The background is a solid teal color. On the left side, there is a large, dark teal DNA double helix that runs vertically. To its left and right, there are fainter, lighter teal DNA helices and molecular structures, including a small cluster of three nodes connected by lines in the top left. On the right side, there is a network of nodes and lines, resembling a molecular or neural network structure, in a lighter teal shade.

A Classification Model for lncRNA and mRNA Based on K- mers and a Convolutional Neural Network



Index

01 Introduction & Dataset

02 Preprocessing

03 Activity Diagram

04 Model Convolution

05 Comparison

06 Conclusion

Dataset

GENCODE database,
gencode.v 26

The datasets used to perform the
analysis are publicly available at
[https:// www.genecodegenes.org/mouse/](https://www.genecodegenes.org/mouse/)

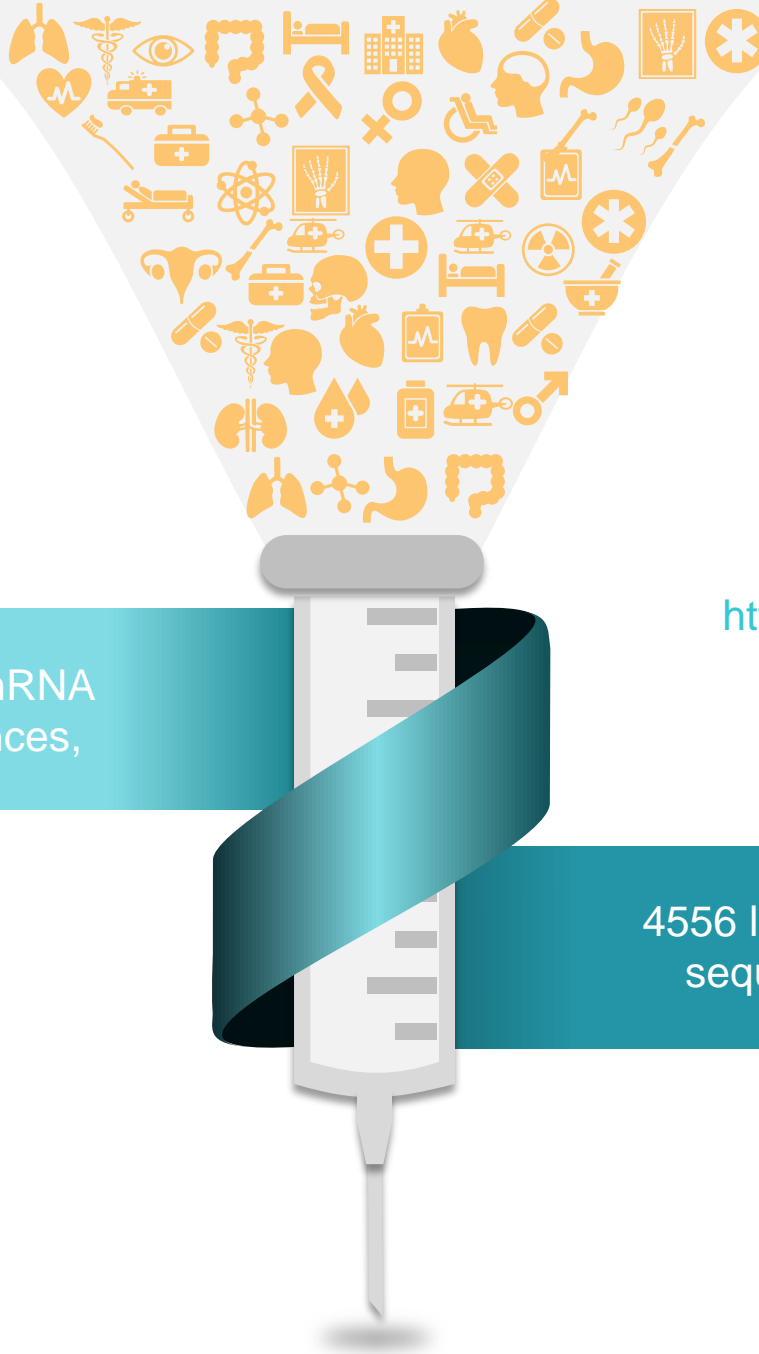
01

4556 mRNA
sequences,

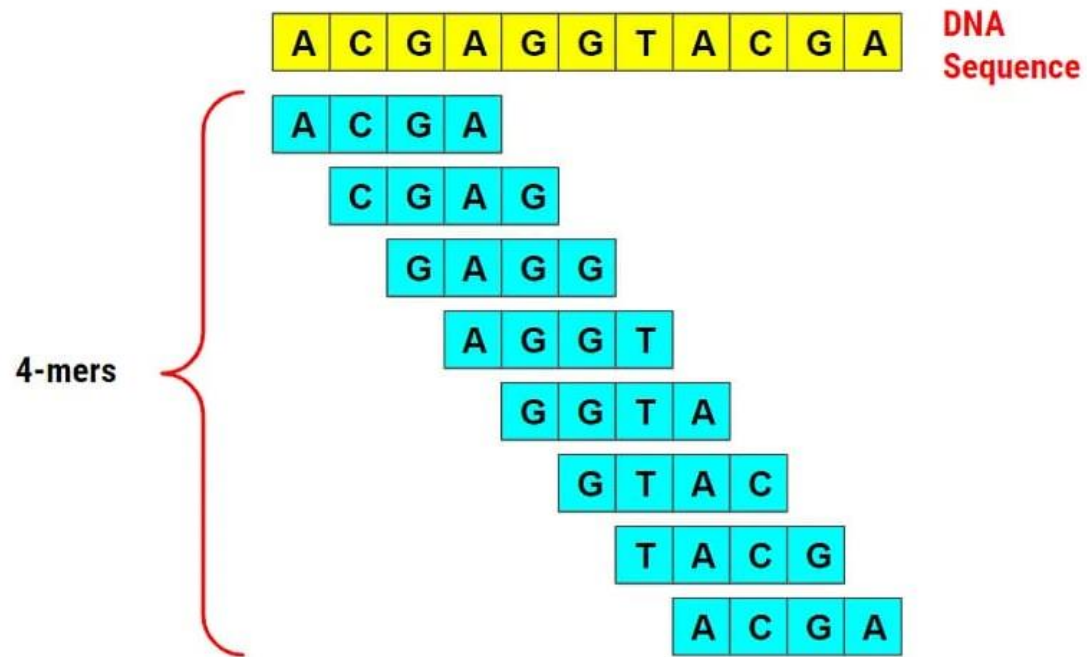
4556 lncRNA
sequences

02

Of which 7745 sequences are
selected as training samples and
the remaining 1367 sequences are
used as test samples.

