

---

# **plantdeepsea Documentation**

---

**Jan 19, 2023**



# OVERVIEW

<b>1 Home Page</b>	<b>3</b>
<b>2 Model Selection</b>	<b>5</b>
<b>3 Submit Task</b>	<b>7</b>
<b>4 Query Results</b>	<b>11</b>
<b>5 Interpretation of Results</b>	<b>13</b>
5.1 Brief Description . . . . .	13
5.2 Main Functions . . . . .	16
5.3 About the result files . . . . .	17
<b>6 Blast</b>	<b>19</b>
<b>7 Case Studies</b>	<b>23</b>
7.1 Case1: The DEP1 gene of rice . . . . .	23
7.2 Case2: The UPA2 gene of maize . . . . .	23
<b>8 Resources</b>	<b>25</b>
8.1 Brief Description . . . . .	25
8.2 Reference Genome Information (Download from Google Drive) . . . . .	25
8.3 Reference Genome Information (Download from Baidu Cloud Disk Drive) . . . . .	26
8.4 Trained model . . . . .	26
8.5 Training data . . . . .	26
8.6 Training config files . . . . .	27
8.7 Model structure . . . . .	27
8.8 A guide to applying the model locally . . . . .	27
<b>9 Statistics</b>	<b>29</b>
<b>10 Citation</b>	<b>45</b>



Welcome! This is the documentation for PlantDeepSEA, a deep learning web service based on Selene Python SDK and Django. The web sever address is located [here](#). PlantDeepSEA is freely accessible for all users.

**Citation :** Hu Zhao#, Zhuo Tu#, Yimeng Liu, Zhanxiang Zong, Jiacheng Li, Hao Liu, Feng Xiong, Jinling Zhan, Xuehai Hu, and Weibo Xie\* (2021). PlantDeepSEA, a deep learning-based web service to predict the regulatory effects of genomic variants in plants. *Nucleic Acids Research*, doi: 10.1093/nar/gkab383



---

**CHAPTER  
ONE**

---

**HOME PAGE**

The appearance of the home page is shown in the figure below.



PlantDeepSEA is a webserver based on deep learning models of chromatin accessibility for multiple plant species. It can predict the impact of genomic variants on chromatin accessibility in multiple tissues. Therefore, it can be used to prioritize genomic variants and discover high-impact cis-regulatory sites within a sequence.

### Select models

### Trained models



#### Arabidopsis

Col-0

A small flowering plant that is widely used as a model organism in plant biology.



#### Rice

Zhenshan 97

Minghui 63

A vital staple crop and a model organism for monocotyledons.



#### Brachypodium

Bd21

A model plant for temperate grasses and herbaceous energy crops.



#### Foxtail millet

Yuguyihao

An important grain crop and forage.



#### Sorghum

JDXBR

A biofuel crop and potential cellulosic feedstock.



#### Maize

B73

A staple food in many parts of the world.

### Query your results

Please input a task id



example task IDs:  
2566\_16148611193615174 or  
2021\_16084538926592038.

### What's New

2021-04-02

We have updated the presentation of the results page.

2021-03-02

The rice model supports coordinates and intervals of multiple reference genomes.

2021-01-19

'Sequence Profiler' can accept input of a custom sequence for analysis.

2020-11-18

 The online service of the Plant Multi-species Deep Learning Model is now available!

Here, you can select the model, go to result, statistics, tutorial, blast Interface.

You can click the home button to return to the home page.

---

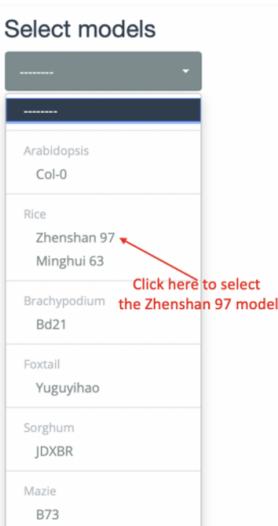
## CHAPTER TWO

---

### MODEL SELECTION

You can click the drop-down box on the left side of the home page or the species reference genome name in the middle column to select the appropriate model for prediction (as shown in the figure below).

Select models



Arabidopsis	Col-0
Rice	Zhenshan 97 Minghui 63
Brachypodium	Bd21
Foxtail	Yuguyihao
Sorghum	JDXBR
Maize	B73

Trained models

 **Arabidopsis**  
**Col-0**  
A small flowering plant that is widely used as a model organism in plant biology.

 **Rice**  
**Zhenshan 97**  
**Minghui 63**  
A vital staple crop and a model organism for monocotyledons.

Click here to select the Zhenshan 97 model

Then, jump to the prediction interface, where you will see the model name you have chosen (as shown in the following figure).

ZS97 Model

**Variant Effector**

(1) Upload a VCF file  
Please select the reference genome corresponding to the VCF file

选择文件 未选择任何文件  
Please upload one file for prediction.

**Sequence Profiler**

(2) Upload a BED file  
Please select the reference genome corresponding to the BED file

选择文件 未选择任何文件  
Please upload one file for prediction.

(3) Select a Chromosome and input a Position  
Please select the reference genome corresponding to the position

Do not know the position? You can go to [Blast Page](#).

(4) Input a sequence (>20-bp is allowed, ≥1-kb is recommended)

**Notice:**

Here is the name of the model you have clicked

This model was constructed based on ATAC-seq data from multiple tissues of the Rice (*Oryza sativa L.*) variety Zhenshan 97 using the deep learning-based algorithmic framework DeepSEA (Zhou et al., *Nature Methods*, 2015) and was implemented using the Selene SDK (Chen et al., *Nature Methods*, 2019).

We provide two services based on this model.

1. Variant Effector, aims to predict the effects of sequence variants on chromatin accessibility. The accepted input is a VCF file containing information on the sequence variants. The results contain information on the effect of variants on chromatin accessibility in each tissue.
2. Sequence Profiler, is a utility that performs "*in silico* saturated mutagenesis" analysis for discovering high-impact sites within a 1-kb sequence. Specifically it performs computational mutation for every base of the input sequence and predicts the effect of every mutation on chromatin accessibility. The accepted inputs are a chromosome and a position, a BED file containing multiple coordinates of genomic regions or a custom sequence. For the BED file, due to the high computational intensity, our service only runs the first five regions. For the custom sequence, a sequence with an effective length greater than 20-bp is accepted. However, since the input of the DNN model is a 1-kb sequence, if the input sequence is less than 1-kb in length, N will be added to both ends of the sequence until the length is equal to 1-kb, which may cause bias in the prediction. When the sequence is greater than 200-bp in length, PlantDeepSEA will perform *in silico* saturation mutagenesis analysis on the middle 200-bp of the input sequence.

Please remember your task ID and check the result page in a few minutes. In addition, your prediction results will be stored on our server for 7 days only.

Examples of Upload Files(Based on Minghui63 Genome of Rice): [VCF](#) [BED](#) [FASTA](#)

## SUBMIT TASK

There are three ways to submit a task.

For Variant Effector, you can upload a VCF for query. The uploaded file must be a .vcf file ending. The format of the file can be viewed by downloading the sample file in the lower right corner.

For Sequence Profiler, you can:

- (1) Upload a BED for query. The uploaded file must be a BED file ending. The format of the file can be viewed by downloading the sample file in the lower right corner. For BED files, we only predict the first 5 lines of the file due to limited GPU resources.
- (2) Using chromosome and position queries. The results will output the effect on the open region of chromatin before and after the 200 bp sequence mutation around the position.

ZS97 Model

The screenshot shows the ZS97 Model interface with three main sections:

- Variant Effector:** A green header with a dropdown for "Reference genome" set to "Nipponbare (MSU7.0)". Below it is a red box labeled "Click here to submit VEF file" pointing to a blue "Upload" button and a "Choose file" input field.
- Sequence Profiler:** A green header with a dropdown for "Reference genome" set to "Nipponbare (MSU7.0)". Below it is a red box labeled "Click here to submit BED file" pointing to a blue "Upload" button and a "Choose file" input field.
- Combined Section:** A green header with a dropdown for "Reference genome" set to "Nipponbare (MSU7.0)". It contains:
  - "(3) Select a Chromosome and input a Position": A dropdown for "Choose Chromosome" and a text input for "Position (e.g: 2018365)".
  - "(4) Input a sequence (>20-bp is allowed, ≥1-kb is recommended)": A text input for "sequence (e.g: ATCG)" and a blue "Submit" button.

