**Project by:** Anand Vibhuti(15MI401), Vertika Sharma(15MI402)
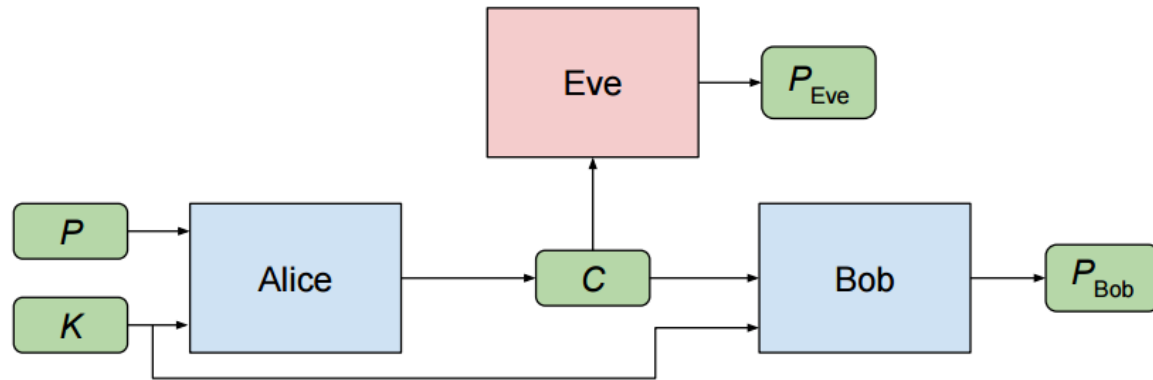
**Objective:** Learning to protect communications by adversarial neural cryptography

**Scenario:** A scenario in security involving three parties: Alice, Bob, and Eve. Typically, Alice and Bob wish to communicate securely, and Eve wishes to eavesdrop on their communications. Eve's goal is simple: to reconstruct P accurately.



**Methodology:**

**1)**In this scenario:

- Alice is the encryption algorithm
- Bob is the decryption algorithm
- Eve is the attacker

  We use 16-bit messages, secret keys and ciphertexts.

**2)**The Alice network will need to have two input vectors: the message to be encrypted, and the secret key. These get concatenated and fed to a dense layer. The signal is then passed through a sequence of four Conv1D layers to create the output.

**3)** The Eve network will have Two dense layers acting on the inputs. The idea is to give it a better chance at decrypting the message, since it doesn't have access to the secret key - it only sees the ciphertext.

**4)** Bob's network is identical to Alice's, except input0 now represents the ciphertext instead of the plaintext.

**5)** The loss for Eve is just the L1 distance between ainput0 and eoutput. Instead of doing an average, sum is taken over all the bits in the message, and its value represents the average number of bits Eve guesses incorrectly. We then take the average across the entire mini-batch with K.mean(). The minimum value of the loss is 0 (Eve guesses all the bits correctly), while the maximum is 16 (Eve is wrong about all the bits).

**6)** First, we want Bob to successfully decrypt the ciphertext.

**7)** We want Alice to learn an encryption scheme which Eve can't break. In an ideal situation, Eve should do random guessing, in which case she would correctly guess half the bits, or m_bits/2 correctly (corresponding to a loss value of 8).

**8)** As per optimizer, we use RMSprop with a default learning rate of 0.001.

**9)** Finally, we create two training models with the loss functions defined earlier.

**10)** Training of the networks.

**11)** Evaluating with providing values.


**Network layers:**

**1)**Alice's Network:

```
[ ]  ainput0 = Input(shape=(m_bits,)) # the message
     ainput1 = Input(shape=(k_bits,)) # the key
     ainput = concatenate([ainput0, ainput1], axis=1)

     adense1 = Dense(units=(m_bits + k_bits))(ainput)
     adense1a = Activation('tanh')(adense1)
     areshape = Reshape((m_bits + k_bits, 1,))(adense1a)
     #output of the Dense layer will be 2-dimensional. This needs to be reshaped to (batch_size, m_bits + k_bits, 1) before

     aconv1 = Conv1D(filters=2, kernel_size=4, strides=1, padding=pad)(areshape)
     aconv1a = Activation('tanh')(aconv1)
     aconv2 = Conv1D(filters=4, kernel_size=2, strides=2, padding=pad)(aconv1a)
     aconv2a = Activation('tanh')(aconv2)
     aconv3 = Conv1D(filters=4, kernel_size=1, strides=1, padding=pad)(aconv2a)
     aconv3a = Activation('tanh')(aconv3)
     aconv4 = Conv1D(filters=1, kernel_size=1, strides=1, padding=pad)(aconv3a)
     aconv4a = Activation('sigmoid')(aconv4)

     aoutput = Flatten()(aconv4a)

     alice = Model([ainput0, ainput1], aoutput, name='alice')
```

