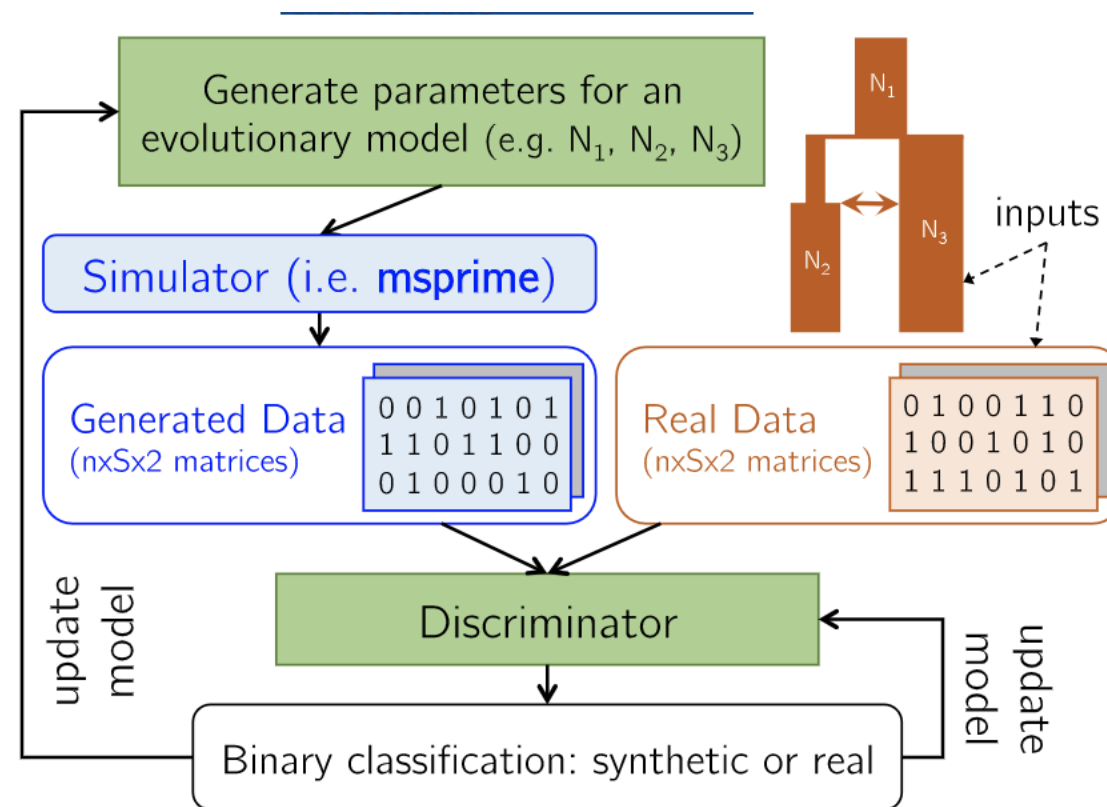
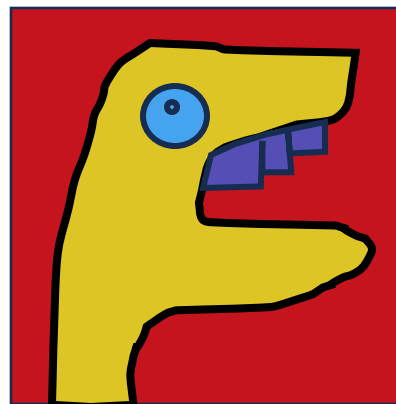
Genetics and population analysis **dnadna: a deep learning framework for population genetics inference**

Théophile Sanchez^{1†}, Erik Madison Bray^{1†}, Pierre Jobic^{1,2}, Jérémy Guez^{1,3}, Anne-Catherine Letournel¹, Guillaume Charpiat¹, Jean Cury ^{1,4*‡} and Flora Jay ^{1*‡}



