# Inference Graphs for CNN Interpretation

Yael Konforti, Alon Shpigler, Boaz Lerner, Aharon Bar Hillel

Ben-Gurion University of the Negev, Beer Sheva, Israel

{yaelkonf, alonshp}@post.bgu.ac.il; {boaz, barhille}@bgu.ac.il

Code: https://github.com/yaelkon/GMM-CNN

# Motivation for CNN interpretation
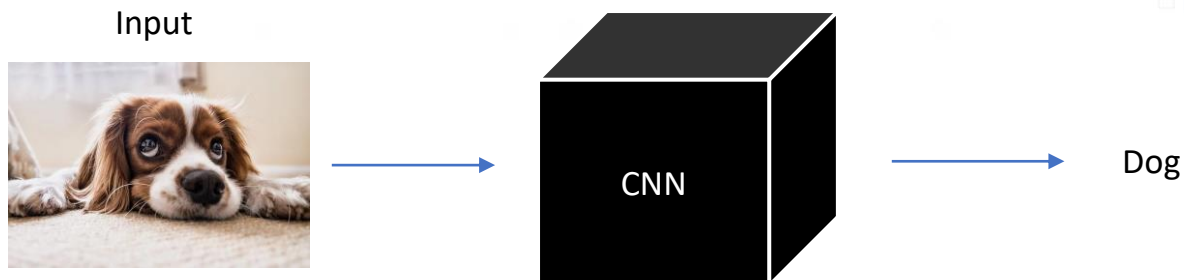
The current state –

      State-of-the-art results for a variety of computer domains.

The problem –

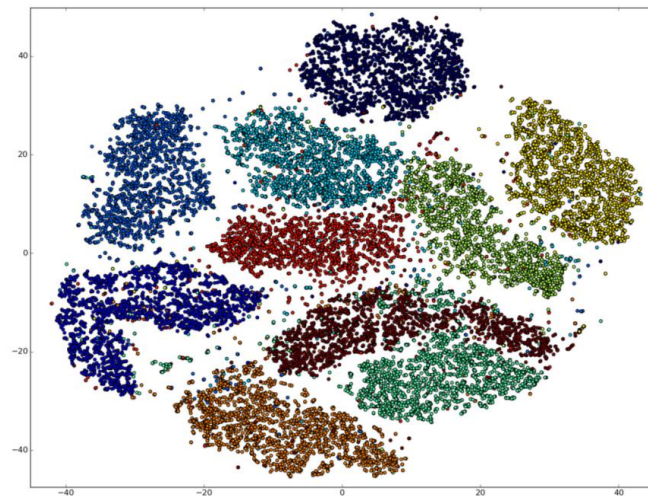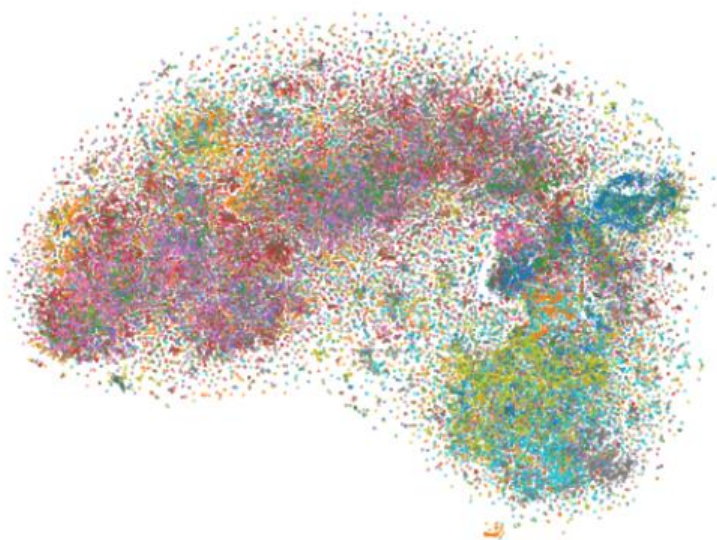      Reasoning of their decision-making process is lacking.