**Notice:** This model was constructed based on ATAC-seq data from multiple tissues of the Rice (*Oryza sativa L.*) variety Zhenshan 97 using the deep learning-based algorithmic framework DeepSEA (Zhou et al., *Nature Methods*, 2015) and was implemented using the Selene SDK (Chen et al., *Nature Methods*, 2019). We provide two services based on this model.  
1. Variant Effector, aims to predict the effects of sequence variants on chromatin accessibility. The accepted input is a VCF file containing information on the sequence variants. The results contain information on the effect of variants on chromatin accessibility in each tissue.  
2. Sequence Profiler, is a utility that performs *in silico* saturated mutagenesis analysis for discovering high-impact sites within a 1-kb sequence. Specifically it performs computational mutation for every base of the input sequence and predicts the effect of every mutation on chromatin accessibility. The accepted inputs are a chromosome and a position, a BED file containing multiple coordinates of genomic regions or a custom sequence. For the BED file, due to the high computational intensity, our service only runs the first five regions. For the custom sequence, a sequence with an effective length greater than 20-bp is accepted. However, since the input of the DNN model is a 1-kb sequence, if the input sequence is less than 1-kb in length, N will be added to both ends of the sequence until the length is equal to 1-kb, which may cause bias in the prediction. When the sequence is greater than 200-bp in length, PlantDeepSEA will perform *in silico* saturation mutagenesis analysis on the middle 200-bp of the input sequence.  
Please remember your task ID and check the result page in a few minutes. In addition, your prediction results will be stored on our server for 7 days only.

Examples of Upload Files(Based on Minghui63 Genome of Rice): VCF BED FASTA

Click here to select reference genome

Click here to download the sample file

Click here to enter the FASTA sequence

Click here to select chromosome and position queries

Click here to submit VEF file

Click here to submit BED file

Note: Only rice model supports coordinates and intervals of multiple reference genomes.

After the task is submitted, it will return to Task ID and run in the background queue. After the task is submitted successfully, it will jump to the waiting interface as shown in the figure below.

For sequence or position input:

2803_16175216338581753	Your TaskID : 2803_16175216338581753 0 sub-mission completed. Still running. Please wait ~ ● ● ● ● ● ● ● ●
------------------------	---

**Tips :**

- 1) Your prediction results will be stored on our server for 7 days only, results beyond this time will not be retrieved.
- 2) Make sure you submit the correct task ID. After your submitting, this page will automatically refresh in every 20 seconds until your result shows up.
- 3) The task will only return results after it is run, so please wait for a while after submitting.
- 4) example task IDs: **2566\_16148611193615174** or **2021\_16084538926592038** .

For BED input:

# Search Your Result

2809\_16175219867692175

Submit

The following prediction requests were received.

Show 10 entries Search:

Chromosome	Start	End
chr09	8998909	8998947
chr10	8461182	8461732
chr11	18836169	18836959
chr12	2794339	2794477
chr12	22076166	22076893

Showing 1 to 5 of 5 entries      Previous 1 Next

Your TaskID : 2809\_16175219867692175  
 0 sub-mission completed.  
 Still running. Please wait ~

Tips :

- 1) Your prediction results will be stored on our server for 7 days only, results beyond this time will not be retrieved.
- 2) Make sure you submit the correct task ID. After your submitting, this page will automatically refresh in every 20 seconds until your result shows up.
- 3) The task will only return results after it is run, so please wait for a while after submitting.
- 4) example task IDs: [2566\\_16148611193615174](#) or [2021\\_16084538926592038](#).

For VCF input:

**Search Your Result**

Submit

The following prediction requests were received.

Chromosome	Position	Variation ID	Ref	Alt
chr09	17599142	vg0916410213	C	A
chr09	17599229	vg0916410299	C	G
chr09	17599532	vg0916410603	G	C
chr09	17599577	vg0916410648	G	T
chr09	17599855	vg0916410925	T	G
chr09	17599855	vg0916410925	T	G
chr09	17599991	vg0916411061	C	G
chr09	17599998	vg0916411068	G	T

Showing 1 to 8 of 8 entriesPrevious1Next

Your TaskID : 2808\_161752196048525  
0 sub-mission completed.  
Still running. Please wait ~

• • • • •

Tips :

- 1) Your prediction results will be stored on our server for 7 days only, results beyond this time will not be retrieved.
- 2) Make sure you submit the correct task ID. After your submitting, this page will automatically refresh in every 20 seconds until your result shows up.
- 3) The task will only return results after it is run, so please wait for a while after submitting.
- 4) example task IDs: [2566\\_16148611193615174](#) or [2021\\_16084538926592038](#).

Remember your task ID and query the results after a while (a few minutes), and your prediction result will store on our server for only 7 days.



---

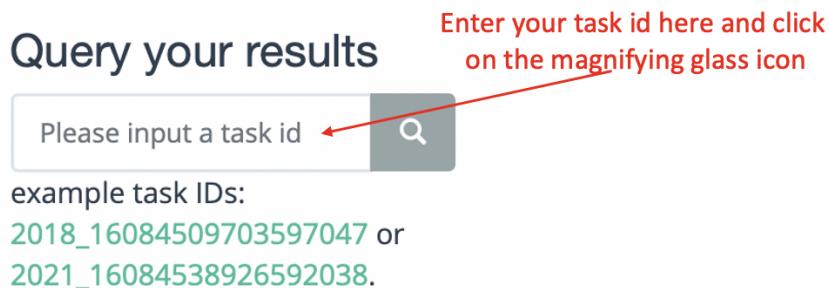
## CHAPTER FOUR

---

### QUERY RESULTS

There are two ways to query the results of submitted tasks.

You can enter your task ID in the query your results column on the right side of the home page and click on the magnifying glass icon (as shown in the figure below).



You can also click the results button on the home page to enter the result query page, and then enter the task ID in the interface and submit.

Enter task ID here (for example, 2021\_16084538926592038) to query the results.

Search Your Result

Input your task ID here

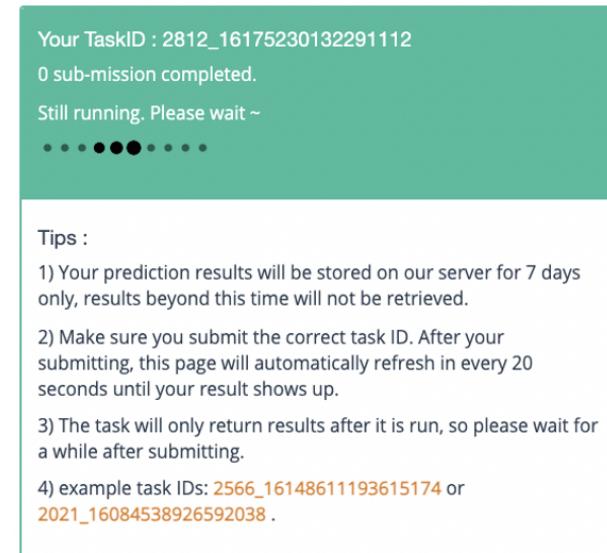
Submit

Please input one task id.

Tips :

1) Your prediction results will be stored on our server for 7 days only, results beyond this time will not be retrieved.  
2) Make sure you submit the correct task ID. After your submitting, this page will automatically refresh in every 20 seconds until your result shows up.  
3) The task will only return results after it is run, so please wait for a while after submitting.  
4) example task IDs: 2566\_16148611193615174 or 2021\_16084538926592038 .

After you submit the query, the information as shown in the following figure will be displayed on the right side of the interface.



The interface will refresh every 10 seconds until the result is returned.

## INTERPRETATION OF RESULTS

After the task is completed, enter the task ID query to obtain the model prediction results.

### 5.1 Brief Description

The left side of the upper part of the result interface is the task ID query column, and the right side shows the download link of the forecast result compression package and task information (as shown in the figure below, take chromosome number and position input as an example).

Search Your Result

2566\_16148611193615174

Submit

Task ID : 2566\_16148611193615174

Mission Info : The MH63 model was used to obtain the input information according to the chromosome number and position.