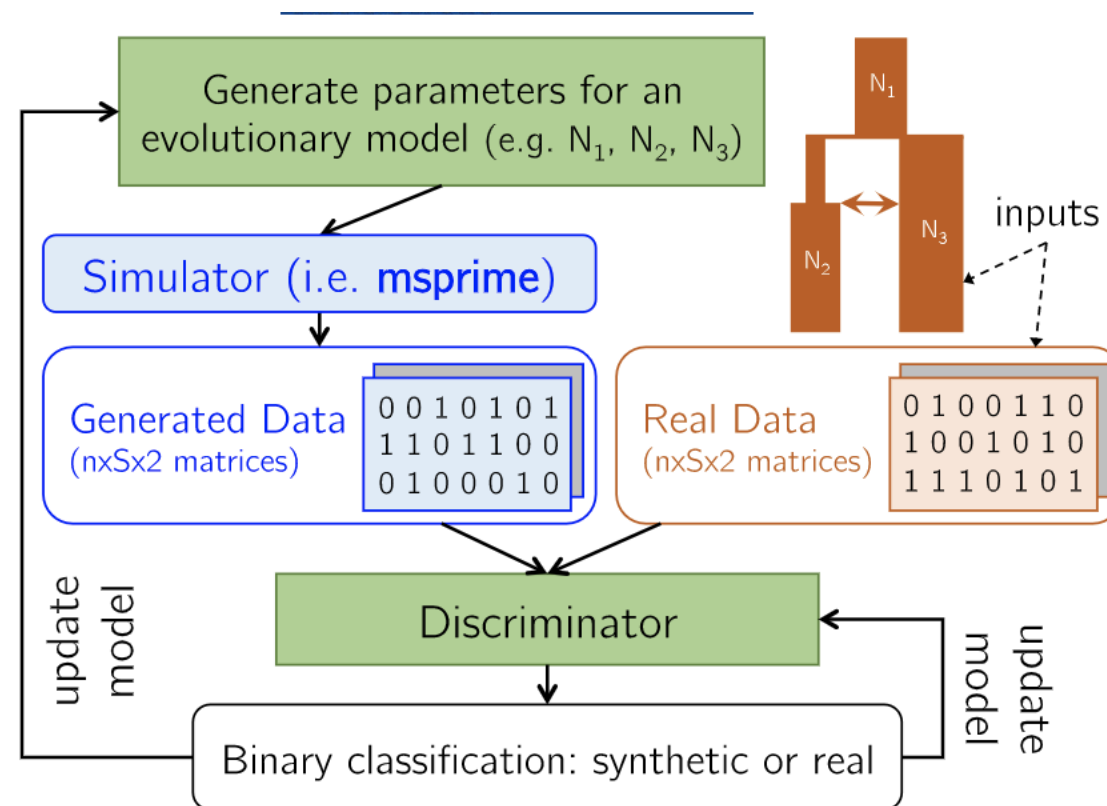
Automatic inference of demographic parameters using generative adversarial networks

Zhanpeng Wang¹ | Jiaping Wang¹ | Michael Kourakos² | Nhung Hoang² |
Hyong Hark Lee² | Iain Mathieson³ | Sara Mathieson¹ 

