# Automating tropical pollen counts using convolutional neural nets: from image acquisition to identification

Derek S. Haselhorst[1], Shu Kong[2], Charless Fowlkes[2], J. Enrique Moreno[3], David K. Tcheng[4] and Surangi W. Punyasena[1]

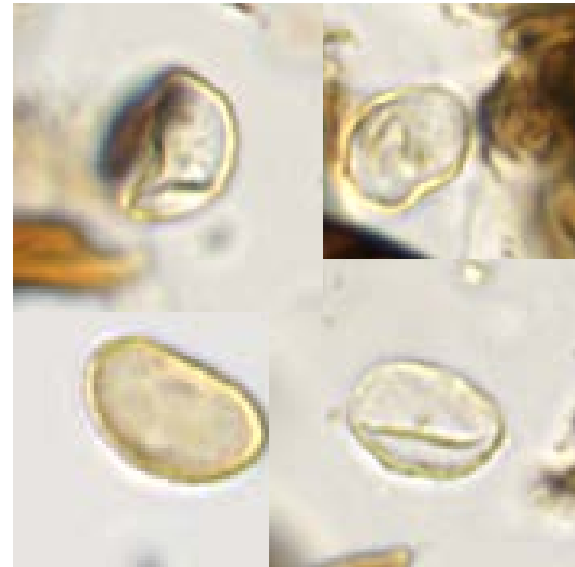[1]*University of Illinois at Urbana-Champaign*
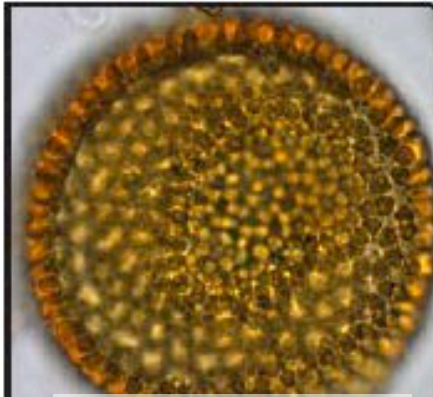[2]*University of California at Irvine*
[3]*Smithsonian Tropical Research Institute, Panama City, Republic of Panama*
[4]*National Center for Supercomputing Applications,*
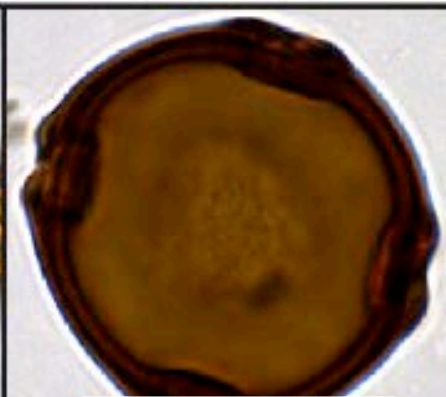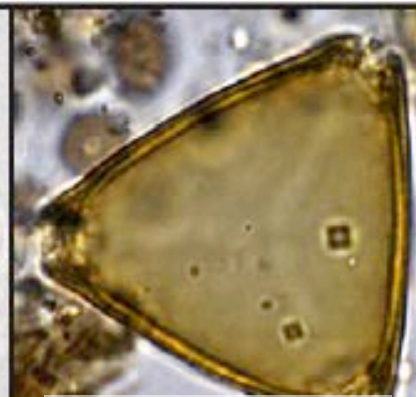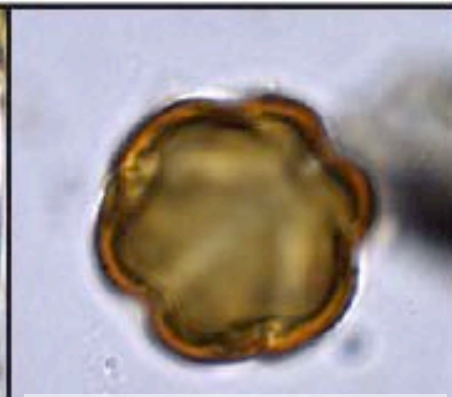*University of Illinois at Urbana-Champaign*

*Cecropia* (Urticaceae)