Our main challenge –

      Improving interpretability of the network inference process.
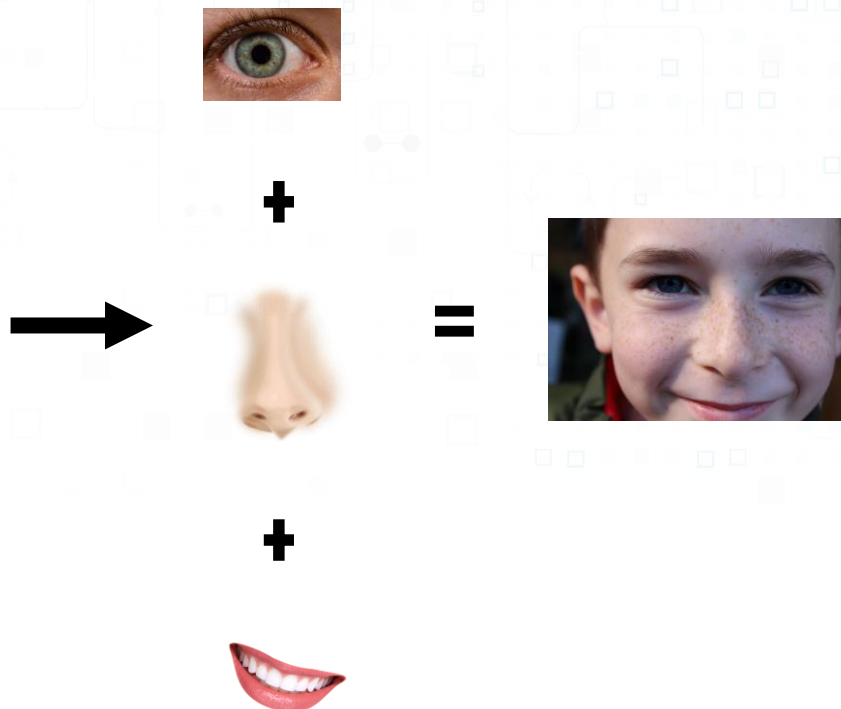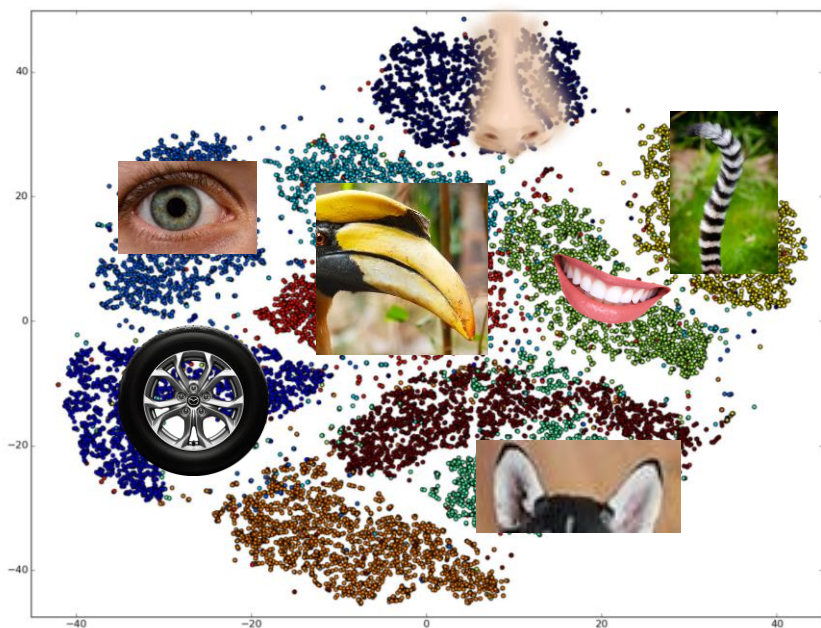
Input

CNN

Dog

# Our challenges

*Can we convert* *distributed representations* *into a* *human-oriented language*?