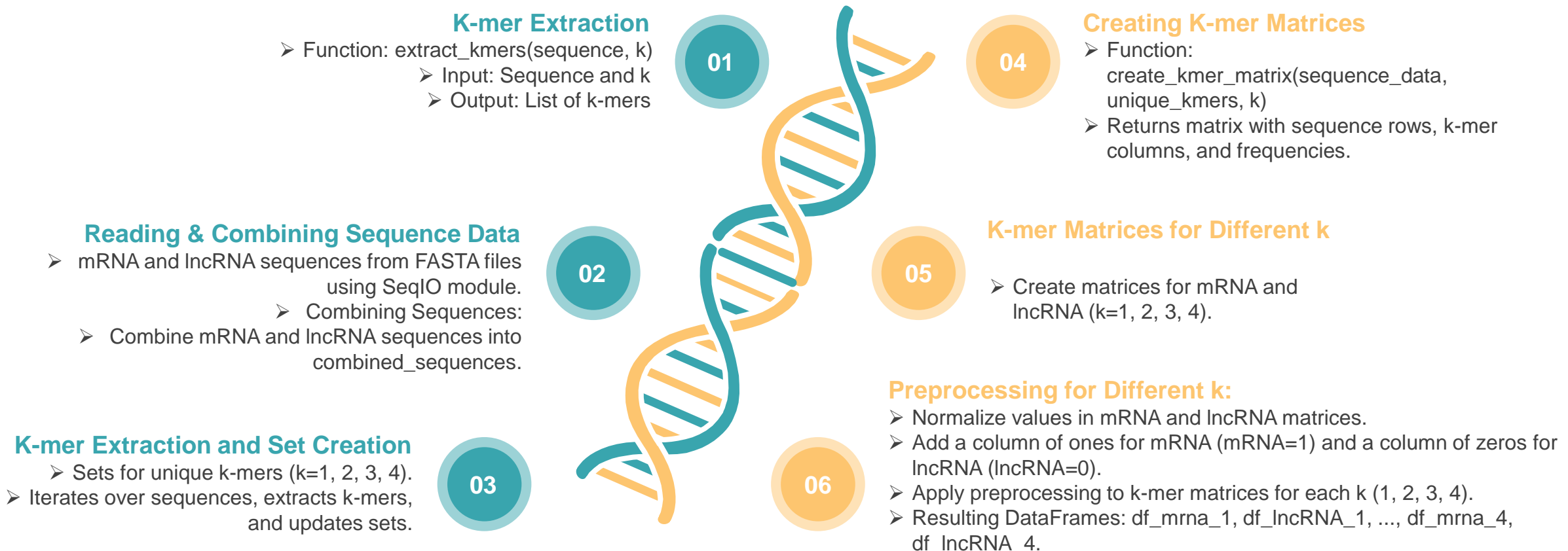


K-mers



K-mers are contiguous subsequences of length k that are derived from a longer sequence, typically in genomics or bioinformatics.

