

Beyond Speech Processing

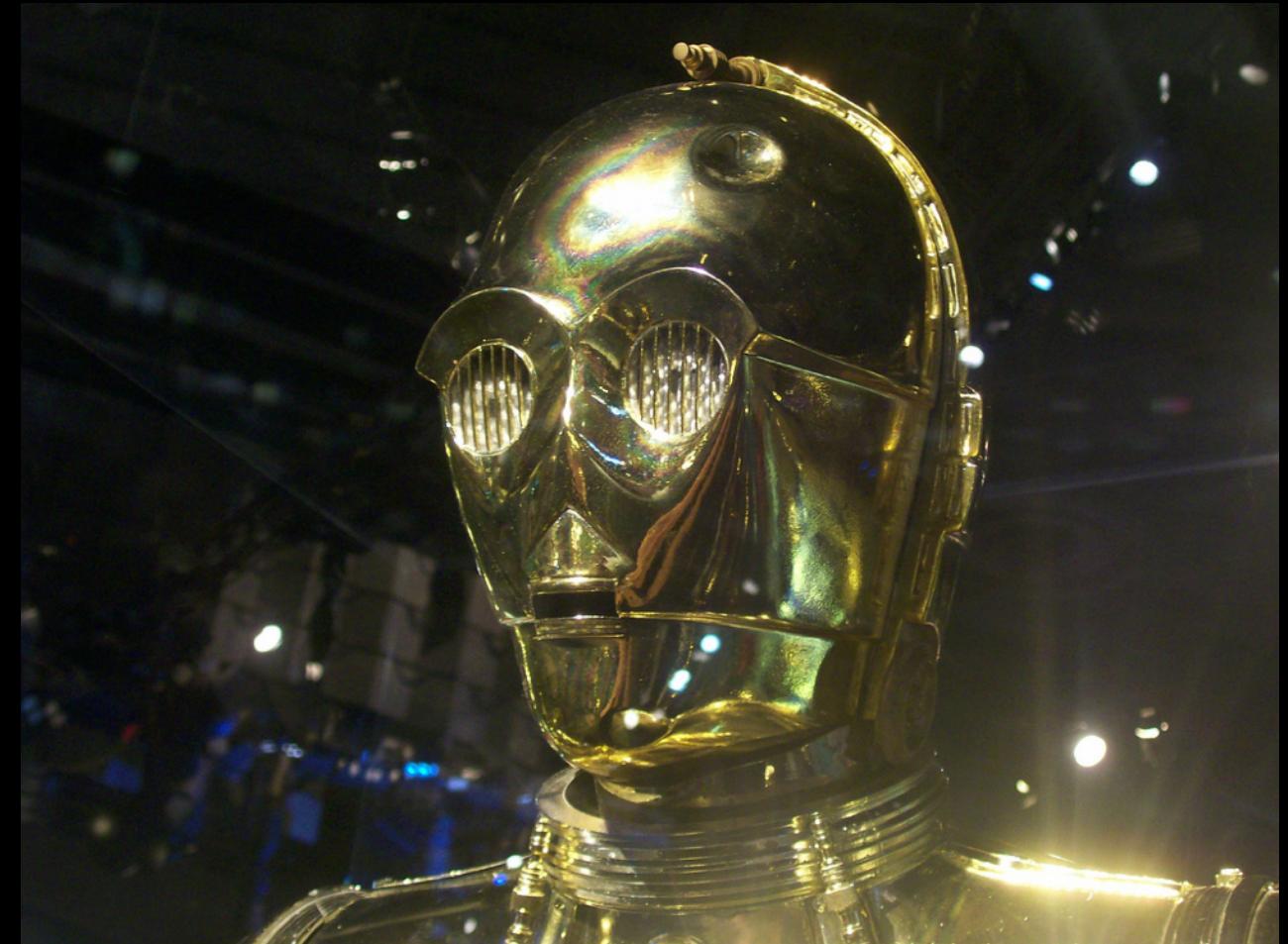
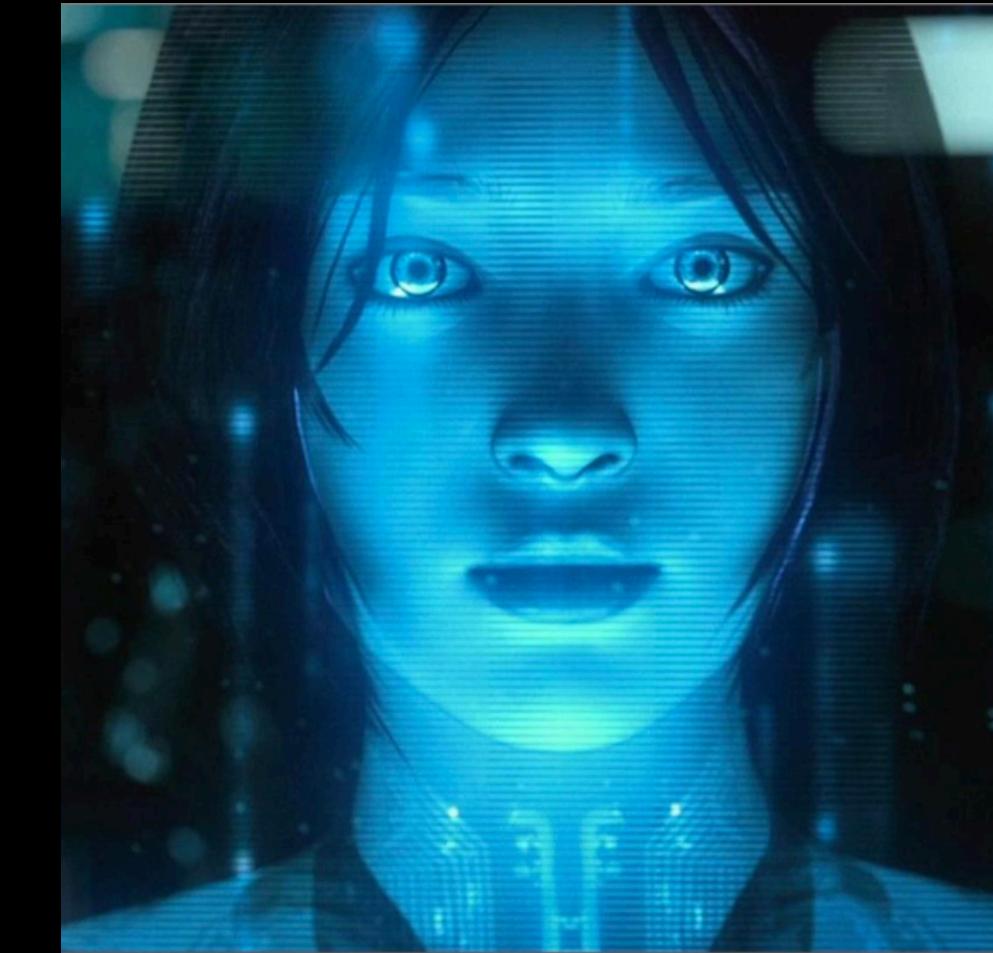
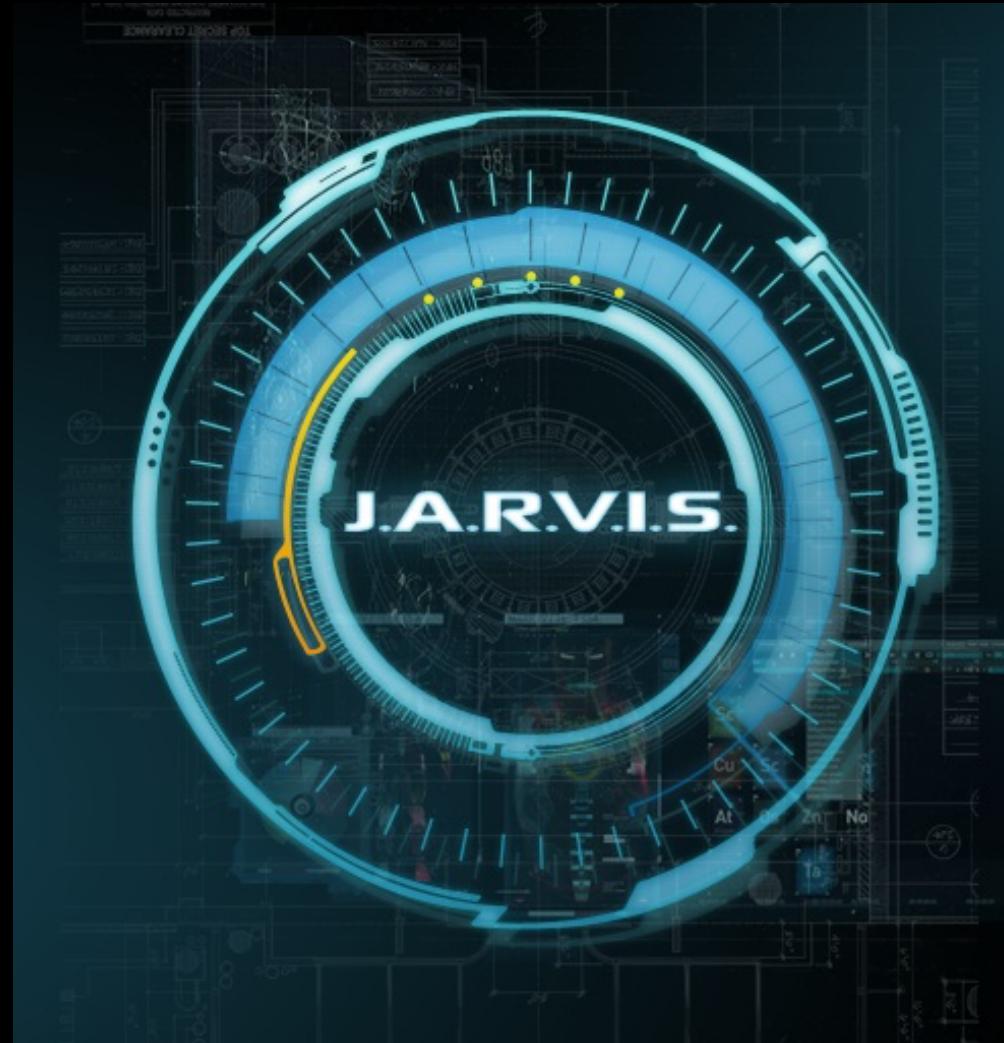
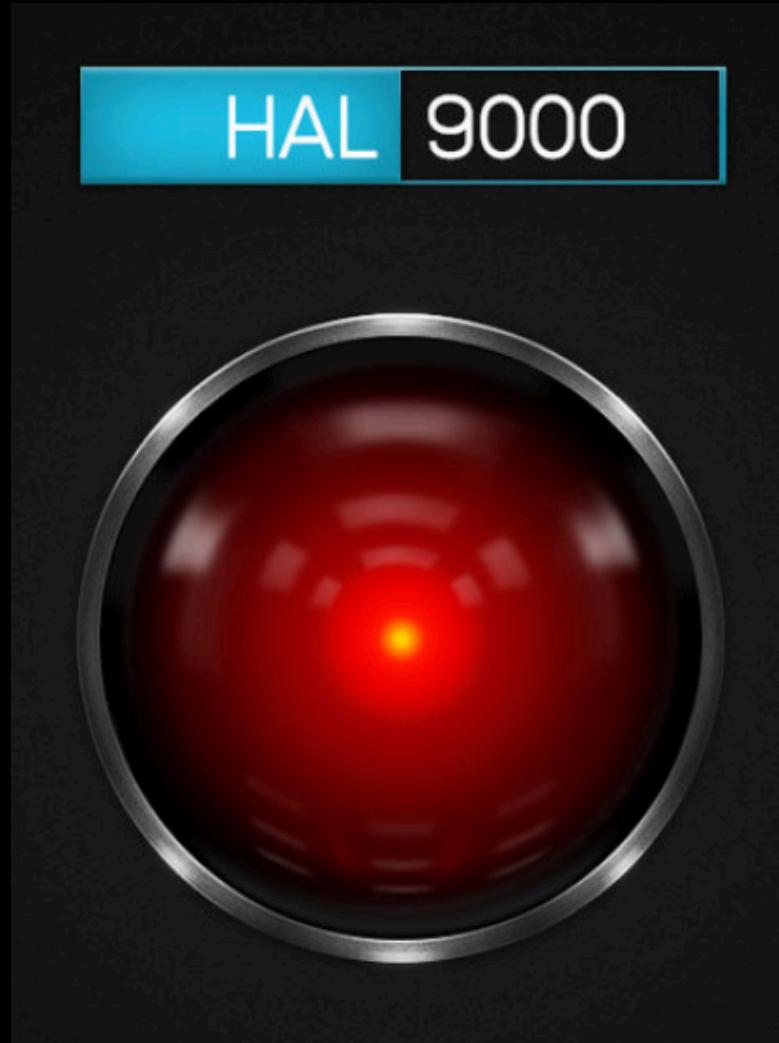
Toolkits for Accelerated Speech Processing and
Machine Learning