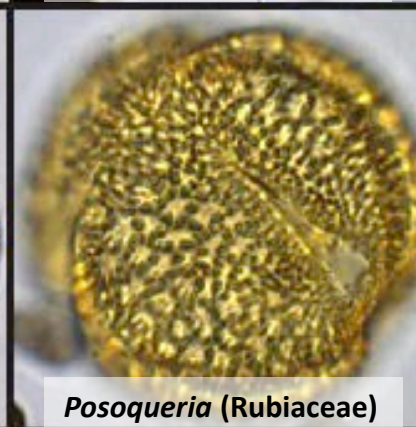*Croton* (Euphorbiaceae)    *Trichilia* (Meliaceae)    *Serjania* (Sapindaceae)    *Combretum* (Combretaceae)
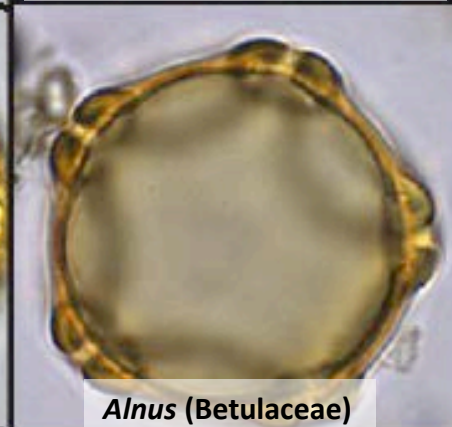
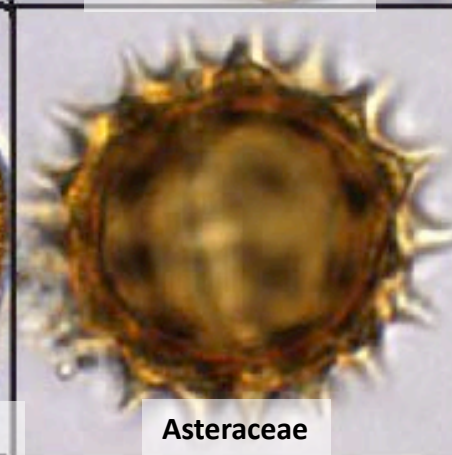*Quararibea* (Malvaceae)    *Pseudobombax* (Malvaceae)    *Posoqueria* (Rubiaceae)    *Alnus* (Betulaceae)

*Quassia* (Simaroubaceae)    *Dalechampia* (Euphorbiaceae)    *Anacardium* (Anacardiaceae)    Asteraceae

Haselhorst , Moreno and Punyasena, 2013, *PLoS ONE*
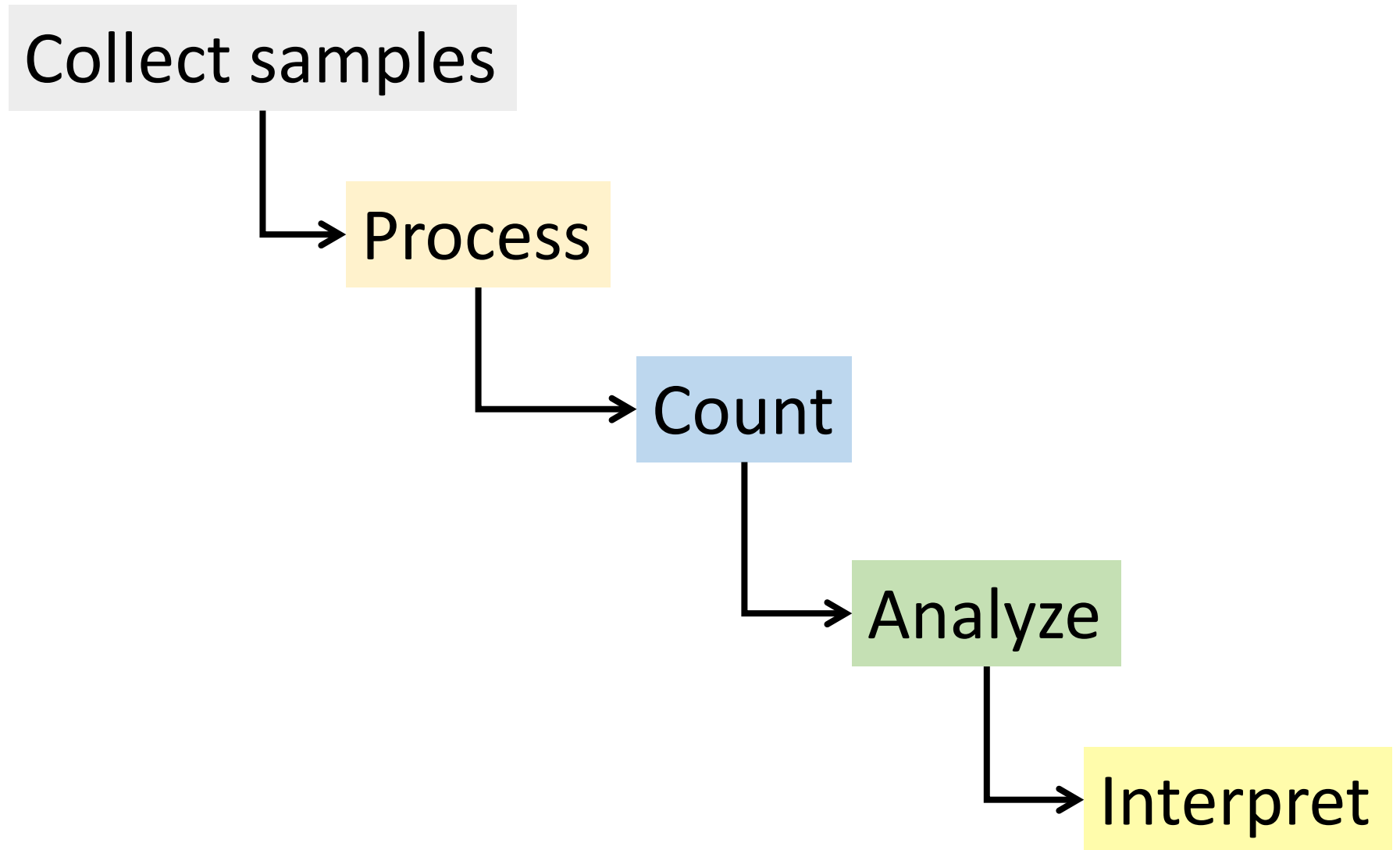
# POLLEN AS "BIG DATA"

