

# AI for Genomics: from CNNs and LSTMs to Transformers

Nikolay Oskolkov, Group Leader (PI) at LIOS, Riga, Latvia

Physalia course, 09.09.2025

Session 3a: CNNs and RNNs applications for genomics and ancient genomics



@NikolayOskolkov



@oskolkov.bsky.social



Personal homepage:  
<https://nikolay-oskolkov.com>

Topics we'll cover in this session:

- 1) Introduction to ancient genomics: contamination problem
- 2) Decontaminating ancient data with neural networks
- 3) Recognizing introgressed regions with neural networks
- 4) One-hot-encoding of biological nucleotide sequences
- 5) Training CNNs and RNNs on genomics data

## PERSPECTIVE

<https://doi.org/10.1038/s41588-018-0295-5>

## A primer on deep learning in genomics

James Zou<sup>1,2,3\*</sup>, Mikael Huss<sup>4,5</sup>, Abubakar Abid<sup>3</sup>, Pejman Mohammadi<sup>6,7</sup>, Ali Torkamani<sup>8,9</sup> and Amalio Teleni<sup>10,6,\*</sup>

**Deep learning methods are a class of machine learning techniques capable of identifying highly complex patterns in large datasets.** Here, we provide a perspective and primer on deep learning applications for genome analysis. We discuss successful applications in the fields of regulatory genomics, variant calling and pathogenicity scores. We include general guidance for how to effectively use deep learning methods as well as a practical guide to tools and resources. This primer is accompanied by an interactive online tutorial.

Deep learning has made impressive recent advances in applications ranging from computer vision to natural-language processing. This primer discusses the main categories of deep learning methods and provides suggestions for how to effectively use deep learning in genomics. The primer is intended for bioinformaticians who are interested in applying deep learning approaches, and for genomists and general biomedical researchers who seek a high-level understanding of this rapidly evolving field. Computer scientists may also use the primer as an introduction to the exciting applications of deep learning in genomics. However, we do not provide a survey of deep learning in the biomedical field, which has been broadly covered in recent reviews<sup>1–3</sup>. This paper is accompanied by an interactive tutorial that we have created for interested readers to build a convolutional neural network to discover DNA-binding motifs (see URLs).

## Deep learning as a class of machine learning methods

Machine learning techniques have been extensively used in genomics research<sup>4–6</sup>. Machine learning tasks fall within two major categories: supervised and unsupervised. In supervised learning, the goal is predicting the label (classification) or response (regression) of each data point by using a provided set of labeled training examples. In unsupervised learning, such as clustering and principal component analysis, the goal is learning inherent patterns within the data themselves.

The ultimate goal in many machine learning tasks is to optimize model performance not on the available data (training performance), but instead on independent datasets (generalization performance). With this goal, data are randomly split into at least three subsets: training, validation and test sets. The training set is used for learning the model parameters (detailed discussion on parameter optimization in ref.<sup>7</sup>), the validation set is used to select the best model, and the test set is kept aside to estimate the generalization performance (Fig. 1). Machine learning must reach an appropriate balance between model flexibility and the amount of training data. An overly simple model will underfit and fail to let the data ‘speak’; an overly flexible model will overfit to spurious patterns in the training data and will not generalize.

Large neural networks, a main form of deep learning, are a class of machine learning algorithms that can make predictions and perform dimensionality reduction. The key difference between deep

learning and standard machine learning methods used in genomics—e.g., support vector machine and logistic regression—is that deep learning models have a higher capacity and are much more flexible. Typical deep learning models have millions of trainable parameters. However, this flexibility is a double-edged sword. With appropriately curated training data, deep learning can automatically learn features and patterns with less expert handcrafting. It also requires greater care to train on and to interpret the underlying biology. Box 1 summarizes the main messages of this primer on how to effectively use deep learning in genomics.

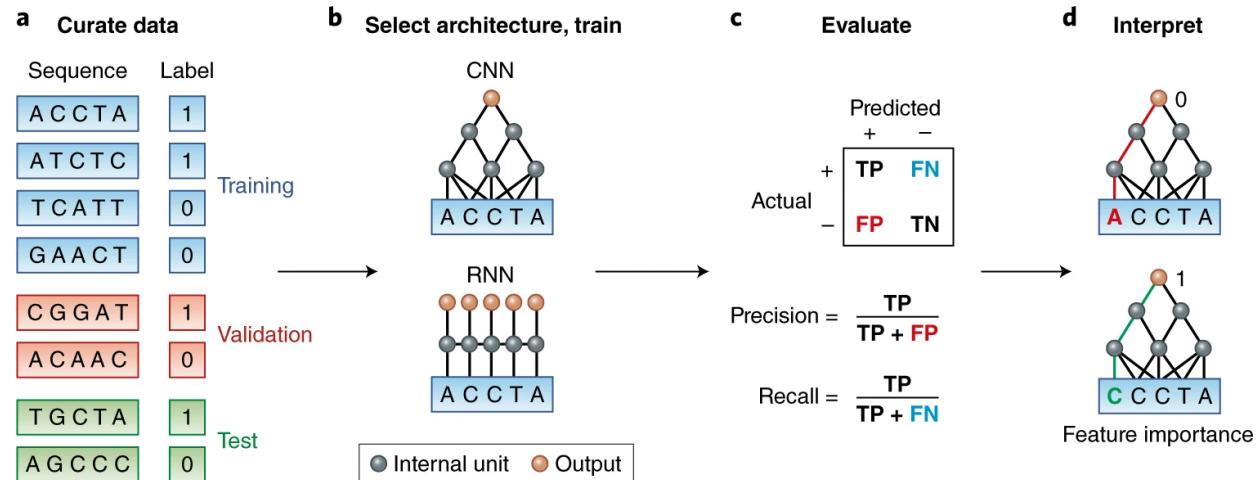
## Setting up deep learning

Deep learning is an umbrella term that refers to the recent advances in neural networks and the corresponding training platforms (e.g., TensorFlow and PyTorch). The starting point of a neural network is an artificial neuron, which takes as input a vector of real values and computes the weighted average of these values followed by a nonlinear transformation, which can be a simple threshold<sup>8</sup>. The weights are the parameters of the model that are learned during training. The power of neural networks stems from individual neurons being highly modular and composable, despite their simplicity<sup>9</sup>. The output of one neuron can be directly fed as input into other neurons. By composing neurons together, a neural network is created.

The input into a neural network is typically a matrix of real values. In genomics, the input might be a DNA sequence, in which the nucleotides A, C, T and G are encoded as [1,0,0,0], [0,1,0,0], [0,0,1,0] and [0,0,0,1]. Neurons that directly read in the data input are called the first, or input, layer. Layer two consists of neurons that read the outputs of layer one, and so on for deeper layers, which are also referred to as hidden layers. The output of the neural network is the prediction of interest, e.g., whether the input DNA is an enhancer. Box 2 describes key terms and concepts in deep learning.

There are three common families of architectures for connecting neurons into a network: feed-forward, convolutional and recurrent. Feed-forward is the simplest architecture<sup>10</sup>. Every neuron of layer  $i$  is connected only to neurons of layer  $i + 1$ , and all the connection edges can have different weights. Feed-forward architecture is suitable for generic prediction problems when there are no special relations among the input data features.