Trung Ngo Trong
trung@imito.ai



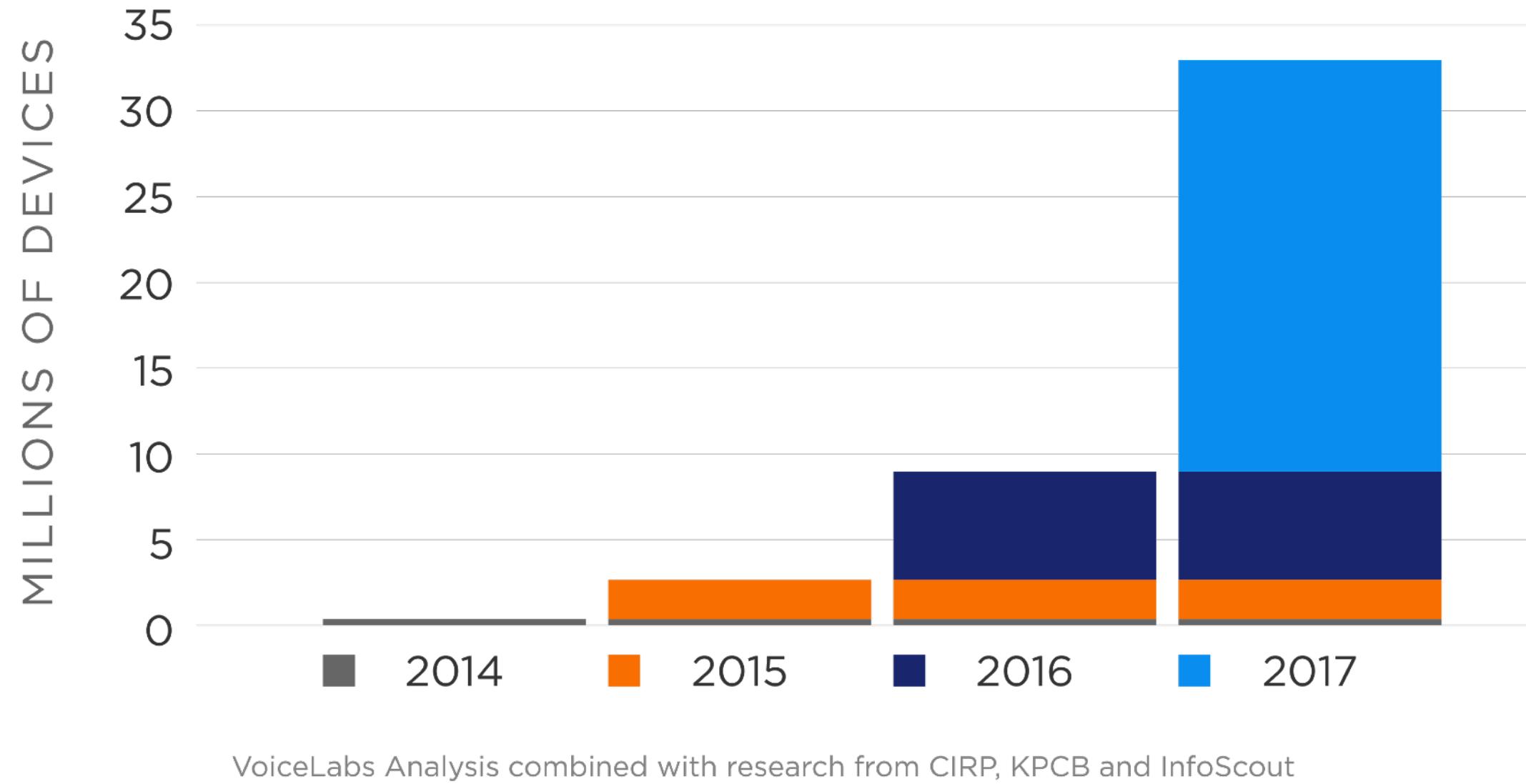
UNIVERSITY OF
EASTERN FINLAND

Computer Speech is the Dream

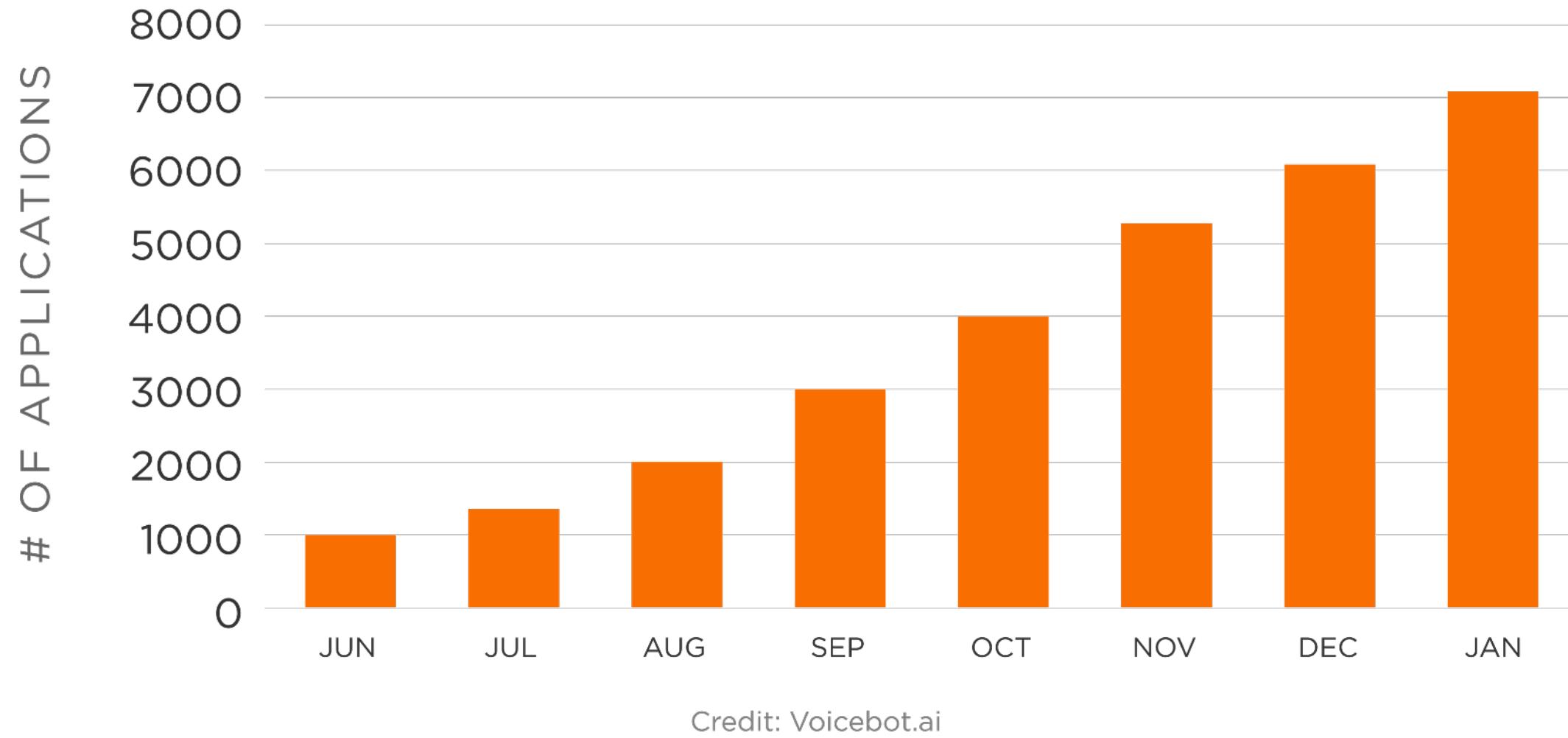


Industrial Response

Voice-First Device Footprint



Alexa Skills June 2016 to January 2017



THE PRIMARY REASON ENTERPRISES CURRENTLY USE AI IS FOR:



48.5%

Automated communications that give business audiences data they can use to make effective business decisions



13.6%

Automated communications that give consumer audiences data they can use to make effective decisions



6.1%

Automation that eliminates manual and repetitive tasks

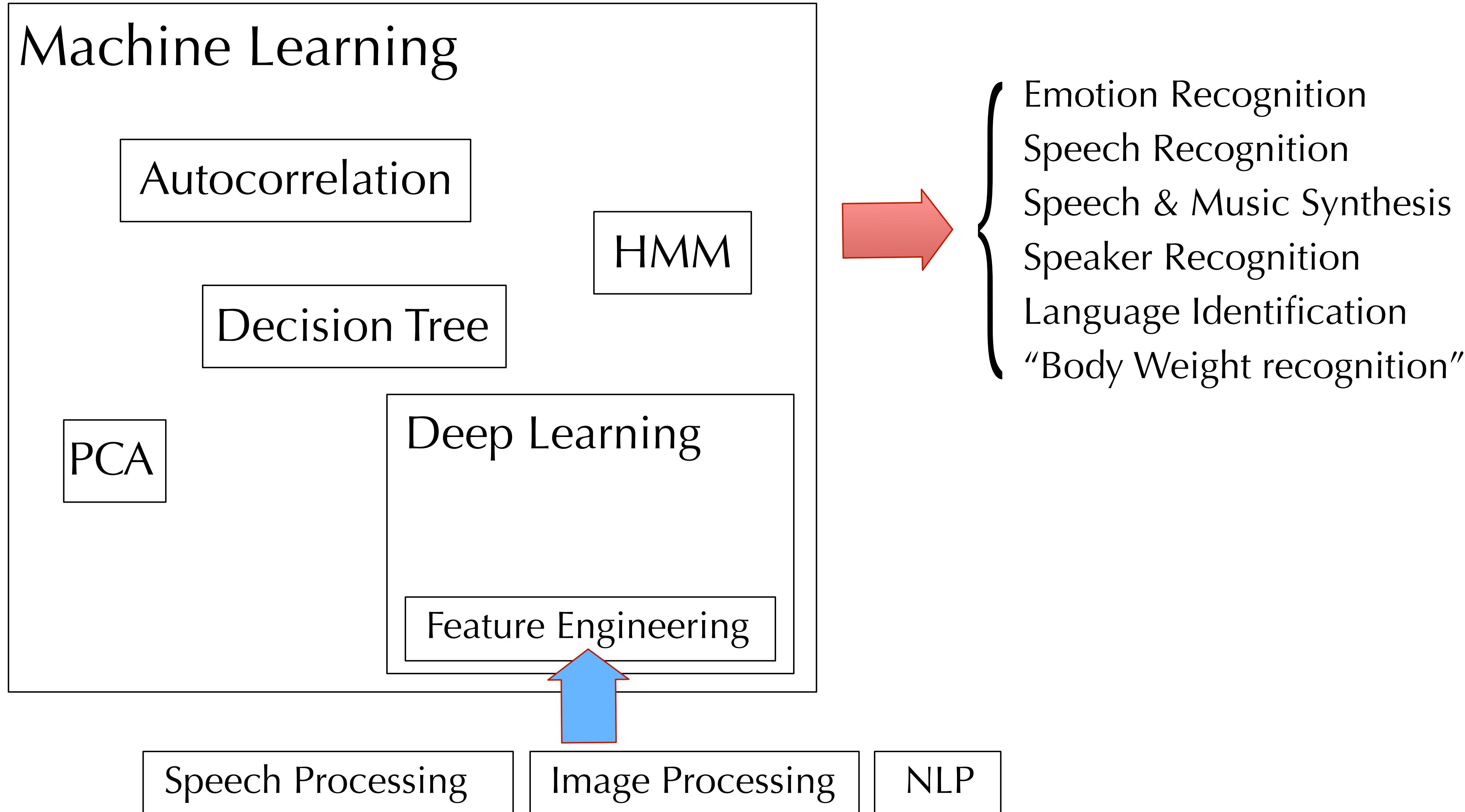


4.6%

Monitoring and alerts about the health of the business. Automated data-driven reporting

