

Reliable Real-time Lip Reading

with no sound and minimal ado

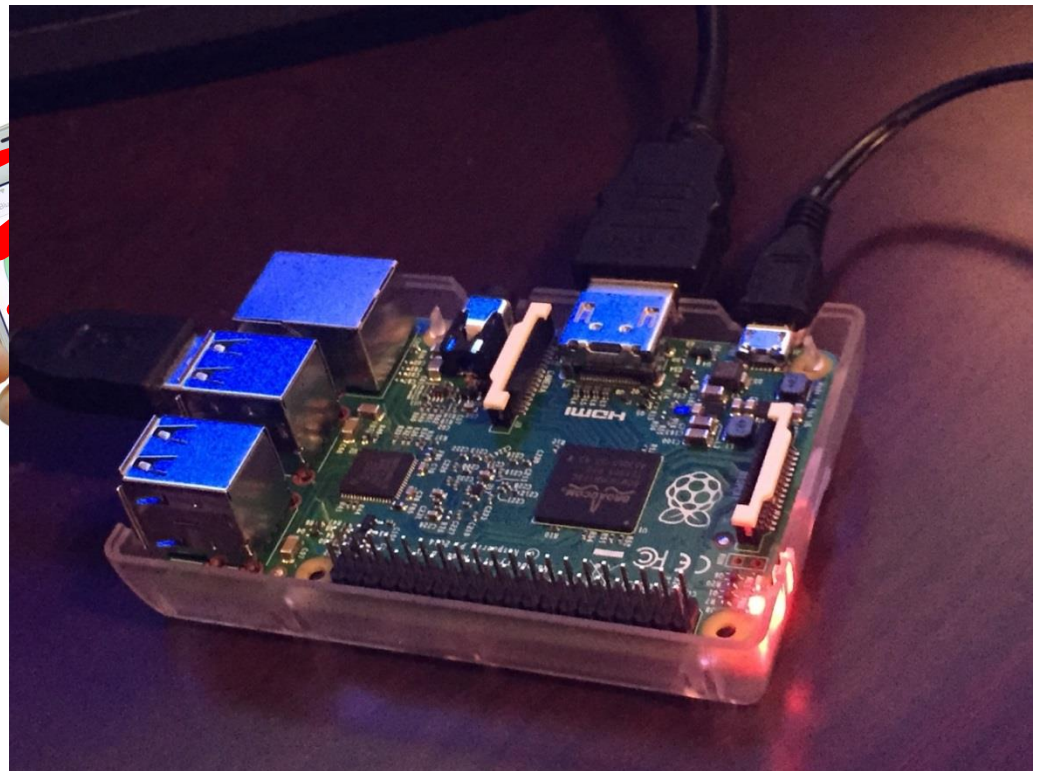
wasn't quite sure what to do with
this birthday present



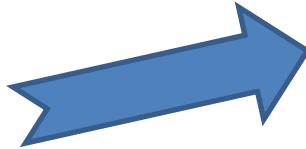
Apps take too long to open



I am not a computer person



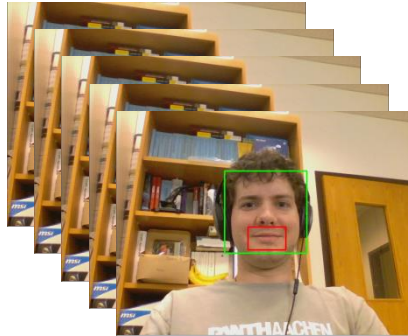
Long Term Plan



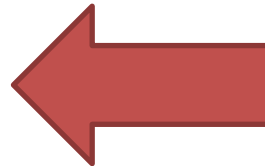
Why Machine Learning?

- Project requirements might change in future
 - Include audio
 - Add more color options
 - Add functionality for other devices
- Data may change
 - New camera
 - New desk location ➔ new lighting conditions

Implementation



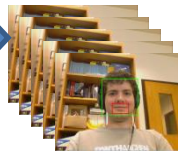
$$\begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \\ a_4 \\ a_5 \\ a_6 \\ a_7 \\ a_8 \\ a_9 \\ a_{10} \\ a_{11} \\ a_{12} \\ a_{13} \\ a_{14} \\ a_{15} \\ a_{16} \\ a_{17} \\ a_{18} \\ a_{19} \\ a_{20} \\ a_{21} \\ a_{22} \\ a_{23} \\ a_{24} \\ a_{25} \\ a_{26} \\ a_{27} \\ a_{28} \\ a_{29} \\ a_{30} \\ a_{31} \\ a_{32} \\ a_{33} \\ a_{34} \\ a_{35} \\ a_{36} \\ a_{37} \\ a_{38} \\ a_{39} \\ a_{40} \\ a_{41} \\ a_{42} \\ a_{43} \\ a_{44} \\ a_{45} \\ a_{46} \\ a_{47} \\ a_{48} \\ a_{49} \\ a_{50} \\ a_{51} \\ a_{52} \\ a_{53} \\ a_{54} \\ a_{55} \\ a_{56} \\ a_{57} \\ a_{58} \\ a_{59} \\ a_{60} \\ a_{61} \\ a_{62} \\ a_{63} \\ a_{64} \\ a_{65} \\ a_{66} \\ a_{67} \\ a_{68} \\ a_{69} \\ a_{70} \\ a_{71} \\ a_{72} \\ a_{73} \\ a_{74} \\ a_{75} \\ a_{76} \\ a_{77} \\ a_{78} \\ a_{79} \\ a_{80} \\ a_{81} \\ a_{82} \\ a_{83} \\ a_{84} \\ a_{85} \\ a_{86} \\ a_{87} \\ a_{88} \\ a_{89} \\ a_{90} \\ a_{91} \\ a_{92} \\ a_{93} \\ a_{94} \\ a_{95} \\ a_{96} \\ a_{97} \\ a_{98} \\ a_{99} \end{bmatrix}$$



Implementation



Viola Jones
Boosted Haar
Cascade



Crop



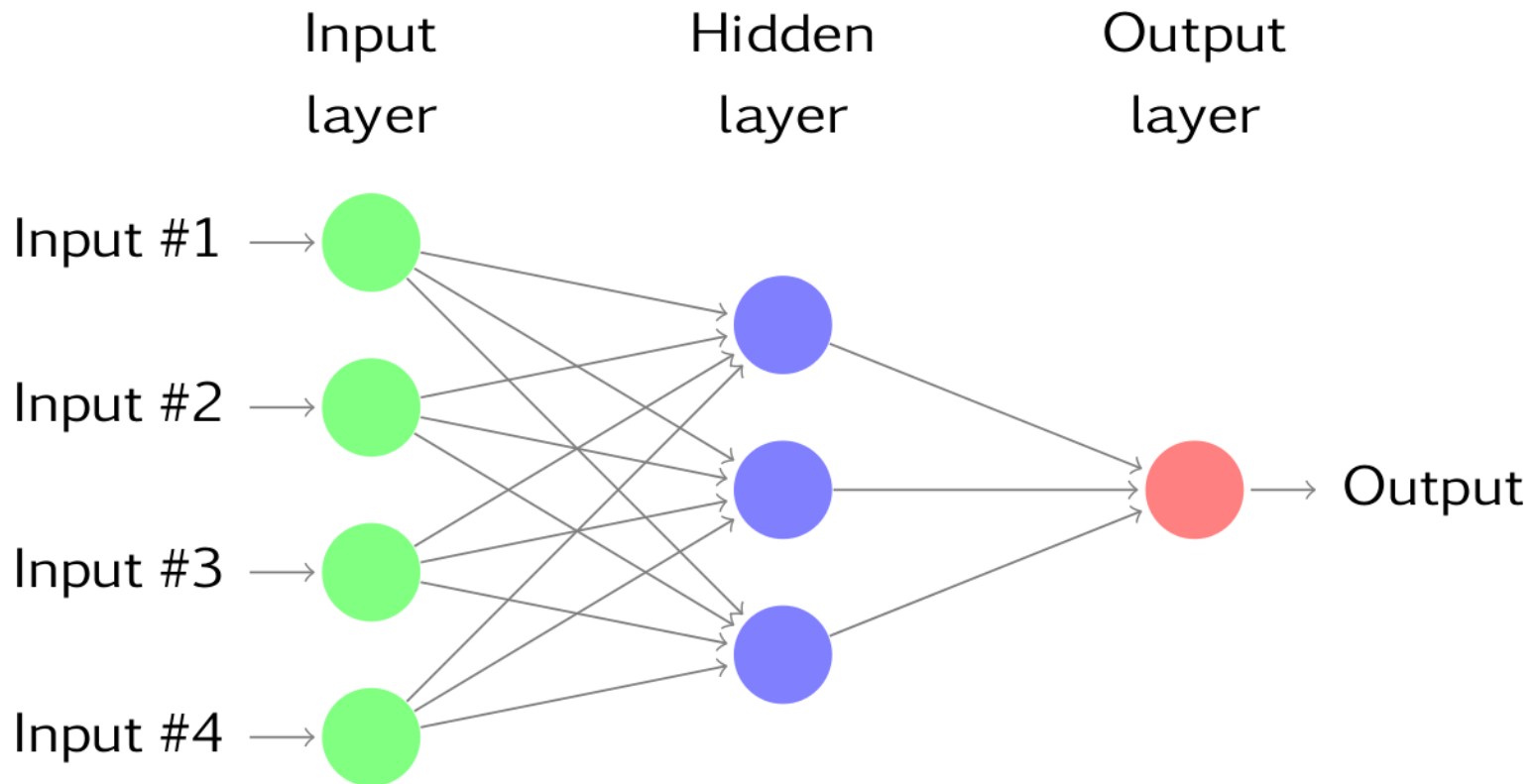
(Unsupervised)
Variational
Convolution
Neural Net
Autoencoder