**2)** Bob's Network:

```
[ ]  binput0 = Input(shape=(c_bits,)) #message
     binput1 = Input(shape=(k_bits,)) #key
     binput = concatenate([binput0, binput1], axis=1)

     bdense1 = Dense(units=(c_bits + k_bits))(binput)
     bdense1b = Activation('tanh')(bdense1)
     breshape = Reshape((c_bits + k_bits, 1,))(bdense1b)

     bconv1 = Conv1D(filters=2, kernel_size=4, strides=1, padding=pad)(breshape)
     bconv1b = Activation('tanh')(bconv1)
     bconv2 = Conv1D(filters=4, kernel_size=2, strides=2, padding=pad)(bconv1b)
     bconv2b = Activation('tanh')(bconv2)
     bconv3 = Conv1D(filters=4, kernel_size=1, strides=1, padding=pad)(bconv2b)
     bconv3b = Activation('tanh')(bconv3)
     bconv4 = Conv1D(filters=1, kernel_size=1, strides=1, padding=pad)(bconv3b)
     bconv4b = Activation('sigmoid')(bconv4)

     boutput = Flatten()(bconv4b)

     bob = Model([binput0, binput1], boutput, name='bob')
```

**3)** Eve's Network:

```
[ ]  einput = Input(shape=(c_bits,)) # the ciphertext

     edense1 = Dense(units=(c_bits + k_bits))(einput)
     edense1a = Activation('tanh')(edense1)
     edense2 = Dense(units=(c_bits + k_bits))(edense1a)
     edense2a = Activation('tanh')(edense2)
     ereshape = Reshape((c_bits + k_bits, 1,))(edense2a)

     econv1 = Conv1D(filters=2, kernel_size=4, strides=1, padding=pad)(ereshape)
     econv1a = Activation('tanh')(econv1)
     econv2 = Conv1D(filters=4, kernel_size=2, strides=2, padding=pad)(econv1a)
     econv2a = Activation('tanh')(econv2)
     econv3 = Conv1D(filters=4, kernel_size=1, strides=1, padding=pad)(econv2a)
     econv3a = Activation('tanh')(econv3)
     econv4 = Conv1D(filters=1, kernel_size=1, strides=1, padding=pad)(econv3a)
     econv4a = Activation('sigmoid')(econv4)

     eoutput = Flatten()(econv4a)# Eve's attempt at guessing the plaintext

     eve = Model(einput, eoutput, name='eve')
```
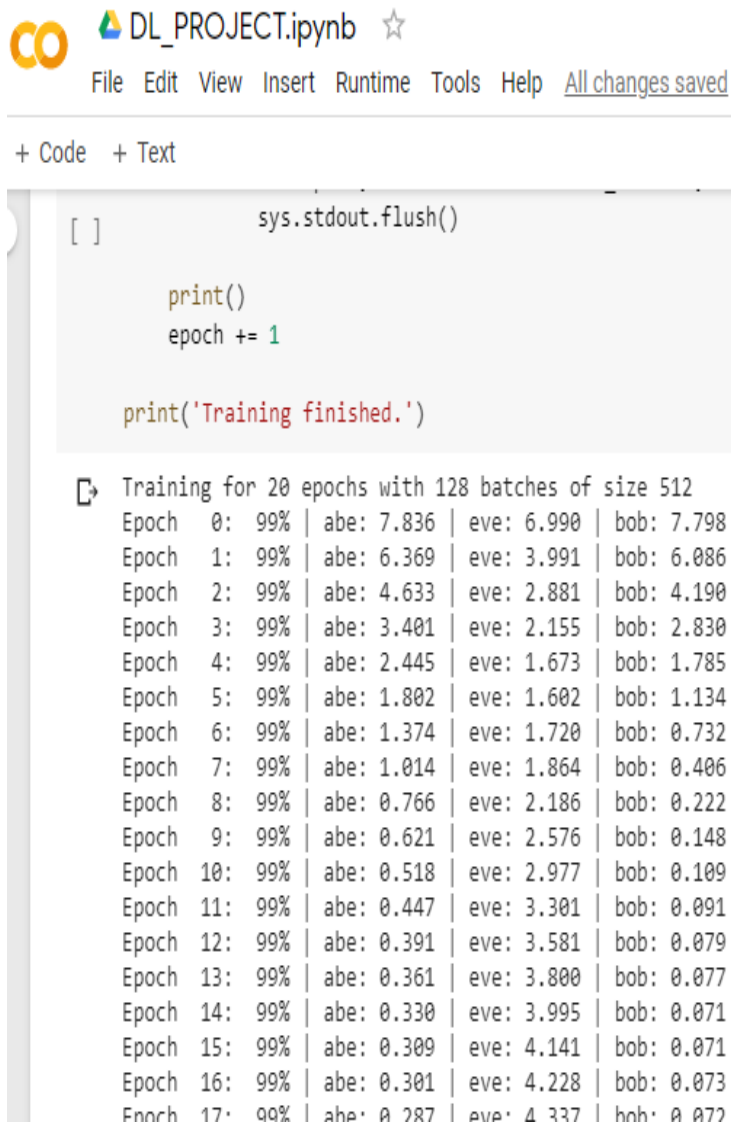
**Results:**



```
CO  DL_PROJECT.ipynb  ☆
    File  Edit  View  Insert  Runtime  Tools  Help   All changes saved

+ Code   + Text

[ ]                    sys.stdout.flush()

         print()
         epoch += 1


    print('Training finished.')


 ⊡  Training for 20 epochs with 128 batches of size 512
    Epoch   0:  99% | abe: 7.836 | eve: 6.990 | bob: 7.798
    Epoch   1:  99% | abe: 6.369 | eve: 3.991 | bob: 6.086
    Epoch   2:  99% | abe: 4.633 | eve: 2.881 | bob: 4.190
    Epoch   3:  99% | abe: 3.401 | eve: 2.155 | bob: 2.830
    Epoch   4:  99% | abe: 2.445 | eve: 1.673 | bob: 1.785
    Epoch   5:  99% | abe: 1.802 | eve: 1.602 | bob: 1.134
    Epoch   6:  99% | abe: 1.374 | eve: 1.720 | bob: 0.732
    Epoch   7:  99% | abe: 1.014 | eve: 1.864 | bob: 0.406
    Epoch   8:  99% | abe: 0.766 | eve: 2.186 | bob: 0.222
    Epoch   9:  99% | abe: 0.621 | eve: 2.576 | bob: 0.148
    Epoch  10:  99% | abe: 0.518 | eve: 2.977 | bob: 0.109
    Epoch  11:  99% | abe: 0.447 | eve: 3.301 | bob: 0.091
    Epoch  12:  99% | abe: 0.391 | eve: 3.581 | bob: 0.079
    Epoch  13:  99% | abe: 0.361 | eve: 3.800 | bob: 0.077
    Epoch  14:  99% | abe: 0.330 | eve: 3.995 | bob: 0.071
    Epoch  15:  99% | abe: 0.309 | eve: 4.141 | bob: 0.071
    Epoch  16:  99% | abe: 0.301 | eve: 4.228 | bob: 0.073
    Epoch  17:  99% | abe: 0.287 | eve: 4.337 | bob: 0.072
```

