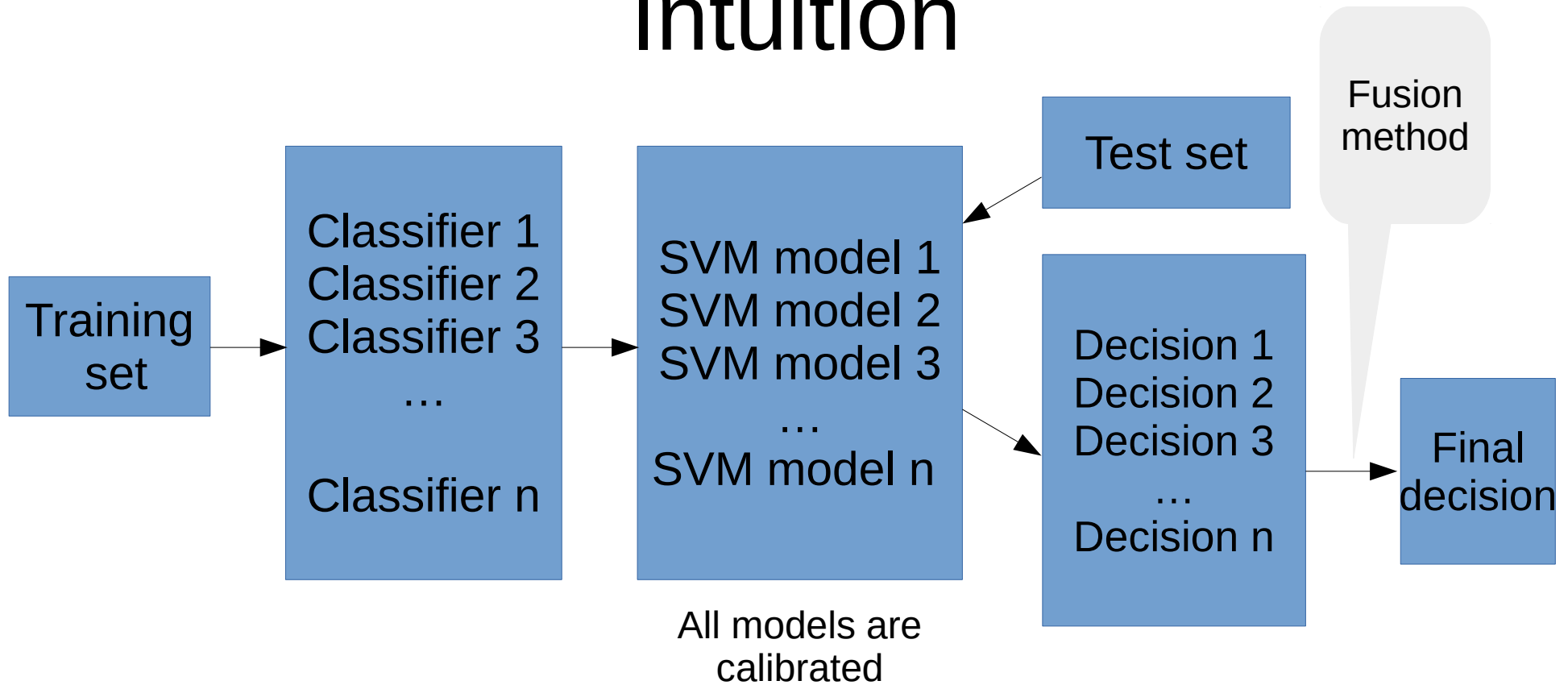


# SVM ensembles approach intro

Nianheng Wu

# Intuition



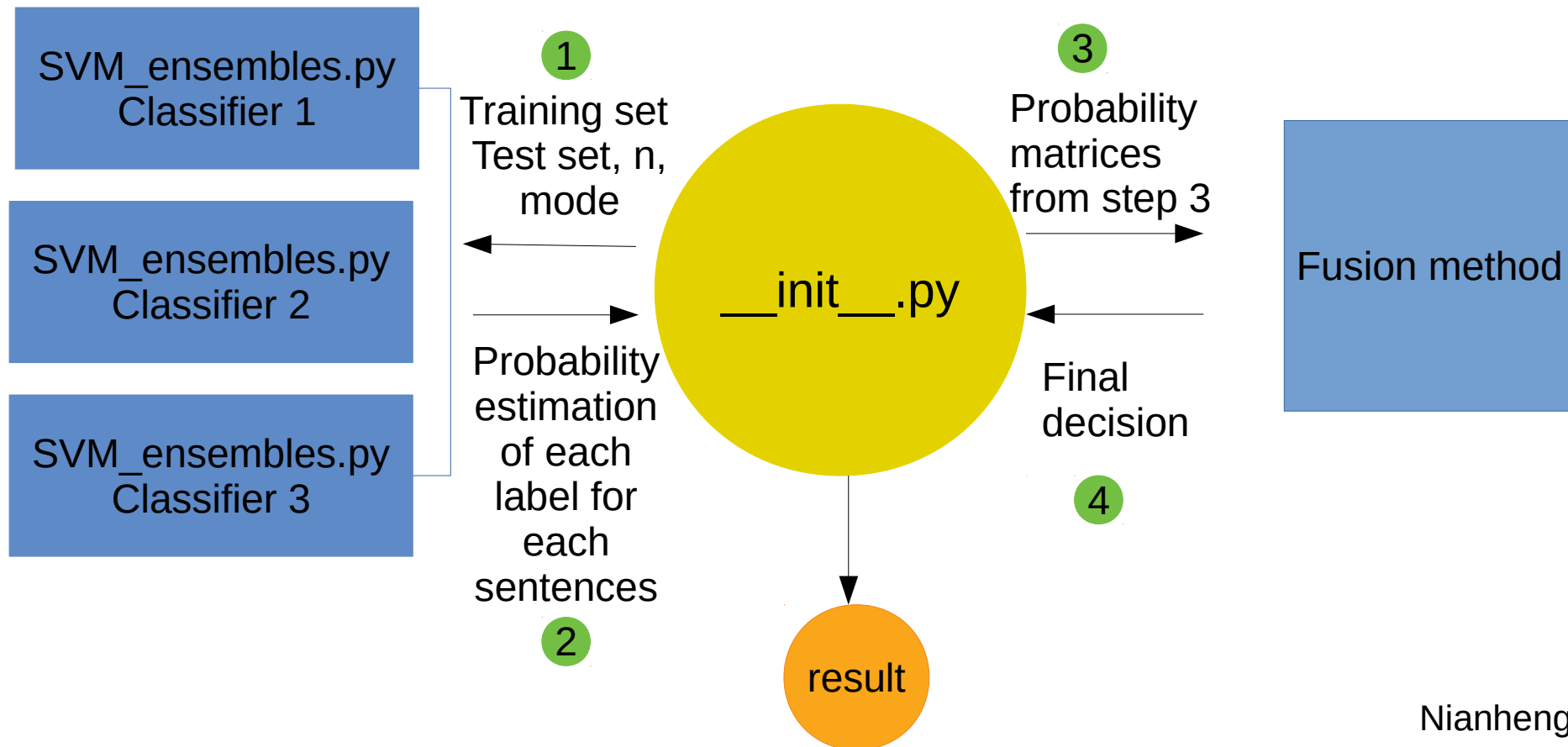
# Our program structure

- `__init__.py`
- `SVM_ensembles.py`
- `Fusion method.py`
- `Other testing functions.py`

# Our program structure

- `__init__.py`
  - An interface, also like a control center
  - Determine how many classifiers we want to generate
  - Get the path to training set and test set, parse the paths to `SVM_ensembles`
  - Get the decision result back from `SVM_ensembles`. Then send the results to fusion method, and get the result from it.