In a convolutional neural network (CNN), a neuron is scanned across the input matrix, and at each position of the input, the CNN



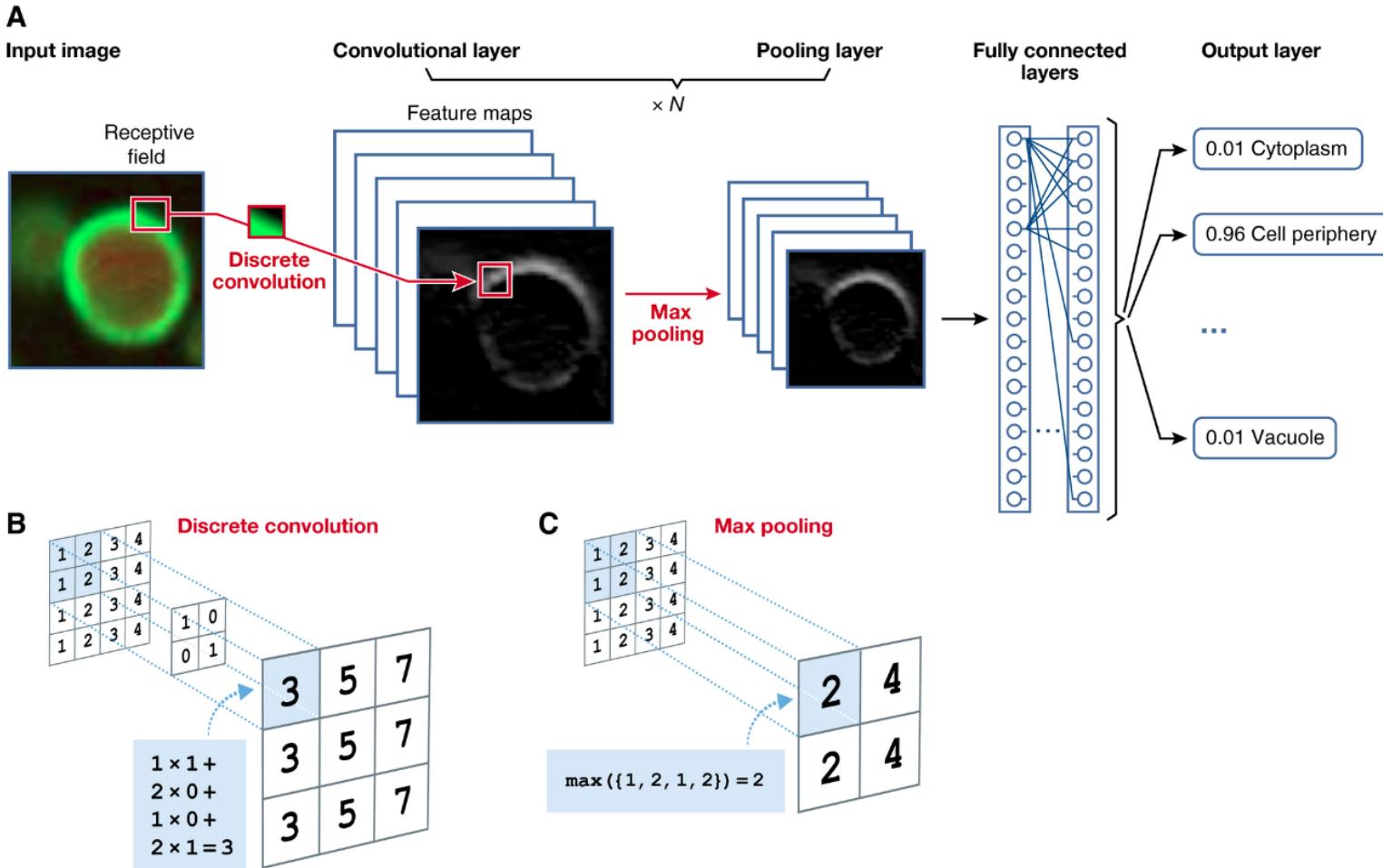
## One-dimensional convolutional neural networks (CNNs)

and Recurrent Neural Networks (example: LSTMs) are

suitable for working with sequential data. Therefore, can

be used for encoding DNA nucleotide sequences.

<sup>1</sup>Department of Biomedical Data Science, Stanford University, Palo Alto, CA, USA. <sup>2</sup>Chan-Zuckerberg Biohub, San Francisco, CA, USA. <sup>3</sup>Department of Electrical Engineering, Stanford University, Palo Alto, CA, USA. <sup>4</sup>Peltarion, Stockholm, Sweden. <sup>5</sup>Department of Learning, Informatics, Management and Ethics, Karolinska Institutet, Stockholm, Sweden. <sup>6</sup>Scripps Research Translational Institute, La Jolla, CA, USA. <sup>7</sup>Department of Integrative Structural and Computational Biology, The Scripps Research Institute, La Jolla, CA, USA. \*e-mail: [jamesz@stanford.edu](mailto:jamesz@stanford.edu); [ateleni@scripps.edu](mailto:ateleni@scripps.edu)



# Introduction to ancient genomics



TTCAGGAAAACCGATAAAAATTCTTATTGGGGGAGGGGCTCAAACAAGAAAATAATCAACAAGTGGTGTCCAGAGTGGAGGCCAGGG  
CCTCCTGGGGACAGCAAGACTGCCTGGGGAGCGGGAAAGCAGCTCCCCGTCTGGGGGAGGCCTGGAGGGGAGGCTGGACCA  
CAGGAGGGCAGCGCCCTGGGCAACCCTATGTAGATGAAGCTGCCGGAGAGGATCAAAGAACCAAGACAGGAGGAAAGAGGGCGGTGAAG  
TCTCCTGTCTCATCCCTTAGGAAGCCTGAGGAGATGGTAAGGGCATTAGAAGCCTCGAACCCCAGGGCAGAGGATAGTTGTAGGCA

Peculiarity: very few samples with little amounts of DNA, hence hard for statistical analysis



[About us](#) [Contact](#) [Events](#) [News](#) [Press](#) [For staff](#) [Vacancies](#) [På svenska](#) [Search](#)

**SciLifeLab** INFRASTRUCTURE SERVICES ▾ RESEARCH ▾ EDUCATION ▾ COLLABORATION ▾ DATA & TOOLS ▾

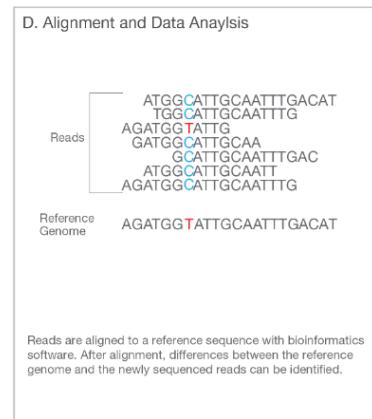
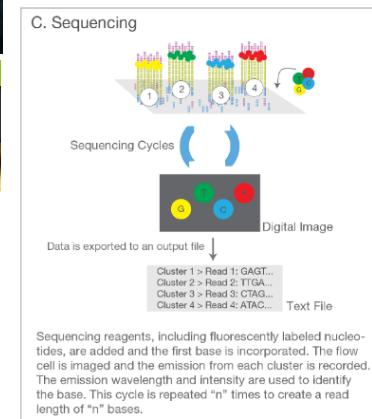
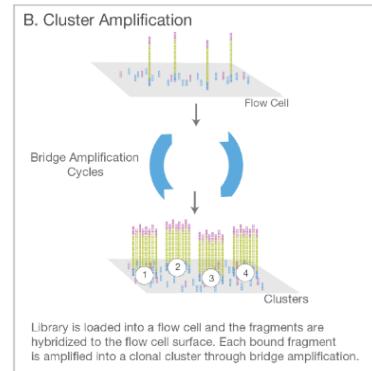
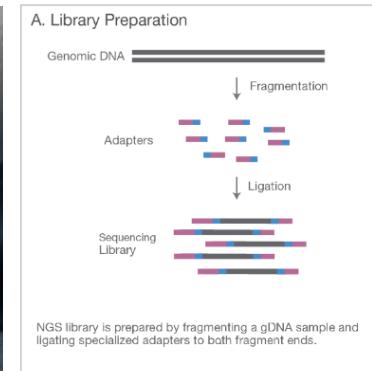
[START](#) [INFRASTRUCTURE SERVICES](#) [ANCIENT DNA](#)