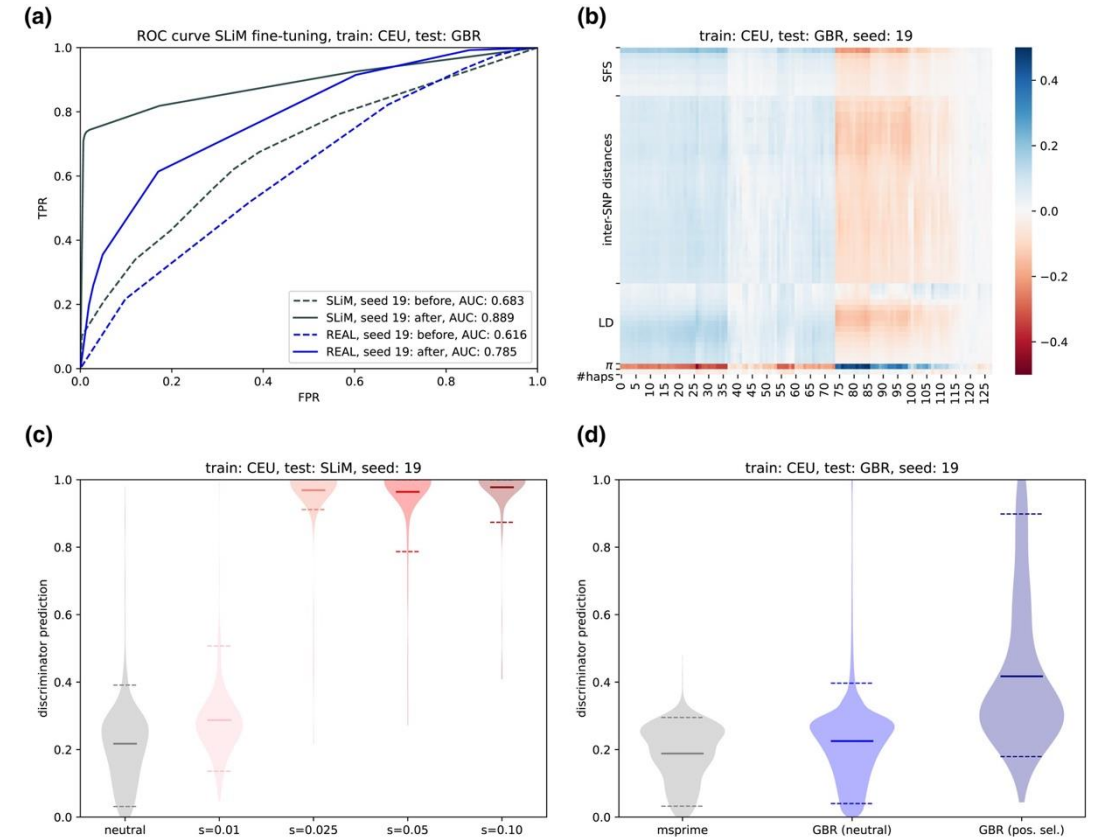
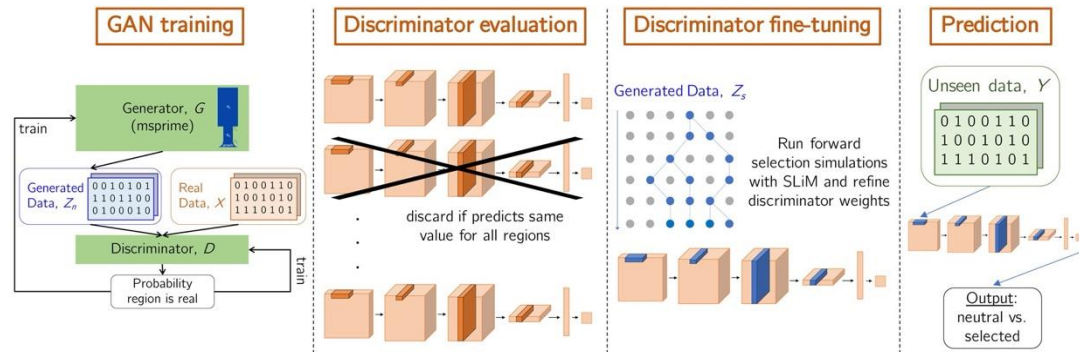


Automatic inference of demographic parameters using generative adversarial networks

Zhanpeng Wang¹ | Jiaping Wang¹ | Michael Kourakos² | Nhung Hoang² |
Hyong Hark Lee² | Iain Mathieson³ | Sara Mathieson¹ 