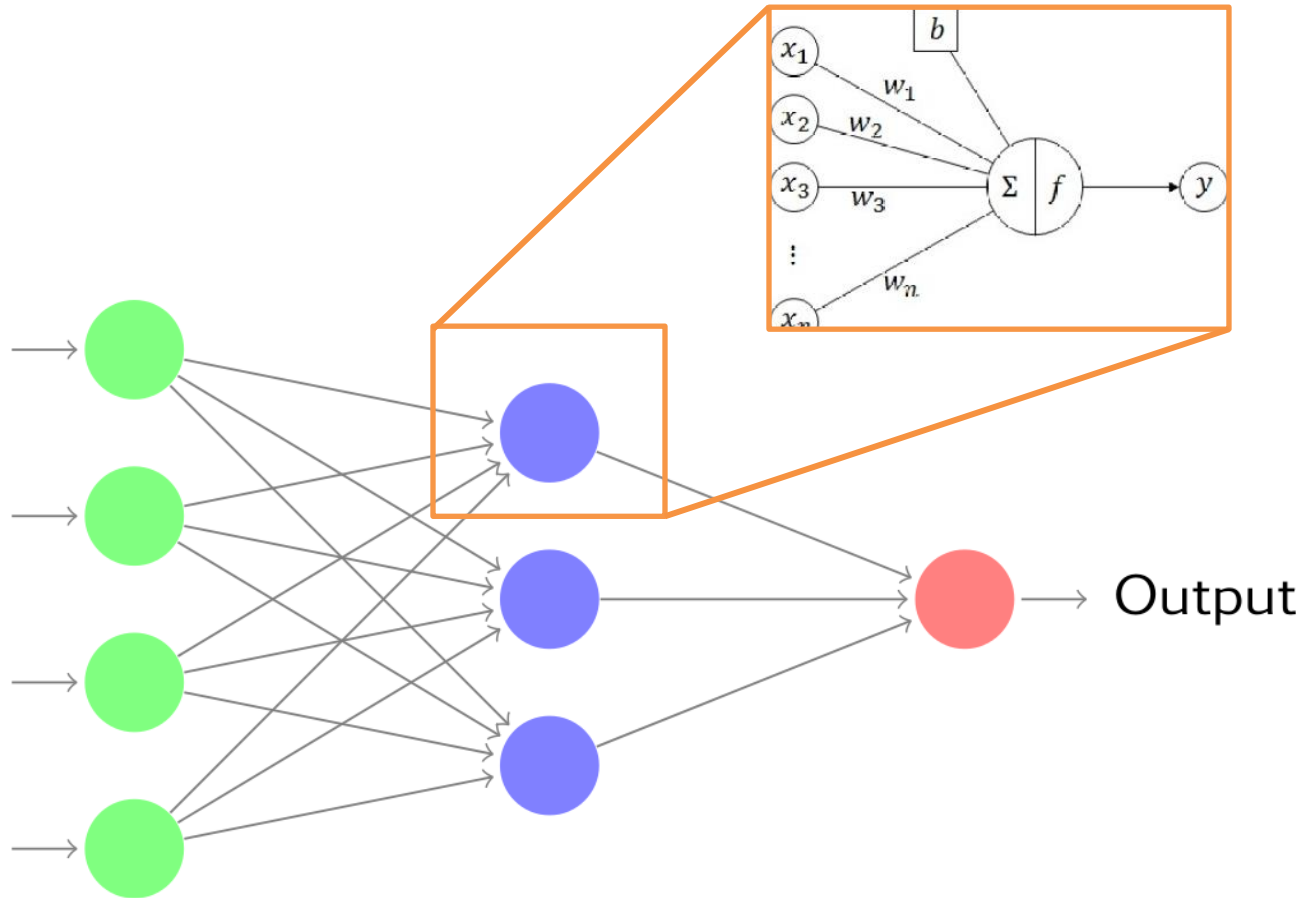
(Supervised) Recurrent
Neural Net

$$\begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$$

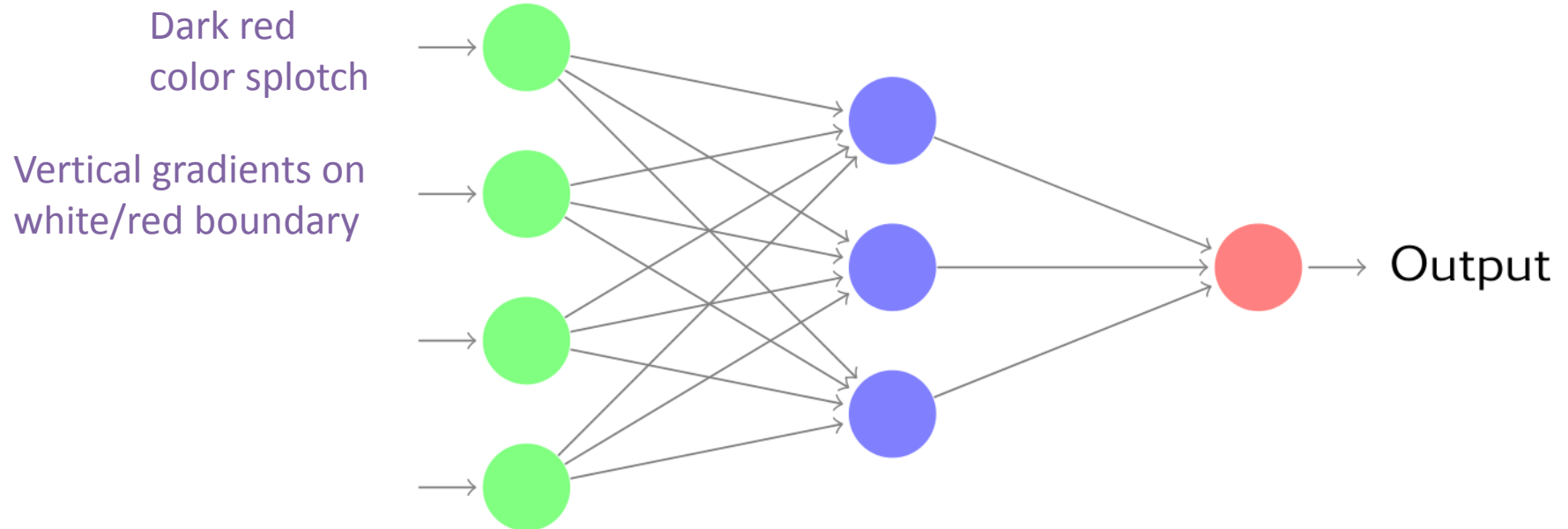
Variational Convolution Neural Net Autoencoder