# Our program structure



# Our program structure

- SVM\_ensembles.py
  - Contains a class, and two major functions: *training* and *testing*
  - Function *training* applies matrix and tf-idf and generate a matrix of features (n value and mode depend on parameters that parsed into it)
  - SVM model = LinearSVC( the matrix of features)
  - Because SVM model do not work based on probabilities, but in order to apply Mean Probability Rule as fusion method, we have to use a trick called *calibration*
  - So: `clf = CalibrationClassifierCV(SVM model)`
  - Now we have a classifier, we put the feature matrix of test set into the classifier and get a matrix of probabilities.

# Our program structure

- SVM\_ensembles.py

Probability of being 'T'

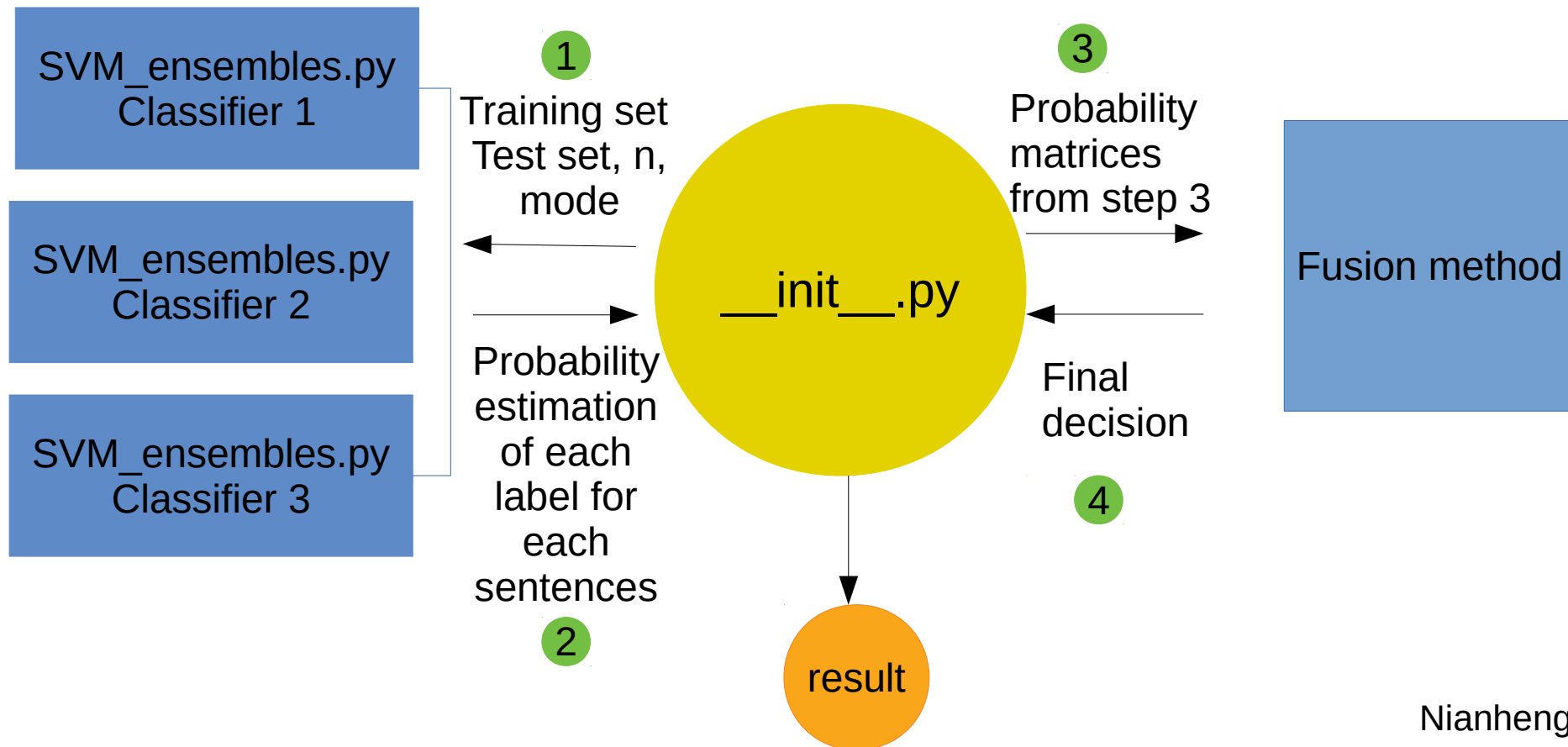
Probability of being 'M'