**Ancient DNA** National facility

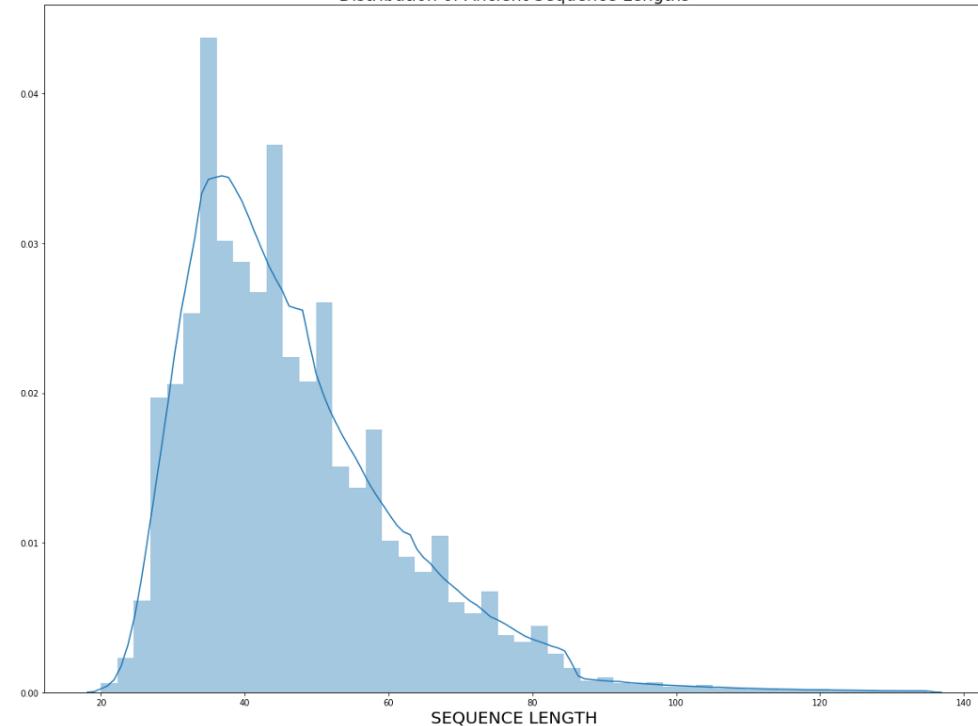
[Share with others](#)

**CONTACT**  
Magnus Lundgren  
[magnus.lundgren@scilifelab.uu.se](mailto:magnus.lundgren@scilifelab.uu.se)  
+46 18 471 2696

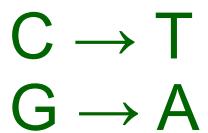
**PERSONNEL**  
Anders Götherström, director  
Mattias Jakobsson, director  
Magnus Lundgren, head of facility  
E-Jean Tan, research engineer  
Thijssen Naidoo, bioinformatician



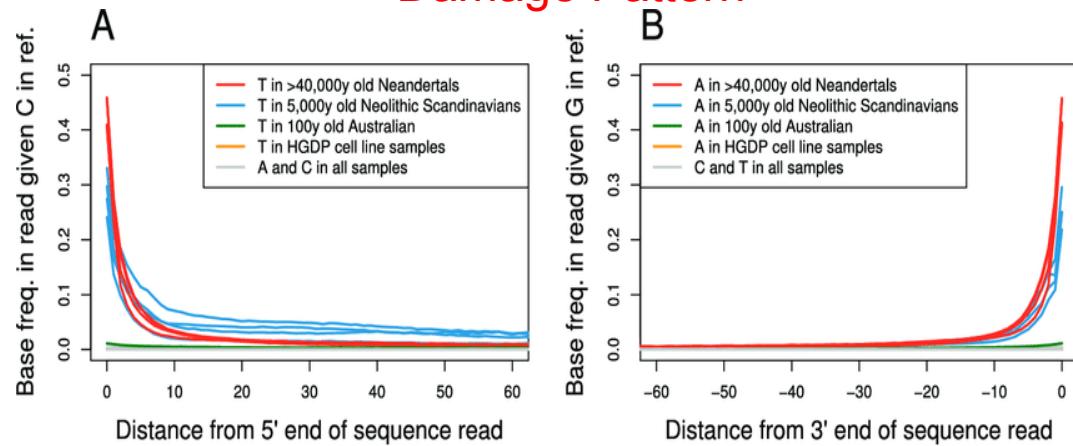
Distribution of Ancient Sequence Lengths



Deamination process:



### Damage Pattern



between enzymes encoded by structural genes originating from two different streptomycetes, rather than the results of the activation of latent genetic information in the recipient strains. (The latter mechanism is the most probable explanation for the three reports of the discovery of novel compounds through natural inter-strain matings (refs 18-20; discussed in ref. 21).) The evidence is most compelling for mederrhodin A, since, of the clones tested, pIJ2301, pIJ2315 and pIJ2316 (which led to mederrhodin A synthesis by AM-7161) all contain a complete transcription unit (absent from pIJ2304, pIJ2308, pIJ2312 and pIJ2317) which complements class V *act* mutants of *S. coelicolor* (F.M., unpublished); such mutants have recently been shown to be blocked in the corresponding hydroxylation involved in actinorhodin biosynthesis (S. P. Cole and H.G.F., unpublished). A true metabolic cooperation between the actinorhodin and gramicidin biosynthetic enzymes is indicated also by complementation of the B1140 mutant by pIJ2308, but not by the other clones tested. B1140 can act as a secretor in co-synthesis tests with *S. coelicolor* *act* mutants of classes I, III and VII, but not of classes IV, V and VI<sup>22</sup>, suggesting that its block is equivalent to that of class IV *act* mutants. Significantly, pIJ2308

## Molecular cloning of Ancient Egyptian mummy DNA

Svante Pääbo

Department of Cell Research, The Wallenberg Laboratory,  
University of Uppsala, Box 562, S-75122 Uppsala, Sweden and  
Institute of Egyptology, Gustavianum, University of Uppsala,  
S-75120 Uppsala, Sweden

**Artificial mummification** was practised in Egypt from ~2600 BC until the fourth century AD. Because of the dry Egyptian climate, however, there are also many natural mummies preserved from earlier as well as later times. To elucidate whether this unique source of ancient human remains can be used for molecular genetic analyses, 23 mummies were investigated for DNA content. One 2,400-yr-old mummy of a child was found to contain DNA that could be molecularly cloned in a plasmid vector. I report here that one such clone contains two members of the *Alu* family of human

**Fig. 2** DNA sequence of part of the mummy clone pMUM2:9. The 9-bp direct repeats flanking the *Alu* sequence are underlined. In the *Alu* consensus sequence<sup>8</sup> only nucleotides differing from the pMUM2:9 sequence are specified.

**Methods.** DNA was prepared from 1.6 g mummy tissue, essentially according to the protocol described previously<sup>5</sup>. The DNA was made blunt-ended with *Escherichia coli* DNA polymerase (Klenow fragment) (Pharmacia P-L Biochemicals) according to the manufacturer's instructions and in the presence of trace amounts of radioactively labelled nucleotides. The sample was then size-fractionated by gel filtration on a G-50 column (Pharmacia, Uppsala). Excluded fractions were pooled and ethanol-precipitated and 25 ng of this material cloned in a *Sma*I-digested to nitrocellulose filters and screened<sup>20</sup> with a *Alu* repeat. The strongly hybridizing clone p hybridization<sup>9</sup>. One of the *Alu* repeats as well after labelling of the *Sph*I and *Nde*I restriction

and alkaline phosphatase-treated pUC8 plasmid<sup>6</sup>. Then, 700 of the white clones were transferred to nitrocellulose filters and probed with a 32P-labeled 550-bp *Bgl*II/*Sph*I fragment from a HLA-DR pseudogene<sup>21</sup>, which contains an *MUM2.9* gene. A clone containing the *MUM2.9* gene was isolated and restriction-mapped. Two *Alu* repeats were identified by Southern blotting analysis. The *MUM2.9* gene was 500 bp of flanking DNA were sequenced according to the Maxam and Gilbert procedure<sup>22</sup> at the sites indicated.