Preprocessing



Activity Diagram

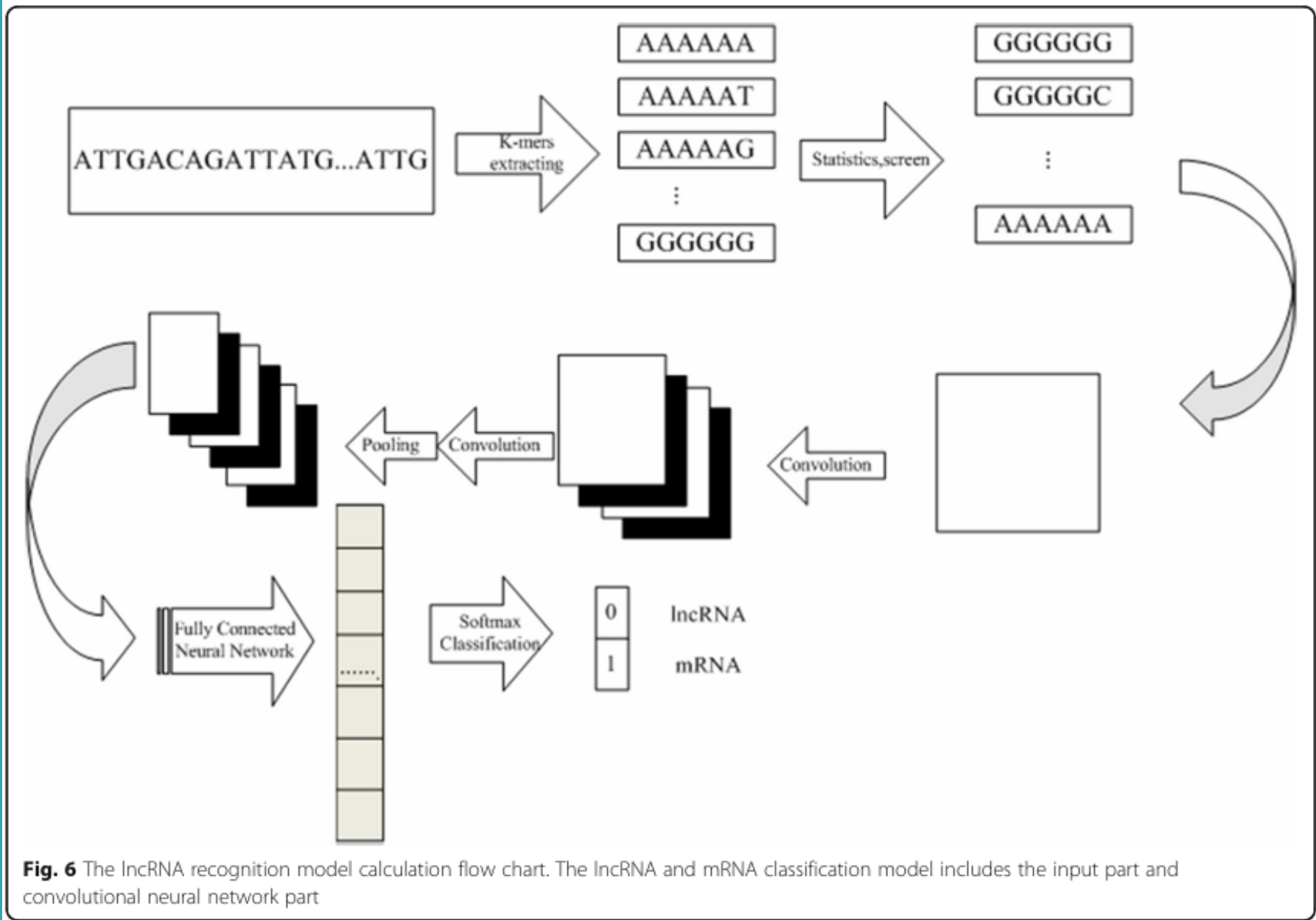