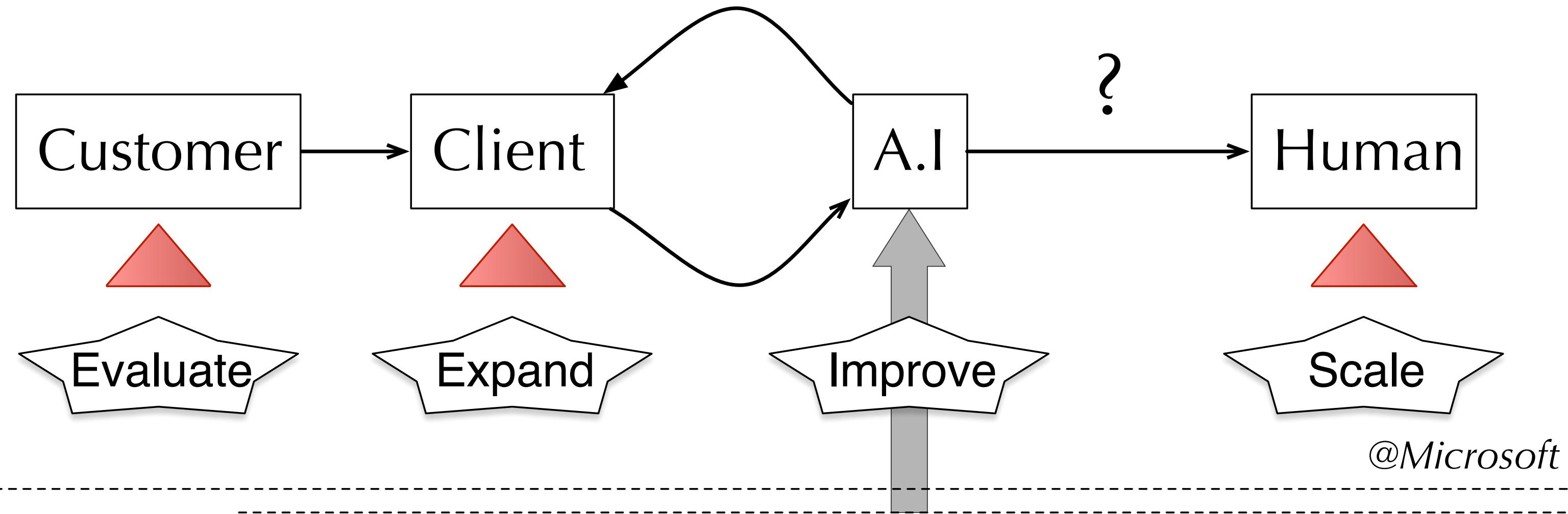
Statistics provided by: @campaignuk, @voicelabs, @microsoft, @reuters

Why Machine Learning !

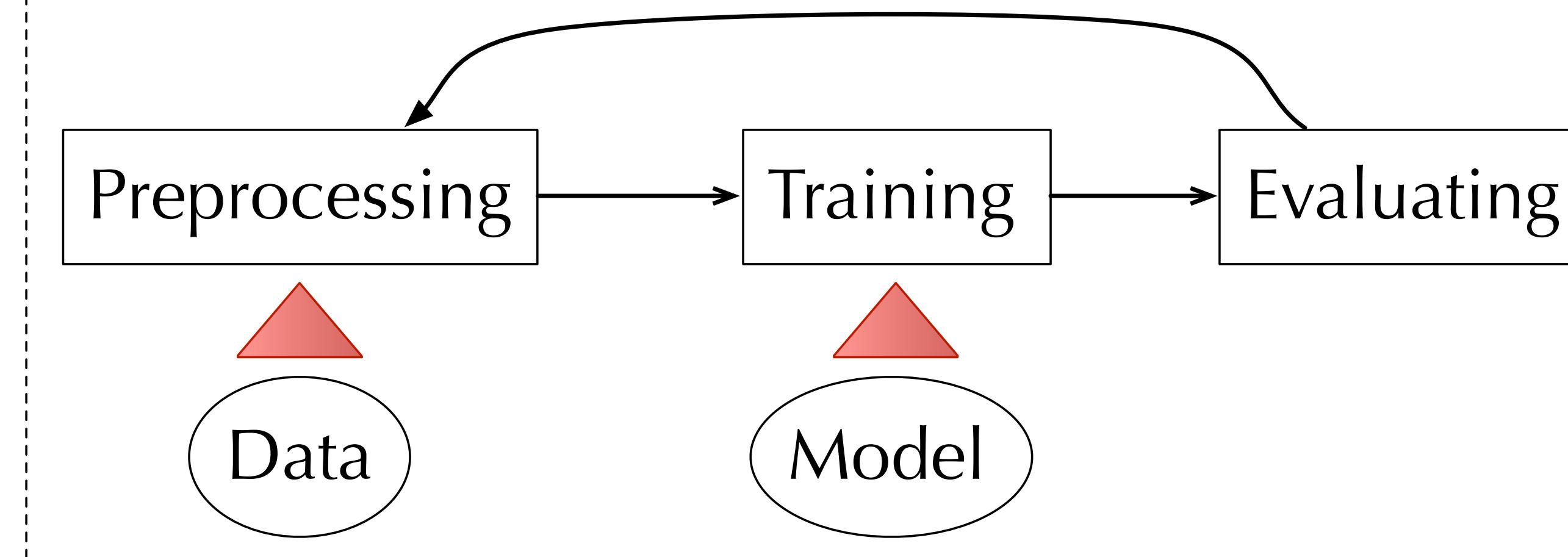


The Pipeline

Industrial deployment pipeline



Research pipeline



- Understand data collection is critical
- It is almost never work for the first time
- Speed iteration is more important than quality iteration

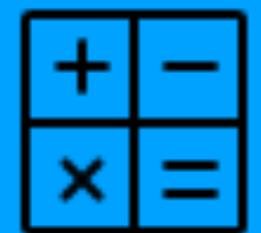
Best Languages ?

	Matlab	Python	R	C/C++
Setup and environment	✓✓	✓✓	✓✓✓	✓
Speed	✓✓	✓✓	✓✓	✓✓✓
Versatile	✓	✓✓✓	✓✓	✓✓
Syntax friendly	✓✓✓	✓✓✓	✓✓✓	✓
Development time	✓✓✓	✓✓	✓✓✓	✓
Libraries	✓✓	✓✓✓	✓✓	✓✓✓
Industrial application	✓	✓✓✓	✓✓✓	✓✓✓
Remark	Research standard, most of speech processing libraries initially built on Matlab	Industrial standard, Amazon, Google, Facebook, Microsoft ...	Extremely easy and efficient for statistical inference	Fast, be able to access full features of ML libraries