# Deep Learning for Ancient Genomics

## Modern Sequences



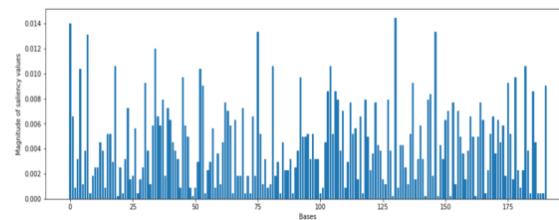
AGCCATATGCATGCATCGTAGTC  
AGCCATATGCATGCATCGTAGTC  
AGCCATATGCATGCATCGTAGTC  
AGCCATATGCATGCATCGTAGTC  
AGCCATATGCATGCATCGTAGTC  
AGCCATATGCATGCATCGTAGTC  
AGCCATATGCATGCATCGTAGTC

## Ancient Sequences



AGCCATATGCATGCATCGTAGTC  
AGCCATATGCATGCATCGTAGTC  
AGCCATATGCATGCATCGTAGTC  
AGCCATATGCATGCATCGTAGTC  
AGCCATATGCATGCATCGTAGTC  
AGCCATATGCATGCATCGTAGTC  
AGCCATATGCATGCATCGTAGTC

Compare and learn  
DNA patterns (k-mers)  
that separate ancient  
and modern sequences



```

from sklearn.preprocessing import LabelEncoder, OneHotEncoder

# The LabelEncoder encodes a sequence of bases as a sequence of integers: 0, 1, 2 and 3
integer_encoder = LabelEncoder()
# The OneHotEncoder converts an array of integers to a sparse matrix where
# each row corresponds to one possible value of each feature, i.e. only 01 and 1 are present in the matrix
one_hot_encoder = OneHotEncoder()
input_features = []

for sequence in sequences:
    integer_encoded = integer_encoder.fit_transform(list(sequence))
    integer_encoded = np.array(integer_encoded).reshape(-1, 1)
    one_hot_encoded = one_hot_encoder.fit_transform(integer_encoded)
    input_features.append(one_hot_encoded.toarray())

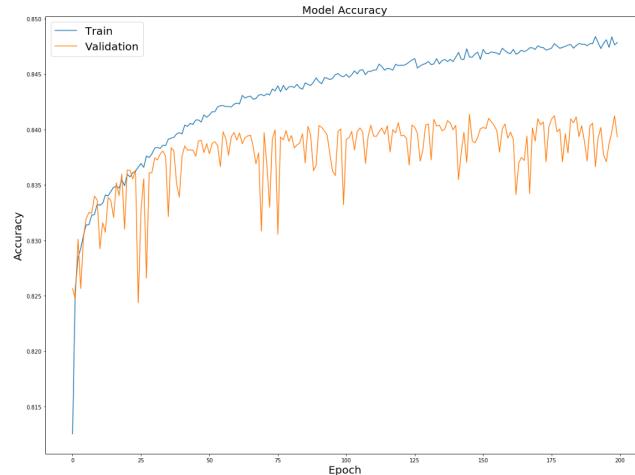
np.set_printoptions(threshold=40)
#print(input_features.shape)
input_features = np.stack(input_features)
print("Example sequence\n-----")
print("DNA Sequence #1:\n", sequences[0][:10], '...', sequences[0][-10:])
print("One hot encoding of Sequence #1:\n", input_features[0].T)

```

Example sequence

-----

DNA Sequence #1:  
AAATCATCAG ... CTCCAAAAAAT  
One hot encoding of Sequence #1:  
[[ 1. 1. 1. ..., 1. 0. 0.]  
 [ 0. 0. 0. ..., 0. 1. 0.]  
 [ 0. 0. 0. ..., 0. 0. 0.]  
 [ 0. 0. 0. ..., 0. 0. 1.]]



```

from keras.optimizers import SGD, Adam, Adadelta
from keras.layers import Conv1D, Dense, MaxPooling1D, Flatten, Dropout, BatchNormalization, Activation
from keras.models import Sequential
from keras.regularizers import l1

model = Sequential()

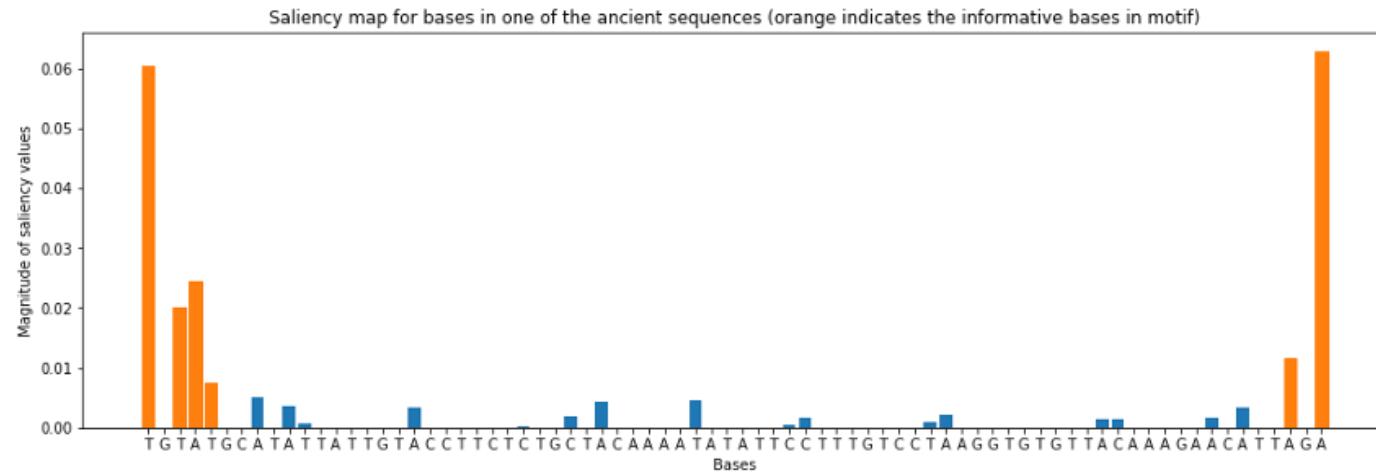
model.add(Conv1D(filters=32, kernel_size=5, padding='same', kernel_initializer= 'he_uniform',
                 input_shape=(train_features.shape[1], 4)))
model.add(Activation("relu"))
model.add(Conv1D(filters=32, kernel_size=5, padding='same', kernel_initializer= 'he_uniform'))
model.add(Activation("relu"))
model.add(MaxPooling1D(pool_size=2))

model.add(Flatten())
model.add(Dense(16, kernel_initializer= 'he_uniform'))
model.add(Activation("relu"))
model.add(Dense(2, activation='softmax'))

epochs = 20
lrate = 0.01
decay = lrate / epochs
sgd = SGD(lr = lrate, momentum = 0.9, decay = decay, nesterov = False)
model.compile(loss='binary_crossentropy', optimizer=sgd, metrics=['binary_accuracy'])
#model.compile(loss='binary_crossentropy', optimizer=Adam(lr = lrate), metrics=[f'binary_accuracy'])
model.summary()

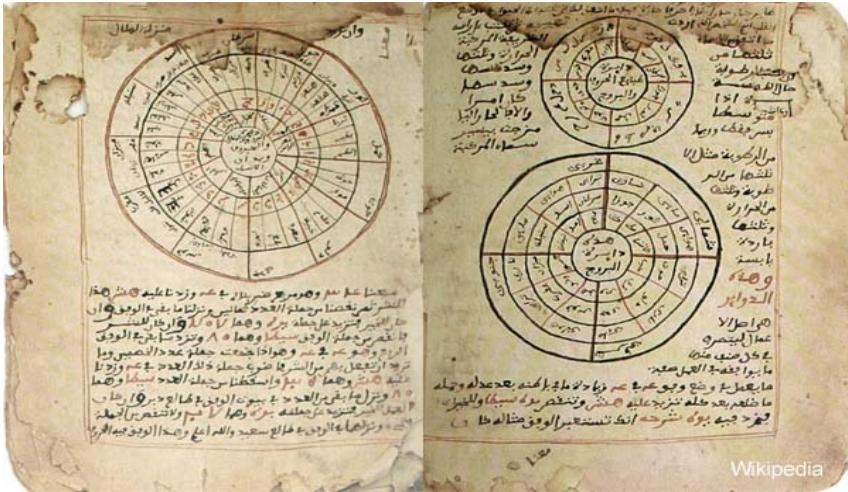
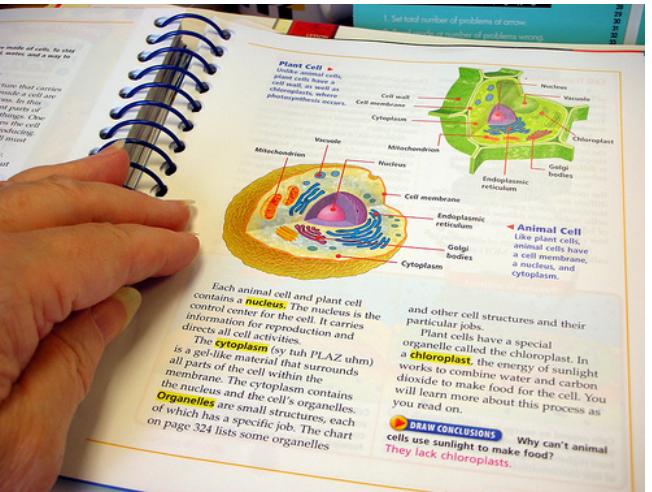
Using TensorFlow backend.

