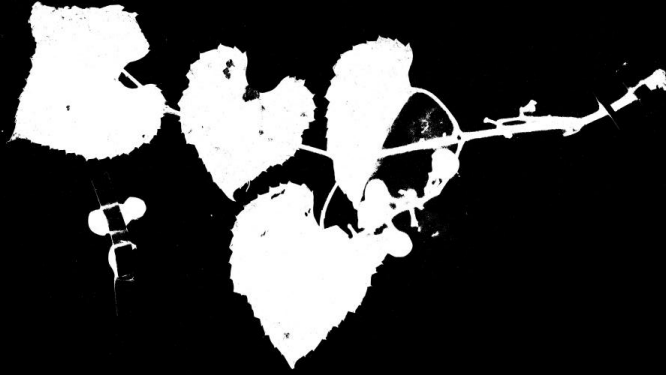


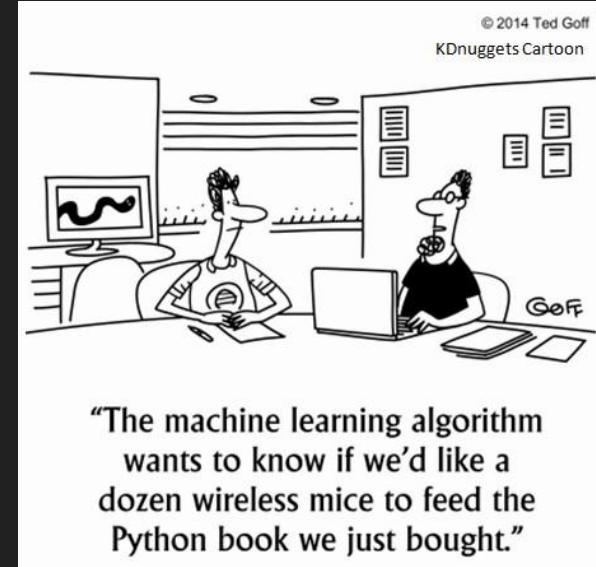
Deep Learning Image Clustering to Aid Species Delimitation Within the *Vitis* *arizonica* complex.



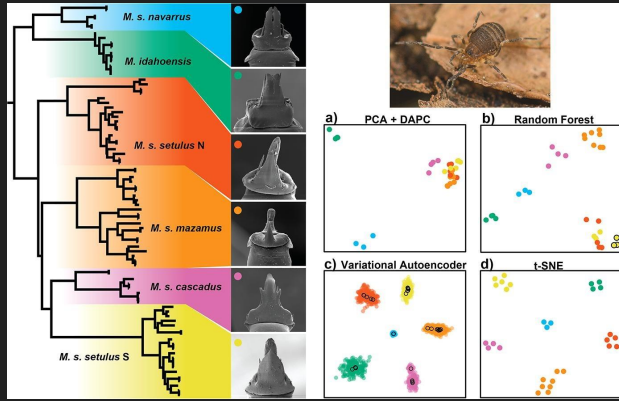
Stefano Fochesatto, Matthew Shavlik,
Steffi Ickert-Bond, Richard Hodel, Jun Wen

Outline

- Machine Learning and Deep Learning
- Research and Goals
- Our Data
- ML and DL Biases
- Preprocessing
- Deep Convolutional Embedded Clustering
- Results
- Further Work/Conclusions
- Plug for OSS Project