Speech and Data Processing



Step 1



Step 2



Step 3



Step 4

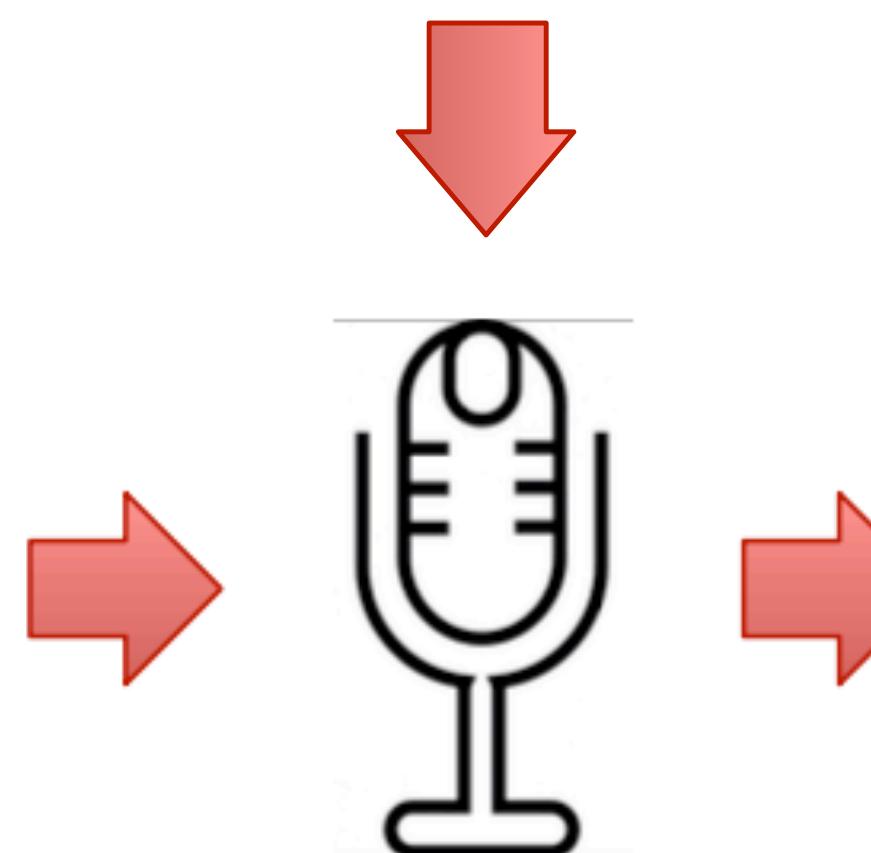


Step 5

Human



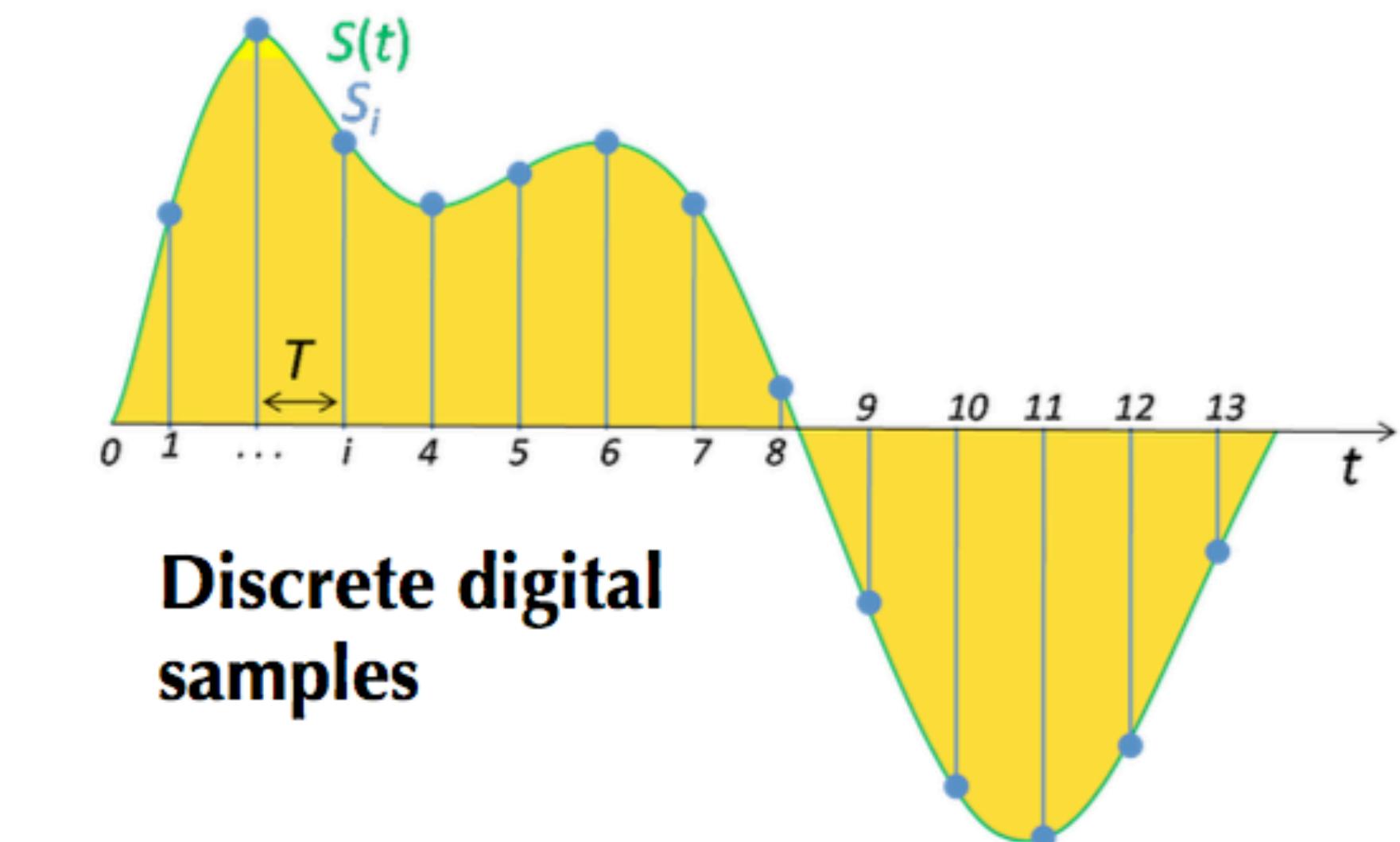
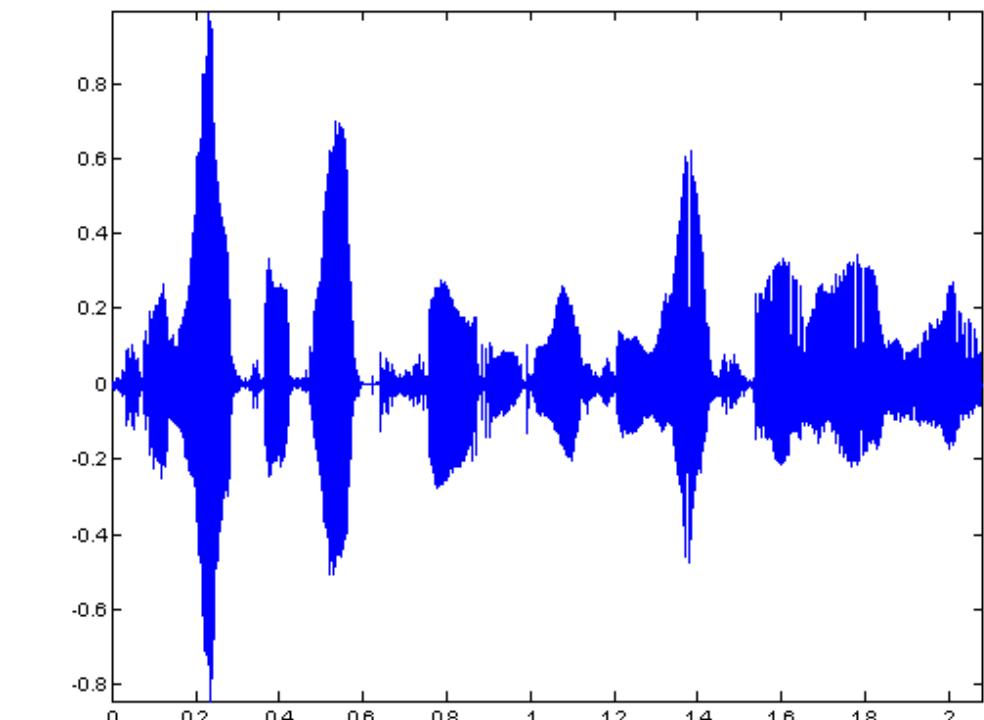
Sampling
rate Bit depth



Microphone

Continuous
sound wave

Machine

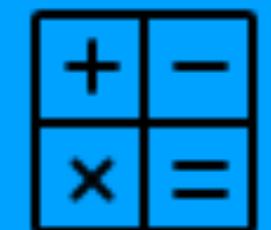


Discrete digital
samples

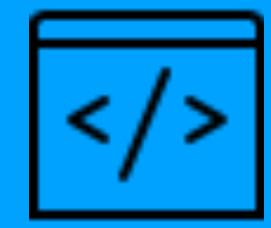
Speech and Data Processing



Step 1



Step 2



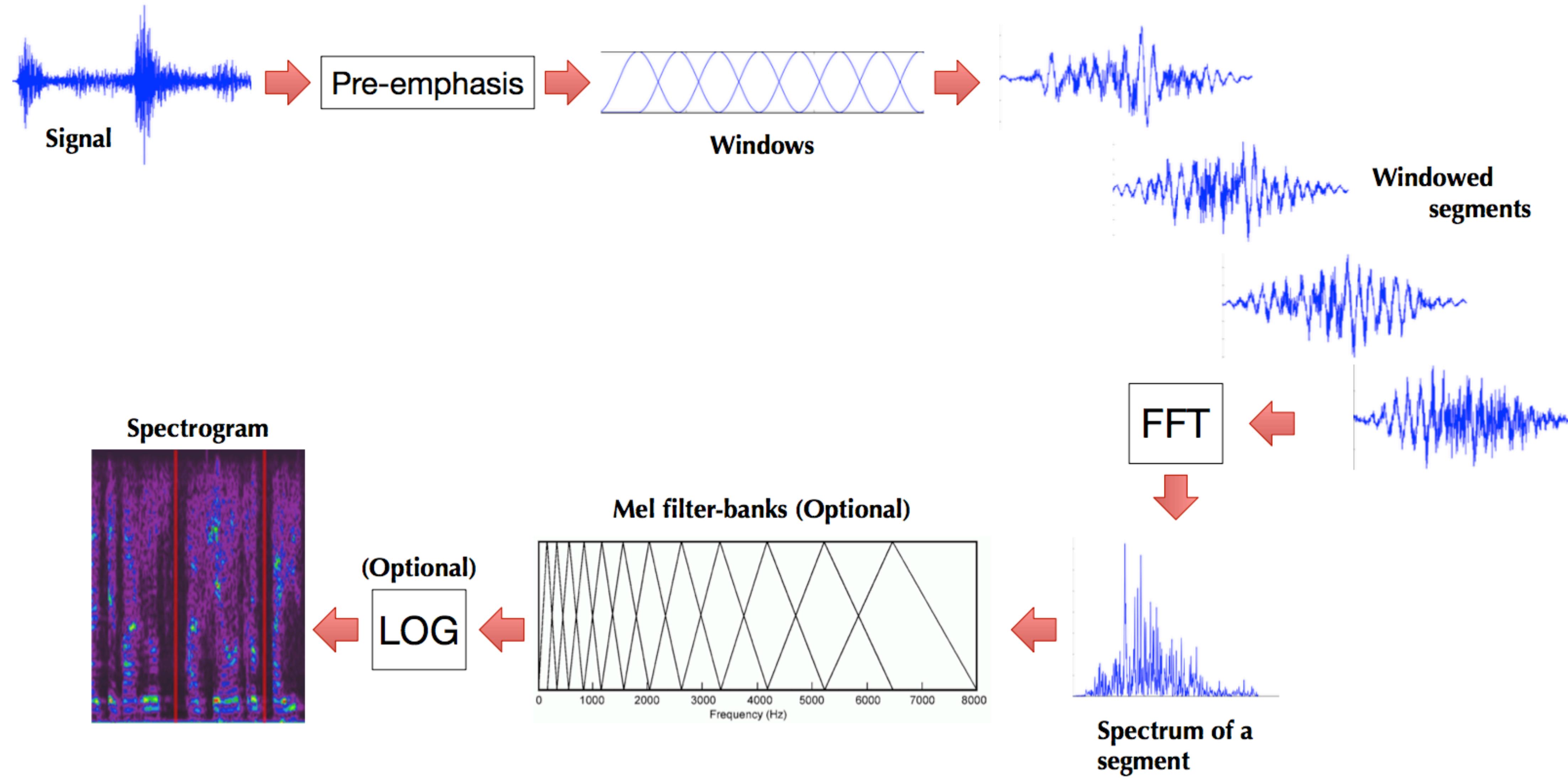
Step 3



Step 4



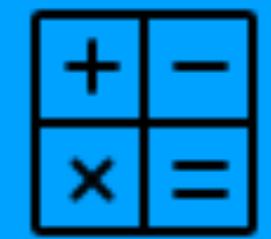
Step 5



Speech and Data Processing



Step 1



Step 2



Step 3



Step 4



Step 5

```
import librosa
import numpy as np
y, sr = librosa.load('path/to/audio')
fft = librosa.stft(y, n_fft=256)
spectrogram = np.abs(fft) ** 2
mel_spec = librosa.feature.melspectrogram(S=spectrogram, sr=sr,
                                            n_mels=40, fmax=8000)
S = librosa.feature.mfcc(S=librosa.power_to_db(mel_spec), n_mfcc=13)
```

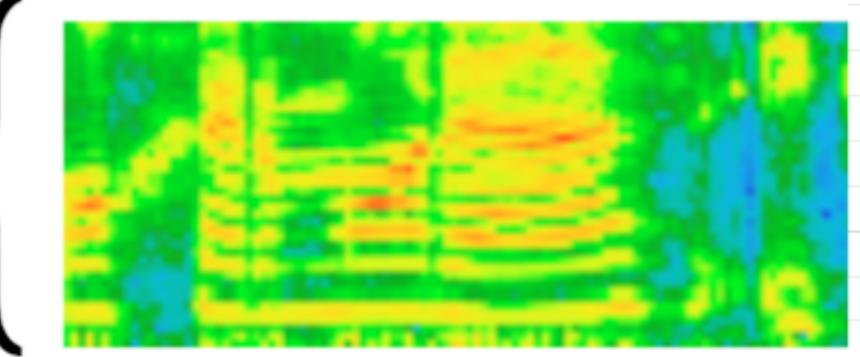
FFT spectrogram

129



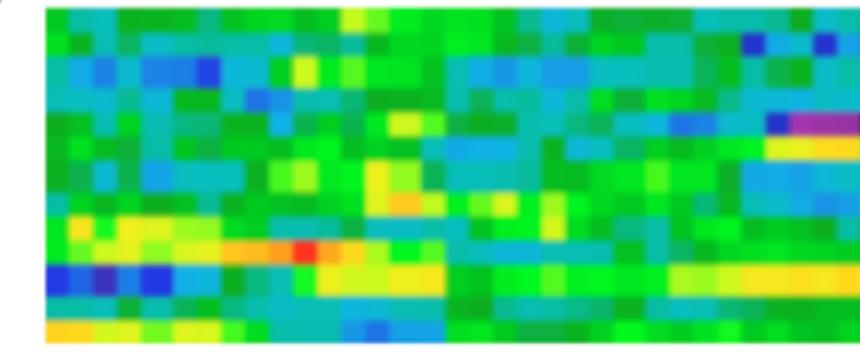
Mel-filterbank

40



MFCC

13

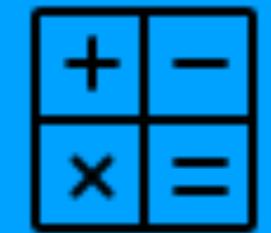


Libraries	sidekit	librosa	pysptk
Highlights	Speaker and Language recognition	Music and audio analysis	Speech analysis and synthesis (wrapper of C/C++ framework SPTK)