```



# DNA Sequence as a Text: Natural Language Processing Approach

<h1>Editing Wikipedia articles on Medicine</h1>	<h2>Engage with editors</h2> <p>Part of the Wikipedia experience is receiving and responding to feedback from other editors. Don't be afraid to ask them to review your edits. They leave <b>feedback</b>. Real human volunteers from the Wikipedia community will likely read and respond to you, and it would be polite for you to respond back to them when they volunteer to polish your work! Everything submitted to Wikipedia is reviewed by multiple, real human users who can get you back on track, but if you do, please acknowledge it.</p>
<h2>Be accurate</h2> <p>You can help increase millions of people use to make medical decisions, so it's vitally important to make sure that the information you provide for medical information than the websites for WebMD, NIH, and the WHO. But with great power comes great responsibility!</p>	<h2>Watch out for close paraphrasing</h2> <p>Plagiarism or close paraphrasing is never okay on Wikipedia and is a violation of your university's academic honor code. It's even worse if you copy and paste someone else's work that could be used to create good content instead used to clean up plagiarized work.</p> <p>If you plagiarize or no closely paraphrase on Wikipedia, it is extremely likely that you'll be caught. If you are caught there will be an online record of your plagiarism not to your permanent online record.</p> <p>Now that even educational materials from organizations like the WHO and abstracts of articles from peer-reviewed journals cannot and cannot be copied. Write them in your own words whenever possible. If you aren't sure about paraphrasing as, visit your university's writing center.</p>
<h2>Scared? Don't be!</h2> <p>Everybody on Wikipedia wants to make the best encyclopedia they can. Take the time to understand the rules, and soon you'll be contributing to a valuable resource you use on a daily basis!</p>	



Tim?

Mat.

Science says that ending texts with a period makes people look like they're jerks. Does that text make me look like a jerk?

I did feel like I'd upset you

Are you okay?

I'm talking to you now, so I'm much better 

Read 9:11 PM

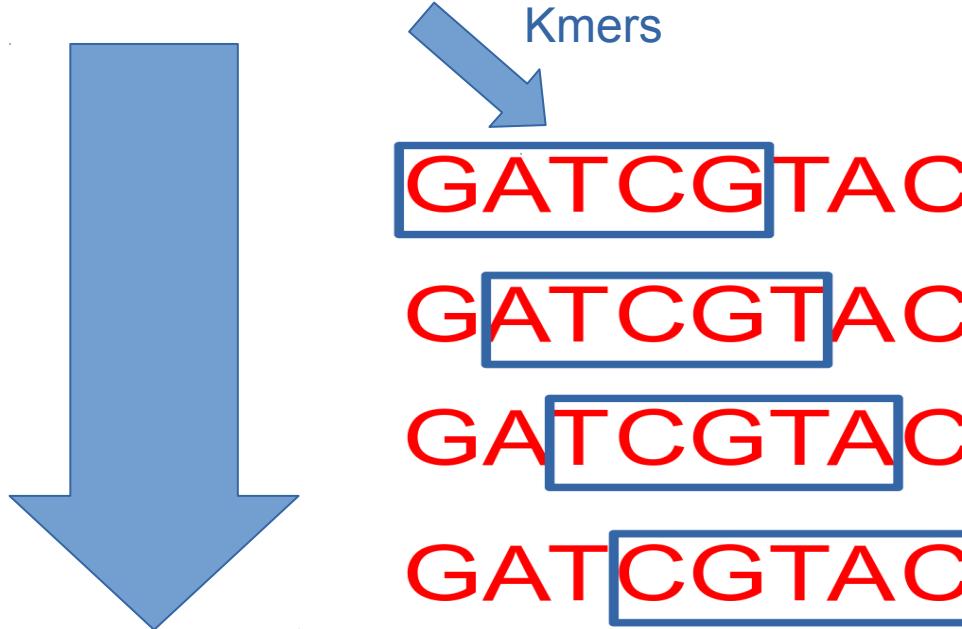
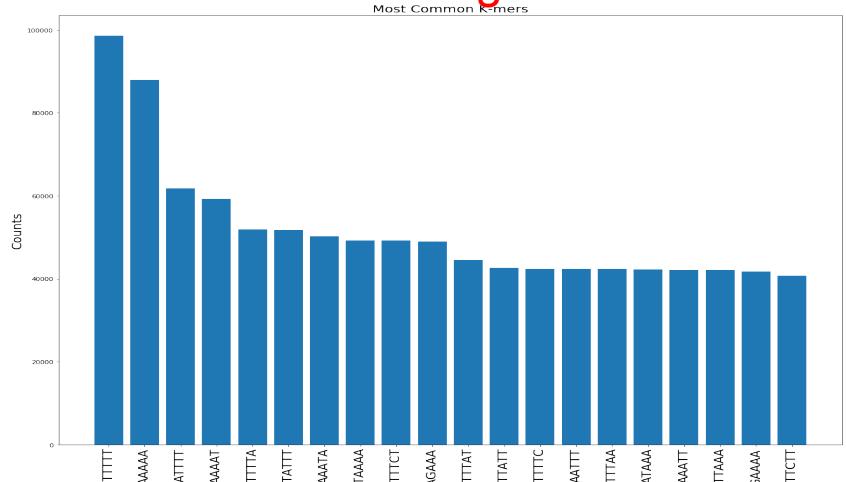
TTTCAAGGAAAACCCGATAAAAATTCTTATTGGGGGAGGGGCTCAAACAAGAAAATAATCAACAAGTGGTGTCCAGAGTGGAGGCCAGGG  
CCTCCTGGGGACAGCAAGACTGCCTGGGAGCGGAAAGCAGCTCCCCCGTCTGGGGGAGGGCGTGGCTGGAGGGAGGGCTGGACCA  
CAGGAGGGCAGGCCCTGGCAACCCTATGTAGATGAAGCTGCCGGAGAGGGATCAAAGAACCAAGACAGGAGGAAAGAGGGCGGTGAAG  
TCTCCTGTCTCATCCCTAGGAAGCCTGAGGAGATGGCTAAGGGCTTAAAGGCTCGAACCCCAGGGCAGAGGATAAGTGTAGGCA  
GATCTGTGAGCTCAGGTCTTACCTGACAGGGGAGGGCCAGGAGCCCCCGAGGCTCATCAGCATCATCACCAAGCCCTGTGGCCTGCA  
CCACACCTCCGACGCCACCAGCCGGGAGTACACCGAAAACCAGCGGGCGTAACCTCCCAGCCTAGCCCAGGCCACAGCCGC  
GGCCAGCAGGGACCCCCCAGCTCTTCGCCGCCACAGGGCACCCACAGCCCCAGCCCAGTCAGAGCCCC  
AATACCGCCAGCAGCGCGGGAGGGCACCCAGGCTTGTCTGCCAGCCACCCAGGAGACCAGCCGGCGCCATCCCCATCCAG  
CCACGGCCACCACCAAGCGCTGCTCGTATCCCCCAGGCCCGGTCTAAAGCGTGGGGAGCCAAGTGCACGTAAGGAACGAAGTACCC  
GCCCCAACCAGGGCTGTGCCAGAGAAGGCCCCGGCTGTGAAACAGACTCAGGCCAGGGCAGCTAGGGACTACGC  
GGTGGGGCTGGGGGTCTCAGGAAGGACCAGGGTAGCAGCAGGGCGCCACAGGGGTGAGGTGGAGGGTATCGGCCGAGGAG  
GAGCAGAGCGCCCCGCCAGCCAAAGTATCGAGAAGAACGCTGCAAGGCCGGCCAGGAGCAGCGAGGAGGCCGTTGCCGGTGAG  
GCCAGCCCCACCAGCCAAGACTCGACGGCGGGAGAAGTAACGCCAGGGGTGCTTAGGGCGGGCGAACACCAGGGCCAACCAA  
GCCTGCGAATGAATAGGAGGGATGGGGCCGGACTGGGACGCCCGCCAGCATTCCAGGCCGGCTCTCCGACCAGGCCCG  
CCTCGTTCGCTACCCAGATCCAAACAGCTCTCACCTCCTCACCTGAATGACCCGGCATCCACTCCCTCACAGCGAGGAG  
GCCCAAGGCCAGGTAGAGATGCAGCAGATCGCTGGCGAAAGCCAGAACGAAAGCCCAGCGAGGCAGGCCAACCACATCAC  
CACGGGGCGGGCCCCCAGCGCTGCTCAGGGCGCTGCCACGGGGCTGAAAGGGGGCGGAGTCACGGAAGAACGCCCGGG  
CCCCAAACTCTCCAACACTCTCATGGCTTTGCCCTCGAACTAATGGTCTTCTGCCACGTTCAAGATGCACTCTTCTAG  
AGCCGGTTGCCCTTCTCAGGTTGCGGTAAATCAGGTGGGTCTCCACTCACGGCTCCTCCCCGTTCTGTCTCCCGCC  
TCAGGGCCCCCTCACTGGCTGCCGCTGCACGCCAGGGCCAGGGCGCTGATCCACGCAGTGTCTGGCGCTCGGTCAAAGTGTCTC  
GGCAAGGTCAAGGAAGGCAAGGCCAGCGAGCGCAGCAGCCCGTAGGACAGCCCCTTATCGCGAAGGCTGCCGCCACACCA  
GCCCAAGCCCCCATCCGGGGTCCGGCGGGCTGGGGGCTATGCCGTCTGCCGGGTGGGAAACATCTGTGAGAGAACGCTCCACGC  
CTGTGCTTCCGCTGGGAGCTGGCATCCCTGAGATCCAGCCTCTGGCTGCTGCCGGGGTGGGCCACTCCAGGGCGATGGCG  
GGGGCGGAGGGAGCCGGAGCCTGGCTAGGGCTGGAACTCCCGGGTCCCGCGCAGGTACGGGGACGGGACAGCCAGATCCCCAGGC  
CCGAGGTTCCCCGTCCCACCGCACACTCCTCCCCCTTACCCGACCTCTCCGGGGCGGTGCCGGGAGGGGAAGGGTGAGGAAGGGCTGG  
CCCGGCTTCTCTGCTTCCAGGCAGGGGGGGGGGAAGGGGAGCGAGGGCAGCGATGGAGCCAACTTGGACGGCTCTCG  
GTAAACAGAGATCACCACAGGGGCTGAGCCACCTGGGACGCAGGGTGTACGGAGGGCGGGGAGCTGAAACAGCCAATCCGGCA  
AGCCCGCGTGAGGCCAGACGCCAGTCCCACGGGATTGAATTATGTACACACACAACCCCCAGGCCGGGGAGGGGGAGCCGGA  
GAAATG