Position :

Chromosome : chr09

Start : 17599129

End : 17599329

Download the predicted results:

[chr09\\_17599129\\_17599329\\_16148611193615174.zip](#)

Click here to download the result file package

The lower part of the result interface is a visual display of the prediction results generated. For different task types, the returned results are slightly different, see the **Main Functions** part of this section for details. Here are two examples (see the figure below) for reference.

Enter the chromosome number and position, and return the following result. The visualization result is generated by bokeh.

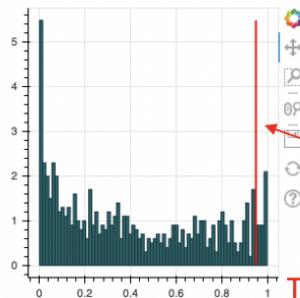
## Results:

You can download the predicted results on this page or select one of the sample you would like to display:

ATAC-seq, flag leaf, rep1

**Click here to select the sample to display**

The currently selected model was trained using the data "ATAC-seq, flag leaf, rep1". The distribution of the prediction scores for 744 sequences labeled as open chromatin regions (OCRs) in the test set is shown below. The prediction score (0.95) of the submitted sequence "chr09:17599129-17599329" is at the 93.7% quantile of this distribution, as shown by the red line.



**The more right the red line is, the higher the confidence is.**

### Sequence used for prediction

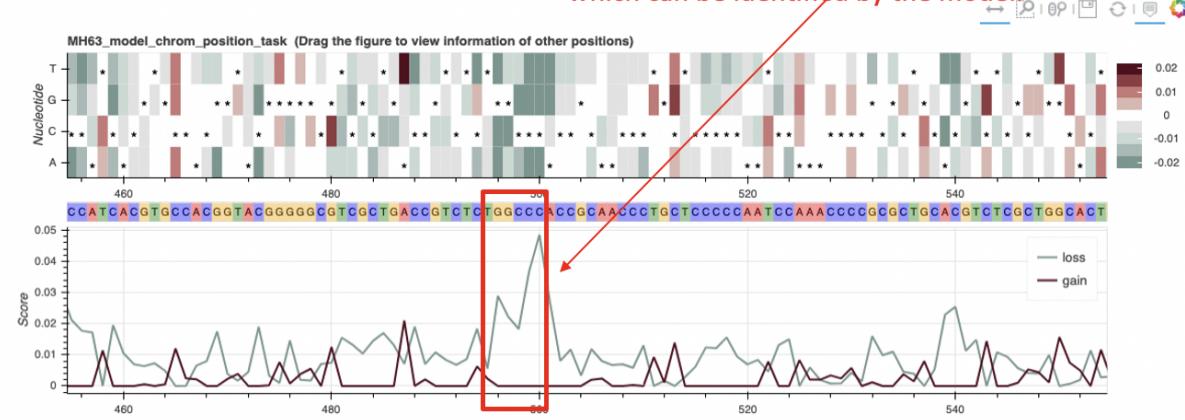
The bases shown in the '*in silico* saturation mutagenesis map' section are highlighted (with line width: 100 bp).

```
ACACAGCTAAATGAATCGCATATGCAGTGAATGAGACTCTGAGAGCCAGTTAAGGTAGGAGTAGCATCAAATTCACTAGTCCTATAACCGGGCACGT  
CGTCGAGTGAATTTATACGATCTCATGTCAGTCATCAGTGCATGCTGTTGAGAAGGGAGTTGGCCCTGCCGCTGCTGAGCTCAACTGAACGC  
TGGCTCCTGCATGGCCAGTTCGACTACAGTTACTACGGCATACCCGGGGTAGCACGGCCGGCTACAGCTACCCGGAGTAGGAGTAGCAGC  
AGTGGTGGTAGTAGAAAGCGCGCAGTGGAAAGGGATAGAGGAGTCCCACATTAATAACTCGTTCAAATCGAACGGCGTGGCTACTCTAGTACACC  
CATTCAACGGGCCGTTCCACTCCACTCCACCGTAACTGCGCGCTGGGAGCCCATCACGTGCCACGGTACGGGGCGTCGCTGACCGTCTCGCC  
CACCGCAACCCCTGCTCCCCAATCCAACCCCGCGCTGACGTCTCGCTGGCACTGGGCAATTGATCCATCGCGTGCCTGGGGGGGGGGGGGGGG  
CGCGCGGGGGCGCCGGCACGCGAGACCGATGTAGACGTGTCACCGGGGAAGCTGTCCGCTCGGGTGCCCTGCGCGGGGGGGGGGGGGGGGG  
GAGCCGTACGTGCGTGCATACCTCGGTGCGTCCCTGCAAGCCGGCATCGCTGCCGCTGCTCAATTATTCCCTGCTGTTCATTCGTACG  
TAGTCGCGCTGGGATGCGCCATAGCCATATCTCGGCAGGACCCTGACGCTCGCTGGCAACTGTACGTGCCGCTCAGGACCGG  
CAGTCATGCGGAGTTGGTACAAACTTGCATGAAAGGAGTGGTACTGTGGTACCCACTGCCCTGCTTGACACACATCAGTCTCCGGAGTTAGC
```

### *in silico* saturated mutagenesis map

In the heatmap, each row represents mutations in the corresponding nucleotide, each column represents a position in the sequence. The color of the heatmap presents the predicted change in probability that an allele belongs to open chromatin. The loss score refers to the maximum decrease in probability that an allele belongs to open chromatin compared to the reference nucleotide in all mutations at each site. And the gain score refers to the maximum increase.

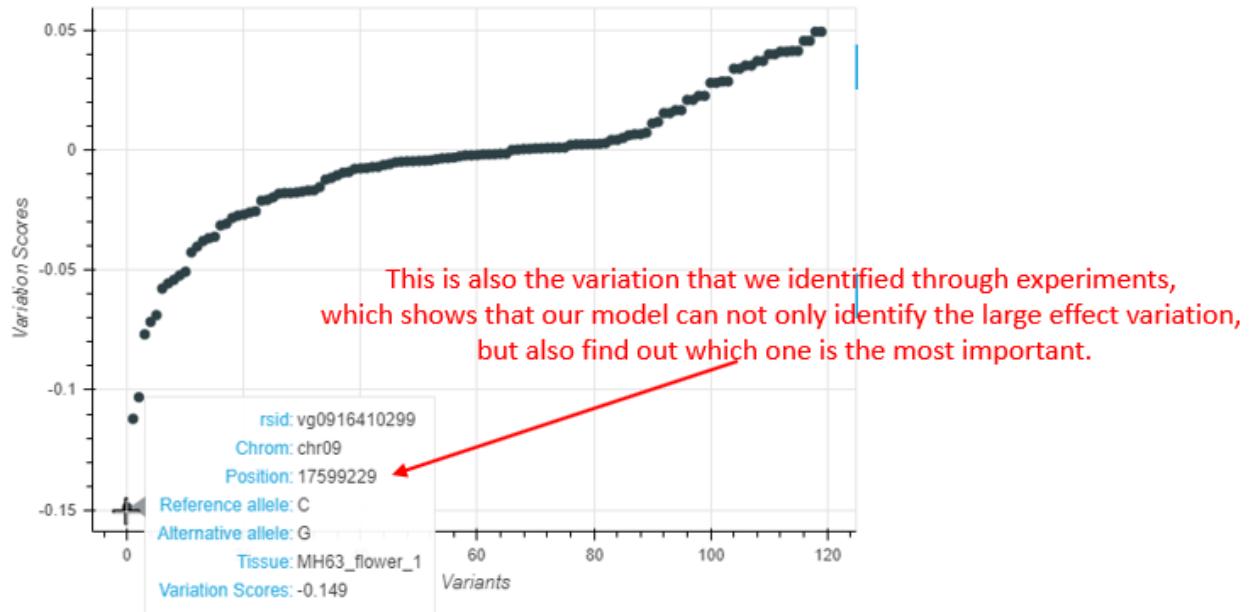
**We have proved that an important motif, which can be identified by the model.**



For VCF format input, the interactive scatter plot is returned as follows.

## Results:

[display.html](#)



In addition, at the bottom of the results interface, you may also get a fimo scan result, if fimo does scan the result(as shown below).