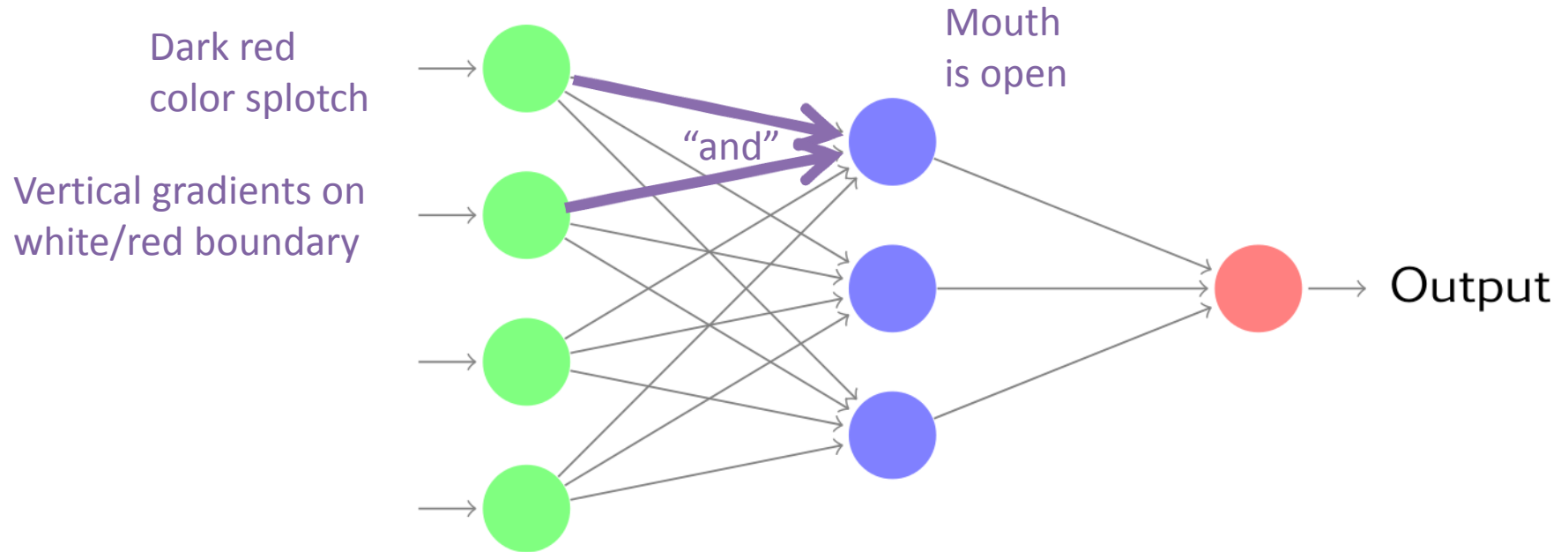
Neural Net



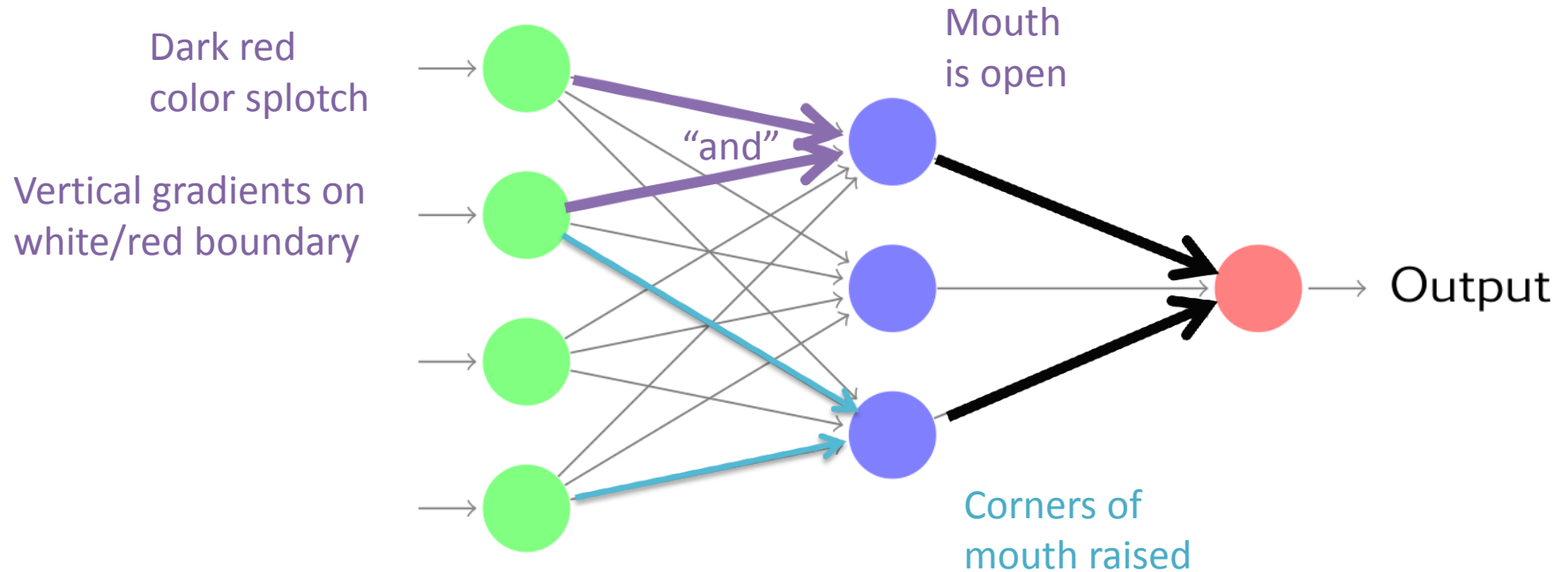
Neural Net can Implement Logic



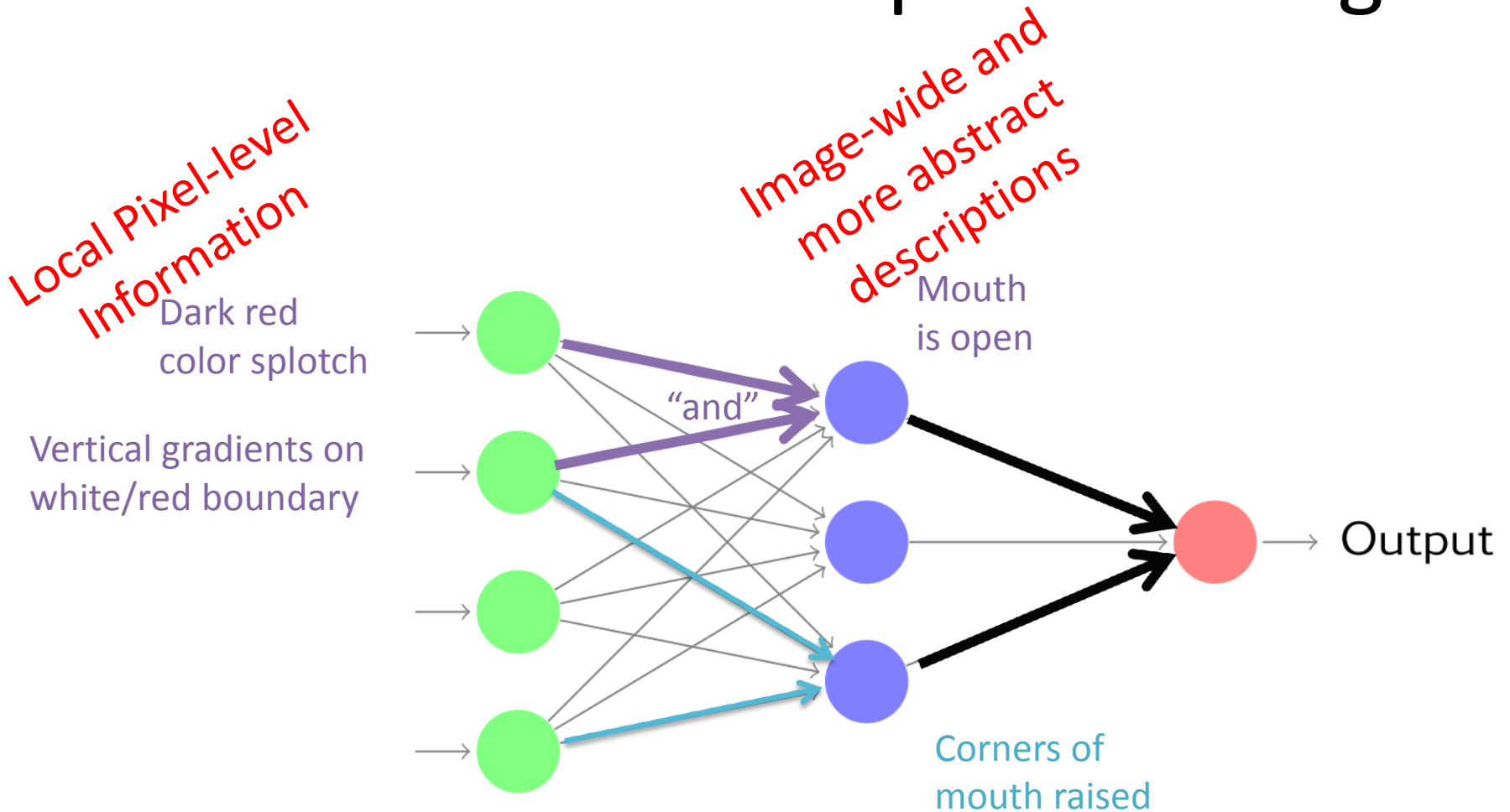
Neural Net can Implement Logic



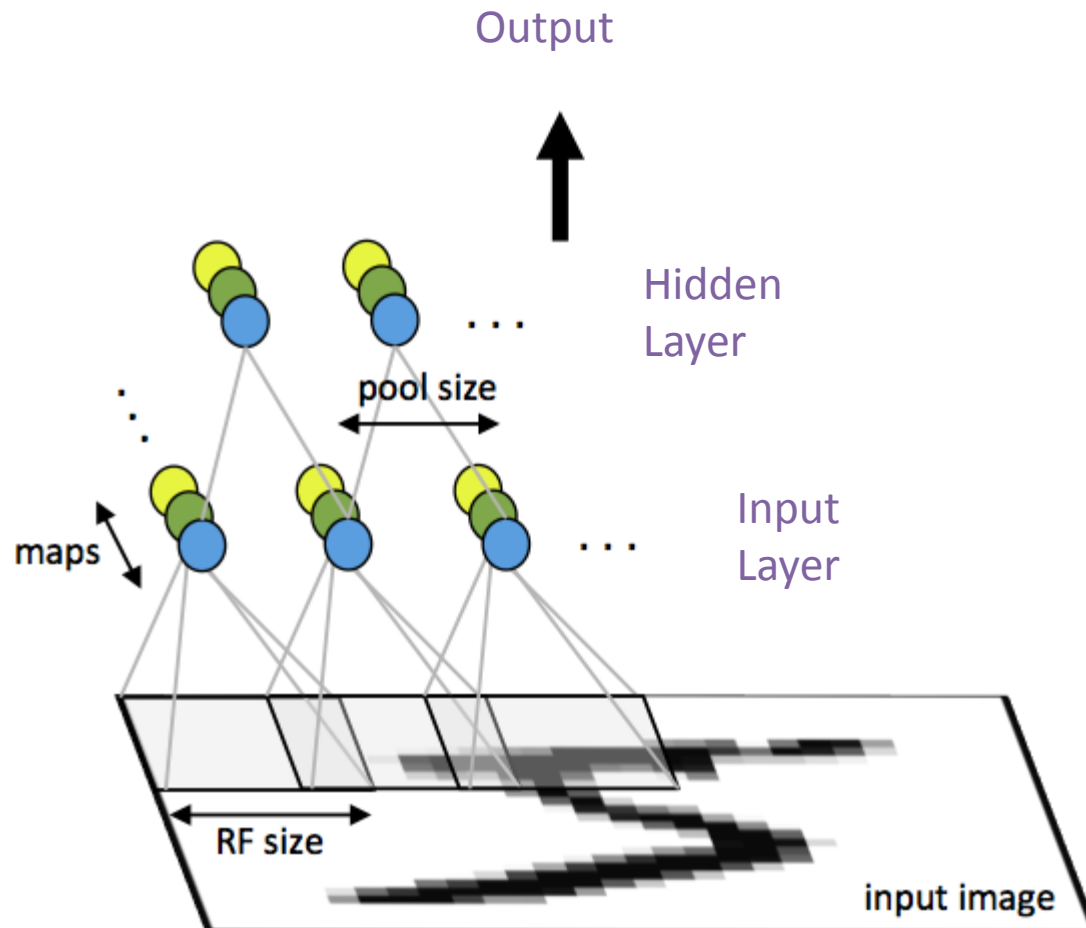
Neural Net can Implement Logic



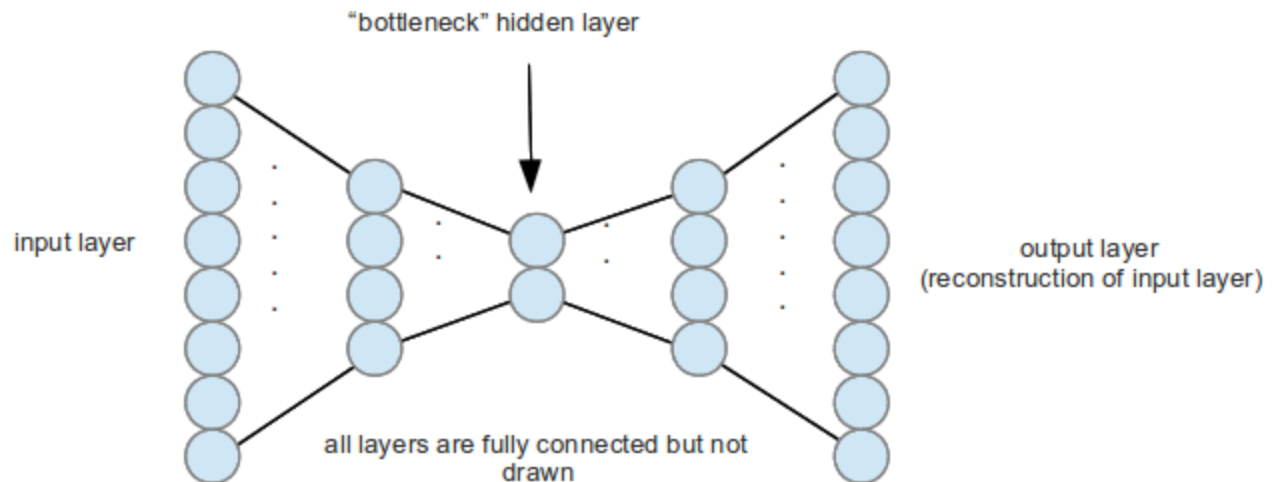