Neural Network for Genomic data



Peer Community Journal

Section: Evolutionary Biology

Research article

Published
2024-03-18

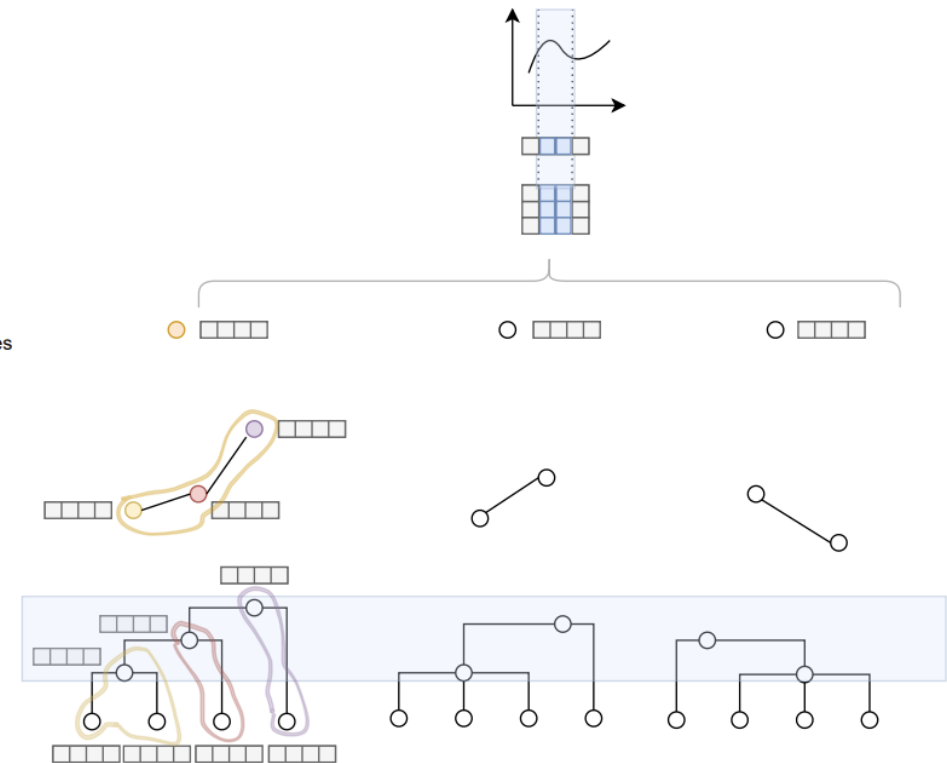
Cite as

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.

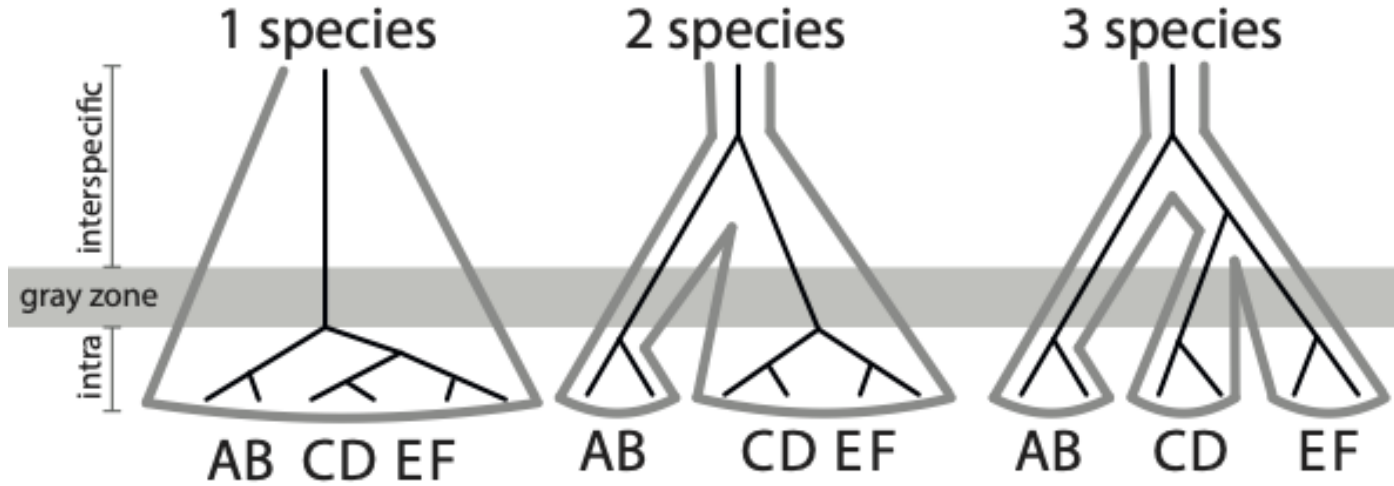
Simultaneous Inference of Past Demography and Selection from the Ancestral Recombination Graph under the Beta Coalescent