The length of the matrix is the size of testing set

And the amount of this matrix is equal to the amount of classifier we generated at the first place

Pt 1	Pm 1
Pt 2	Pm 2
Pt 3	Pm 3
...	...
Pt n	Pm n

# Our program structure





# Our program structure

- `fusion_method.py`
  - We have already gotten  $x$  of probability matrices, the `__init__.py` will parse them to the fusion method. Here I'll just explain how Mean Probability Rule works.

# Our program structure

- fusion\_method.py

$$\begin{bmatrix} \text{Pt1} & \text{Pm1} \\ \text{Pt2} & \text{Pm2} \\ \text{Pt3} & \text{Pm3} \\ \vdots & \vdots \\ \text{Ptn} & \text{Pmn} \end{bmatrix} + \begin{bmatrix} \text{Pt1}' & \text{Pm1}' \\ \text{Pt2}' & \text{Pm2}' \\ \text{Pt3}' & \text{Pm3}' \\ \vdots & \vdots \\ \text{Ptn}' & \text{Pmn}' \end{bmatrix} + \begin{bmatrix} \text{Pt1}'' & \text{Pm1}'' \\ \text{Pt2}'' & \text{Pm2}'' \\ \text{Pt3}'' & \text{Pm3}'' \\ \vdots & \vdots \\ \text{Ptn}'' & \text{Pmn}'' \end{bmatrix}$$

$$= \begin{bmatrix} \text{Pt1} + \text{Pt1}' + \text{Pt1}'' & \text{Pm1} + \text{Pm1}' + \text{Pm1}'' \\ \text{Pt2} + \text{Pt2}' + \text{Pt2}'' & \text{Pm2} + \text{Pm2}' + \text{Pm2}'' \\ \text{Pt3} + \text{Pt3}' + \text{Pt3}'' & \text{Pm3} + \text{Pm3}' + \text{Pm3}'' \\ \vdots & \vdots \\ \text{Ptn} + \text{Ptn}' + \text{Ptn}'' & \text{Pmn} + \text{Pmn}' + \text{Pmn}'' \end{bmatrix}$$

# Our program structure

$$\left[ \begin{array}{cc} Pt1 + Pt1' + Pt1'' & Pm1 + Pm1' + Pm1'' \\ Pt2 + Pt2' + Pt2'' & Pm2 + Pm2' + Pm2'' \\ Pt3 + Pt3' + Pt3'' & Pm3 + Pm3' + Pm3'' \\ \vdots & \vdots \\ Ptn + Ptn' + Ptn'' & Pmn + Pmn' + Pmn'' \end{array} \right] / 3$$

$$= \left[ \begin{array}{cc} Pt1 \text{ avg} & Pm1 \text{ avg} \\ Pt2 \text{ avg} & Pm2 \text{ avg} \\ Pt3 \text{ avg} & Pm3 \text{ avg} \\ \vdots & \vdots \\ Ptn \text{ avg} & Pmn \text{ avg} \end{array} \right]$$

And we compare each row in this matrix, if  $Ptn \text{ avg} > Pmn \text{ avg}$  then this sentence is Taiwan Guoyu. And vice versa.

# Our program structure

$$\begin{bmatrix} T \\ M \\ T \\ \dots \\ T \end{bmatrix}$$

And we get this ←

Then we compare this to the label list  
we got from test set, and calculate the  
f1 score