Neural Net can Implement Logic



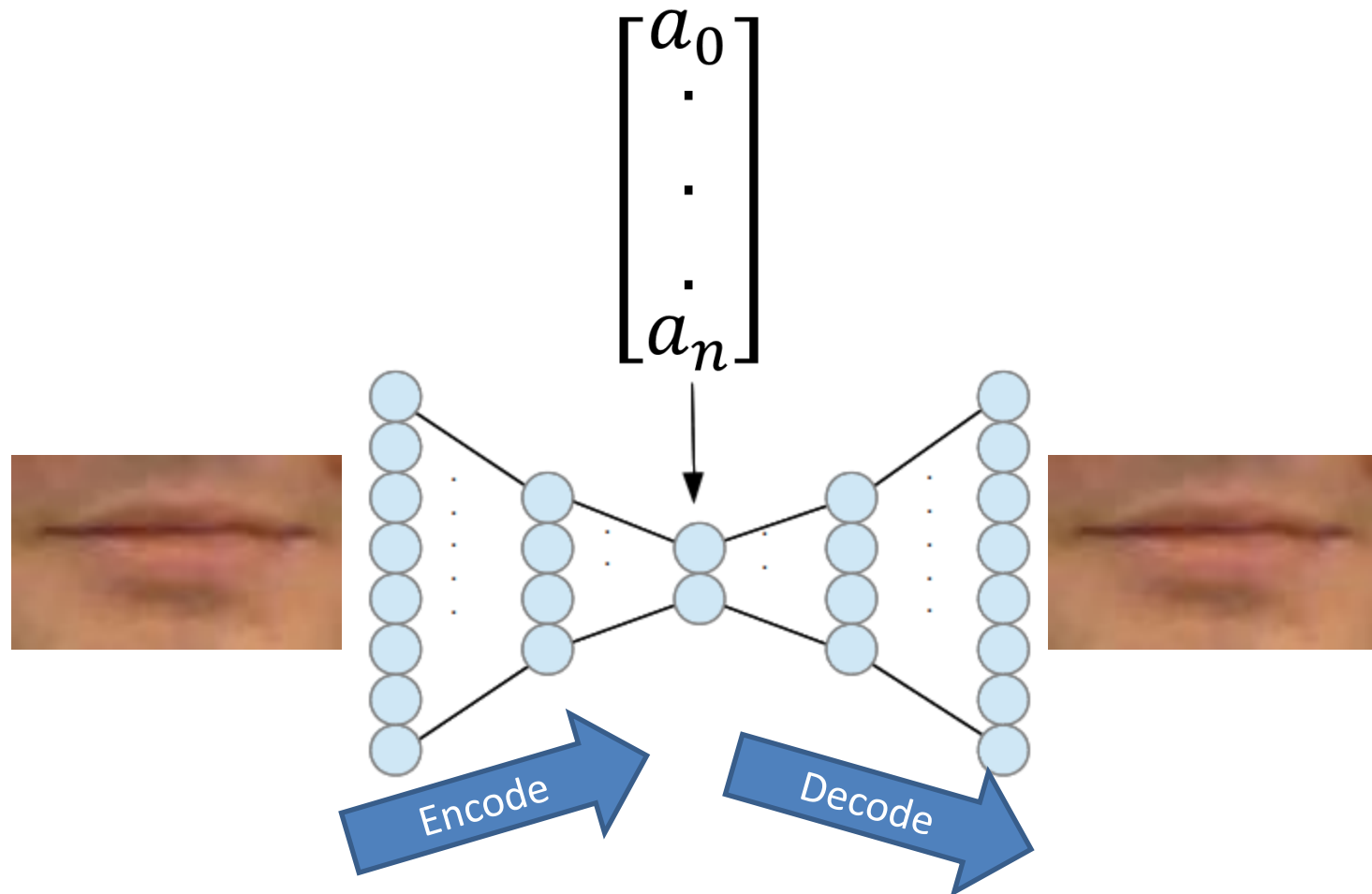
Variational Convolution Neural Net Autoencoder



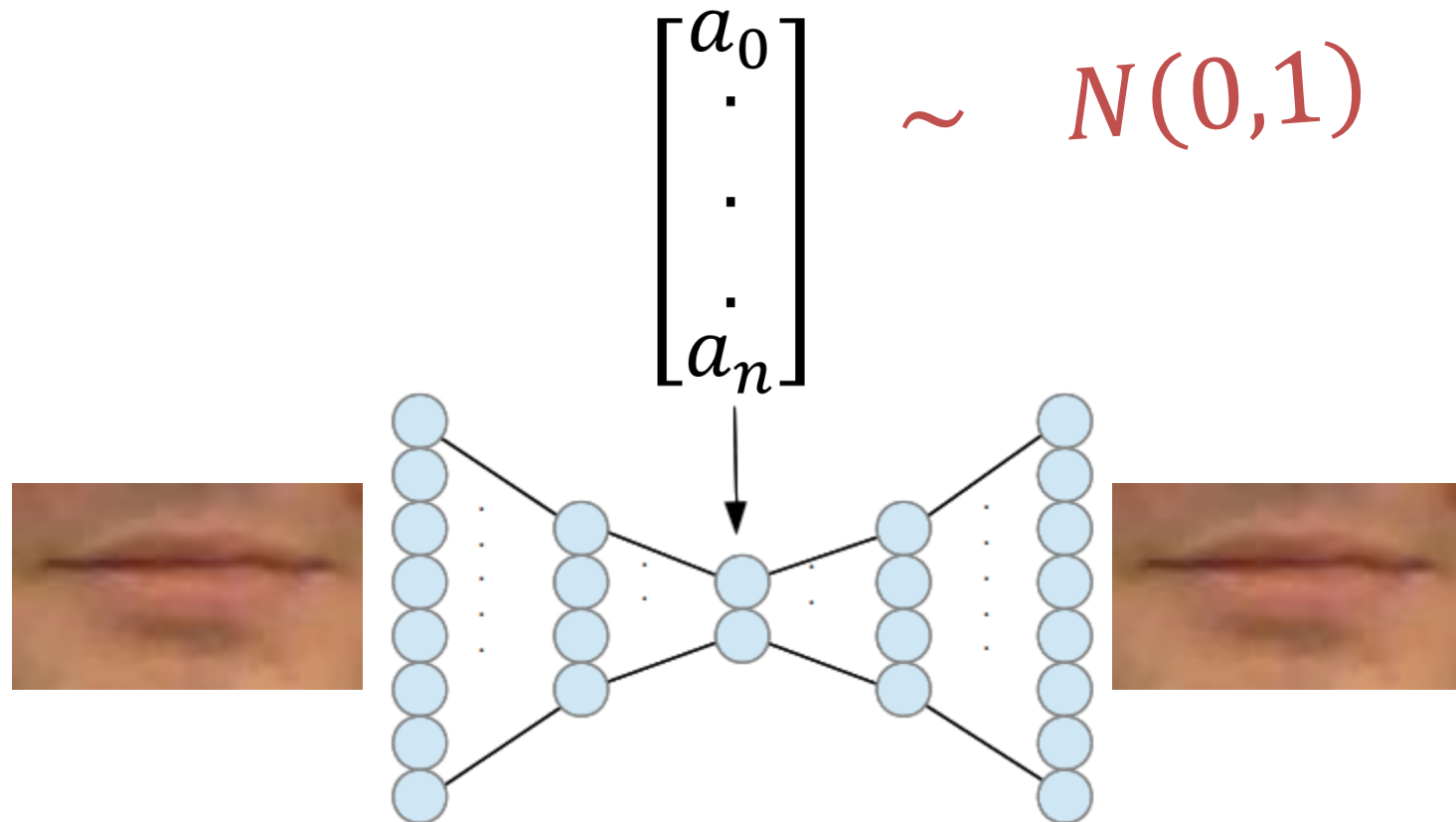
Variational Convolution Neural Net **Autoencoder**



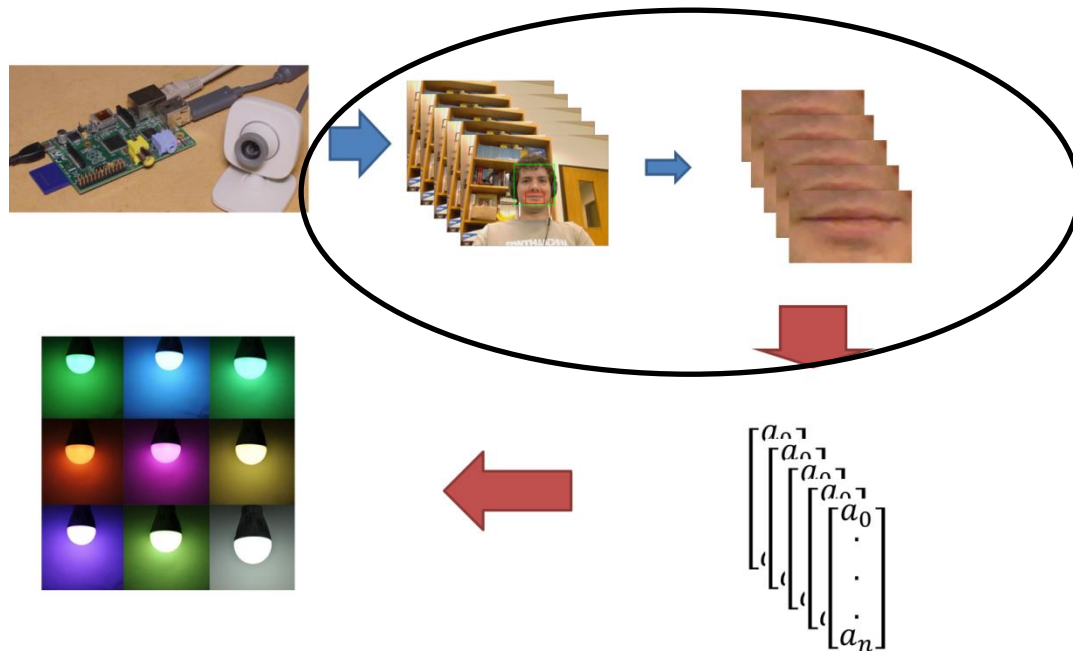
Variational Convolution Neural Net Autoencoder