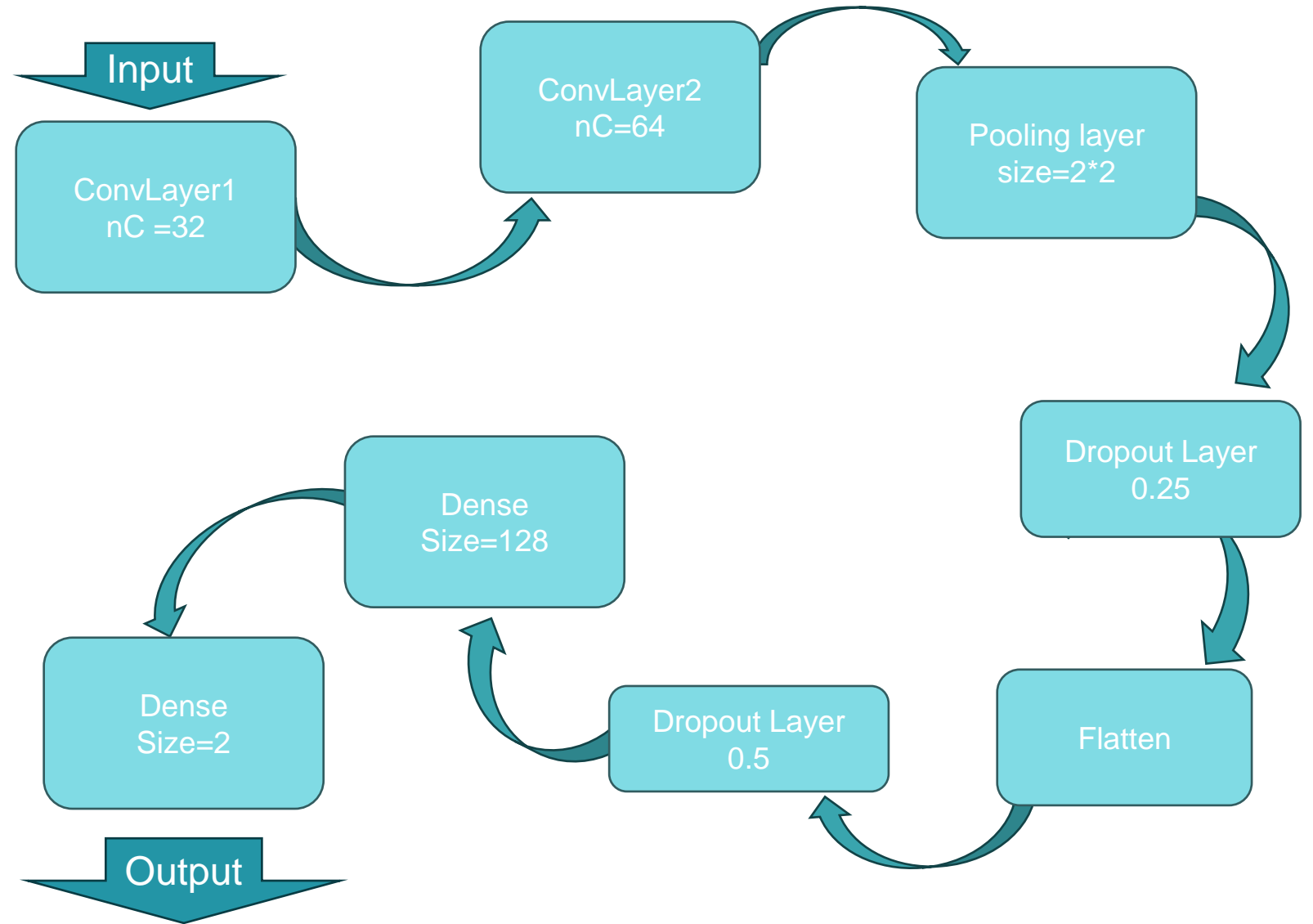
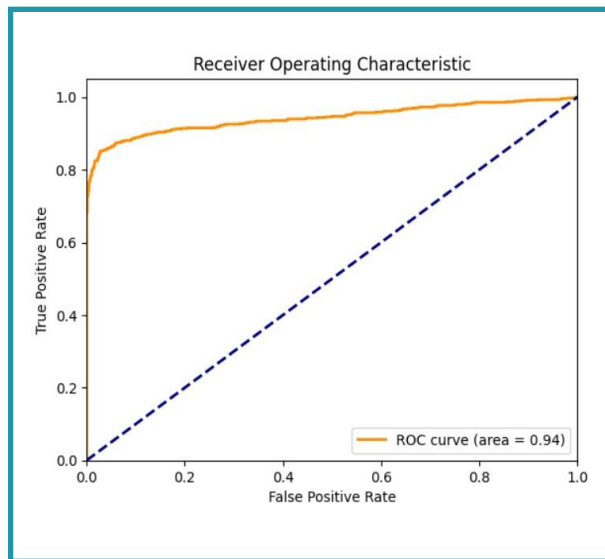