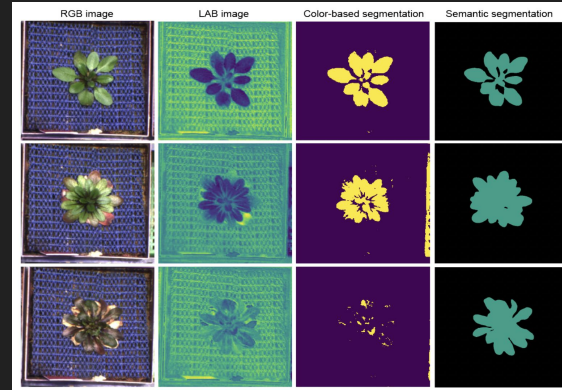


Machine Learning and Deep Learning



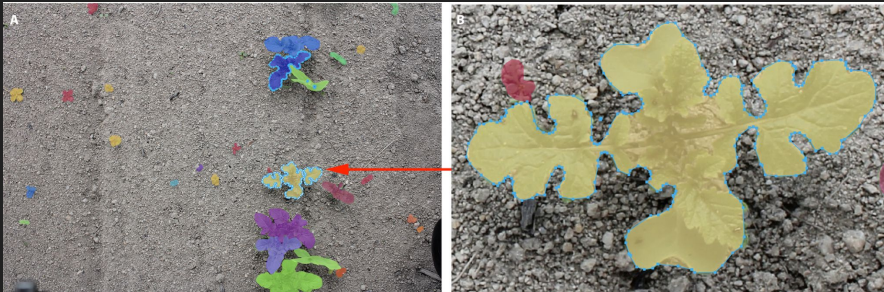
Shahan Derkarabetian et al. 2020

A demonstration of unsupervised machine learning in species delimitation



Hüther et al. 2020

araDEEPopsis: From images to phenotypic traits using deep transfer learning

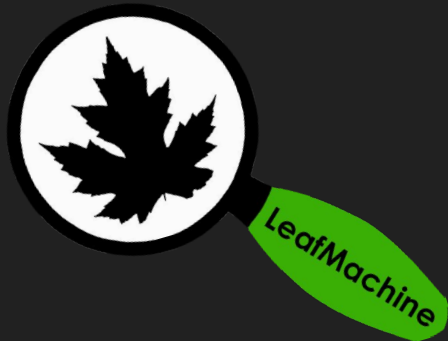


Champ et al. 2020

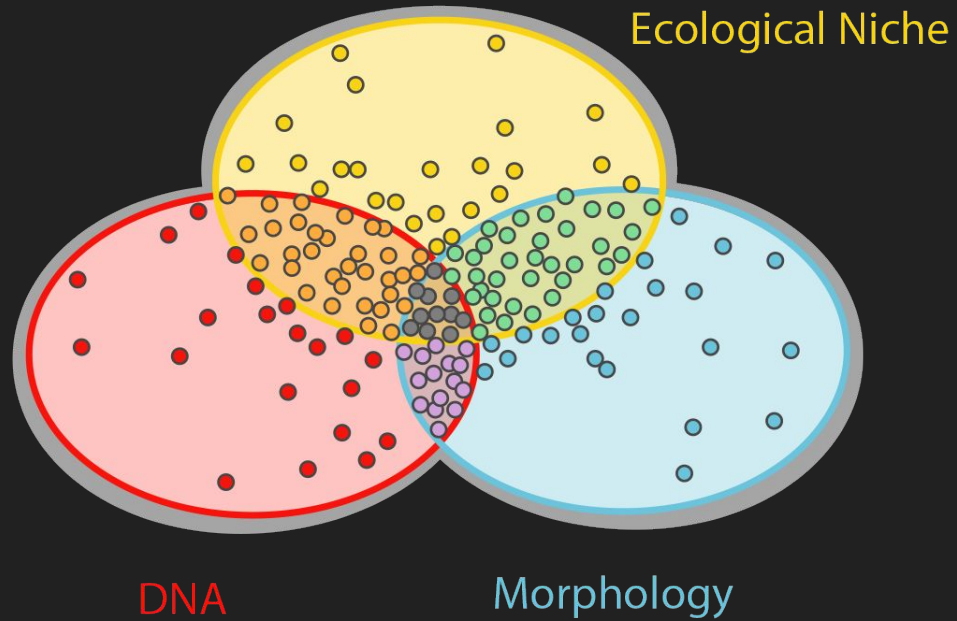
Instance segmentation for the fine detection of crop and weed plants by precision agricultural robots