Variational Convolution Neural Net Autoencoder

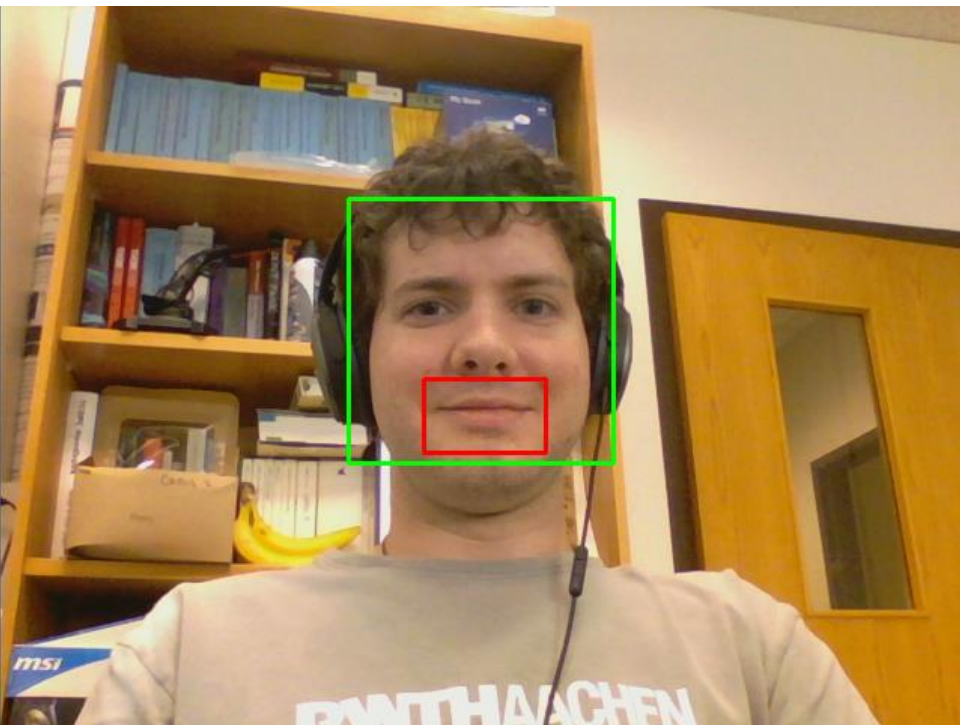


Step 1: Detect and Crop

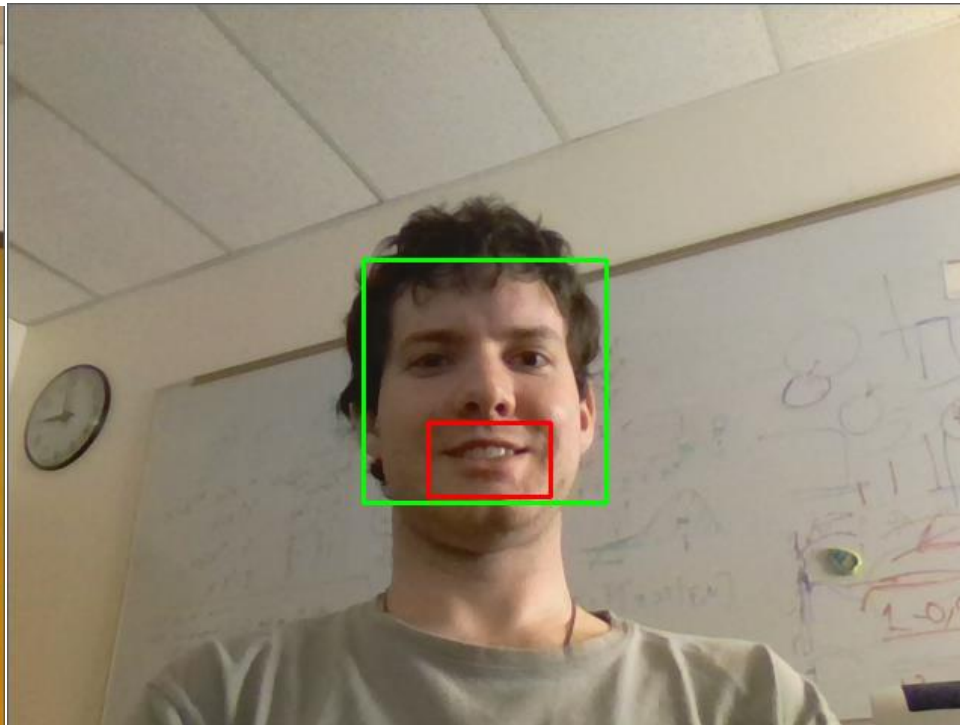


Step 1: Detect and Crop

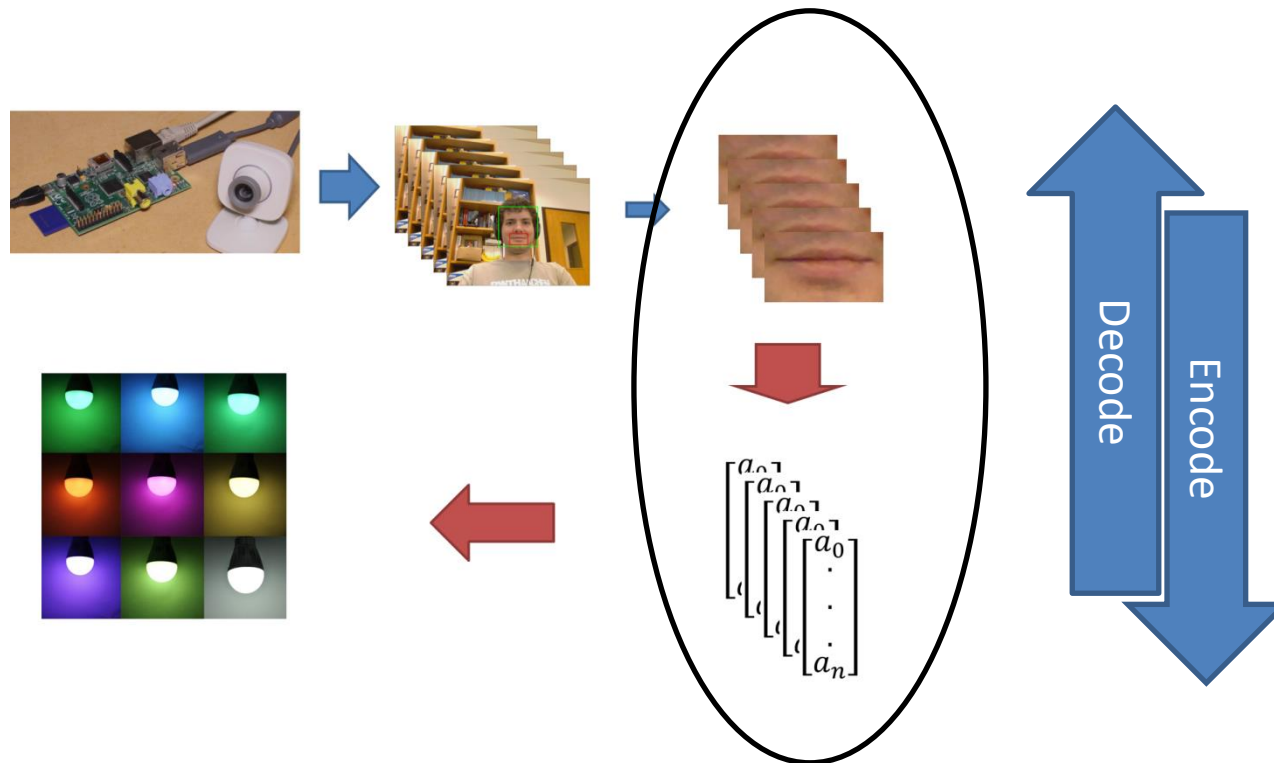
Even Lighting
(easy)



Uneven Lighting
(difficult)



Step 2: Encode



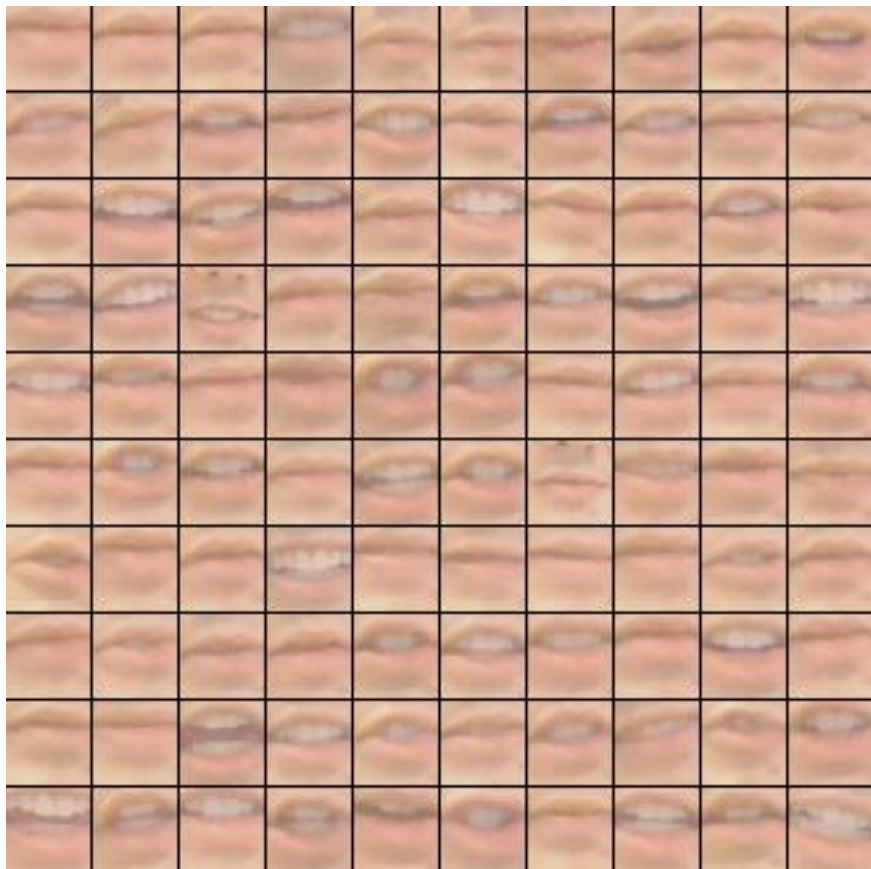
Training

(Variational Autoencoder)