### Regulatory elements identified by FIMO

Detailed Fimo results and the parameters used can be viewed at [fimo.html](#).

For the motifs identified by FIMO (with the *P*-value < 1e-4), we calculated the average loss and gain values in the *in silico* saturated mutagenesis map corresponding to the motif regions, which are listed in the table below as 'Loss value' and 'Gain value', respectively.

Motif ID	Name	Start	End	Strand	p-value	q-value	Matched Sequence	Loss value	Gain value
<a href="#">MP00534</a>	AT5G38860	449	466	+	2.13e-06	5.81e-04	GCGACCCATCACGTGCC	1.13e-02	3.45e-03
<a href="#">MP00594</a>	AT5G67000	573	587	+	1.04e-05	3.07e-03	CGGTCGCGGCCGCGG	5.76e-03	5.85e-03
<a href="#">MA1261.1</a>	AT5G67000	573	587	+	1.04e-05	3.07e-03	CGGTCGCGGCCGCGG	5.76e-03	5.85e-03
<a href="#">MP00078</a>	AT3G22170	437	448	-	1.16e-05	3.63e-03	CCCACGCGCGCA	9.72e-03	1.16e-03
<a href="#">MA0557.1</a>	FHY3	437	448	-	1.16e-05	3.63e-03	CCCACGCGCGCA	9.72e-03	1.16e-03
<a href="#">MA1011.1</a>	PHYPADRAFT_72483	457	466	+	1.22e-05	2.30e-03	ATCACGTGCC	6.18e-03	2.99e-03
<a href="#">MA1011.1</a>	PHYPADRAFT_72483	457	466	-	1.22e-05	2.30e-03	GGCACGTGAT	6.18e-03	2.99e-03
<a href="#">MP00102</a>	AT4G14410	457	466	+	1.41e-05	2.15e-03	ATCACGTGCC	6.18e-03	2.99e-03
<a href="#">MP00102</a>	AT4G14410	457	466	-	1.41e-05	2.15e-03	GGCACGTGAT	6.18e-03	2.99e-03
<a href="#">MA0960.1</a>	BHLH104	457	466	+	1.41e-05	2.15e-03	ATCACGTGCC	6.18e-03	2.99e-03

Showing 1 to 10 of 119 entries

Previous 1 2 3 4 5 ... 12 Next

## 5.2 Main Functions

### 1. Variant Effector

We predicted changes in chromatin features in a 1-kb region around the mutation (upstream 500 bp and downstream 500 bp). You are supposed to submit a vcf file to generate a background predictiton task, note that there is no limit for the num of lines of your VCF file by now. For each mutation in your VCF file, we predict the scores of changes (may be positive or negative) in every sample of this model, then display all results via scatterplot.

**Input:** A VCF file in any length.

**Output:** After your submission, our system will return a task id, which is essential, go to the result page to find out your predictiton result via the task id. Input your task id to search your result stored on our server (contains the vcf file you submitted and all scores file and result scatterplot in html format)

### 2.Sequence Profiler

Sequence Profiler is used to analyse possible motif or important mutation in a 1-kb DNA sequence from reference genome file. You can generate a predictiton request through 2 ways:

(1)By select a chromomsome and input a position:

After receiving the chromosome and position information, we predict the mutation of the middle 200-bp sequence in the range of the middle 1-kb region (for each position of the 200-bp, our server simulate the difference mutation effect of the reference allele and the other 3 alter allele, thus return a 200\*4 matrix afterwards, and display the matrix in the html result page through a interactive bokeh heatmap).

**Input:** Select a chromosome and input a position

**Output:** The web page returns a heatmap, a line plot and fimo result, the heatmap shows the 200\*4 matrix and line plot shows continuous positive or negative scores of a possible motif (named gain line or loss line respectively). Fimo result shows possible motif get in meme database of this DNA sequence.

(2) By upload a BED file:

Unlike predicting VCF file in Variant Effector, we only calculate the top 5 lines of your BED file due to our limited GPU resource. For each line in the top 5 lines of your BED file, the server does the same work described in (1).

**Input:** A BED file in any length.

**Output:** Prediction result of the top 5 lines of your BED file, web result description same as (1).

(3) By inputting the FASTA sequence:

Input any length of FASTA sequence you are interested in, and the model will intercept the middle 1000bp (if the input length is greater than 1000bp) or fill 1000bp at both ends (if the input length is less than 1000bp) for prediction.

**Input:** Any length of FASTA sequence you are interested in.

**Output:** The web page returns a heatmap, a line plot and fimo result, the heatmap shows the 200\*4 matrix and line plot shows continuous positive or negative scores of a possible motif (named gain line or loss line respectively). Fimo result shows possible motif get in meme database of this DNA sequence.

### 5.3 About the result files

Got task ID : 2021\_16084538926592038

Mission Info : MH63\_model\_vcf\_task

[Click here to download the compressed file of the results.](#)

The VCF File You Submit :

[dep1\\_MH63\\_ref.vcf](#)

Download the predicted results:

[vcf\\_16084531203644767.zip](#)

1. The task of predicting sequence (For example, ID 2847\_16177604444503944)

The contents of the compressed package are generally as follows:

File Name	Note
cisml.xml	Motif information (CisML format) matched by 200 bp segment scanned by Fimo
fimo.gff	Motif information (GFF format) matched by 200 bp segment scanned by Fimo
fimo.html	Motif information (HTML format) matched by 200 bp segment scanned by Fimo
fimo.tsv	Motif information (TSV format) matched by 200 bp segment scanned by Fimo
fimo.xml	Motif information (in XML format) matched by 200 bp segment scanned by Fimo
fimo_cut.tsv	Optimized(by PlantDeepSEA) version of fimo.tsv
do_fimo.bed	Record the position information of 200 bp section
do_fimo.fasta	The FASTA file of the sequence
ism_logits.tsv	$\log(P_{\text{mut}} / (1-P_{\text{mut}})) - \log(P_{\text{ref}} / (1-P_{\text{ref}}))$ , the difference between <i>logit(mut)</i> and <i>logit(ref)</i> predictions
ism_diffs.tsv	$P_{\text{mut}} - P_{\text{ref}}$ , the difference between <i>mut</i> and <i>ref</i> predictions
ism_abs_diffs.tsv	$ P_{\text{mut}} - P_{\text{ref}} $ , the absolute difference between <i>mut</i> and <i>ref</i> predictions
ism_predictions.tsv	After Selene SDK prediction, the input sequence was scored three times per mutation

The mut in the formula in the table refers to mutagenesis.

The interpretation of various formats of Fimo results is detailed [here](#).

ism\_predictions.tsv is the score of each mutation of the input 1000bp sequence three times after the prediction of Selene SDK. It records the absolute value of the difference, the true value of the difference, the Logits value and the true score of each mutation sequence.

Each PDF file in the compressed package is the drawing of each organization.

## 2.VCF variation annotation task (For example,[ID2870\\_16178001127678075](#))

The contents of the compressed package are generally as follows:

File Name	Note
dep1_MH63_ref.vcf	VCF files submitted by users
dep1_MH63_ref_abs_diffs.tsv	$ P_{\text{mut}} - P_{\text{ref}} $ , Absolute value of difference before and after mutation of each locus
dep1_MH63_ref.alt_predictions.tsv	Scoring value of each locus variation
dep1_MH63_ref_diffs.tsv	$P_{\text{mut}} - P_{\text{ref}}$ , value of difference before and after mutation of each locus
dep1_MH63_ref_logits.tsv	Logits value before mutation of each locus
dep1_MH63_ref.ref_predictions.tsv	Score before variation of each locus

The above TSV files are the scoring files predicted by [selene SDK](#)

---

**CHAPTER  
SIX**

---

**BLAST**

Click the blast button in the main interface to enter the blast (for conservative sequence analysis) page. On this page, you can enter the sequence of interest and select the appropriate database and parameters to query (as shown in the figure below).

**PlantDeepSEA-Multi-Genome-BLAST**

Input your sequence below:

GC GGAGGAGGCGGAAGCGGCGC

**Input query sequence** 

Select database —> Database: Z597

E-value: 0.001

Word size: 11

Sensitivity: NORMAL

**Setting parameters** 

Copyright © 2020 National Key Laboratory of Crop Genetic Improvement, Huazhong Agricultural University

In the returned blast result interface, click the corresponding result to jump to the model (same as database) prediction interface, and automatically fill the corresponding position to (3) by select a chromosome and input a position under the sequence profiler. The user can directly click submit to submit the task for prediction.