~470 million years of plant history

Billions of potential specimens

~Continuous deposition across a range of environments

# REIMAGINING THE WORKFLOW

Collect samples

Process

Count

Analyze

Interpret

# IS AUTOMATION THE ANSWER?

**Quantity**: increase the throughput of pollen analysis

**Reproducibility**: improve the consistency and accuracy of pollen identifications

**Resolution**: produce repeatable recognition of *species* from pollen for more precise biome reconstructions
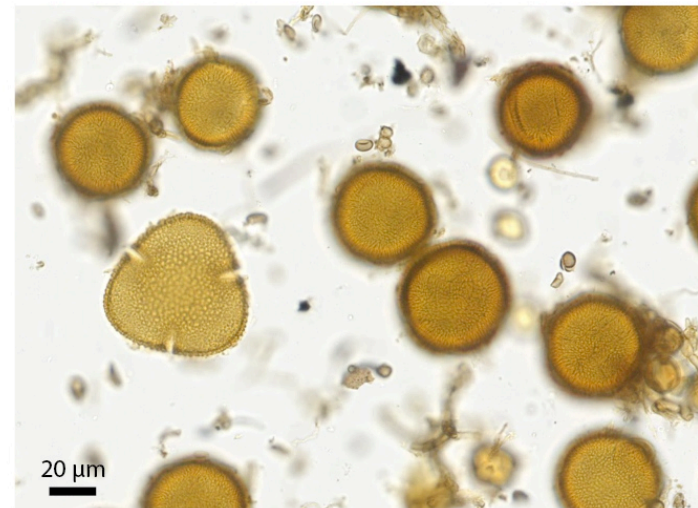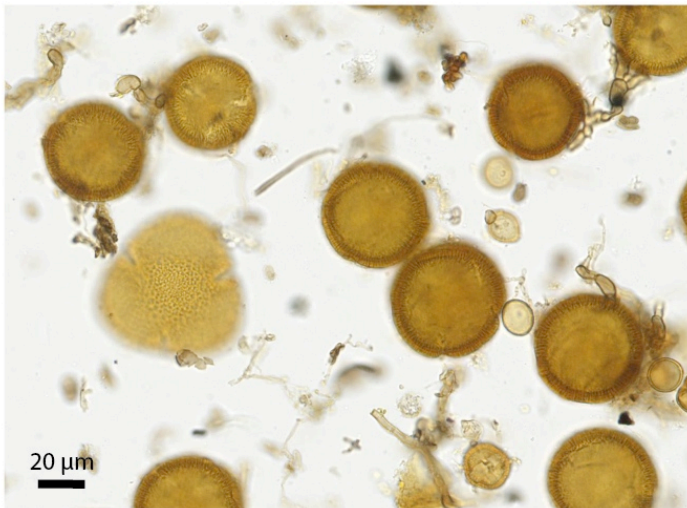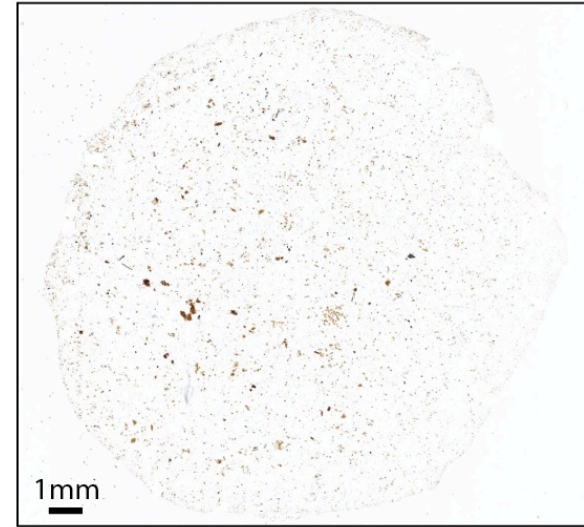
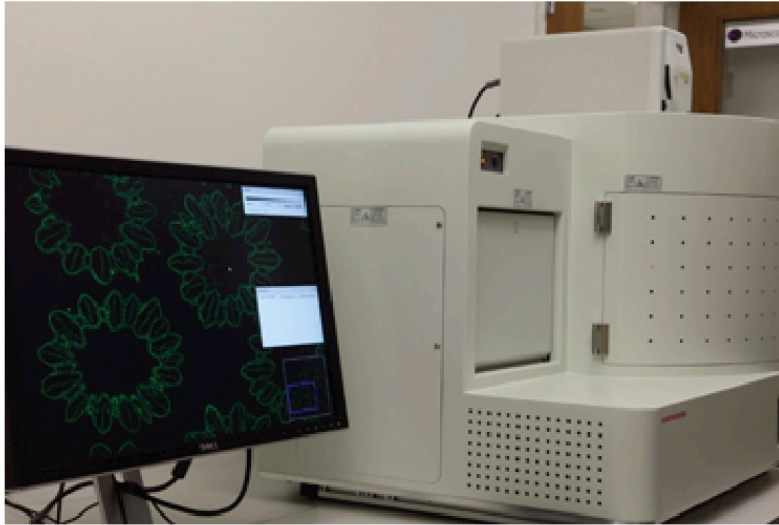# TRAINING ON HYPERDIVERSE SAMPLES

- A 15-year pollen rain record from Barro Colorado Island, Panama
  - Obtained from a series of 20, evenly spaced pollen traps along two parallel transect in the 50 ha CTFS plot