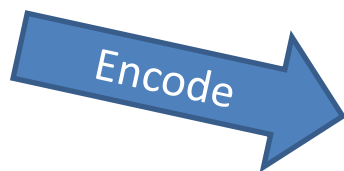
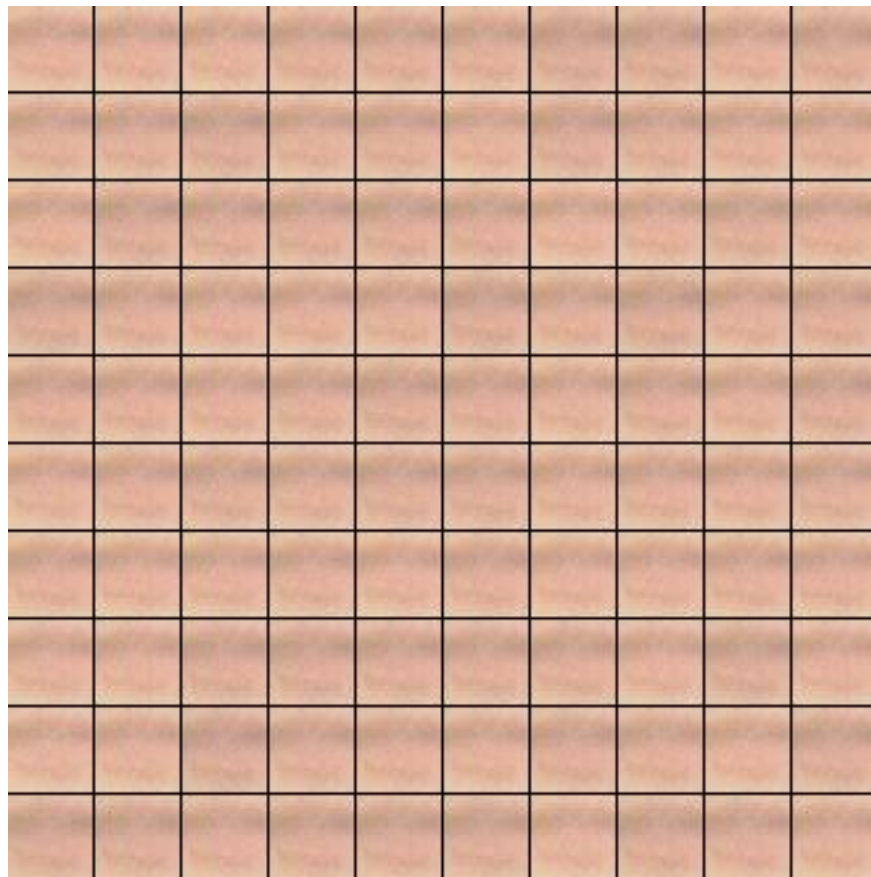
- 20,000 images (32x32x3) taken during Skype conversations
- 11,966,848 parameter model
 - 500x as many parameters as independent inputs!
 - Only 20mil pixels in training set!
- Projected 7 months for model to converge on modest laptop cpu
- ~14 hours on TACC gpu

Training

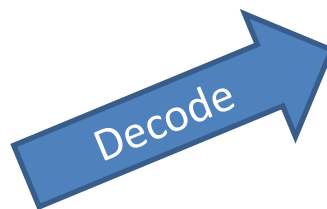
Original



“Reconstructed”

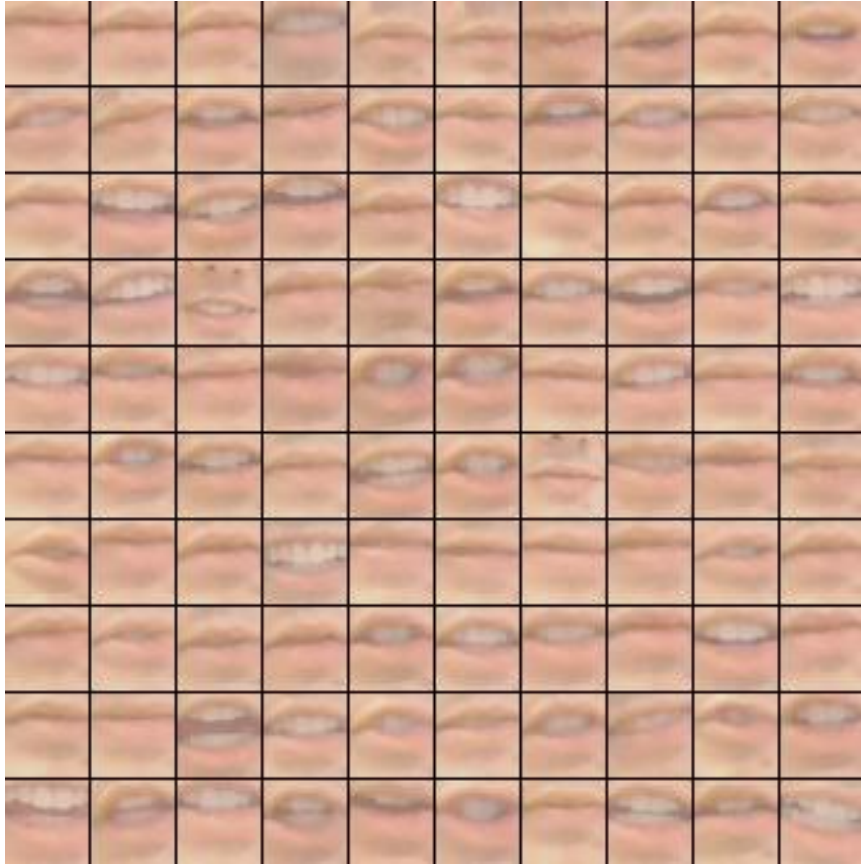


$$\begin{bmatrix} a_0 \\ \vdots \\ \vdots \\ \vdots \\ a_n \end{bmatrix}$$

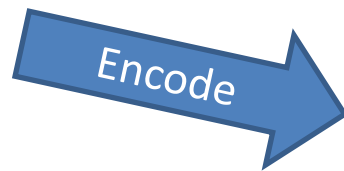
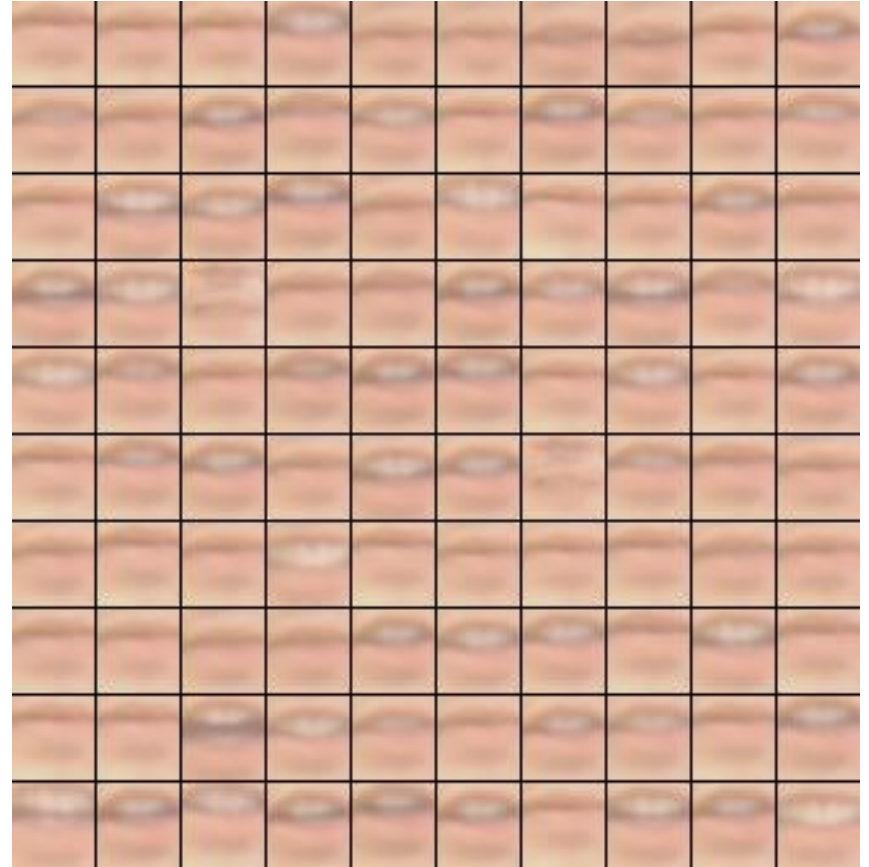


Training: 300epoch (6mil iter) later...