Figure 1: Training the networks

```python
#Plotting
steps = -1

plt.figure(figsize=(7, 4))
plt.plot(abelosses[:steps], label='A-B')
plt.plot(evelosses[:steps], label='Eve')
plt.plot(boblosses[:steps], label='Bob')
plt.xlabel("Iterations", fontsize=13)
plt.ylabel("Loss", fontsize=13)
plt.legend(fontsize=13)
plt.show()
```
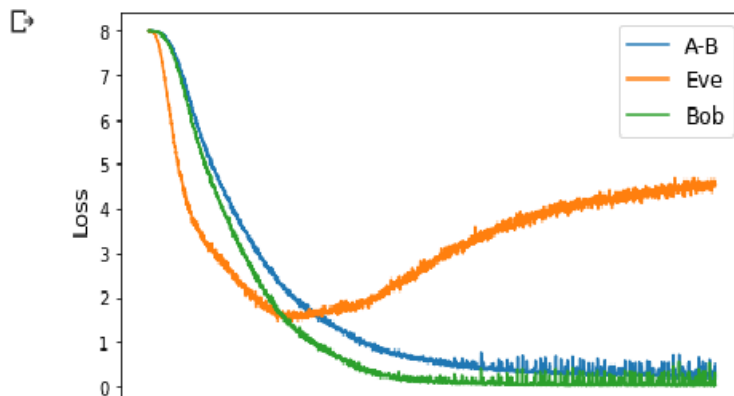


Figure 2: Plots for the computed losses

Reference Paper: **Learning to Protect Communications with Adversarial Neural Cryptography by Martín Abadi, David G. Andersen (Google Brain).**

*(Submitted on 21 Oct 2016)*