Research Goals and Outcomes

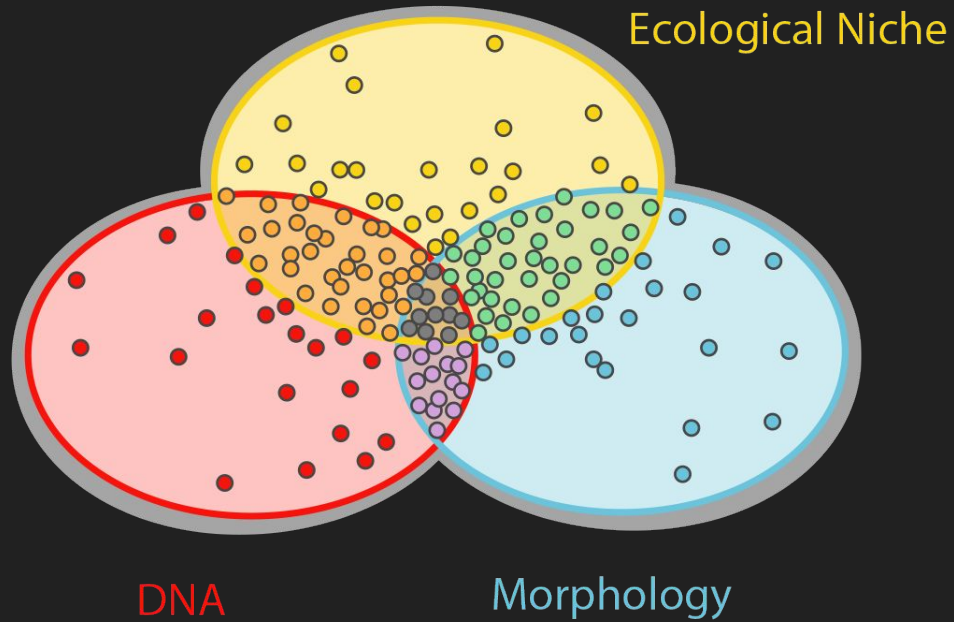
- Identify a deep learning workflow which can cluster herbarium sheet imagery in a way which signals species delimitation.
- Ideally clustering will be able to identify which specimen are most likely to return significant results from DNA sequencing
- Eventually incorporating tools like Leaf Machine to automate, and leverage large amounts of data from sources like iDigBio, GBIF