- A 10-year pollen record from the Lutz weather tower
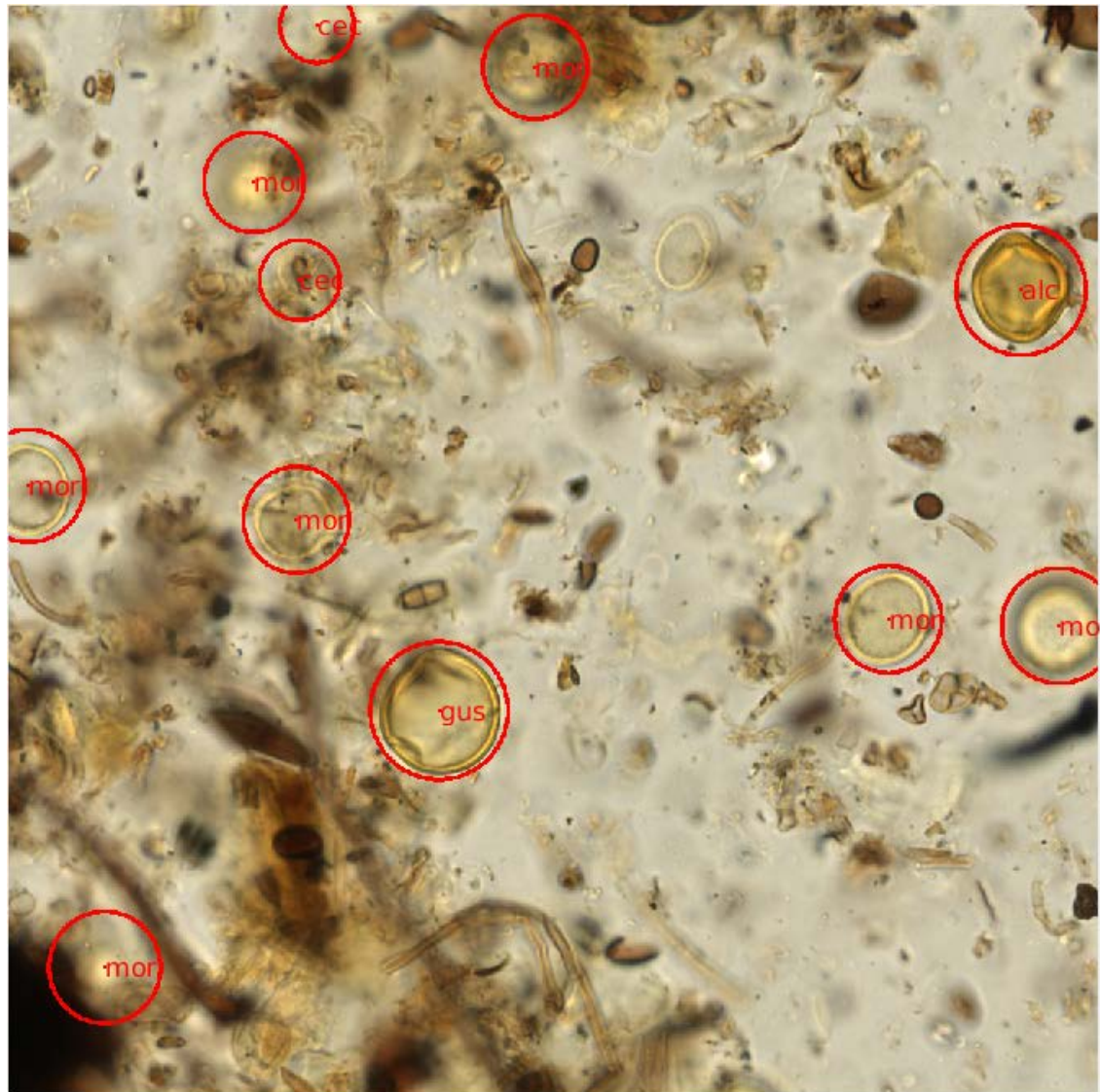  - Images to be analyzed this summer