Speech and Data Processing



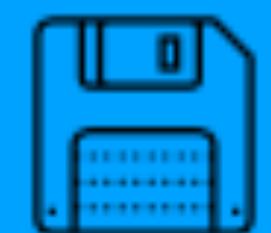
Step 1



Step 2



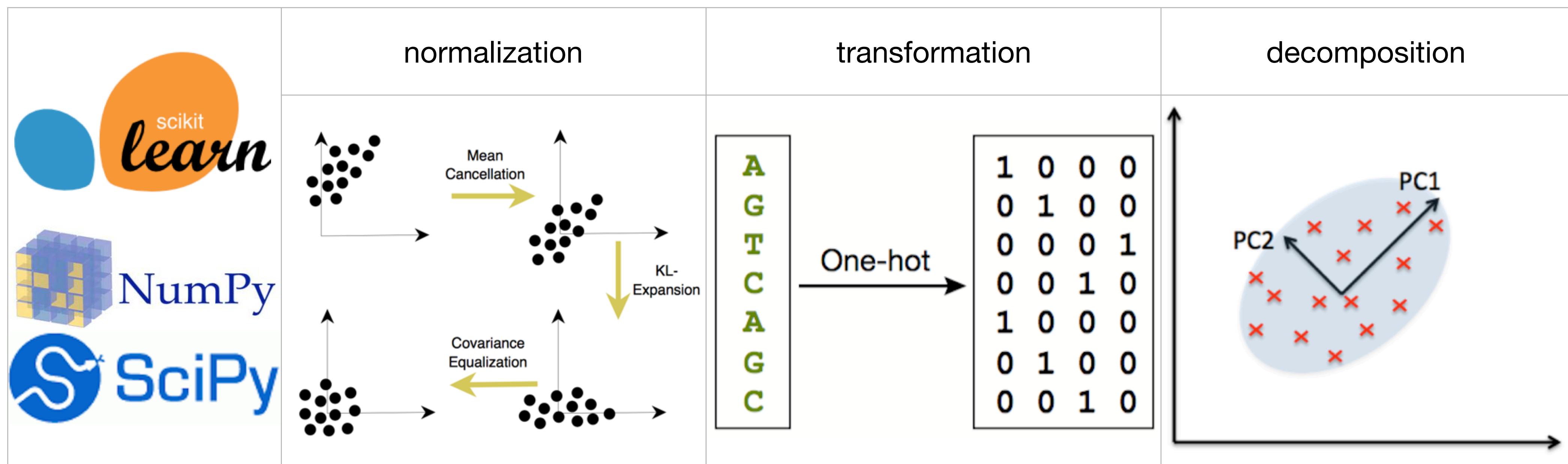
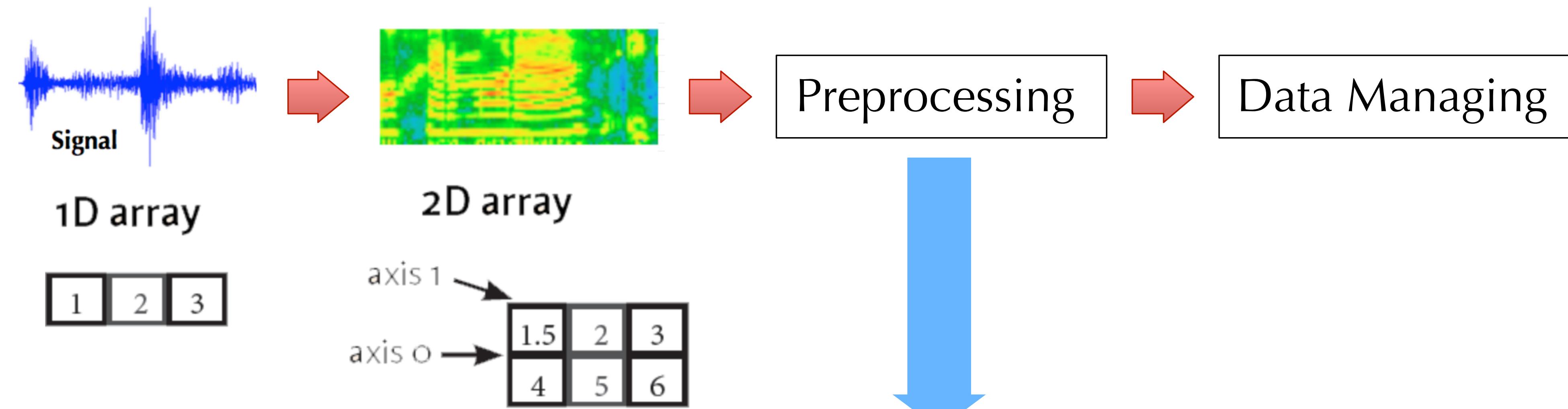
Step 3



Step 4



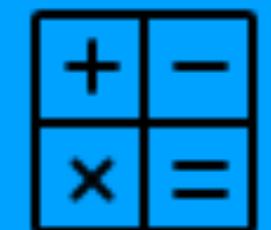
Step 5



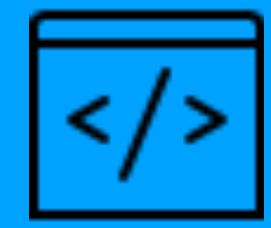
Speech and Data Processing



Step 1



Step 2



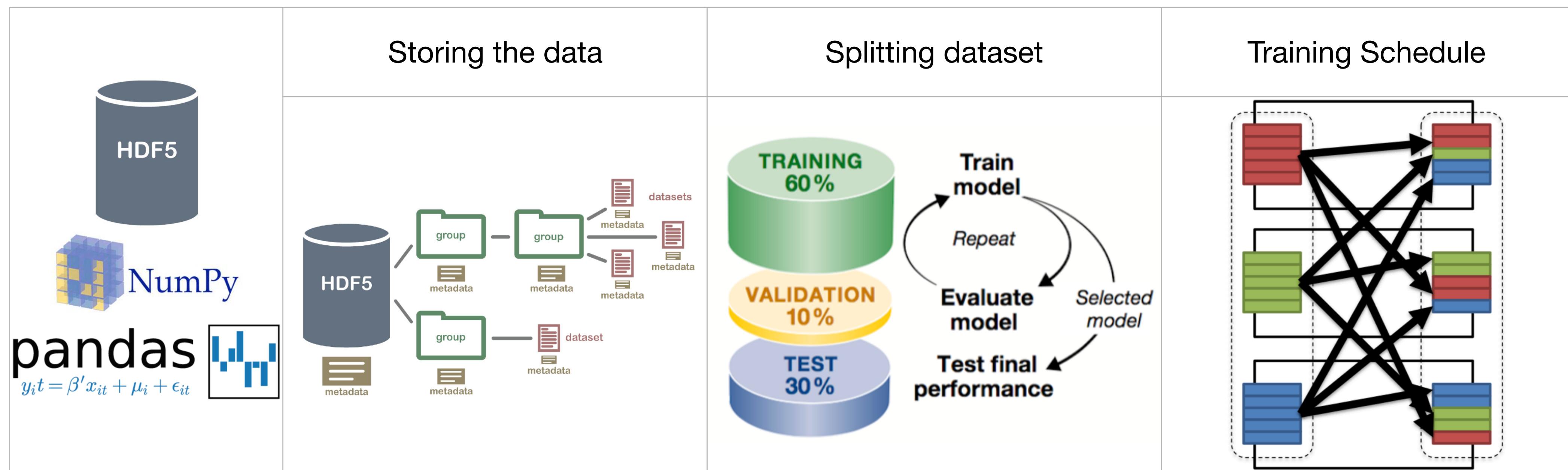
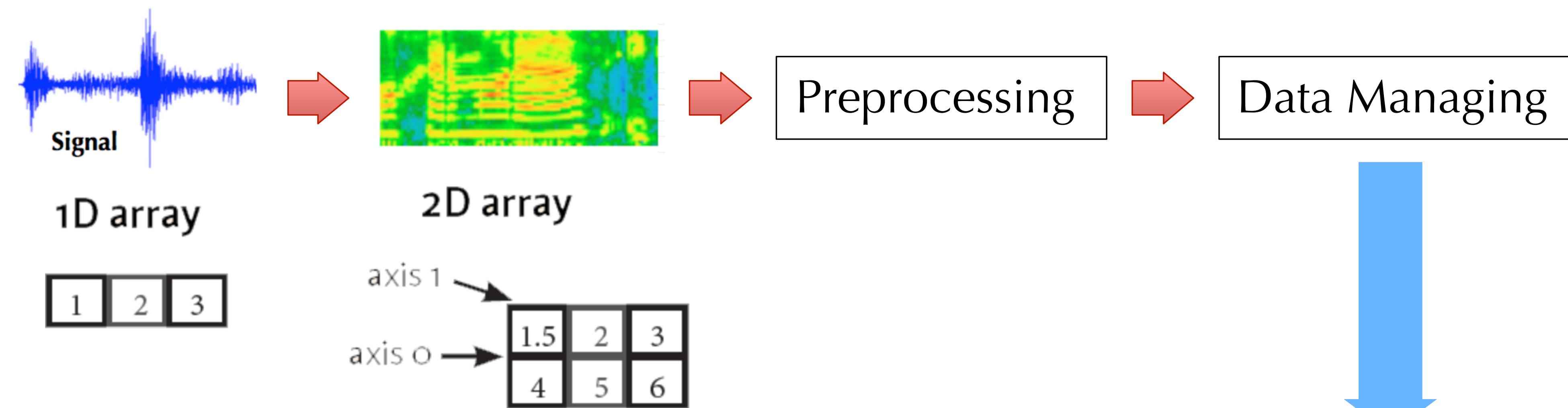
Step 3