## PlantDeepSEA-Multi-Genome-BLAST-Result

Got Zhenshan 97 BLASTN Result !

BLASTN 2.9.0+

Description

#	contig	query	length	e-value	score	ident %
1	chr10	seq1	25619316	0.000155868	46.0	100.0
2	chr08	seq1	30215435	0.000155868	46.0	100.0
3	chr01	seq1	44482313	0.000155868	46.0	100.0

Click here to jump to the model prediction interface,  
and the corresponding location information will be filled in automatically.

Alignments

chr10 length = 25619316

hsp	1
length	23
e-value	0.000155868
score	46.0
identities	23
positives	23
bits	42.8
query start	1
query end	23
subject start	2474308
subject end	2474286

Q: GCGGAGGGAGCCGAAGCGCGC  
M: ||||||| :  
S: GCGGAGGGAGCCGAAGCGCGC

**ZS97 Model**

Variant Effector

(1) By Upload VCF file

**Upload one file:**  未选择文件

**Upload**

Sequence Profiler

(2) By Upload BED file

**Upload one file:**  未选择文件

**Upload**

(3) By select a Chromosome and input a Position

Do not know the position? You can go to [Blast Page](#)

**Submit**

(4) By input a sequence

**Submit**

Notice:  
The model corresponds to the choice of database

This model was constructed based on ATAC-seq data from multiple tissues of the Rice (*Oryza sativa L*) variety Zhenshan 97 using the deep learning-based algorithmic framework DeepSEA (Zhou et al., *Nature Methods*, 2015) and was implemented using the Selene SDK (Chen et al., *Nature Methods*, 2019).

We provide two services based on this model.

1. **Variant Effector**, aims to predict the effects of sequence variants on chromatin accessibility. The accepted input is a VCF file containing information on the sequence variants. The results contain information on the effect of variants on chromatin accessibility in each tissue.
2. **Sequence Profiler**, is a utility that performs "in silico saturated mutagenesis" analysis for discovering high-impact sites within a sequence. Specifically it performs computational mutation for every base of the input sequence and predicts the effect of every mutation on chromatin accessibility. The accepted inputs are a chromosome and a position, or a BED file containing multiple coordinates of genomic regions. For the BED file, due to the high computational intensity, **our service only runs the first five regions**.

Please remember your **task ID** and check the result page in a few minutes. In addition, your prediction results will be stored on our server for **7 days** only.

Example Upload File: [VCF](#) [BED](#)

The corresponding location information has been automatically filled in.



## CASE STUDIES

### 7.1 Case1: The DEP1 gene of rice

A recent study(Fu , Xu etal. ,2019) have shown that nine NCVs in the DEP1 promoter region can regulate the gene expression and leaf-trait variation .

- (1) Mapping these nine variants to RiceVarMap database and constructing the VCF file .
- (2) Using PlantDeepSEA to make predictions for this VCF file .
- (3) One SNP, named vg0916410299 in RiceVarMap, has the greatest effect score (example result ID:[2870\\_16178001127678075](#)).
- (4) Entering the genomic coordinates(Chromosome:chr09 ; position:17599229) of vg0916410299 into ‘Sequence Profiler’ for prediction. The sequence TGGCCC, which overlaps with vg0916410299, has the extreme effect scores. FIMO results also indicate that this sequence overlaps with a binding motif of the TCP transcription factors.(example result ID: [2847\\_16177604444503944](#)) .

### 7.2 Case2: The UPA2 gene of maize

Two validated 240bp regions of the (UPA2) haplotype(Tian , Wang etal. ,2019) were analyzed with “Sequence Profiler”

The in silico saturated mutagenesis map showed the haplotype of CIMMYT 8759 (with AGTGTG) has intensive high effect scores in the C2C2 motif region compared to the haplotype of W22 (with AGTG–) (example result ID: [2867\\_16177642225564935](#) and [2868\\_1617764237554985](#)) .



---

**CHAPTER  
EIGHT**

---

**RESOURCES**

## 8.1 Brief Description

Click on the “Resources” link in the main menu to view training data statistics, quality control information, and model evaluation results. In addition, we provide links to download the reference genomes, the trained models and files for training models.

Please note: *We are still analyzing the ATAC-seq data we generated from multiple tissues of grass species. We have provided lists of OCRs obtained from these ATAC-seq data so that researchers can use it to reproduce the results of PlantDeepSEA. However, please do not use these lists for genome-scale analyses, until the publication of our article on data analysis (by May 2022 at the latest). Or contact us for permission. Users who download data from this site are considered to have accepted this statement.*

## 8.2 Reference Genome Information (Download from Google Drive)

variety	Species	Version	Reference link
Arabidopsis	Arabidopsis thaliana	TAIR10.1	<a href="#">Download link</a>
Rice	Oryza sativa	MSU v7.0	<a href="#">Download link</a>
Rice	Oryza sativa	MH63 RS2	<a href="#">Download link</a>
Rice	Oryza sativa	ZS97 RS2	<a href="#">Download link</a>
Brachypodium	Brachypodium distachyon	Bd21-3 v1.1	<a href="#">Download link</a>
Foxtail millet	Setaria italica	v2.0	<a href="#">Download link</a>
Sorghum	Sorghum bicolor	v3.1.1	<a href="#">Download link</a>
Mazie	Zea mays	AGPv4	<a href="#">Download link</a>

## 8.3 Reference Genome Information (Download from Baidu Cloud Disk Drive)

variety	Species	Version	Reference link
Arabidopsis	Arabidopsis thaliana	TAIR10.1	<a href="https://pan.baidu.com/s/1TvqpN94yE6gq7VHVRNt6kA?pwd=iv43">https://pan.baidu.com/s/1TvqpN94yE6gq7VHVRNt6kA?</a> pwd=iv43 iv43
Rice	Oryza sativa	MSU v7.0	: <a href="https://pan.baidu.com/s/1SYpyzJLd6X_3sTVoTbyxbw">https://pan.baidu.com/s/1SYpyzJLd6X_3sTVoTbyxbw</a> : tdgi
Rice	Oryza sativa	MH63 RS2	<a href="https://pan.baidu.com/s/1uXqWbSdipnNHfckH6xOZoA?pwd=7fed">https://pan.baidu.com/s/1uXqWbSdipnNHfckH6xOZoA?</a> pwd=7fed 7fed
Rice	Oryza sativa	ZS97 RS2	<a href="https://pan.baidu.com/s/1Xm8GoRQ2wCt1aSjU7DEqUQ?pwd=bhtw">https://pan.baidu.com/s/1Xm8GoRQ2wCt1aSjU7DEqUQ?</a> pwd=bhtw bhtw
Brachypodium	Brachypodium distachyon	Bd21-3 v1.1	: <a href="https://pan.baidu.com/s/1hiUYOE9g3KwtDucUkRkPTw">https://pan.baidu.com/s/1hiUYOE9g3KwtDucUkRkPTw</a> : r4ka
Foxtail millet	Setaria italica	v2.0	<a href="https://pan.baidu.com/s/1qjwnowHMEf_psOLpTgEo0g?pwd=4m10">https://pan.baidu.com/s/1qjwnowHMEf_psOLpTgEo0g?</a> pwd=4m10 4m10
Sorghum	Sorghum bicolor	v3.1.1	<a href="https://pan.baidu.com/s/1AqbUPws9AyTtjnBv3vUMmQ?pwd=3j0q">https://pan.baidu.com/s/1AqbUPws9AyTtjnBv3vUMmQ?</a> pwd=3j0q 3j0q
Mazie	Zea mays	AGPv4	<a href="https://pan.baidu.com/s/1Ppexb6mzC_FtfvIzfYb4vw?pwd=1xf8">https://pan.baidu.com/s/1Ppexb6mzC_FtfvIzfYb4vw?pwd=</a> 1xf8 1xf8

## 8.4 Trained model

variety	Reference link
Arabidopsis	<a href="#">Download link</a>
Zhenshan 97	<a href="#">Download link</a>
Minghui 63	<a href="#">Download link</a>
Brachypodium	<a href="#">Download link</a>
Foxtail millet	<a href="#">Download link</a>
Sorghum	<a href="#">Download link</a>
Mazie	<a href="#">Download link</a>

## 8.5 Training data

variety	Reference link
Arabidopsis	<a href="#">Download link</a>
Zhenshan 97	<a href="#">Download link</a>
Minghui 63	<a href="#">Download link</a>
Brachypodium	<a href="#">Download link</a>
Foxtail millet	<a href="#">Download link</a>
Sorghum	<a href="#">Download link</a>
Mazie	<a href="#">Download link</a>

## 8.6 Training config files

variety	Reference link
Arabidopsis	<a href="#">Download link</a>
Zhenshan 97	<a href="#">Download link</a>
Minghui 63	<a href="#">Download link</a>
Brachypodium	<a href="#">Download link</a>
Foxtail millet	<a href="#">Download link</a>
Sorghum	<a href="#">Download link</a>
Mazie	<a href="#">Download link</a>

## 8.7 Model structure

model name	Reference link
DeeperDeepSEA	<a href="#">Download link</a>

This model structure is Selene's DeeperDeepSEA structure .