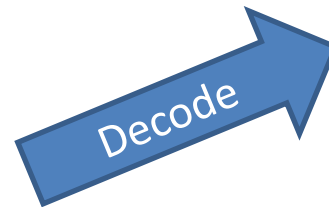
Original



“Reconstructed”



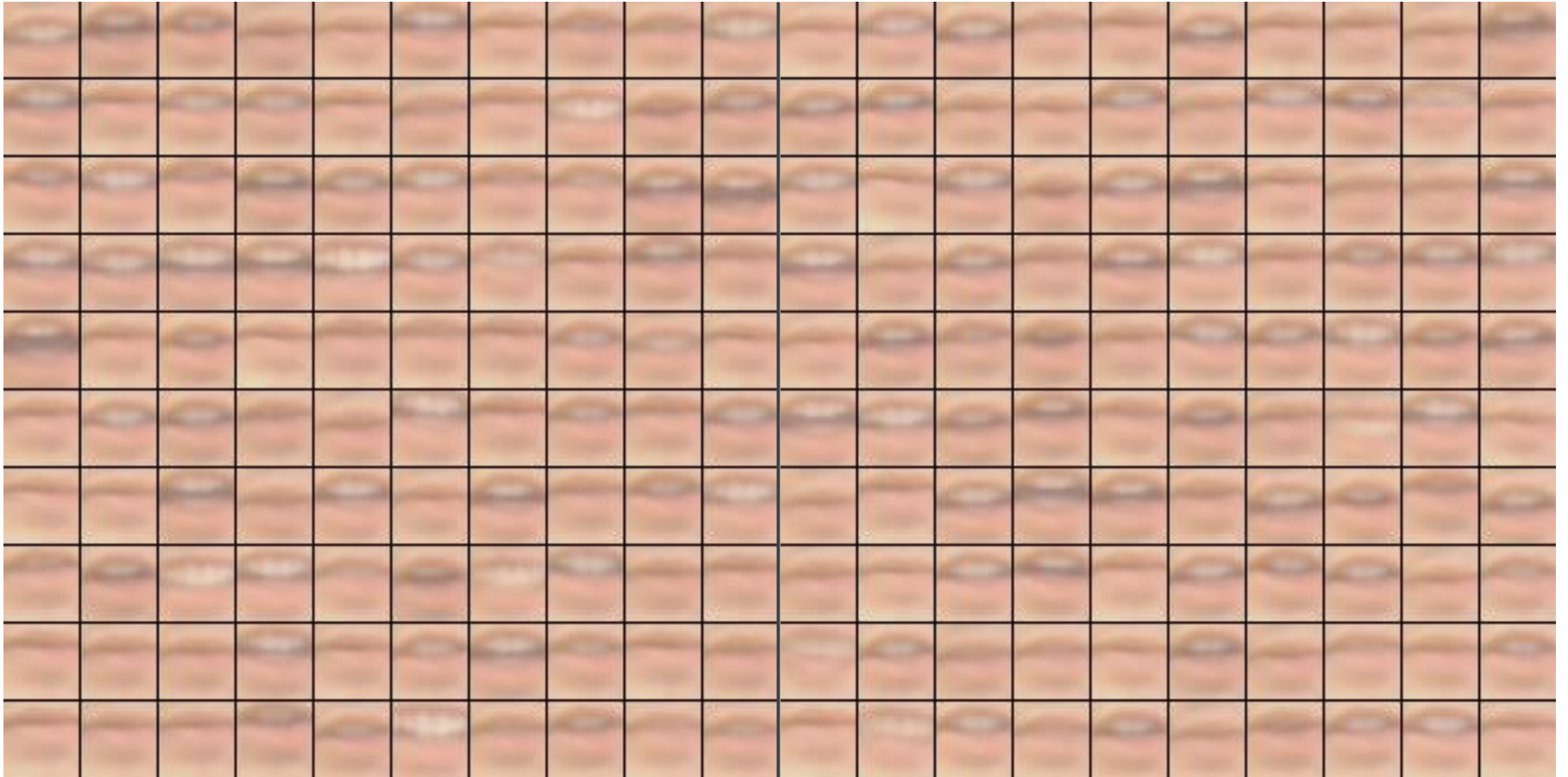
$$\begin{bmatrix} a_0 \\ \vdots \\ a_n \end{bmatrix}$$



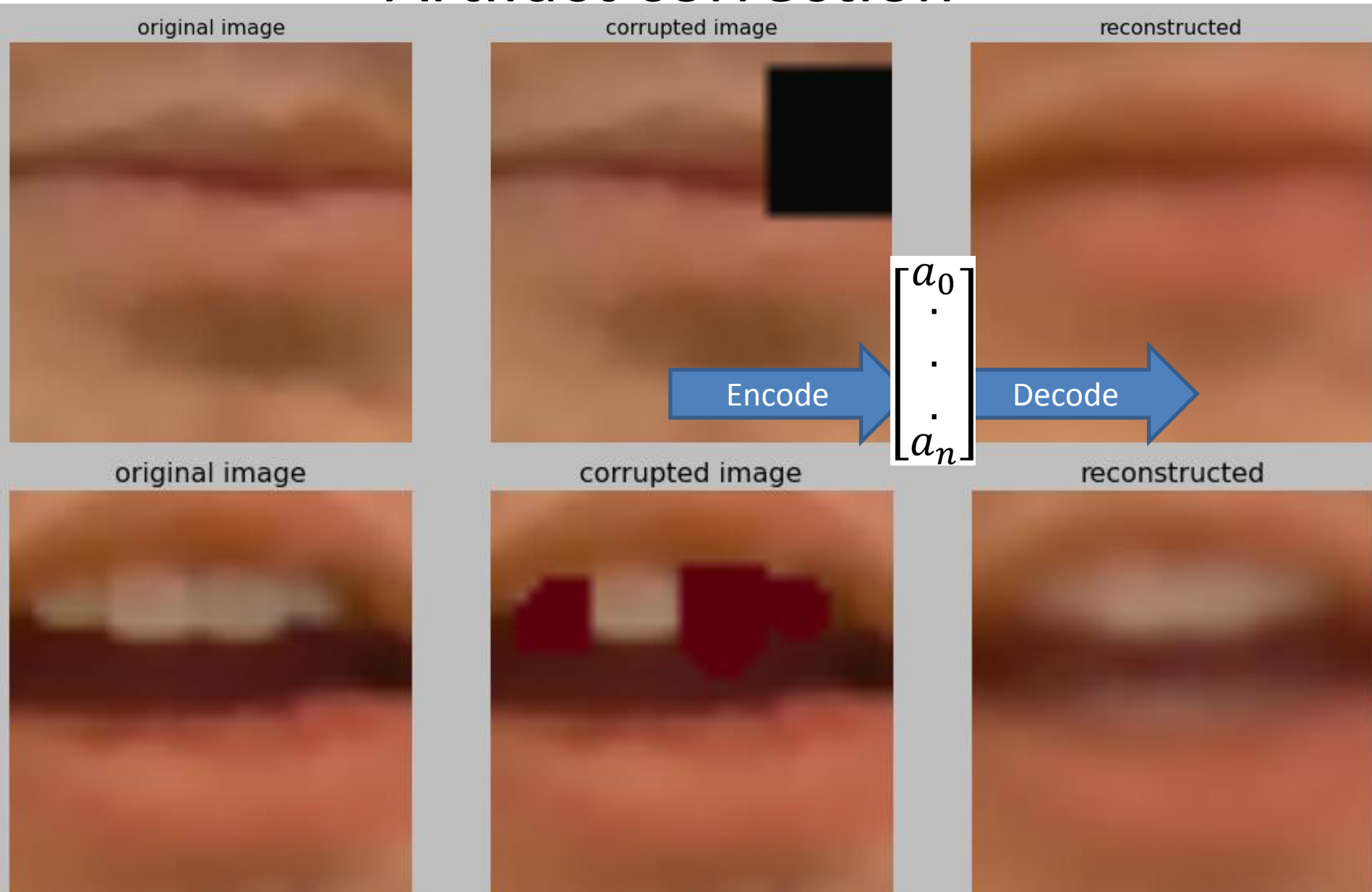
“Variational” Convolution Neural Net Autoencoder:

- Allows for generation of new images, which is nice

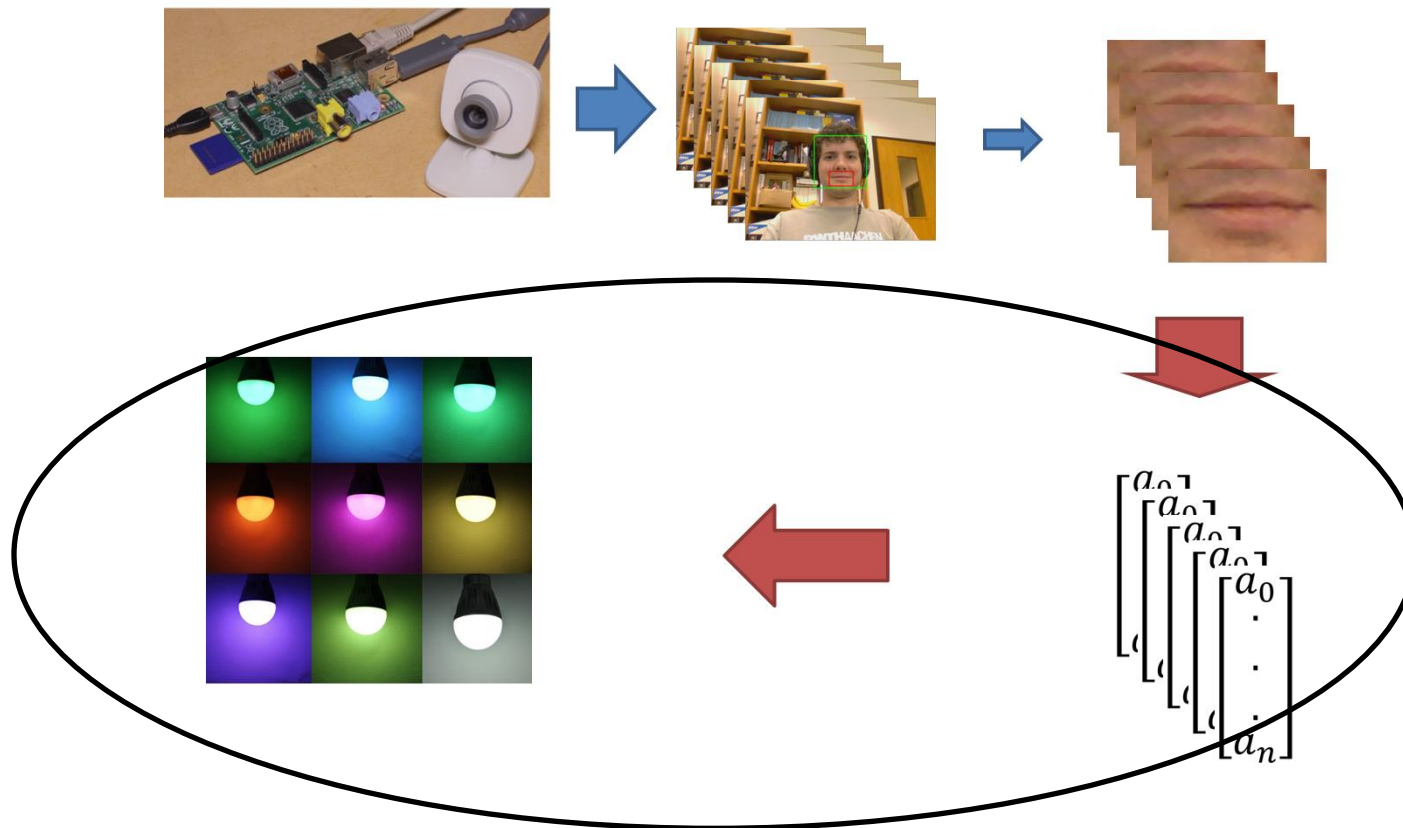
Algorithm Imagined or
“Hallucinated” images