Fig. 6 The lncRNA recognition model calculation flow chart. The lncRNA and mRNA classification model includes the input part and convolutional neural network part

Convolution Model



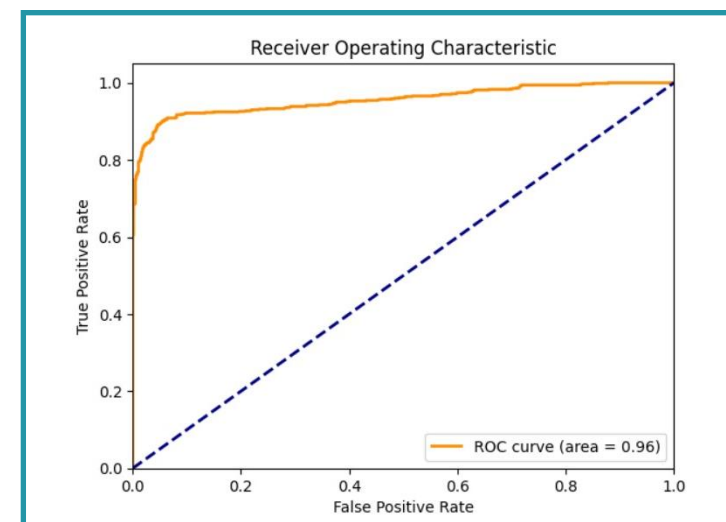
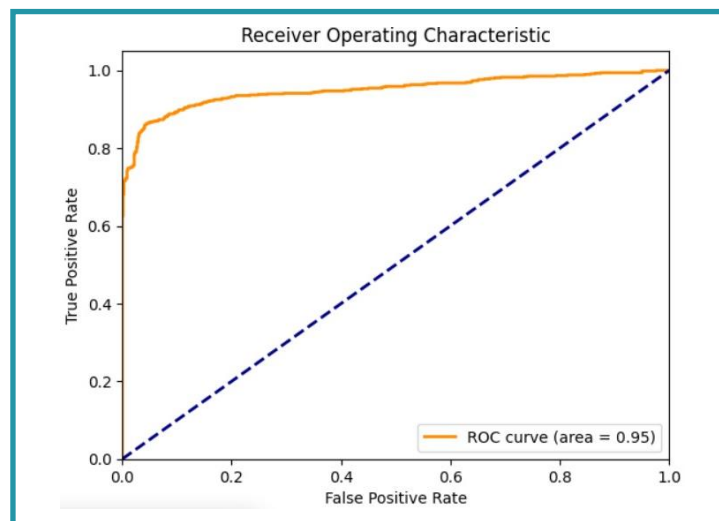


ROC curve for 3-mers
($n=4^3=64$)

N(k-mers)	accuracy	Precision	Recall	F1 score
16	0.9078	0.9644	0.8516	0.9045
64	0.9239	0.9771	0.8691	0.9199
256	0.9005	0.9223	0.8735	0.8972

Results and Output

ROC curve for 4-mers
($n=4^4$)