Species Delimitation



Species Delimitation



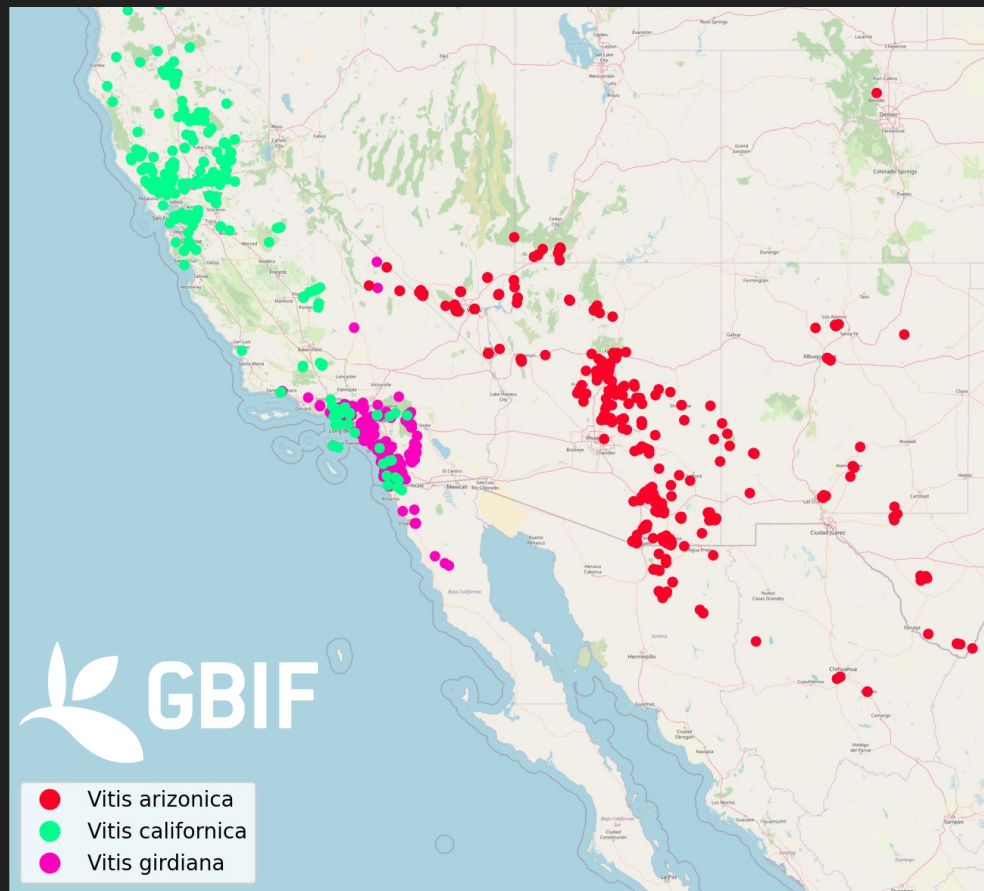
Vitis arizonica complex



Methods: Our Data

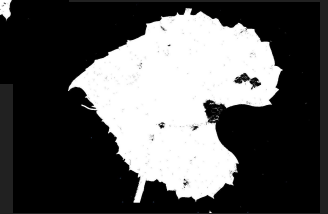
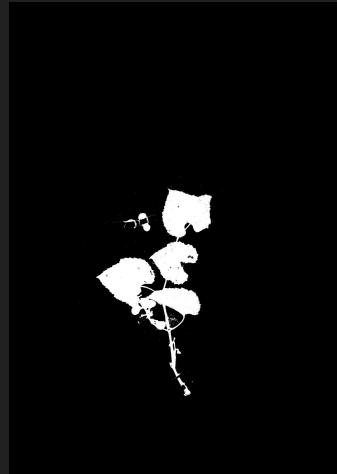
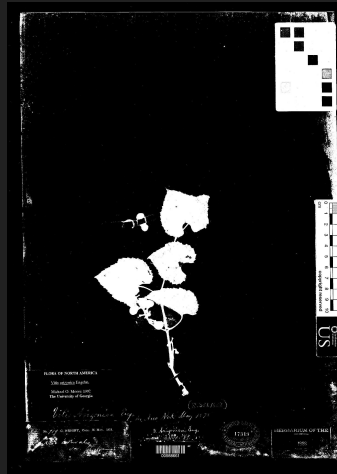
	SI	RSA	ASU	Total
<i>Vitis arizonica</i> Engelm	65	29	117	211
<i>Vitis girdiana</i> Munson	0	39	0	39
<i>Vitis californica</i> Benth	0	27	0	27
Total	65	95	117	277

Species Distribution Map



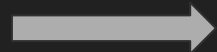
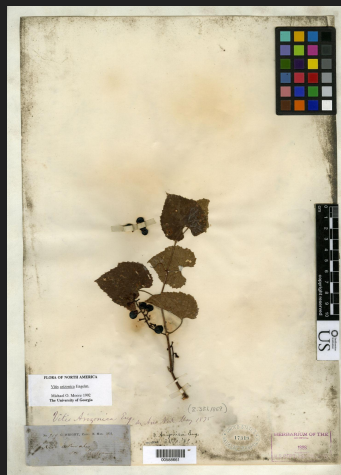
Methods: Preprocessing

- The goal is to remove any biasing information
 - Camera settings, lighting conditions, labels are all features that obstruct morphological signal.
 - Clustering on the segmentation masks captures the leaf morphological traits of the specimen, while removing biasing information.

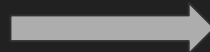
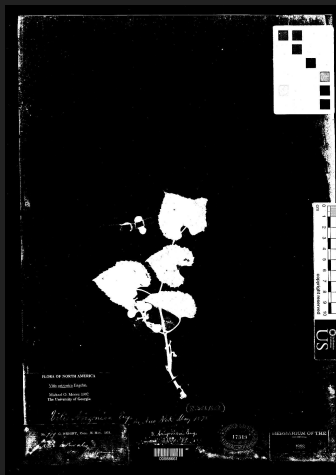