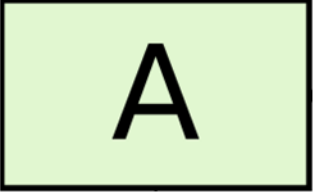
Just for fun: Image Inpainting and Artifact correction

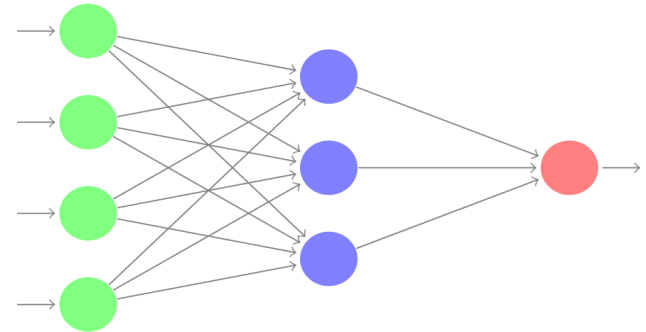


Step 3: Predict Color with vector sequence

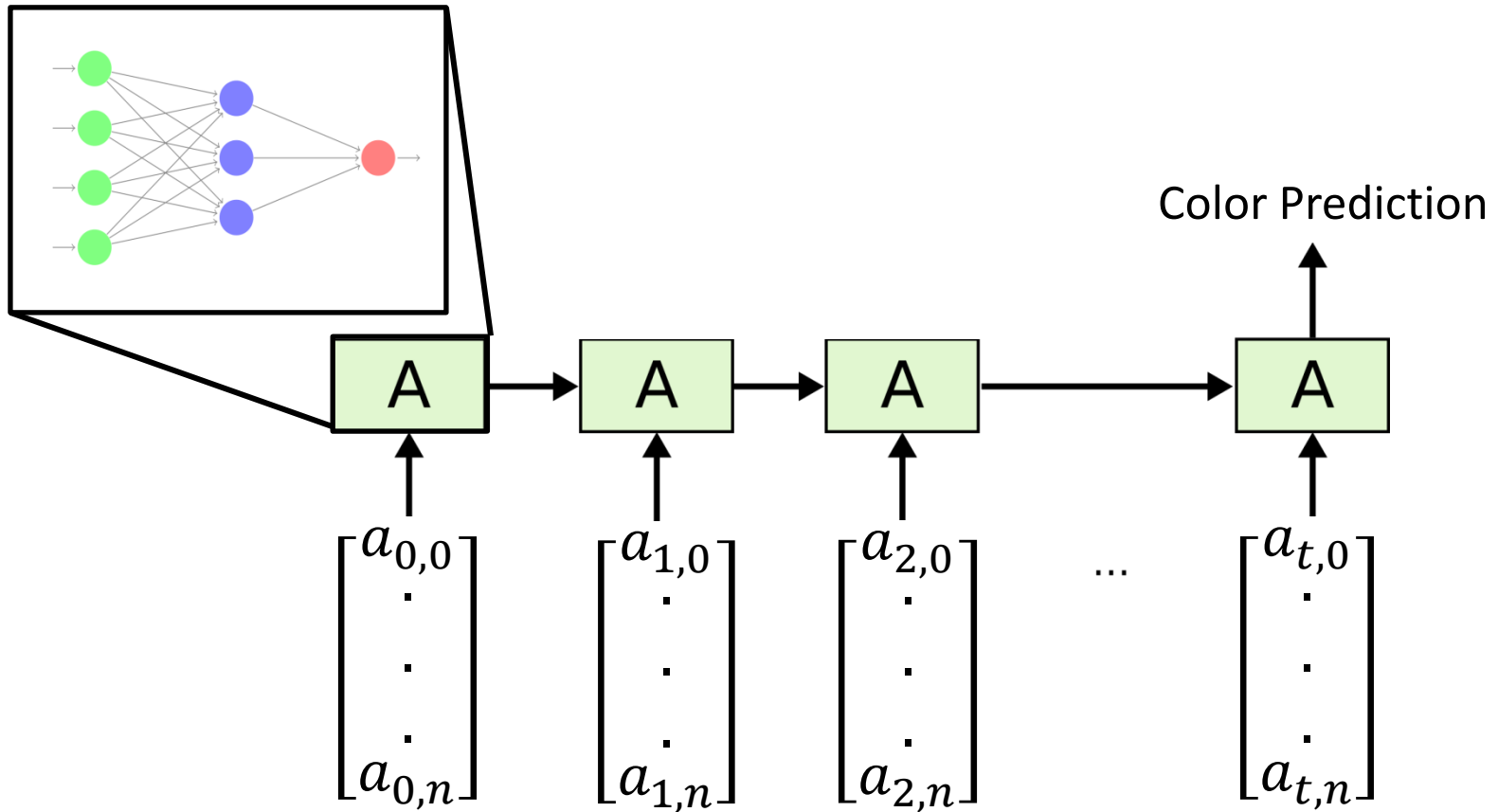


Recurrent Network

Let “” Represent a single layer neural network



Recurrent Network



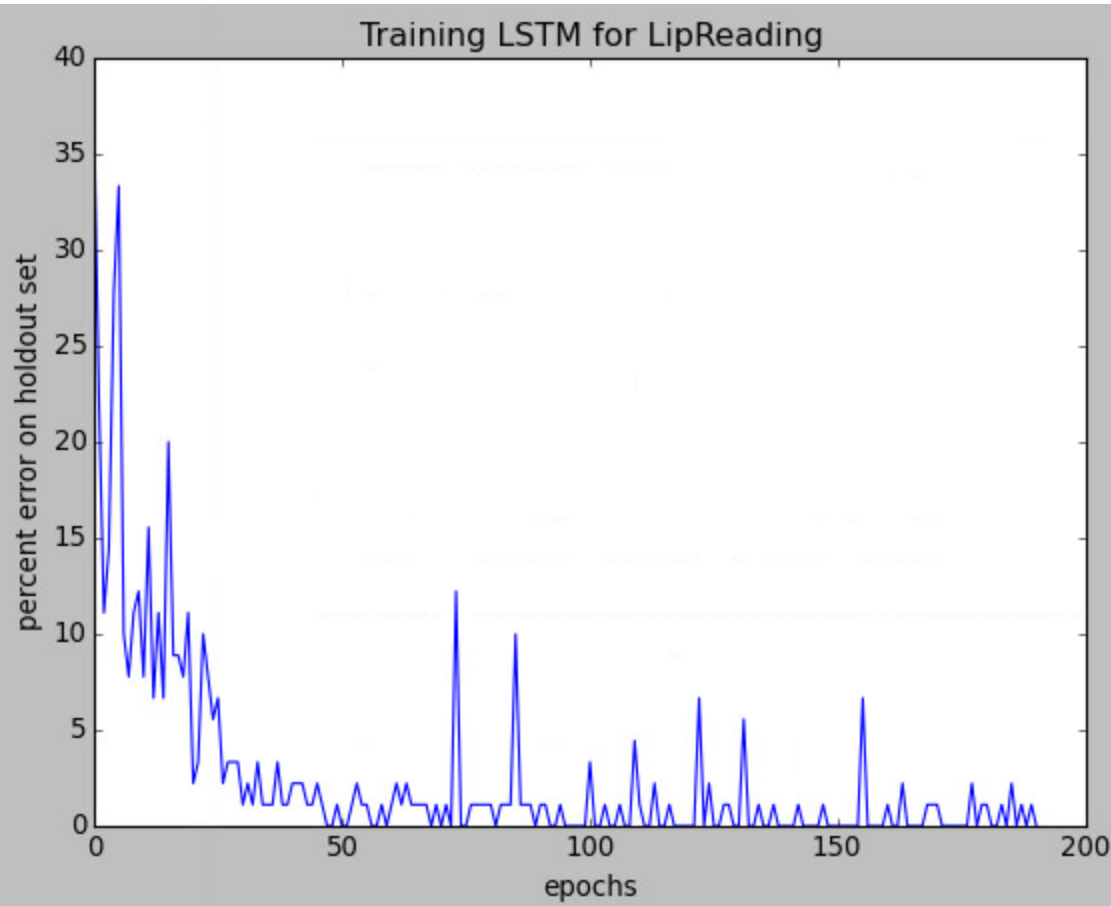
Recurrent Network

- Model motivation: pass on a “memory” vector of output of previous frames to influence how the next frame is interpreted
- Color outputs { red, purple, yellow }
 - Purple is similar to two syllables of red
 - Yellow could potentially be predicted by a single wide-mouth frame alone

Recurrent Network: training

- ~300 training examples of each color
 - 903 total
 - 16-25 time points each example (*variable!*)
 - 512 dim vector “code” for each frame
- A modest 2,099,200 variables
- Training takes ~1hr (gpu)
- Randomly choose 90%/10% split :
803 for training : 90 for testing

A Surprise! A Perfect Classifier



- Quickly achieve 0% error on both training and test set (all colors)
- Rare feat in computer vision
- Probably would not be flawless with 10x more examples