Step 4



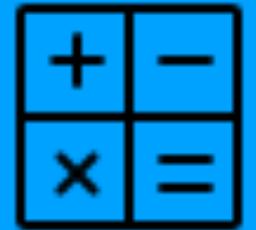
Step 5



Speech and Data Processing



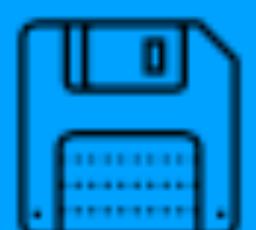
Step 1



Step 2



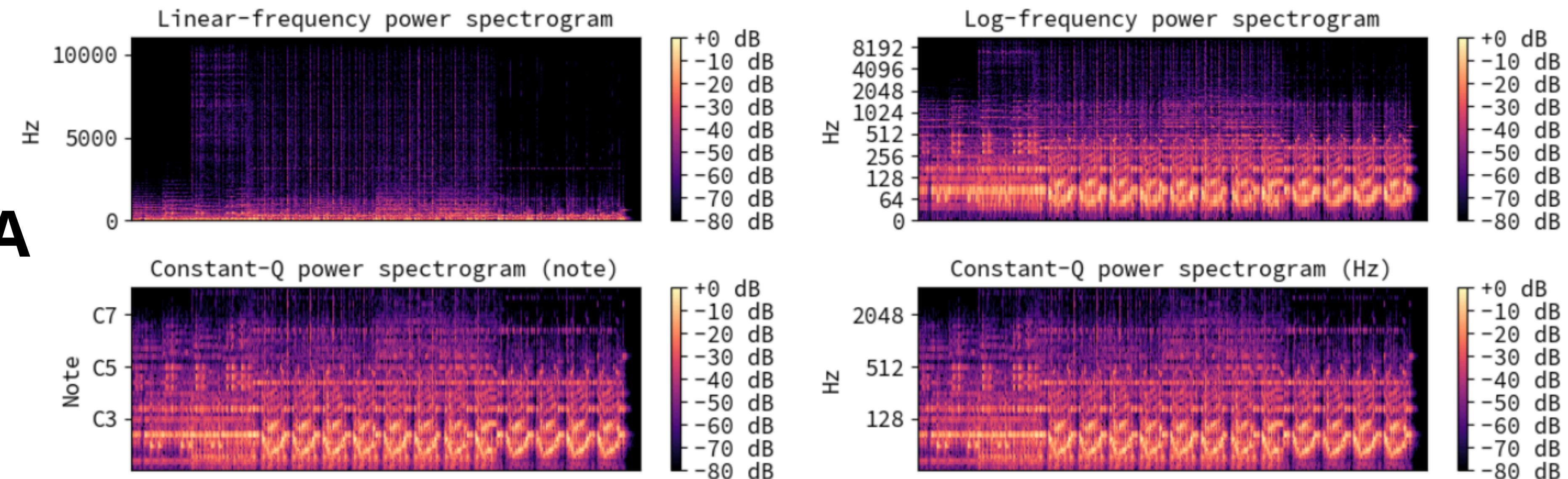
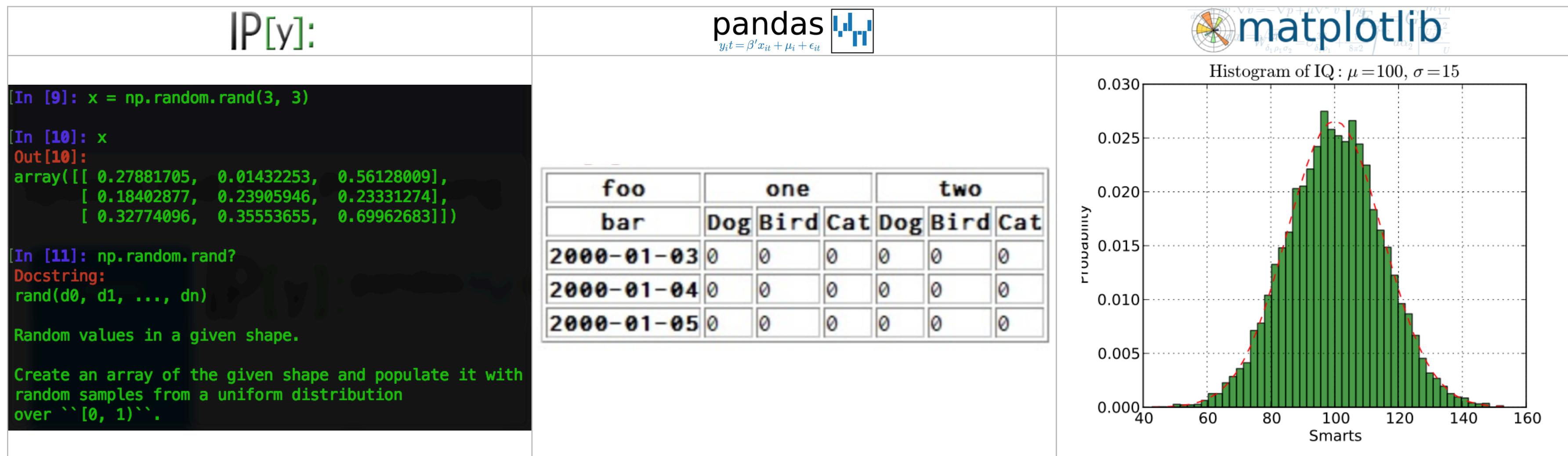
Step 3



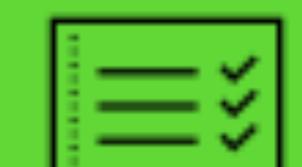
Step 4



Step 5

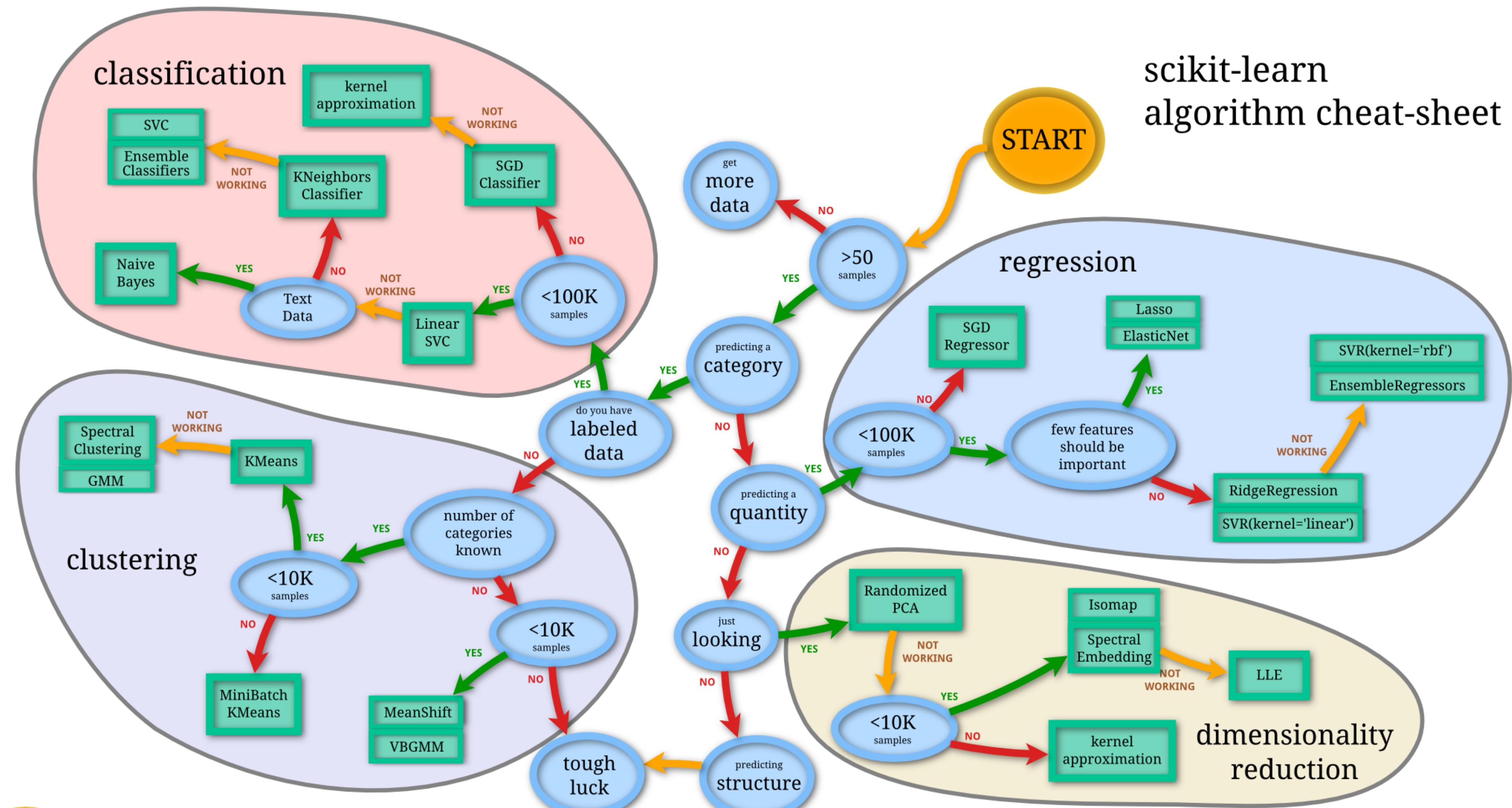


Machine Learning



Back
scikit
learn

scikit-learn
algorithm cheat-sheet



Machine Learning



Step 1



Step 2

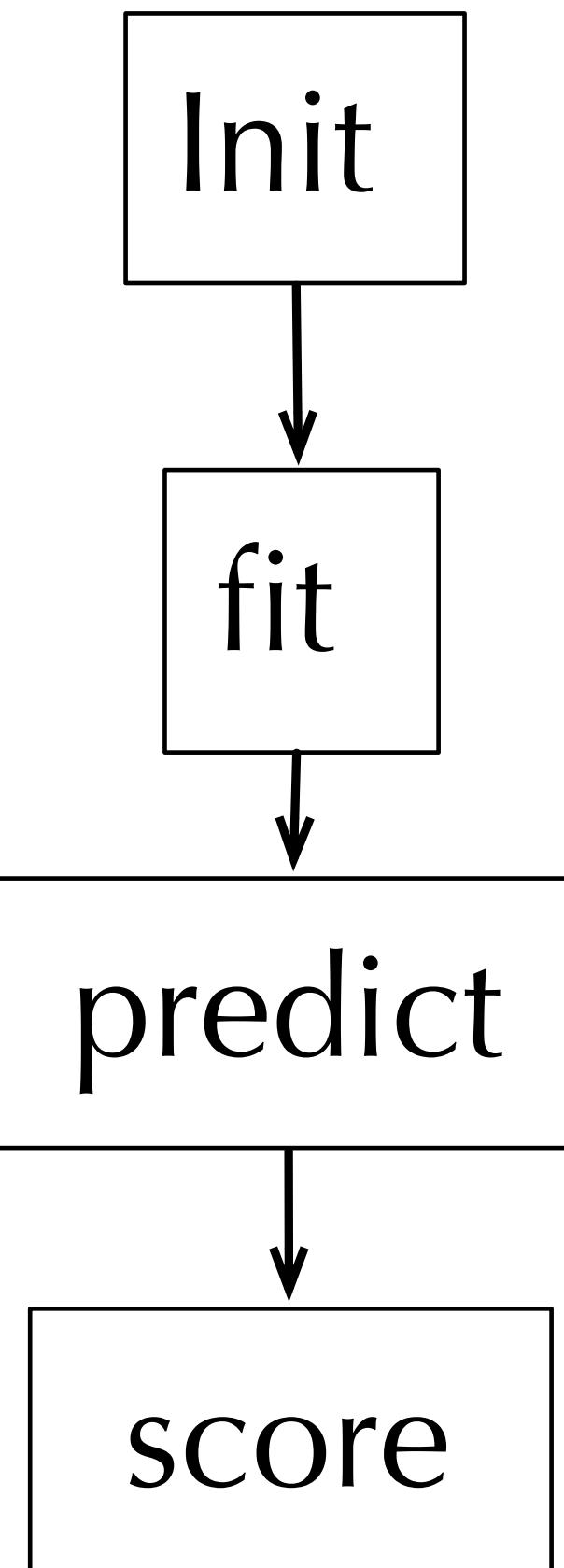


Step 3



Step 4

Workflow



```
# Import datasets, classifiers and performance metrics
from sklearn import datasets, svm, metrics
# The digits dataset
digits = datasets.load_digits()
# Create a classifier: a support vector classifier
classifier = svm.SVC(gamma=0.001)
# We learn the digits on the first half of the digits
classifier.fit(data[:n_samples // 2], digits.target[:n_samples // 2])
# Now predict the value of the digit on the second half:
expected = digits.target[n_samples // 2:]
predicted = classifier.predict(data[n_samples // 2:])
print("Classification report for classifier %s:\n%s\n"
      % (classifier, metrics.classification_report(expected, predicted)))
print("Confusion matrix:\n%s" % metrics.confusion_matrix(expected, predicted))
```