Methods: Preprocessing

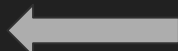
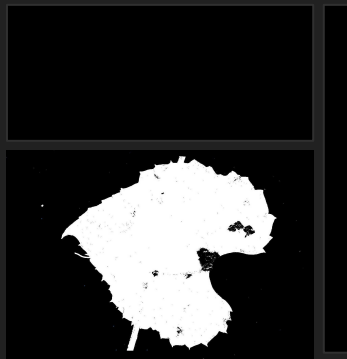
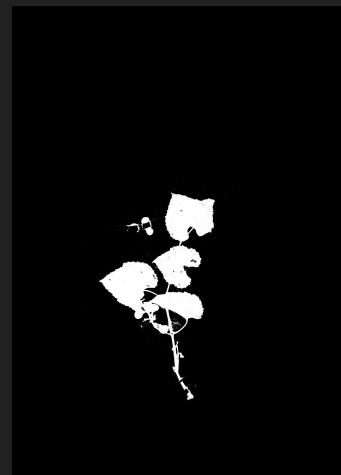
- Tools are in production for generating segmentation masks automatically (Leaf Machine 2)
- Masks for this project were generated adapting a workflow from *Generating segmentation masks of herbarium specimens and a data set for training segmentation models using deep learning* (White and Dikow, et al 2020)



Global Otsu
Binarization



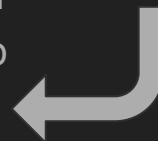
PhotoShop Batch
Processing



Border Padding + Final Resize
(retains relative size)



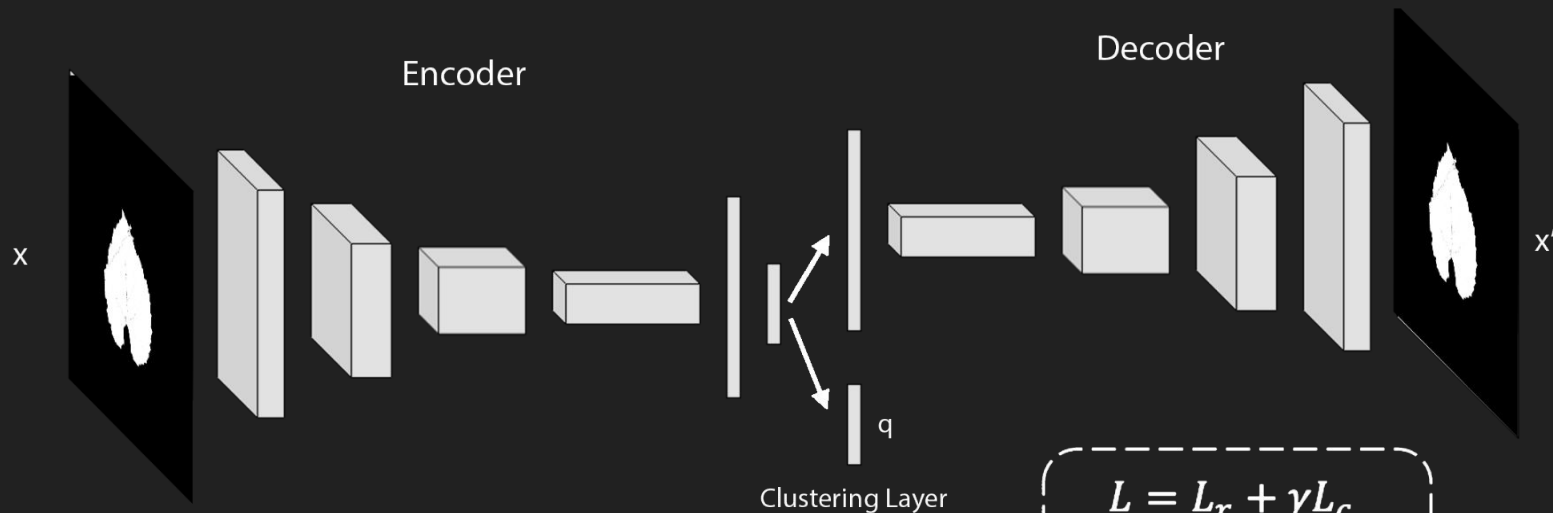
Resize + Leaf
Instance Crop



Methods: Image Clustering

- The Deep Learning Algorithm that we are using for clustering is called Deep Convolutional Embedded Clustering (Guo, Liu, et. al. 2017).
 - (pretraining) A Deep Convolutional AutoEncoder is trained on the data.
 - K-means is used in the latent space to identify n cluster centers.
 - (clustering) A Clustering layer is then incorporated alongside the latent space which maps embedded points in the latent space to a Student's t-distribution with n-dimensions.
 - KL Divergence is added to the loss function.
 - Cluster centers are updated alongside AE weights.
- Autoencoders preserve local structure of data in the latent space.
- Convolutional Layers learn image features.