ROC curve for 2-mers
($n=4^2$)

Comparison

CNN

Convolutional Neural Network (CNN) model can self-learn the characteristics of the sequence through continuous training **without artificial intervention and efficiently calculate large amounts of data**, no domain-expert knowledge or fine-tuning of parameters to increase accuracy are needed

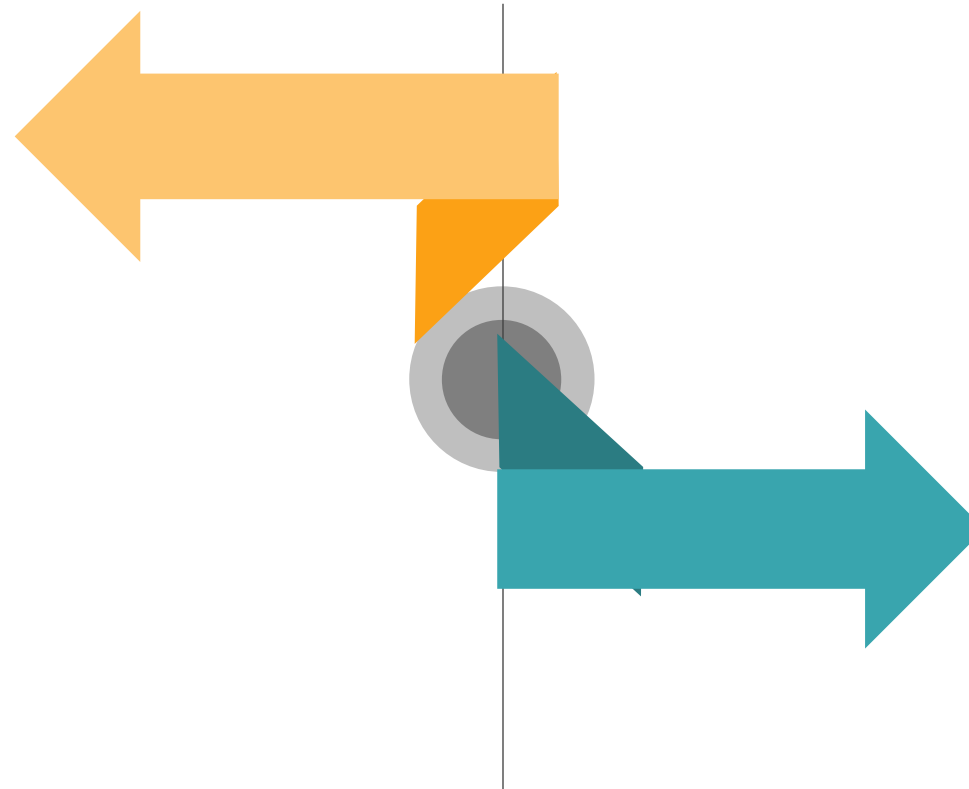
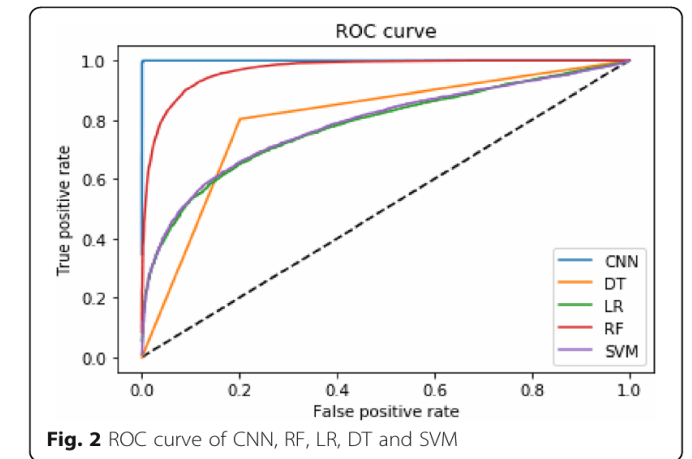