This is a text, so we can use Natural Language Processing (NLP)  
Where are sentences? Where are words?

Sequence  
**GATCGTAC**

Functional elements vs. junk DNA

NLP: Bag of Words



**GATCG** word    **ATCGT** word    **TCGTA** word    **CGTAC** word



```
import warnings
warnings.filterwarnings('ignore')

from gensim.models import Word2Vec
model = Word2Vec(sentences, min_count = 2, workers = 4)
print(model)

Word2Vec(vocab=1024, size=100, alpha=0.025)
```

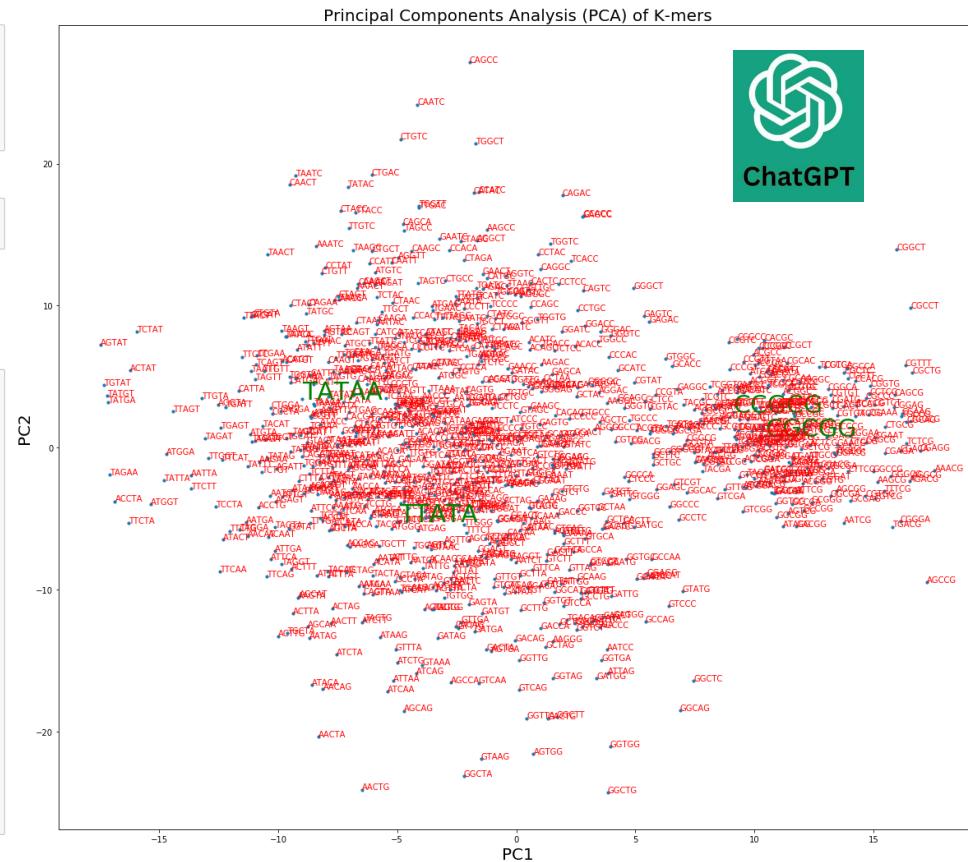
```
X = model[model.wv.vocab]
X.shape

(1024, 100)
```

Now each word is one observation, this observation has 100 coordinates, i.e. the default number of latent variables for word2vec. Next we can try to use the constructed word vectors and visualize the k-mers space using PCA and UMAP.

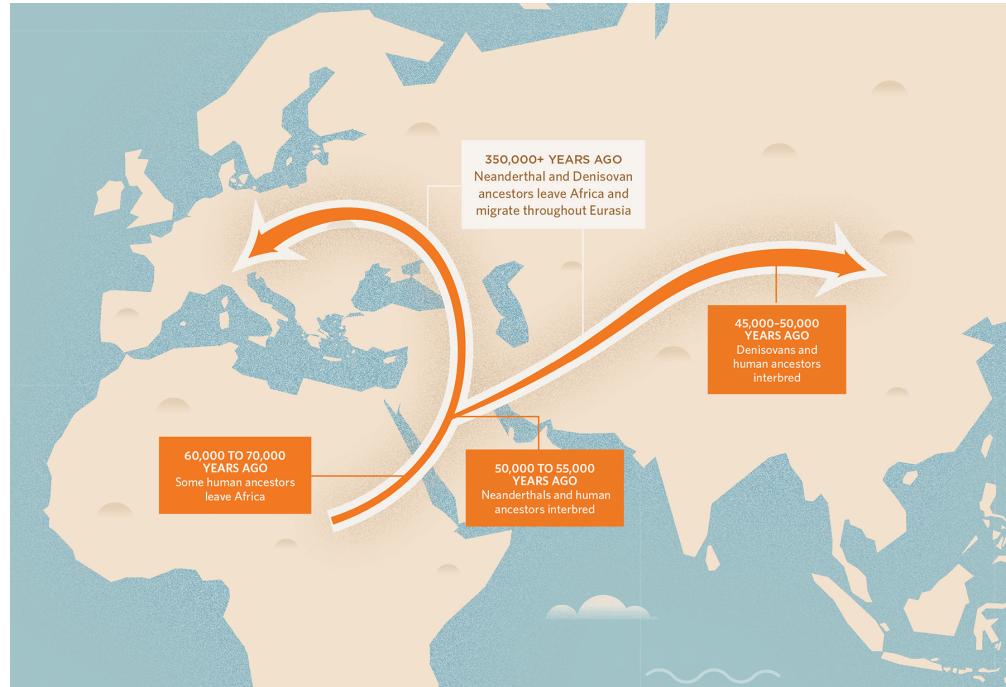
```
import matplotlib.pyplot as plt
from sklearn.decomposition import PCA
X = model[model.wv.vocab]
pca = PCA(n_components = 2)
result = pca.fit_transform(X)

plt.figure(figsize=(20,18))
plt.scatter(result[:, 0], result[:, 1], s = 10, cmap = 'tab10')
plt.title('Principal Components Analysis (PCA) of K-mers', fontsize = 20)
plt.xlabel("PC1", fontsize = 20)
plt.ylabel("PC2", fontsize = 20)
words = list(model.wv.vocab)
for i, word in enumerate(words):
    if word == 'CCGG':
        plt.text(result[i, 0], result[i, 1], word, fontsize = 30, c = 'green')
    elif word == 'CGCG':
        plt.text(result[i, 0], result[i, 1], word, fontsize = 30, c = 'green')
    elif word == 'TTAA':
        plt.text(result[i, 0], result[i, 1], word, fontsize = 30, c = 'green')
    elif word == 'TATAA':
        plt.text(result[i, 0], result[i, 1], word, fontsize = 30, c = 'green')
    else:
        plt.text(result[i, 0], result[i, 1], word, fontsize = 10, c = 'red')
plt.show()
```