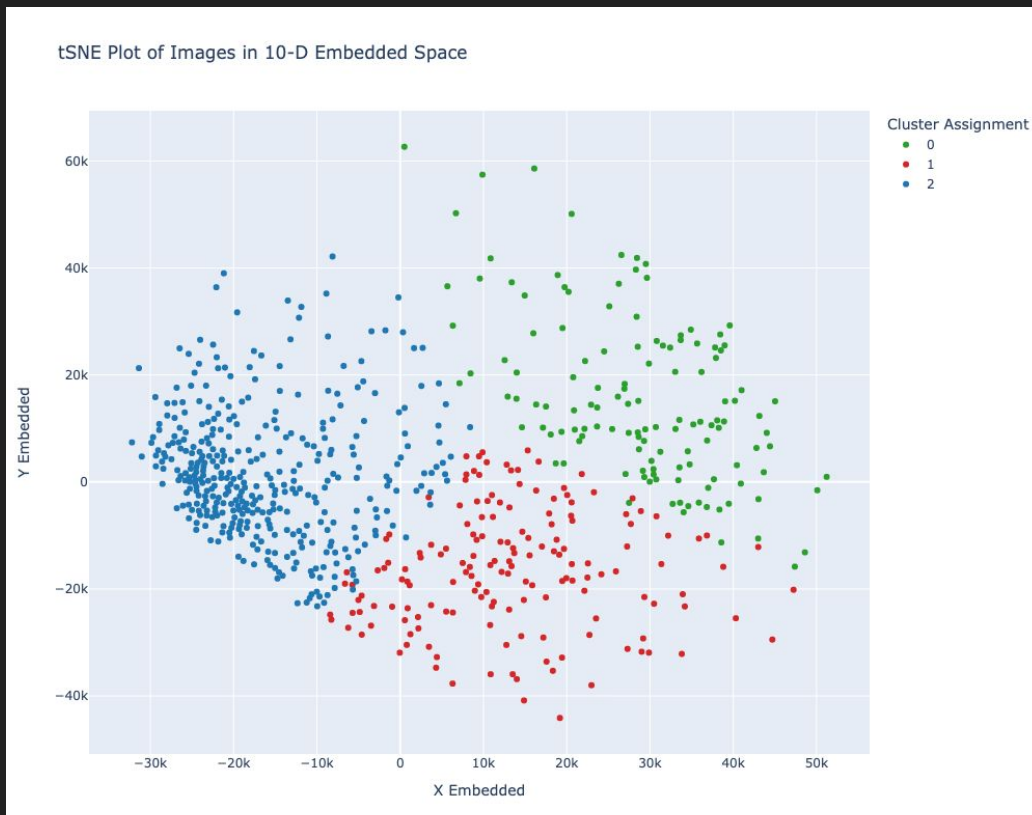
Methods: Image Clustering



Deep Clustering with Convolutional Autoencoders (Xifeng Guo, Xinwang Liu, et.al 2017)