## 8.8 A guide to applying the model locally

You can click [here](#) to see how the models are trained and how to apply the trained models provided in PlantDeepSEA locally.



---

CHAPTER  
NINE

---

**STATISTICS**

You can slide the table to see more information.

Species	Samples	Raw reads number	Mapping rate	Markdoup rate	Q30 reads number	TSS Enrichment	Peak number	Sample name	Number of sequences labeled as OCR in training set	Ratio of sequences labeled as OCR in training set	Number of sequences labeled as OCR in test set	Ratio of sequences labeled as OCR in test set	AUC	AUPRC
A.	Ara-bidop-thalidam_stem_cell_rep3	151660948953	0.25649713306813	0.3019732950966	1985AC seq, stem cells, rep3	-139540	0.0272529962	50.0249309438401538						
A.	Ara-bidop-thalidam_stem_cell_rep2	997336029951	0.203540483857164988736684	0.27226AC	-240981	0.0470666615605043666984979087398								
A.	Ara-bidop-thalidam_stem_non_hair_cell_rep1	9407985029949	0.101730624398530735054894	0.18815AC	-344379	0.067266423430506430288781638252								
A.	Ara-bidop-thalidam_stem_tip_rep1	479528009871	0.060392488436951652722942	0.050015AC	-344701	0.067324534060506540405743689505								

continues on next page

Table 1 – continued from previous page

Species	Samples	Raw reads number	Mapping rate	MarkdU30 reads number	TSS Enrichment	Peak number	Sample name	Number of sequences labeled as OCR in training set	Ratio of sequences labeled as OCR in training set	Number of sequences labeled as OCR in test set	Ratio of sequences labeled as OCR in test set	AUC	AUPRC
A.	<i>Arabidopsis thaliana</i> _stem_cell_rep1	846878049949	0.2180939062	4645169329395	5ATAC	356863	0.069698946805690	10.9562667357					
A.	<i>Arabidopsis thaliana</i> _leaf_7days_leaf_rep1	452255029949	0.304	13461588.32	25571	ATAC	3935800.076871093750.071016.92149085419						
A.	<i>Arabidopsis thaliana</i> _root_tip_rep2	121576499858	0.08497109088859762200784	5ATAC	3953220.077217588126.075920.97435892883								
A.	<i>Arabidopsis thaliana</i> _root_hair_cell_rep1	10916207993	0.10396284278950922633886	8ATAC	4145660.080968981876.081930.98974058947								
A.	<i>Arabidopsis thaliana</i> _mesophyll_cell_rep1	112729519914	0.3276053380926378927465	9ATAC	5098070.0995710096805100660.97990.9562								
A.	<i>Arabidopsis thaliana</i> _mesophyll_cell_rep3	100129049927	0.25777809057304088265024944AC	5466160.1067609375	0.107870.990508446								

continues on next page

Table 1 – continued from previous page

Species	Samples	Raw reads number	Mapping rate	MarkdU@30 reads number	TSS Enrichment	Peak number	Sample name	Number of sequences labeled as OCR in training set	Ratio of sequences labeled as OCR in training set	Number of sequences labeled as OCR in test set	Ratio of sequences labeled as OCR in test set	AUC	AUPRC
A.	Ara-bidop-thaliamot_non_hair_cell_rep2	242969529936	0.057855546678750582696243	3416AC- seq, root non-hair, rep2	557959	0.108976368180510797095681025593							
A.	Ara-bidop-thaliamesophyll_cell_rep2	999383089932	0.2860128268704610974239160251AC-	544592 seq, meso-phyll cells, rep2	544592	0.1063656825	0.108874995610964885						
A.	Ara-bidop-thaliamot_hair_cell_rep2	291715489938	0.064234858303615726267335	ATAC- seq, root hair, rep2	570957	0.111515029060511217048827048039							
A.	Ara-bidop-thaliamdays_leaf_rep2	149379047171	0.3	30345758.54	33398	ATAC- seq, leaf, 7 days, rep2	585868	0.11442782750.11288060897035818					
B.	Bdis-tachyon_root_1	145430284715	0.0551	5170227.4744296244	ATAC- seq, root, rep1	40093	0.00783986406250066002955510225137						
B.	Bdis-tachyon_young_leaf_2	102916469754	0.1146	7932578.1456617113	ATAC- seq, young leaf, rep2	74266	0.014505028126.011584952420137792						
B.	Bdis-tachyon_young_leaf_1	108702049604	0.1215	8569899.7817823358	ATAC- seq, young leaf, rep1	112185	0.0219128028105018850.99842068269						

continues on next page

Table 1 – continued from previous page

Species	Samples	Raw reads number	Mapping rate	MarkdU30	reads number	TSS Enrichment	Peak number	Sample name	Number of sequences labeled as OCR in training set	Ratio of sequences labeled as OCR in training set	Number of sequences labeled as OCR in test set	Ratio of sequences labeled as OCR in test set	AUC	AUPRC
B.	Bdistachyon_flower_2 distachyon	243299529589	0.5249	84039674.747946148				ATAC seq, flower, rep2	162902	0.0318137796876	0.0253909398	0.5461		
B.	Bdistachyon_root_2 distachyon	457011565064	0.0649	173550649749663702				ATAC seq, root, rep2	180246	0.035204296876	0.0286909351586216			
B.	Bdistachyon_flower_1 distachyon	395727069536	0.5504	13233688.728452655				ATAC seq, flower, rep1	215663	0.0421256096805034348960240852589				
B.	Bdistachyon_flag_leaf_2 distachyon	639324969606	0.3272	37316076.099052939				ATAC seq, flag leaf, rep2	439499	0.085839678430507213092842015894				
B.	Bdistachyon_panicle_1 distachyon	116923089716	0.5741	395569215.731589336				ATAC seq, young panicle, rep1	408596	0.079808904250.074350445897086807				
B.	Bdistachyon_flag_leaf_1 distachyon	940677029695	0.3964	49821067.450757838				ATAC seq, flag leaf, rep1	452662	0.0884105476876	0.07522898866896854			
O.	ZS97_sativa	576204869742	0.3419	28421794.377067787				ATAC seq, young panicle, rep1	34913	0.00681894531050084639416606666666				

continues on next page

Table 1 – continued from previous page

Species	Samples	Raw reads number	Mapping rate	MarkdQ30	reads number	TSS Enrichment	Peak number	Sample name	Number of sequences labeled as OCR in training set	Ratio of sequences labeled as OCR in training set	Number of sequences labeled as OCR in test set	Ratio of sequences labeled as OCR in test set	AUC	AUPRC
O.	MH63_sativa	858925029856	0.151	584615	157024	693603	ATAC seq, stamen & pistil, rep1	49527	0.009673421	8050088	20.937466896551			
O.	MH63_sativa	103252060853	0.2364	761475	9522170	926572	ATAC seq, stamen & pistil, rep2	68509	0.01338066406050121	90.98396655626				
O.	ZS97_sativa	p650029069873	0.3547	291508640.718667179			ATAC seq, young panicle, rep5	140295	0.0274048671	8050160	90.45452096349			
O.	ZS97_sativa	p51132549902	0.2823	28807070.307267392			ATAC seq, young panicle, rep4	148634	0.0290360778126	0.017928.055890.45588				
O.	MH63_sativa	647456029861	0.2291	306585	828341	29063	ATAC seq, flag leaf, rep2	119482	0.023338928126	0.019860.955880.65532				
O.	ZS97_sativa	y674438678755	0.5716	45901408.782362667			ATAC seq, young leaf, rep1	234188	0.045738823750.0281	10.04539915805				