Computational Graph

Deep learning is matrix manipulation



Step 1



Step 2



Step 3



Step 4

Tensorflow, Theano

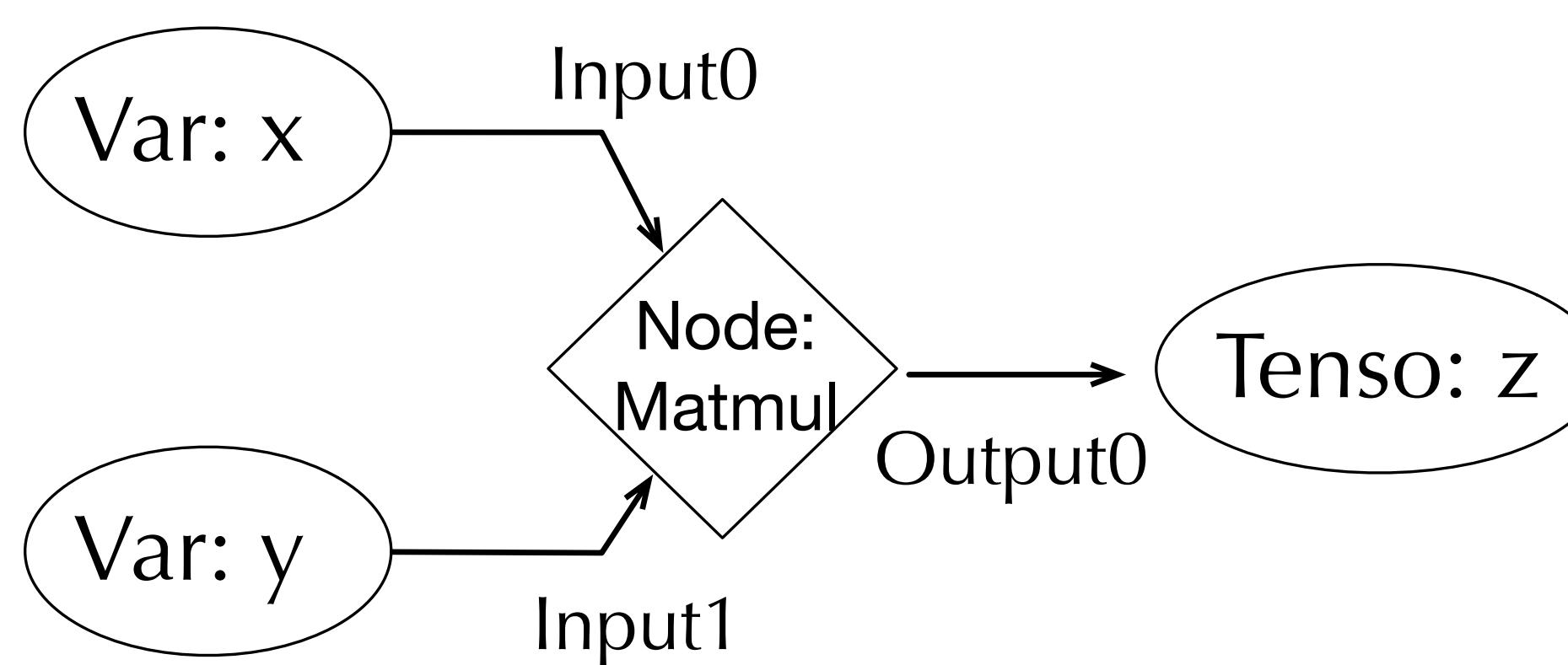
```
import tensorflow as tf  
  
x = tf.placeholder(dtype='float32', shape=(2, 3))  
y = tf.placeholder(dtype='float32', shape=(3, 4))  
z = tf.matmul(x, y)  
print(z)
```

Numpy, PyTorch

```
import numpy as np  
  
x = np.random.rand(2, 3)  
y = np.random.rand(3, 4)  
z = np.dot(x, y)  
print(z)
```

```
Tensor("MatMul:0", shape=(2, 4), dtype=float32)  
[[ 0.78686185  0.61629636  0.59583339  1.04375596]  
 [ 1.48166385  1.08328796  1.10720437  1.58379081]]
```

```
[[ 0.78686185  0.61629636  0.59583339  1.04375596]  
 [ 1.48166385  1.08328796  1.10720437  1.58379081]]
```



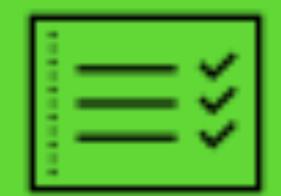
Computational Graph



Step 1



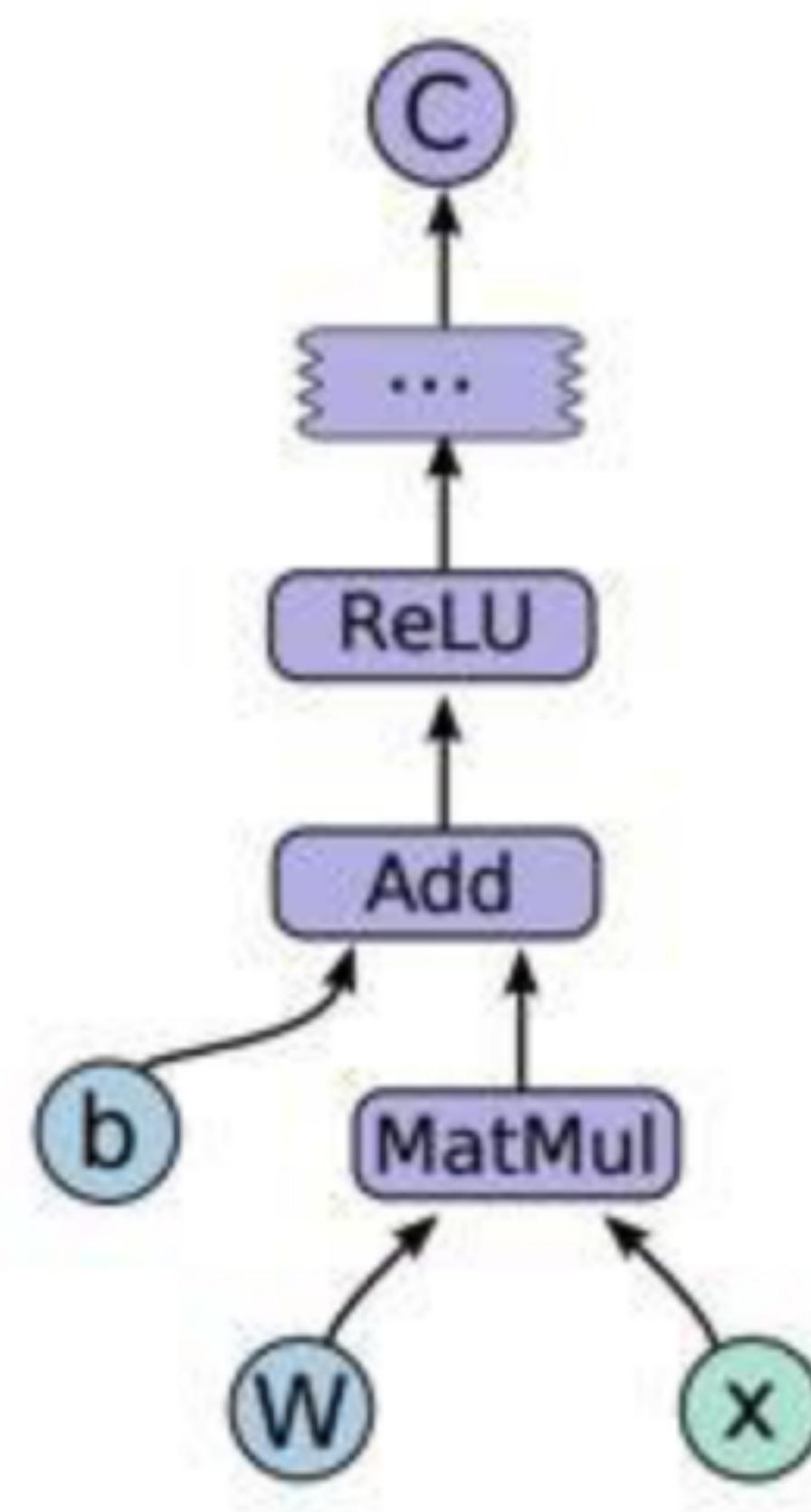
Step 2



Step 3



Step 4



Computational graph of
a feed forward layer

```
import tensorflow as tf

# define two matrices
matrix1 = tf.constant([[3., 3.]])
matrix2 = tf.constant([[2.], [2.]])
# do multiplication
product = tf.matmul(matrix1, matrix2)
# Let's get a new session
sess = tf.Session()
result = sess.run(product)
print(result)
# the output is : [[ 12.]]
```

Tensorflow code

```
from keras.models import Sequential
from keras.layers import Dense, Activation
model = Sequential()
model.add(Dense(2, input_dim=3))
model.add(Activation('relu'))
```

Keras code

Computational Graph



Step 1



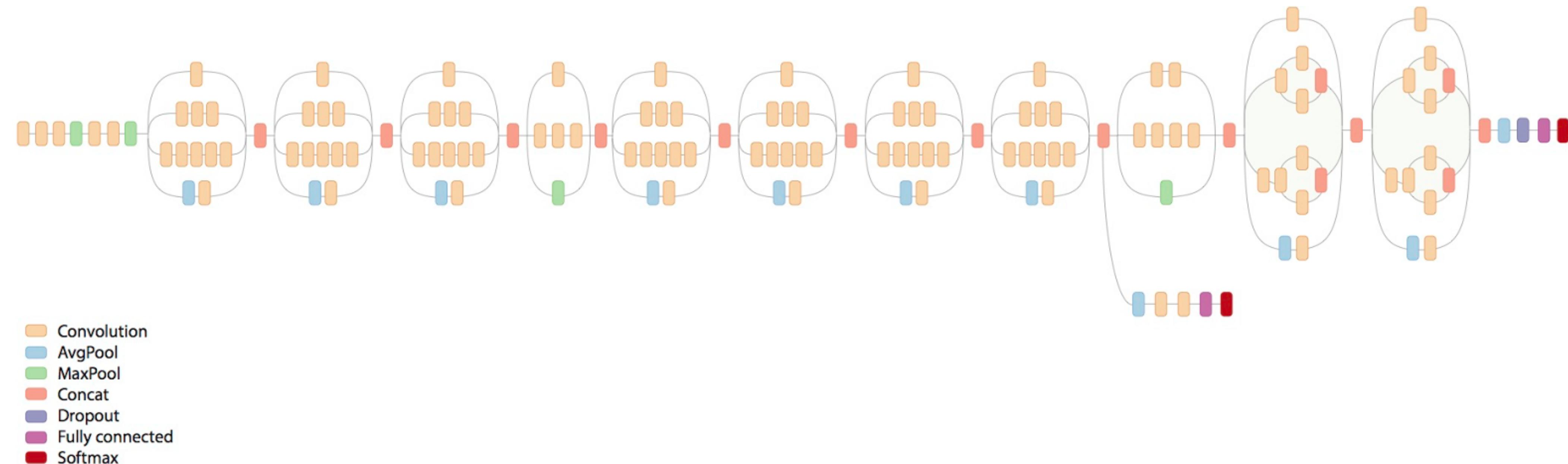
Step 2



Step 3



Step 4



```
from tensorflow.contrib import slim
logits, endpoints = slim.inception.inception_v3(
    images,
    dropout_keep_prob=0.8,
    num_classes=num_classes,
    is_training=for_training,
    restore_logits=restore_logits,
    scope="Inception")
```