- ~ 130 pollen morphotypes



© Christian Ziegler

Photo: STRI
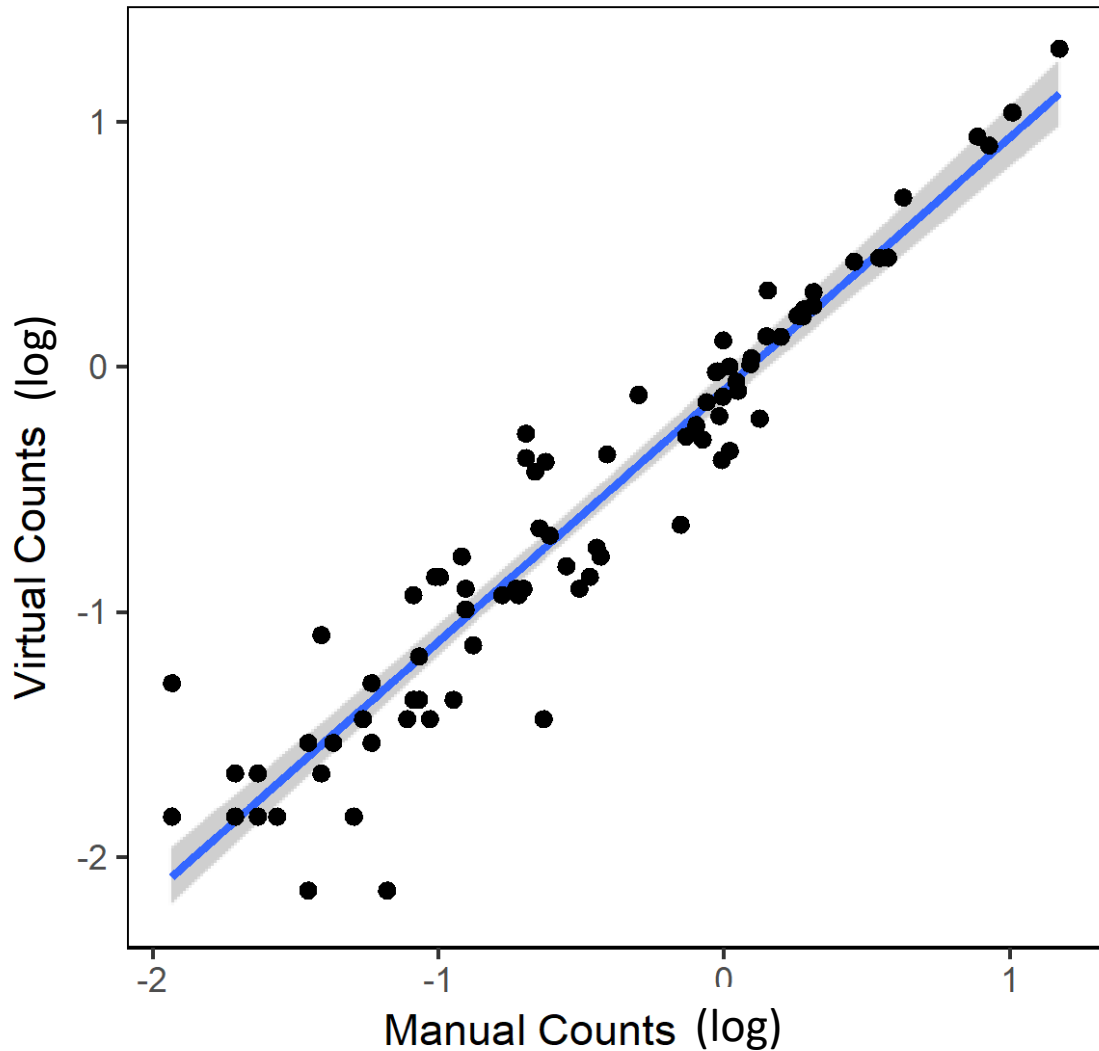
# IMAGING POLLEN SLIDES FOR COUNTING



1mm

400x, 0.23 μm/pixel

20 μm

20 μm

One sample (41 @ 1 μm axial planes) = ~400 GB

# VIRTUAL MICROSCOPE POLLEN IDENTIFICATIONS



- PNG images subsampled using a java script

- Slide images read from Matlab script

- Image metadata recorded for each individual pollen grain
    - Slide ID
    - Pollen coordinate
    - Pollen radius
    - 3-letter ID
    - Confidence level (0-9)

- Images and metadata then shared with UC-Irvine Computer Vision Collaborators

# COMPARING VIRTUAL AND LIGHT MICROSCOPY COUNTS



- Testing the fidelity of the virtual microscope using the 10-year Lutz tower record

- Can pollen can be identified at the same taxonomic resolution and frequency?

- $R^2 = 0.97$

- Observed differences are likely reflective of differences in counting strategy:
  Manual: slide transects
  Virtual: randomized images

(log)

# TRAINING CONVOLUTIONAL NEURAL NETS (CNN) FOR POLLEN IDENTIFICATIONS FROM ANNOTATED IMAGE DATA

- Images were randomly split into training and testing sets

- Annotated training (ground truth) pollen image examples included the pollen id, location coordinate, and pollen grain radius

- CNN searches each image for patterns corresponding to each pollen id morphology

- Non-maximum suppression was used to identify pollen grains according to pollen ornamentation

**Human**

# SIMULTANEOUS POLLEN SEGMENTATION AND IDENTIFICATION

**Machine**

**Human**



- 48-way classification matrices were constructed using the 47 most abundant pollen types and an additional category called "reject" comprised of pollen types not included in the 47 most abundant

# SIMULTANEOUS POLLEN SEGMENTATION AND IDENTIFICATION

**Machine**

**Machine**



- 48-way classification matrices were constructed using the 47 most abundant pollen types and an additional category called "reject" comprised of pollen types not included in the 47 most abundant

# COMPLETE AUTOMATION: CONFUSION MATRICES

**~70% accurate on full 47 pollen type training set, 87.25% on 25 most accurate types**



confusion matrix on test set (acc=87.25%)

# COMPLETE AUTOMATION: CONFUSION MATRICES

**>90% accuracy**



confusion matrix on test set (acc=87.25%)