Table 5 Five model effect comparison table in human

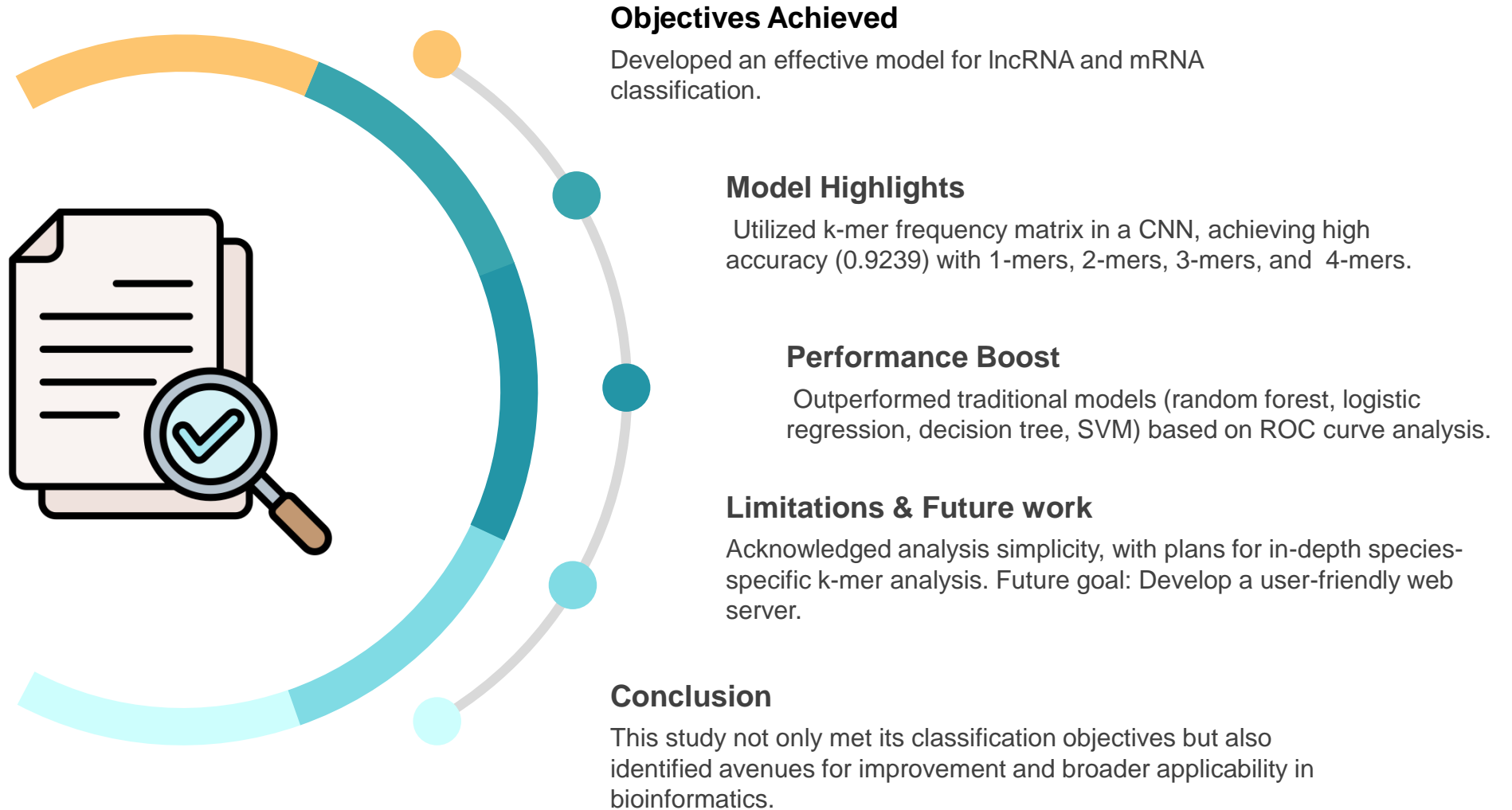
model	model accuracy	precision rate(P)	recall rate(R)	F_1 score
CNN	0.9872	0.9993	0.9955	0.9974
RF	0.8820	0.8949	0.8867	0.8925
LR	0.7020	0.7247	0.7183	0.7218
DT	0.8030	0.7873	0.7852	0.7869
SVM	0.7020	0.7245	0.7158	0.7179



ML methods

- The study used the **maximum entropy algorithm** for k-mer screening and support vector machines for classification, revealing high computational complexity and cost.
- **Expert-driven pre-processing** and feature selection were crucial, emphasizing parameter fine-tuning for accuracy in various conventional machine learning algorithms like **SVM, logistic regression, decision trees, NN, BNs, GAs, HMMs**, etc.

Conclusion





Thank You