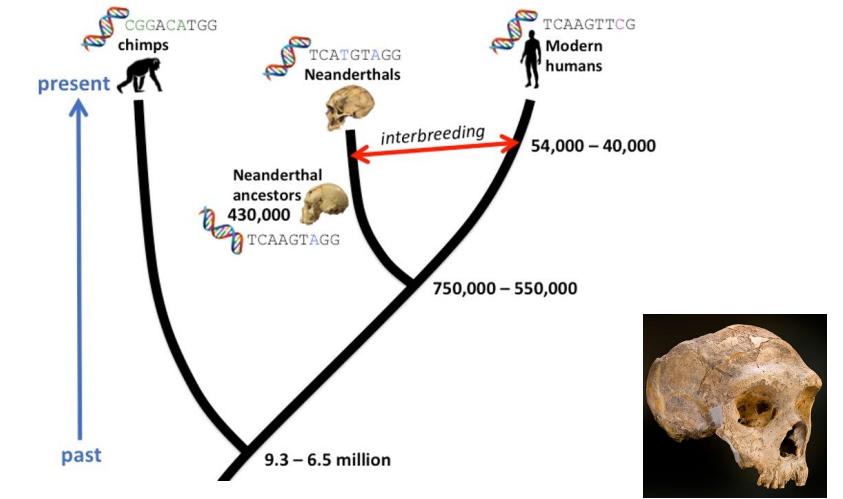
GC-rich and AT-rich 5-mers seem to form separate clusters of nucleotide “words”

# Applications for detecting Neanderthal DNA in modern genomes

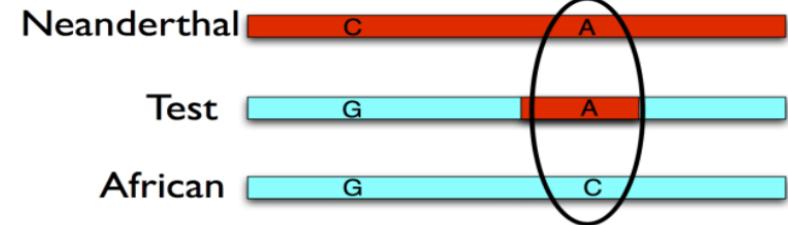
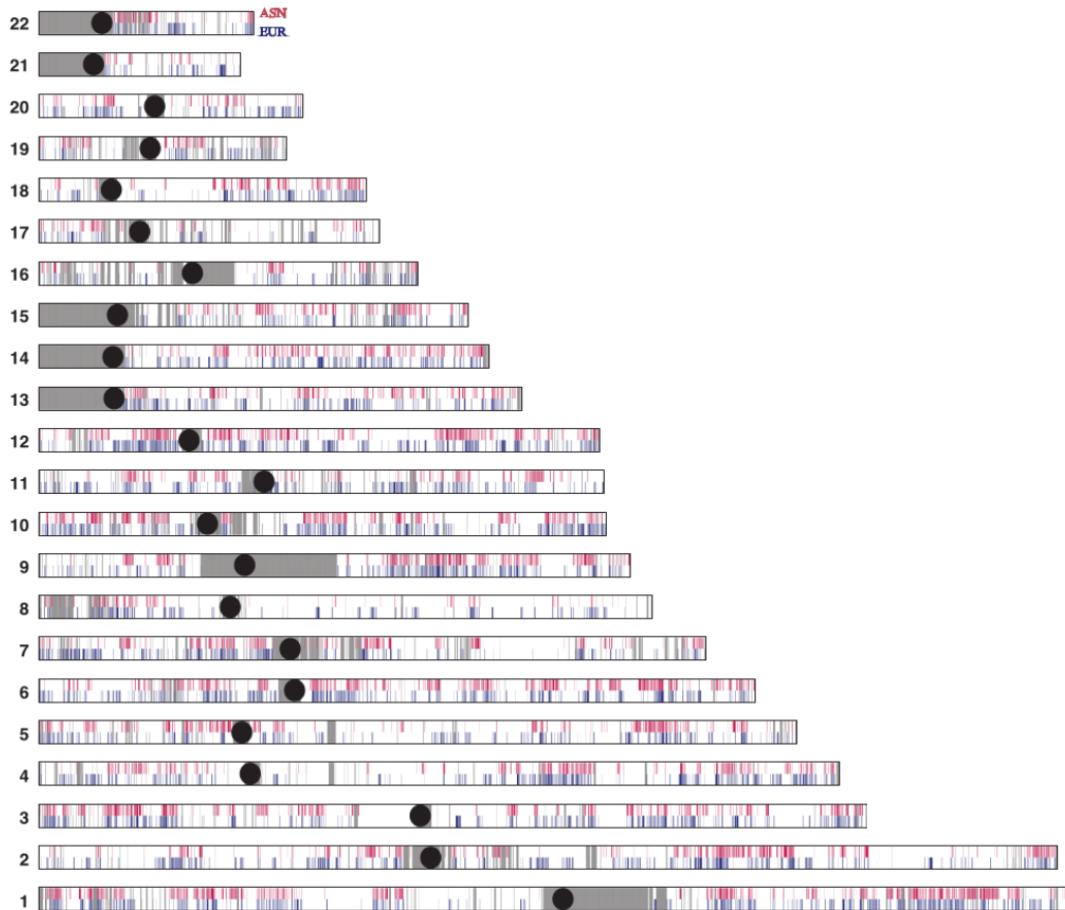


### Neanderthal DNA in Modern Human Genomes Is Not Silent

From skin color to immunity, human biology is linked to our archaic ancestry



2010 Draft Neanderthal Genome Sequenced



# Sentiments Analysis: Bag of Words

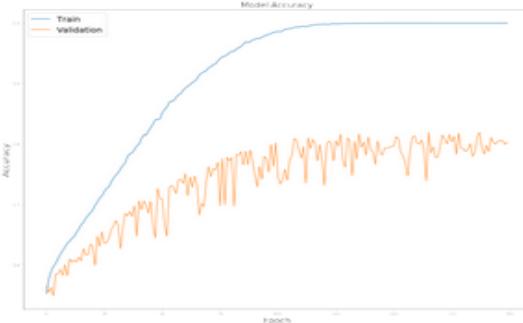
Shallow MLP: 82% accuracy

```
from keras.models import Sequential
from keras.regularizers import l2, l1
from keras.callbacks import ModelCheckpoint
from keras.optimizers import SGD, Adam, Adadelta
from keras.layers import Dense, Flatten, Dropout

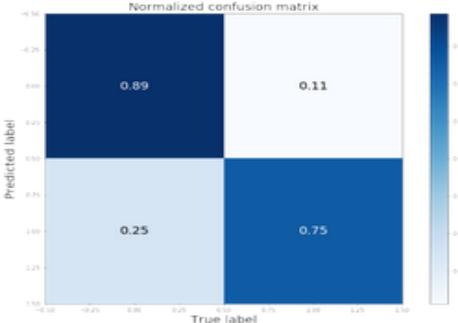
model = Sequential()
model.add(Dense(3000, input_shape = (X.shape[1], ), activation = 'sigmoid',
               kernel_regularizer = l1(0.00001)))
model.add(Dense(1, activation = 'sigmoid'))
sgd = SGD(lr = 0.0001, momentum = 0.9, nesterov = False)
model.compile(loss = 'binary_crossentropy', optimizer = sgd, metrics = ['binary_accuracy'])
checkpoint = ModelCheckpoint("weights.best.hdf5", monitor='val_binary_accuracy', verbose=1,
                             save_best_only = True, mode = 'max')
history = model.fit(x_train, y_train,
                     epochs = 200, verbose = 1, validation_split = 0.2, batch_size = 32,
                     shuffle = True, callbacks = [checkpoint])
```