# COMPLETE AUTOMATION: CONFUSION MATRICES

**>90% accuracy**

# COMPLETE AUTOMATION: CONFUSION MATRICES

# COMPLETE AUTOMATION: CONFUSION MATRICES

**<50% accuracy**

# THE PAIRWISE COMPARISONS IN THE CONFUSION MATRIX SHOWING THE MOST DISAGREEMENT ARE MORPHOLOGICALLY VERY SIMILAR

*Ficus* (**66% accuracy)** was misclassified **29%** of the time as *Brosinum*-type



*Ficus* (Moraceae)

*Brosinum*-type(Moraceae)

# THE PAIRWISE COMPARISONS IN THE CONFUSION MATRIX SHOWING THE MOST DISAGREEMENT ARE MORPHOLOGICALLY VERY SIMILAR

Trema (**47% accuracy**) was misclassified **50%** of the times as *Brosinum*-type



*Trema* (Ulmaceae)

*Brosinum*-type(Moraceae)
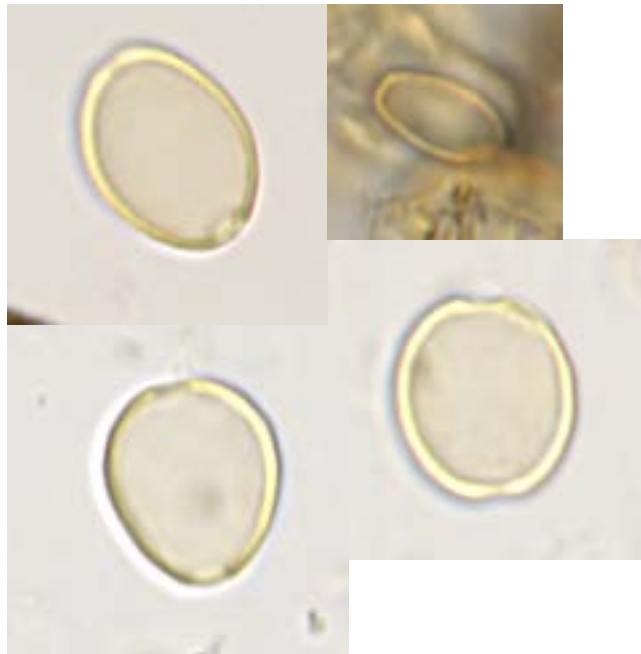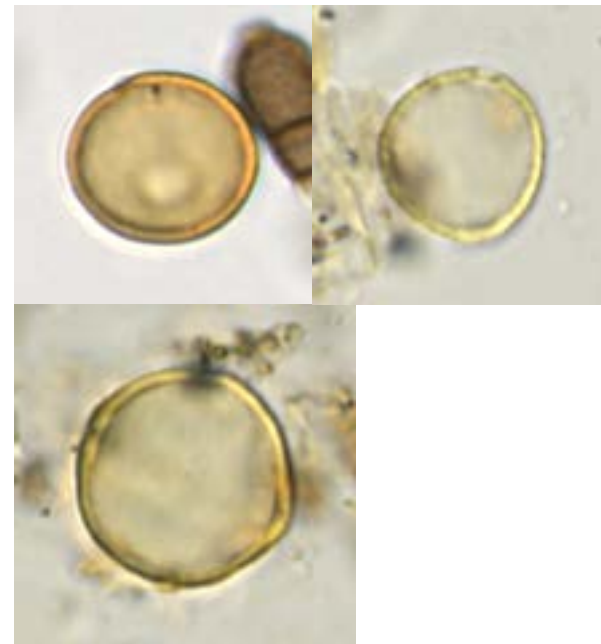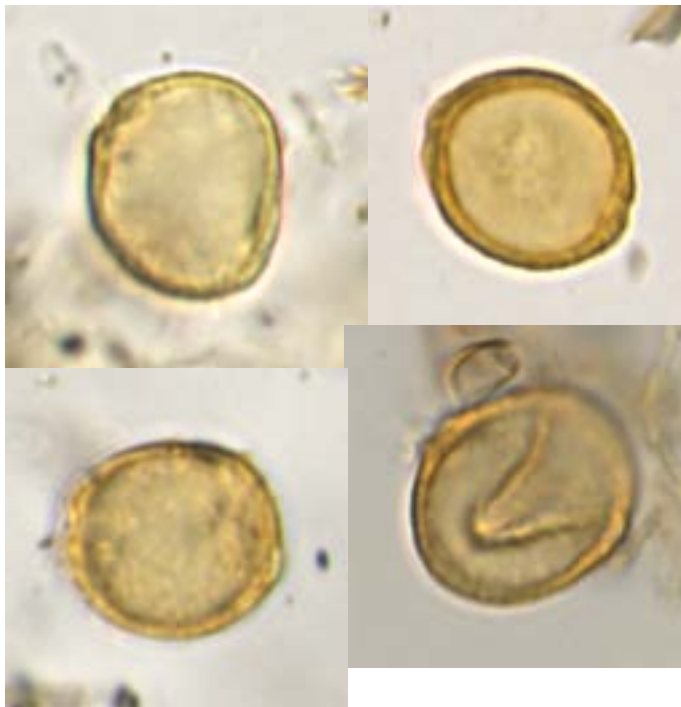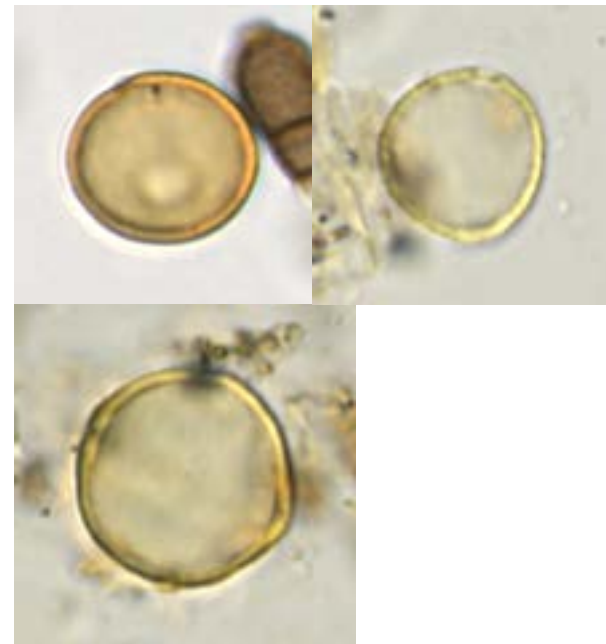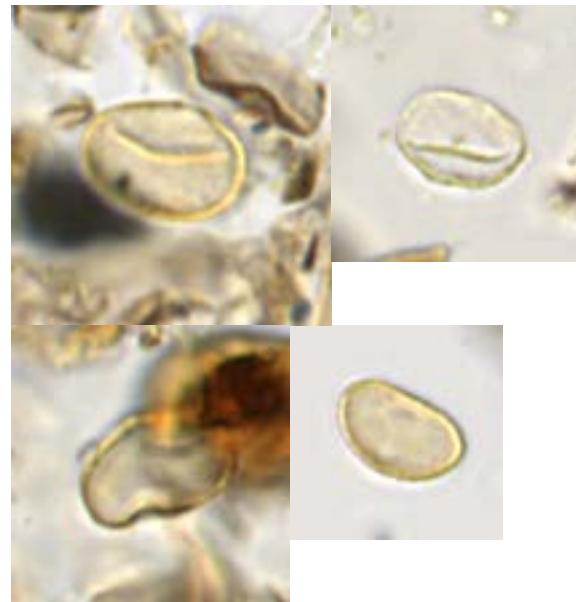
# THE PAIRWISE COMPARISONS IN THE CONFUSION MATRIX SHOWING THE MOST DISAGREEMENT ARE MORPHOLOGICALLY VERY SIMILAR

Piperaceae (**27% accuracy**) was misclassified **57%** of the times as *Cecropia*



Piperaceae



*Cecropia* (Urticaceae)

# CONCLUSIONS

- Overall, our results are very promising given this difficult classification problem

- Our preliminary results show that automated segmentation and classification models can distinguish pollen types from hyper-diverse samples

- Model performance is poorest on pollen types that are morphologically very similar

- Predicted outputs should improve as the neural nets are trained on more tagged examples

# FUTURE DIRECTIONS

- The same system can be implemented using training data from herbarium and reference material collections

- Apply the methodology to fossil records

- Create a collaborative pollen identification database that harnesses the expertise of multiple palynologists

- Expand to other proxies
  - Diatoms
  - Phytoliths
  - Cuticles

# ACKNOWLEDGEMENTS