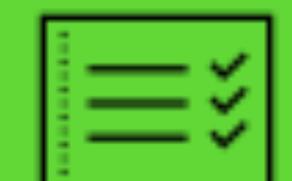
Evaluation



Step 1



Step 2

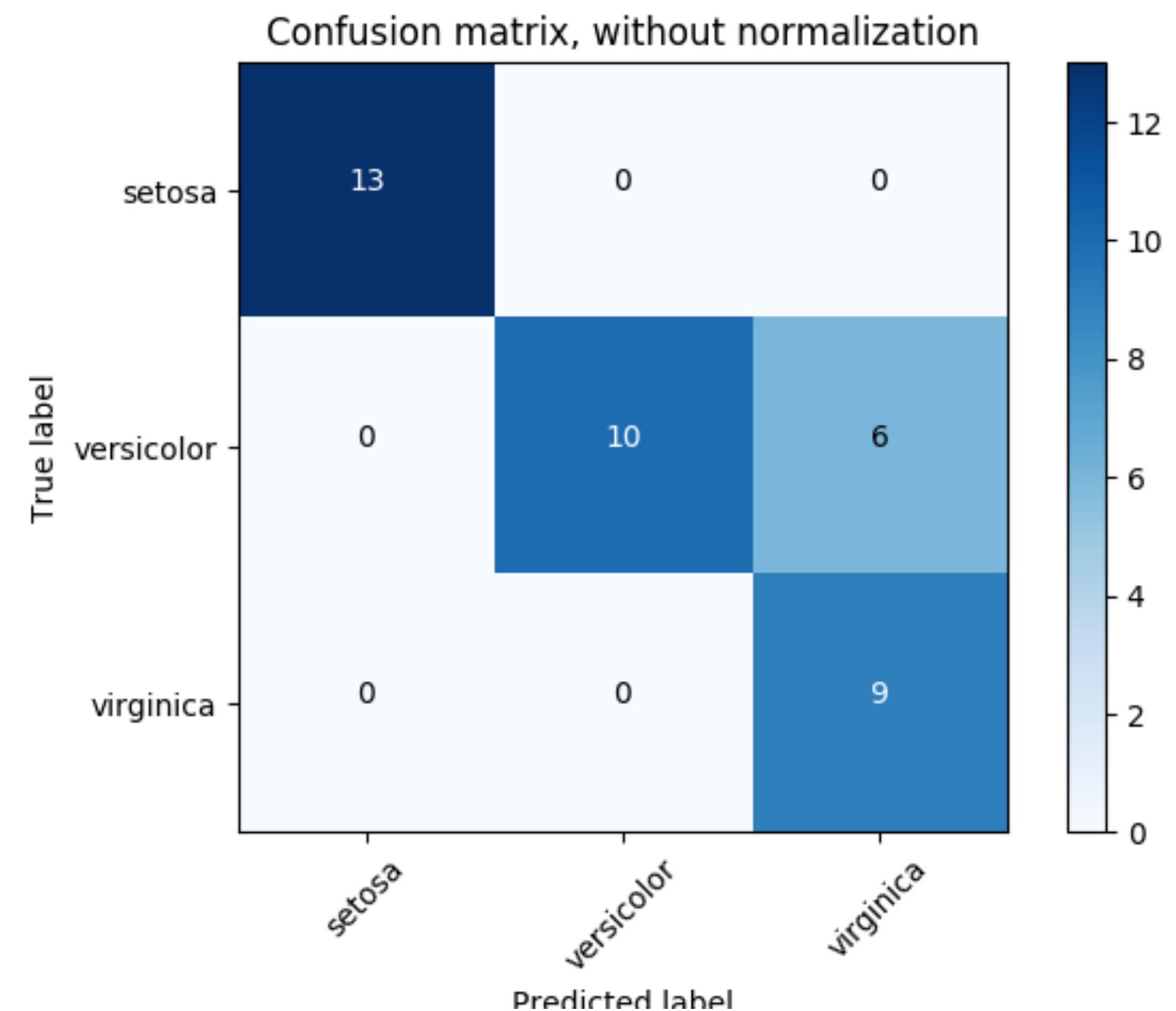


Step 3

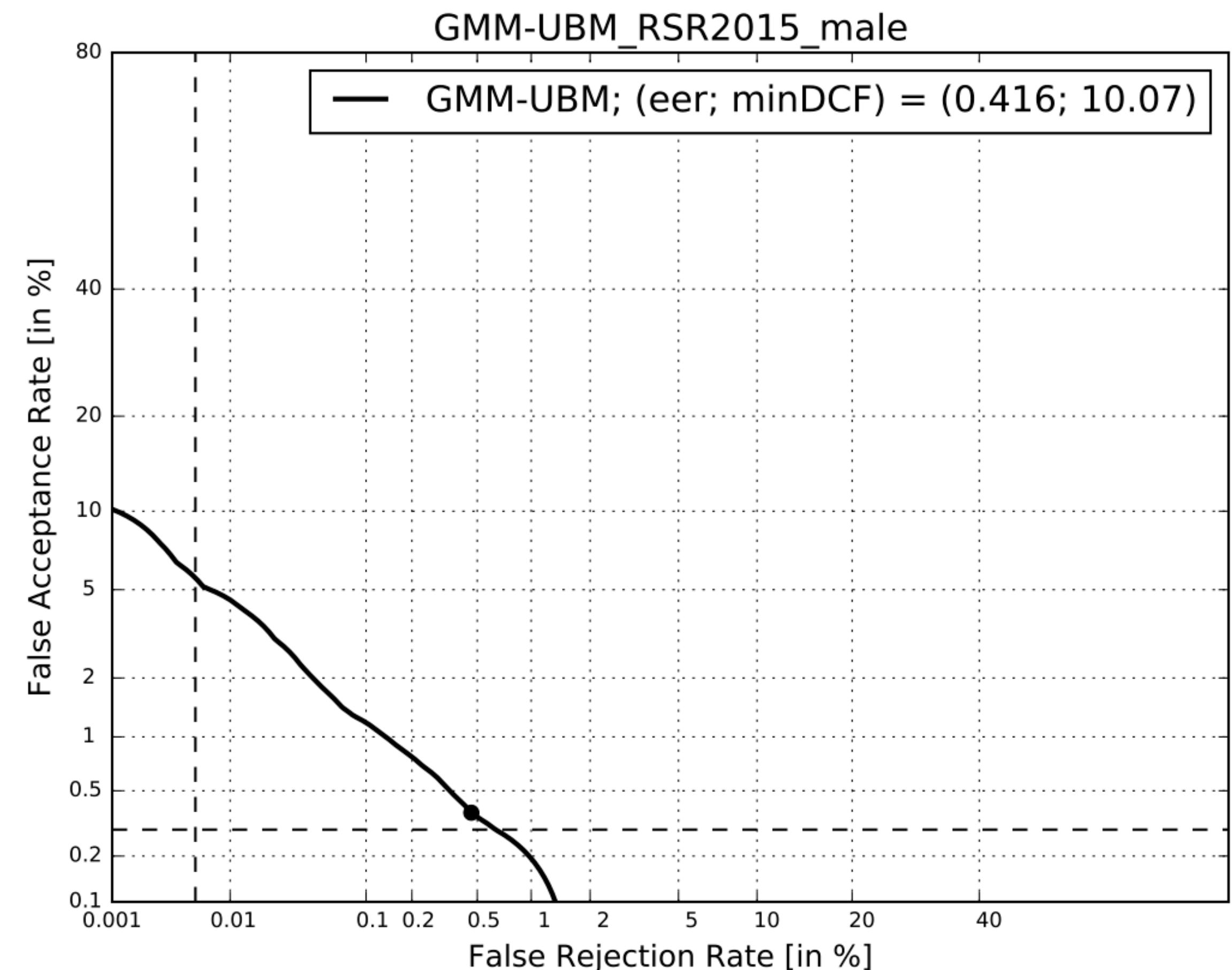


Step 4

Confusion matrix (sklearn)



EER (implemeted in sidekit)



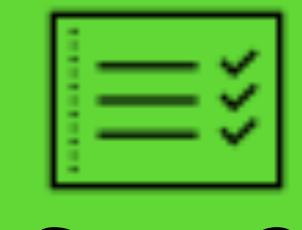
Git: version control



Step 1



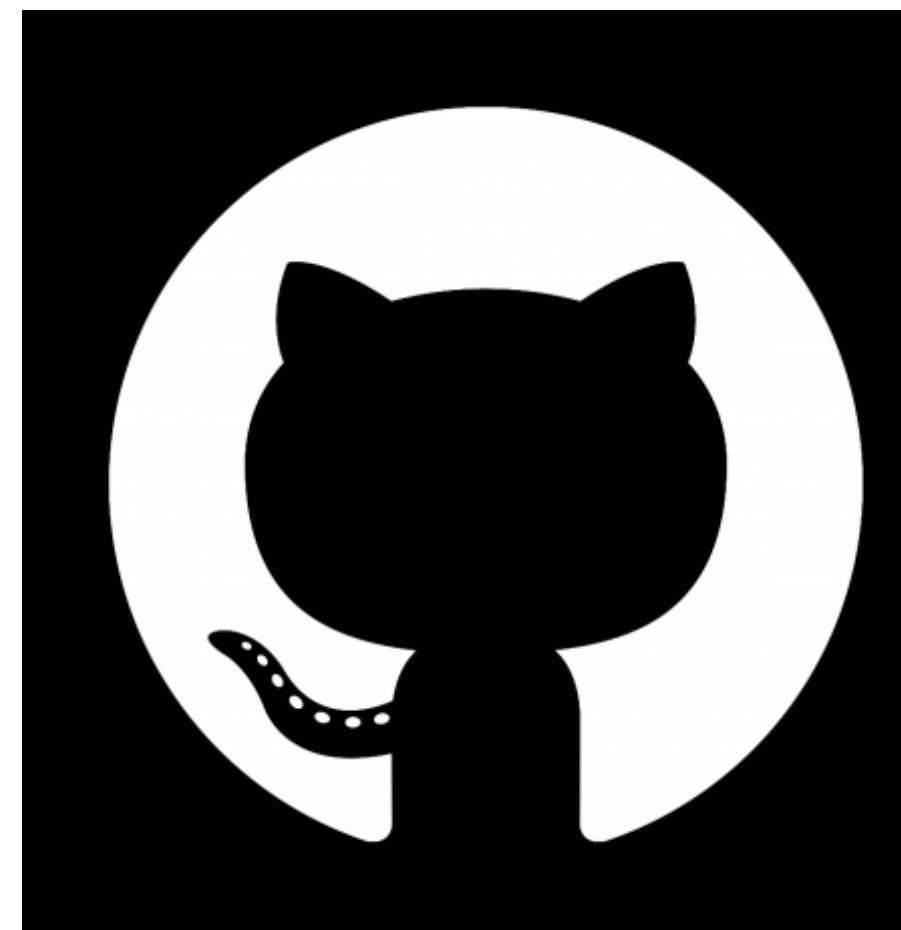
Step 2



Step 3



Step 4



GitLab