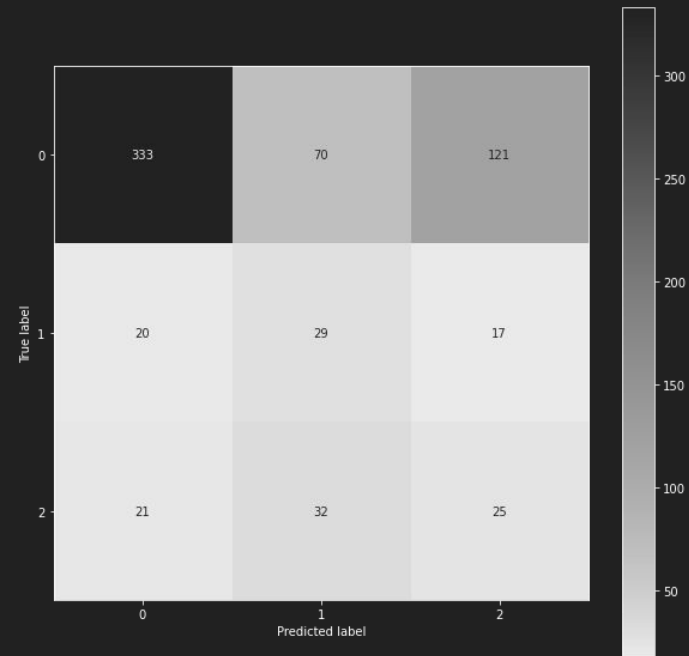
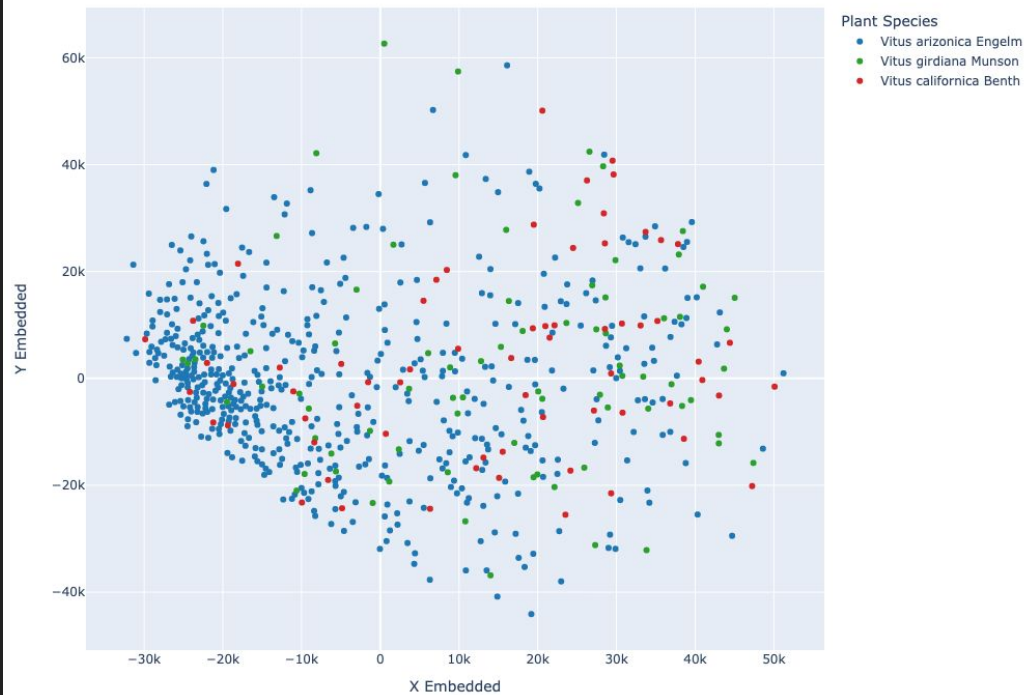
$$L = L_r + \gamma L_c$$
$$L_r = \|x - x'\|_2^2$$
$$L_c = KL(p||q)$$

Results



Results

tSNE Plot of Images in 10-D Embedded Space



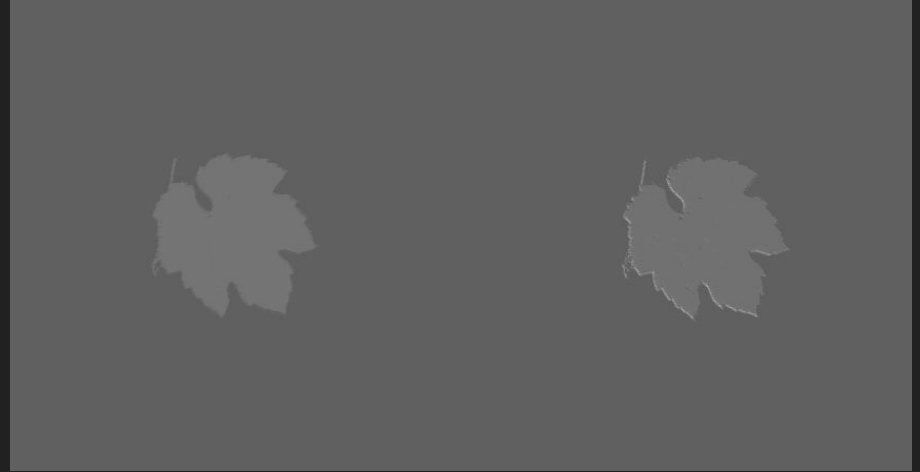
Results

tsNE Plot of Images in 10-D Embedded Space



Results: Feature Maps

- Feature maps show promise
 - Extracting edge morphology
 - Extracting size, area, and texture