continues on next page

Table 1 – continued from previous page

Species	Samples	Raw reads number	Mapping rate	MarkdQ30	reads number	TSS Enrichment	Peak number	Sample name	Number of sequences labeled as OCR in training set	Ratio of sequences labeled as OCR in training set	Number of sequences labeled as OCR in test set	Ratio of sequences labeled as OCR in test set	AUC	AUPRC
O.	ZS97	p2412166049875	0.3067	1104782	B.954281875	ATAC seq, young panicle, rep2	120640	0.02356283	0.0291	96.974384056329				
O.	sativa	MH63	p331095829863	0.3309	14732476.590184837	ATAC seq, young panicle, rep2	167176	0.0326543525	0.0292	96.87686	0.6312			
O.	sativa	MH63	young3029824	0.4114	28029716.212452855	ATAC seq, young leaf, rep2	169796	0.033164281250.0299	70.9662379.5959					
O.	sativa	MH63	p8812559843	0.4198	24409468.836046717	ATAC seq, young panicle, rep1	204168	0.0398755625	0.0340	78.06379016246				
O.	sativa	MH63	p11621459863	0.3469	22157515.72402842	ATAC seq, young panicle, rep3	201314	0.039319130626	0.034684.806830.4627					
O.	sativa	MH63	p62730759874	0.2876	25483264.945473604	ATAC seq, young panicle, rep5	236853	0.0462668515605391	20.84883096203					

continues on next page

Table 1 – continued from previous page

Species	Samples	Raw reads number	Mapping rate	MarkdQ30	reads number	TSS Enrichment	Peak number	Sample name	Number of sequences labeled as OCR in training set	Ratio of sequences labeled as OCR in training set	Number of sequences labeled as OCR in test set	Ratio of sequences labeled as OCR in test set	AUC	AUPRC
O.	MH63_sativa	511583969752	0.3268	26097080	752674025	ATAC seq, young panicle, rep4	234237	0.045749440	60503933	0.89655076266				
O.	MH63_sativa	1389881016864	0.4966	44734544	0.0129d5991	ATAC seq, young leaf, rep1	229037	0.0447359890	605040070	0.9450689655				
O.	MH63_sativa	68284992987	0.4127	32567024	4.212471764	ATAC seq, root, rep2	253463	0.0495044921	80504263	0.9563103408				
O.	ZS97_sativa	08035949843	0.3616	11787013	222785795	ATAC seq, young leaf, rep2	356413	0.0696112970	60504330	0.9605203555				
O.	ZS97_sativa	p36435609844	0.3379	15282506	704685220	ATAC seq, young panicle, rep3	1825570	0.0356589640	605044570	0.96450056259				
O.	MH63_sativa	112634202987	0.273	68545456	1039285130	ATAC seq, flag leaf, rep1	2926020	0.057148828126	0.050100	0.9586206896				
O.	ZS97_sativa	rd669914592781	0.5518	54000696	5.533989269	ATAC seq, root, rep1	4077480	0.0796382871250	0.052310	0.940598206				
O.	ZS97_sativa	r746423946985	0.4001	34476213	763274959	ATAC seq, root, rep2	4269070	0.083380204430556890	0.056440052768					

continues on next page

Table 1 – continued from previous page

Species	Samples	Raw reads number	Mapping rate	MarkdQ30	reads number	TSS Enrichment	Peak number	Sample name	Number of sequences labeled as OCR in training set	Ratio of sequences labeled as OCR in training set	Number of sequences labeled as OCR in test set	Ratio of sequences labeled as OCR in test set	AUC	AUPRC
O.	MH63_sativa	17413919867	0.2213	4873230	16.194293847	ATAC-seq, lemma & palea, rep1	330188	0.06448847437	50.057044.91982056885					
O.	MH63_sativa	1850476019856	0.1906	5687589	280929B00375	ATAC-seq, lemma & palea, rep2	336402	0.065708515626	0.057650.84286896285					
O.	MH63_sativa	18237110592769	0.5957	67402069	5.614483825	ATAC-seq, root, rep1	343090	0.067008755626	0.058930.49556896651					
O.	ZS97_sativa	flag63620198661	0.6583	39551245	215747310	ATAC-seq, flag leaf, rep2	270937	0.0529173828105067808.495588975818						
O.	ZS97_sativa	flag625Rep1019825	0.305	46956084	773185516	ATAC-seq, flag leaf, rep1	324916	0.0634601956250.079720.564100519417						
O.	ZS97_sativa	flb0166019832	0.2821	61616639	2.8360810778	ATAC-seq, stamen & pistil, rep1	449357	0.0877630690605108420.477836016021						

continues on next page

Table 1 – continued from previous page

Species	Samples	Raw reads number	Mapping rate	MarkdQ30	reads number	TSS Enrichment	Peak number	Sample name	Number of sequences labeled as OCR in training set	Ratio of sequences labeled as OCR in training set	Number of sequences labeled as OCR in test set	Ratio of sequences labeled as OCR in test set	AUC	AUPRC
O.	ZS97_sativa	1986068069858	0.2658	61262472.177650387	12 ATAC seq, stamen & pistil, rep2	46294	1.0904	1818306051093	24.98827076344					
O.	ZS97_sativa	1884472049865	0.1542	26758090.020383452	ATAC seq, lemma & palea, rep2	494763	0.09663234843051171	80.9563	0.6211					
O.	ZS97_sativa	1870556609861	0.1391	26193428.661885427	ATAC seq, lemma & palea, rep1	6016670.11751288893051421	20.465940163366							
S.	Sorghum bi-color	18288204987005751	0.05751	530417840126120294	ATAC seq, bottom part of panicle, rep2	14006	0.0021888675	0.002328.975960.53846						
S.	Sorghum bi-color	179710048507	0.2408	162396811746750825	ATAC seq, upper part of panicle, rep1	14920	0.00233203	0.002529.946470.35897						

continues on next page

Table 1 – continued from previous page

Species	Samples	Raw reads number	Mapping rate	MarkdU30 reads number	TSS Enrichment	Peak number	Sample name	Number of sequences labeled as OCR in training set	Ratio of sequences labeled as OCR in training set	Number of sequences labeled as OCR in test set	Ratio of sequences labeled as OCR in test set	AUC	AUPRC
S.	Sorghum bicolor	636398788728	0.4883	10827038.898647897	ATAC seq, root, rep1	37103	0.005797643750.005796.994839754759						
S.	Sorghum bicolor	16914622961	0.1574	1144435806683616149	ATAC seq, lemma & palea, rep2	44558	0.0069628875	0.007360.978846169824					
S.	Sorghum bicolor	8595820.987520.1869	52951313.034692824	ATAC seq, young leaf, rep2	47084	0.0073568375	0.008175.985629453818						
S.	Sorghum bicolor	128870396120.5199	28765923.350760239	ATAC seq, young leaf, rep1	46878	0.0073245875	0.008188.984960.531646						
S.	Sorghum bicolor	8459610980103542	384766746925880589	ATAC seq, bottom part of panicle, rep1	51301	0.008016781250.008410.969538456453							
S.	Sorghum bicolor	3962820290	0.3257	193204888587697830	ATAC seq, upper part of panicle, rep2	54562	0.0085269025	0.008638.8227510.58051					