[mlp.py](#) hosted with ❤ by GitHub

[view raw](#)



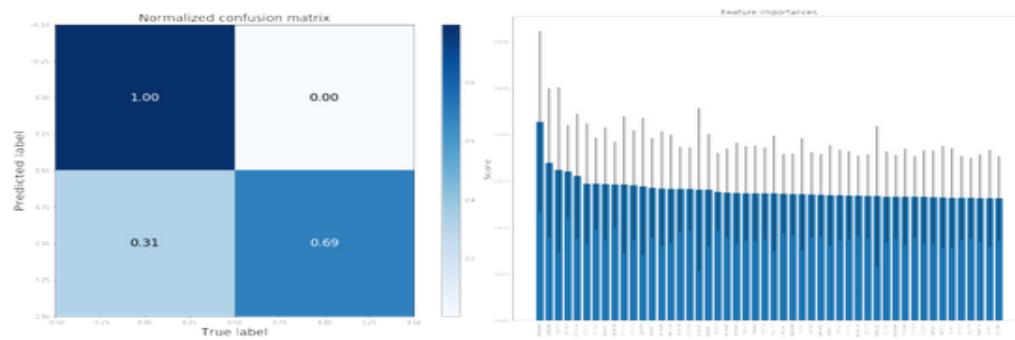
MLP Accuracy training curve (left) and confusion matrix of evaluation on test data set (right)



```
1 import pickle
2 import matplotlib.pyplot as plt
3 from sklearn.ensemble import RandomForestClassifier
4
5 classifier = RandomForestClassifier(n_estimators = 500)
6 classifier.fit(X_train, y_train)
7 y_pred = classifier.predict(X_test)
8 pickle.dump(classifier, open('RF_model_Neand_Intr_vs_Depl.sav', 'wb'))
9
10 importances = classifier.feature_importances_
11 std = np.std([tree.feature_importances_ for tree in classifier.estimators_], axis = 0)
12 indices = np.argsort(importances)[::-1]
13 plt.title("Feature importances", fontsize = 20)
14 plt.bar(range(X_train.shape[1])[0:50], importances[indices][0:50],
15         yerr = std[indices][0:50], align = "center")
16 plt.xticks(rotation = 90); plt.ylabel("Score", fontsize = 20)
17 plt.xticks(range(X_train.shape[1])[0:50], np.array(names)[indices][0:50])
```

[RF.py](#) hosted with ❤ by GitHub

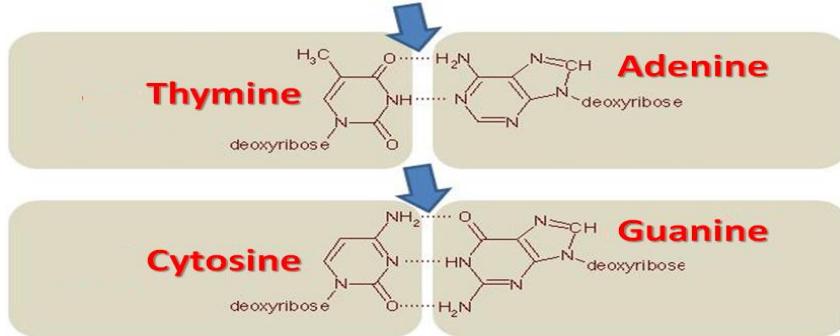
[view raw](#)



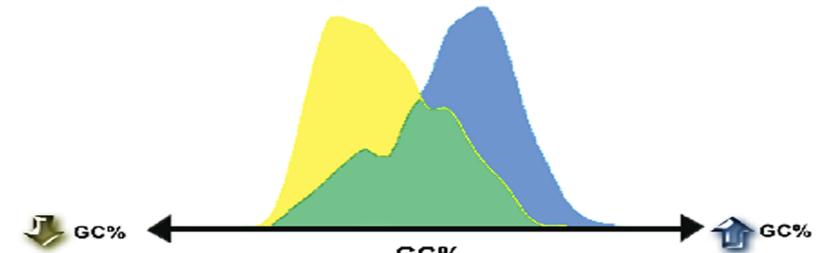
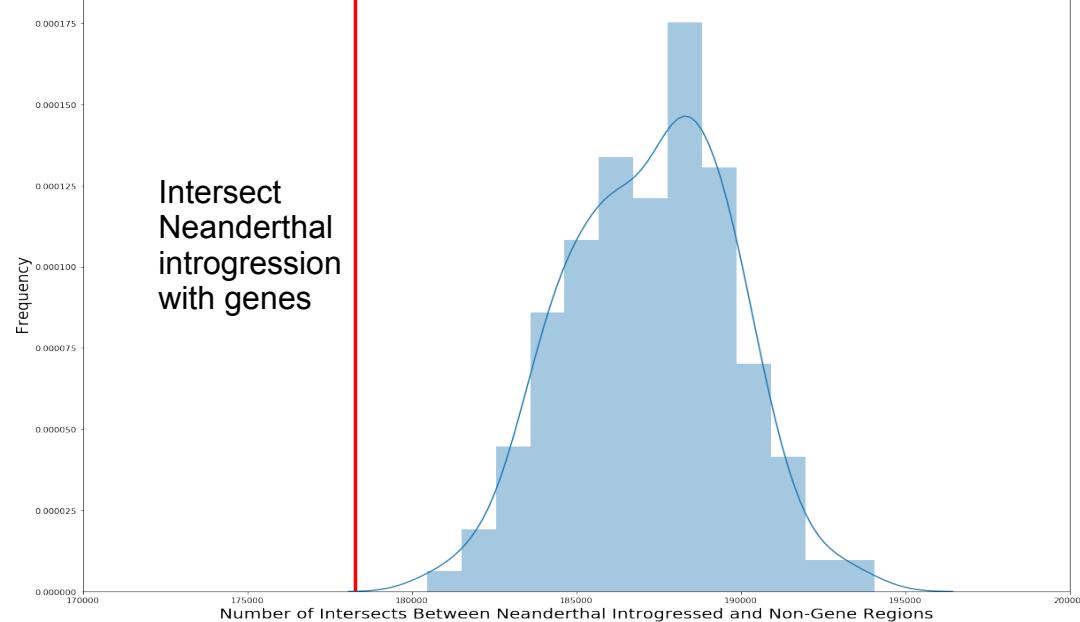
Random Forest classifier confusion matrix (left) and feature importances dominated by AT-rich k-mers (right)

AAAAA, CAAAAA, CATT and TTTT are most predictive, i.e. **AT-rich (gene depletion)**

# Evolution seems to eliminate Neanderthal ancestry from modern human genes



Distribution of Non-Gene Intersects: Vernot and Akey, Science 2016



### AT rich

- high evolutionary rate
- low transcription/repair
- increased methylation/mutation
- increased intron insertion
- chromatin condensation
- pseudogenization

### GC rich

- low evolutionary rate
- high transcription
- decreased methylation/mutation
- decreased intron insertion
- increased duplication
- increased amplification

**It seems that breeding with Neanderthals was not beneficial for modern humans?**

Take home messages of the session:

- 1) One-hot-encoding is essential to train CNN on DNA
- 2) Deep Learning powerful for functional element detection
- 3) Genomics is a potential Big Data with correct setup
- 4) RNNs do not outperform CNNs for genomic applications
- 5) DNA is a text, which is suitable for NLP applications



*Knut och Alice  
Wallenbergs  
Stiftelse*



**LUNDS  
UNIVERSITET**