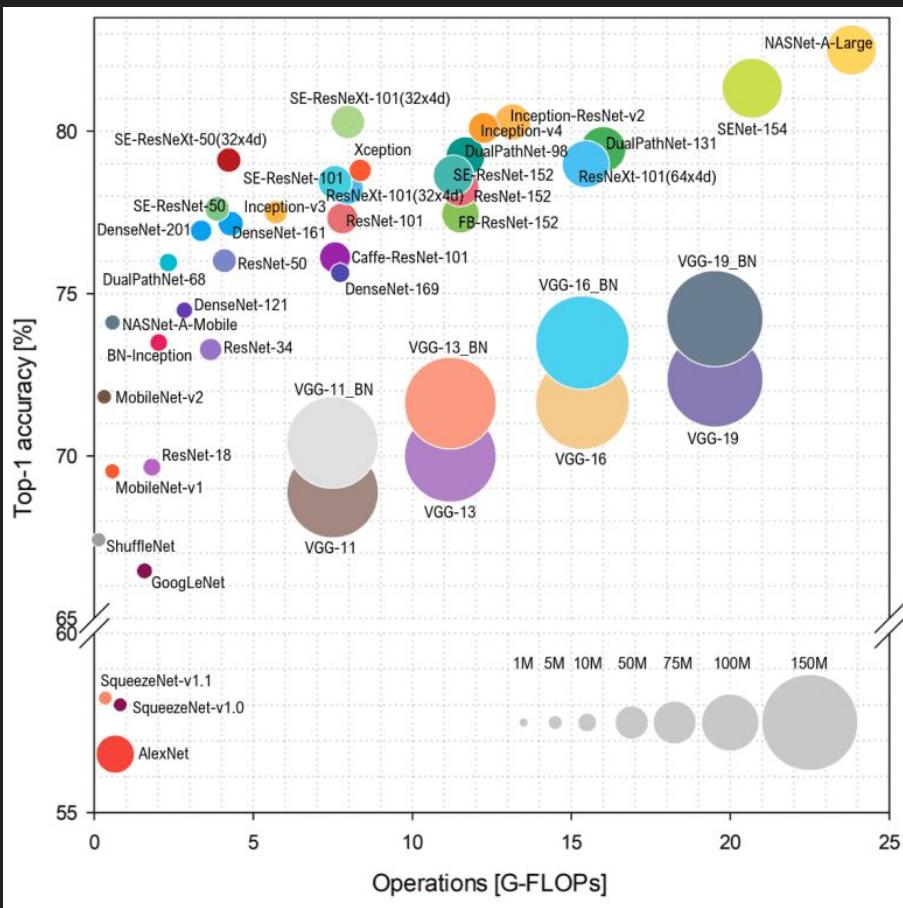


Further Work

More robust network architecture.

- Used poor performing traditional CNN architecture.
- ResNet/VGG are used in segmentation software like LeafMachine.

➔ **Transfer Learning opportunity**



Bianco et al. 2018

Benchmark Analysis of Representative Deep Neural Network Architectures

Further Work

Alternative clustering methodology.

- ClusterGAN
- DAIC (Deep Adaptive Image Clustering)
- ASPC-DA (Adaptive Self-Paced Deep Clustering with Data Augmentation)

Morphology retaining data augmentation for pretraining step.

- Rotations
- Translations

Plug for OSS Project

No high level library for DL Image Clustering.

Current Implementation Workflow

- Read paper
- Hope and pray for author's github link works 🙏
- Refactor (often times) depreciated code/Integration Hell

Goal Workflow

- Read paper
- Scikit-learn esque implementation

Example: Segmentation Models



Segmentation
Models

Acknowledgments



Steffi Ickert-Bond

Matt Shavlik

Cam Webb



Smithsonian

Richie Hodel

Jun Wen

Questions?