continues on next page

Table 1 – continued from previous page

Species	Samples	Raw reads number	Mapping rate	MarkdU30 rate	U30 reads number	TSS Enrichment	Peak number	Sample name	Number of sequences labeled as OCR in training set	Ratio of sequences labeled as OCR in training set	Number of sequences labeled as OCR in test set	Ratio of sequences labeled as OCR in test set	AUC	AUPRC
S.	Sorghum bicolor	2011682918957	0.1605	1335334567	7843324	ATAC-seq, stamen & pistil, rep1	64437	0.010068281250.01021637215084423						
S.	Sorghum bicolor	24504467966	0.1495	1601337732019460069	ATAC-seq, stamen & pistil, rep2	65513	0.010236486250.01049098145897743							
S.	Sorghum bicolor	131234958927	0.4356	52334264.240055991	ATAC-seq, flag leaf, rep1	68154	0.0106480625	0.011030.94823016784						
S.	Sorghum bicolor	1445046921	0.412	60112860.82269930	ATAC-seq, flag leaf, rep2	72738	0.0113692825	0.011590.9496794891						
S.	Sorghum bicolor	951659309952	0.1691	654426909907105032	ATAC-seq, lemma & palea, rep1	75782	0.0118409375	0.012394.8745307692						
S.	Sorghum bicolor	3970022962907	0.7644	60516298.393379326	ATAC-seq, root, rep2	97185	0.015185136250.015474.95961638624							

continues on next page

Table 1 – continued from previous page

Species	Samples	Raw reads number	Mapping rate	MarkdU30 rate	Q30 reads number	TSS Enrichment	Peak number	Sample name	Number of sequences labeled as OCR in training set	Ratio of sequences labeled as OCR in training set	Number of sequences labeled as OCR in test set	Ratio of sequences labeled as OCR in test set	AUC	AUPRC
S.	Sitalica_panicle_1 italica	177674200883	0.2021	1262273889160802264	ATAC seq, young panicle, rep1	48566	0.009485546876	0.005240.924898547688						
S.	Sitalica_young_leaf_2c_1 italica	700671089694	0.4212	29659320260909107	ATAC seq, young leaf, rep1	65424	0.012778145	0.007478032470662547						
S.	Sitalica_flower_1 italica	131086530924	0.23	918095211961096573	ATAC seq, flower, rep1	91008	0.01777318	0.010616989779483772						
S.	Sitalica_root_2c_1 italica	257357508962	0.4955	7709618.4828	ATAC seq, root, rep1	1021400.019949398750	0.011318.98867035896							
S.	Sitalica_young_leaf_2c_2 italica	57397792982	0.5489	196386697398941305	ATAC seq, young leaf, rep2	72681	0.011356166250	0.011969.950640062568						
S.	Sitalica_flower_2 italica	13835702992	0.2077	9927057723444490325	ATAC seq, flower, rep2	1210540.023644359376	0.014790.934929658337							
S.	Sitalica_flag_leaf_1 italica	141996009915	0.2786	8592447690569965882	ATAC seq, flag leaf, rep1	1626850.03177445406052053084256016257								

continues on next page

Table 1 – continued from previous page

Species	Samples	Raw reads number	Mapping rate	MarkdU30 rate	Q30 reads number	TSS Enrichment	Peak number	Sample name	Number of sequences labeled as OCR in training set	Ratio of sequences labeled as OCR in training set	Number of sequences labeled as OCR in test set	Ratio of sequences labeled as OCR in test set	AUC	AUPRC
S.	Sitalica_panicle_2 italica	15014206871	0.2297	100733619785594699				ATAC seq, young panicle, rep2	192586	0.0376144331	2602477097068647491			
S.	Sitalica_flag_leaf_2 italica	10070620926	0.2437	636930602937368992				ATAC seq, flag leaf, rep2	224032	0.04375826	0.02857095982006598			
26.	Zmays_mays	96724598955	0.3522	40851048.561478382				ATAC seq, upper part of ~1cm ear, rep1	39169	0.006125186250.00661057869220057692				
26.	Zmays_mays	965148099390.3728		38466973.245580019				ATAC seq, bottom part of ~1cm ear, rep2	53425	0.008346956250.008710.942300659256				
26.	Zmays_mays	926165609760.8325		40570460.986698983				ATAC seq, bottom part of ~1cm ear, rep1	61812	0.009658125	0.01030398637820826			

continues on next page

Table 1 – continued from previous page

Species	Samples	Raw reads number	Mapping rate	MarkdU30 reads number	TSS Enrichment	Peak number	Sample name	Number of sequences labeled as OCR in training set	Ratio of sequences labeled as OCR in training set	Number of sequences labeled as OCR in test set	Ratio of sequences labeled as OCR in test set	AUC	AUPRC
26.	Zmays_mays	401409800939	0.3889	38784430	0.248695514	ATAC seq, upper part of ~1cm ear, rep2	87961	0.0137490625	0.0137849826804543				
26.	Zmays_mays	106130162727	0.3883	57581514	0.049183836	ATAC seq, young leaf, rep3	101406	0.019803859378	0.0173769815862068				
26.	Zmays_mays	15842062432	0.1117	24442269	7941615869	ATAC seq, root, rep2	111460	0.021769531250	0.0208789862586209				
26.	Zmays_mays	1593420750975	0.3186	80469075	952080430	ATAC seq, root, rep3	158819	0.0310145859305030846981370451318					
26.	Zmays_mays	young308528	0.3596	41587718	617390227	ATAC seq, young leaf, rep1	242864	0.047434875	0.04148798666551354				
26.	Zmays_mays	23672605024	0.1413	42431275	0.0827729812	ATAC seq, root, rep1	226605	0.04425828906050417569851024657					
26.	Zmays_mays	103926088484	0.3064	30621515	5.4010744194	ATAC seq, young leaf, rep2	234570	0.045814253128	0.04182098865015087				

continues on next page

Table 1 – continued from previous page

Species	Samples	Raw reads number	Mapping rate	MarkdU30 reads number	TSS Enrichment	Peak number	Sample name	Number of sequences labeled as OCR in training set	Ratio of sequences labeled as OCR in training set	Number of sequences labeled as OCR in test set	Ratio of sequences labeled as OCR in test set	AUC	AUPRC
26.	Zmays_mays	<u>105904100923</u>	0.2676	538499373660	694274		ATAC seq, young tassel spikelet, rep1	250943	0.049016304680	50420	0.88206892625		
26.	Zmays_mays	<u>043752500963</u>	0.3106	424170858691	656382		ATAC seq, bottom part of ~1cm tassel, rep1	243257	0.0475163782810	5042900	0.9805600589		
26.	Zmays_mays	<u>11b2607990975</u>	0.289	551913315284525693			ATAC seq, young tassel spikelet, rep2	264358	0.051630521870	0.044110	0.98884082758		
26.	Zmays_mays	<u>049875149944</u>	0.3956	36094115046647335			ATAC seq, ~5mm ear, rep1	261472	0.05106883		0.045990.98120685262		
26.	Zmays_mays	<u>861783629955</u>	0.3613	353519437079461634			ATAC seq, bottom part of ~1cm tassel, rep2	254837	0.049770891560504640	0.58063048225			

continues on next page

Table 1 – continued from previous page

Species	Samples	Raw reads number	Mapping rate	MarkdU30 reads number	TSS Enrichment	Peak number	Sample name	Number of sequences labeled as OCR in training set	Ratio of sequences labeled as OCR in training set	Number of sequences labeled as OCR in test set	Ratio of sequences labeled as OCR in test set	AUC	AUPRC
26.	Zmays_mays	<u>0601sh409924</u>	0.325	42612007.053287925	ATAC seq, ~5mm ear, rep2	276548	0.0540172281250.048690.928172015579						
26.	Zmays_mays	<u>180483079963</u>	0.413	715798470775	ATAC seq, upper part of ~1cm tassel, rep1	286296	0.0559171875	0.050983.982410.5279					
26.	Zmays_mays	<u>flag784042951</u>	0.4233	398283974888660508	ATAC seq, flag leaf, rep2	386108	0.075411014750.068290.925360.0539						
26.	Zmays_mays	<u>952493049969</u>	0.3709	435707825500591086	ATAC seq, flag leaf, rep1	389894	0.0761510251876.069030.886370344085						

---

**CHAPTER  
TEN**

---

**CITATION**

If you use this web sever for your research, please cite our paper:

Hu Zhao#, Zhuo Tu#, Yinmeng Liu, Zhanxiang Zong, Jiacheng Li, Hao Liu, Feng Xiong, Jinling Zhan, Xuehai Hu, and Weibo Xie\* (2021). PlantDeepSEA, a deep learning-based web service to predict the regulatory effects of genomic variants in plants. *Nucleic Acids Research*, doi: 10.1093/nar/gkab383