Kevin Korfmann ,^{#,1}, Thibaut Paul Patrick Sellinger ,^{#,2,1}, Fabian Freund ,^{3,4}, Matteo Fumagalli ,^{5,6}, and Aurélien Tellier ,¹

1. Coalescent trees with feature vectors
2. Learned subgraph with updated feature vectors
3. Last pooling step with feature vector containing inferred variables
4. Masking of time-relevant regions and column-wise mean
5. Visualization of inferred variables



Integrative Deep Learning species delimitation

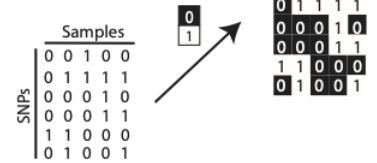


Simulate SNPs
and tree

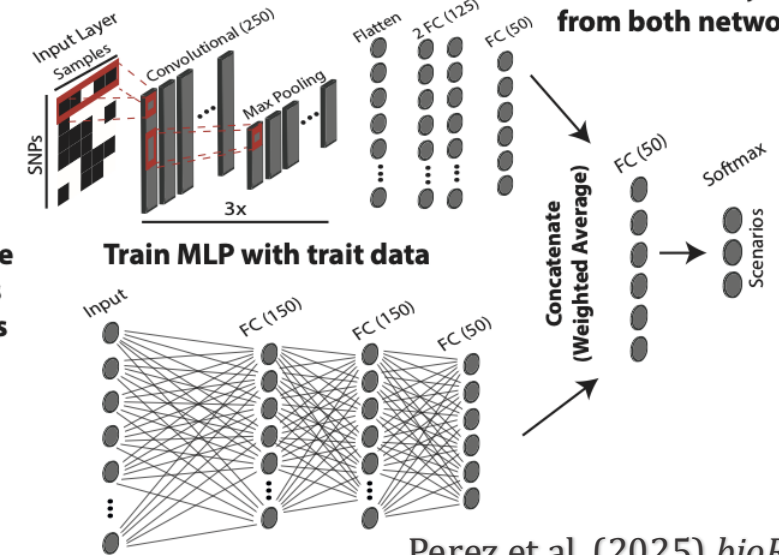
Transform SNPs
to image

Train CNN with SNP data

Combine the
Dense (FC) layers
from both networks



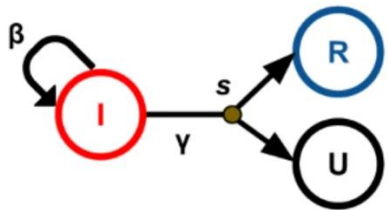
Simulate discrete
and continuous
traits from trees



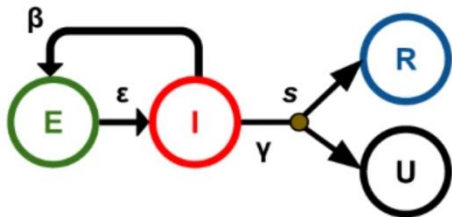
Deep Learning for **phylogenetics** and **macroevolution**

Perez & Gascuel (2025) *bioRxiv*

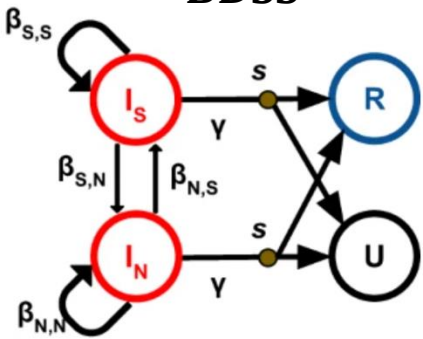
BD



BDEI

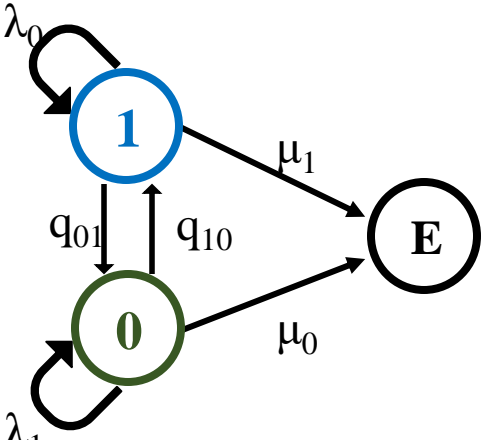
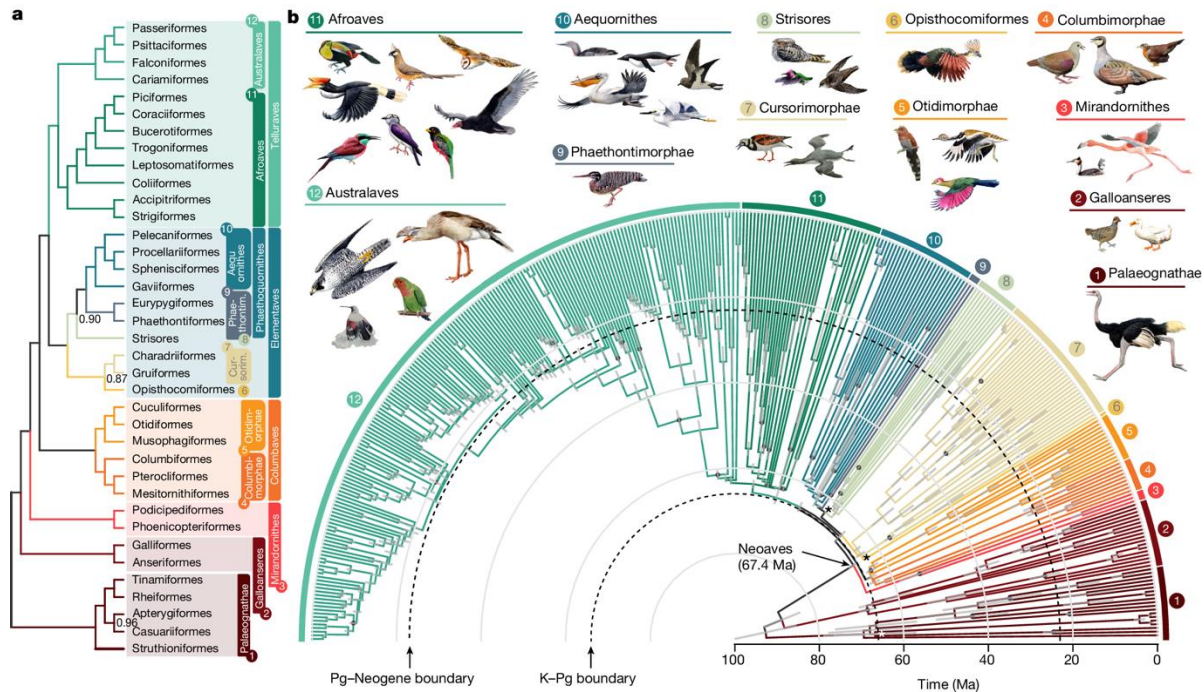


BDSS



Voznica et al. (2022) *Nat Comm*

Stiller et al 2024 Nature



$q = q_{01} = q_{10}$
 $\epsilon = \mu_0 / \lambda_0 = \mu_1 / \lambda_1$

Perez et al. (2025) *bioRxiv*

What's next?

Your
Project
Here

Part III: Wrapup.

Goals

- Conceive and simulate genetic data under competing demographic scenarios
- Understand deep learning background and how a CNN works
- How to use deep learning to compare demographic scenarios