# Our challenges

*Can we learn a dictionary of visual words and model their interrelationships?*

# Generative modeling of MLPs

Trained MLP inference process as a Hidden Markov Model (HMM):

- **Single MLP layer distribution**: a mixture model of multivariate Gaussians (GMM).

- **Connections between layer representations**: conditional probability tables between GMM components of adjacent layers.

- **Post-Relu activations**: rectified Gaussians via additional hidden variables.

- **Optimization:** online Expectation Maximization (EM).

# Layers Dictionaries for CNNs

- Each spatial location example vector $x_p^l$, located at
  $p = (i, j) \in \left\{ \{1, \ldots, H^l\} \times \{1, \ldots, W^l\} \right\}$, is described as
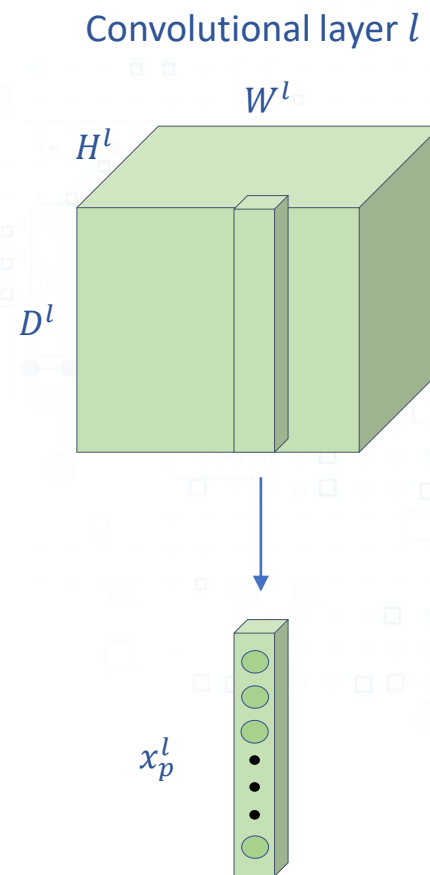  arising from a GMM of $K^l$ components.

- Example $x_p^l$ is assigned to cluster $C_k^l$ iff $P\left(h_p^l = k \big| x_p^l\right)$ is
  the highest.

- Each cluster $C_k^l$ represents a **visual word** (edge, texture, body part etc.), together forming the layer **dictionary**.

Convolutional layer $l$
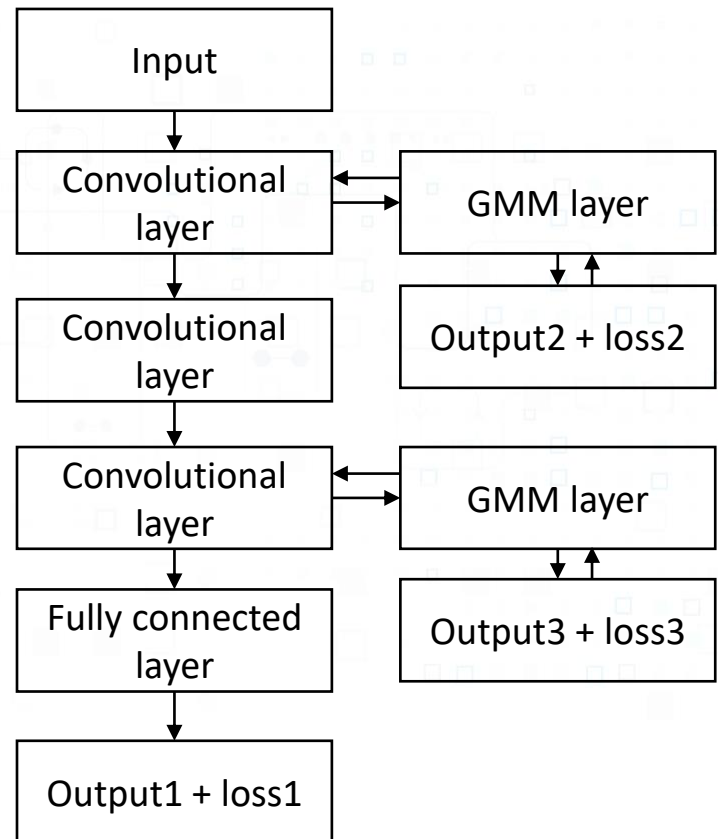
$W^l$

$H^l$

$D^l$

$x_p^l$

# Layers Dictionaries for CNNs

Train a CNN-GMM model:

1. Get a pre-trained CNN.

2. For every modeled layer $l$, append to it a GMM layer with the GMM parameters $\Theta^l = \left\{\pi_k^l, \mu_k^l, \Sigma_k^l\right\}_{k=1}^{K^l}$ as learning weights.

3. Train all GMM layers and estimate their parameters independently, using SGD.

```
        ┌──────────────┐
        │    Input     │
        └──────────────┘
               │
               ▼
   ┌──────────────────┐       ┌──────────────┐
   │  Convolutional   │ ◄────►│  GMM layer   │
   │      layer       │       └──────────────┘
   └──────────────────┘              │ ▲
               │                     ▼ │
   ┌──────────────────┐       ┌──────────────┐
   │  Convolutional   │       │Output2 + loss2│
   │      layer       │       └──────────────┘
   └──────────────────┘
               │
   ┌──────────────────┐       ┌──────────────┐
   │  Convolutional   │ ◄────►│  GMM layer   │
   │      layer       │       └──────────────┘
   └──────────────────┘              │ ▲
               │                     ▼ │
   ┌──────────────────┐       ┌──────────────┐
   │ Fully connected  │       │Output3 + loss3│
   │      layer       │       └──────────────┘
   └──────────────────┘
               │
   ┌──────────────────┐
   │ Output1 + loss1  │
   └──────────────────┘
```

# Graph Node Selection Algorithm

Consider a graph in which visual words $\left\{ C_k^l \right\}_{l=1,k=1}^{L,\,K^l}$ are the nodes, and transition probabilities between visual words of consecutive layers quantify edges.
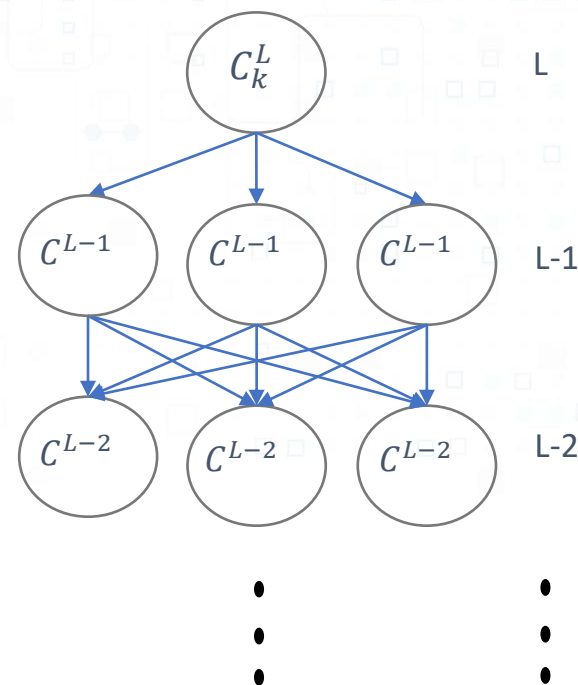
- Given a selected subset of images to be explained $\Omega = \{I_n\}_{n=1}^N$ (e.g., class or a single image), a specific subgraph can have high explanatory value.

  *How can we find the most explanatory visual words?*

- **Node selection algorithm** –
  iterative algorithm starting from "explaining" the classification decision node, then explaining layers backward, until outputting a subgraph.

# Graph Node Selection Algorithm

Given an instance of a single visual word $h_p^l = s$ to "explain", we look for the visual words $T$ in its receptive field $R(p)$ most contributing to its likelihood:

$$\max_{T,|T|=Z} \log P\left(h_p^l = s | h_q^{l'} : q \in R(p), h_q^{l'} \in T\right)$$

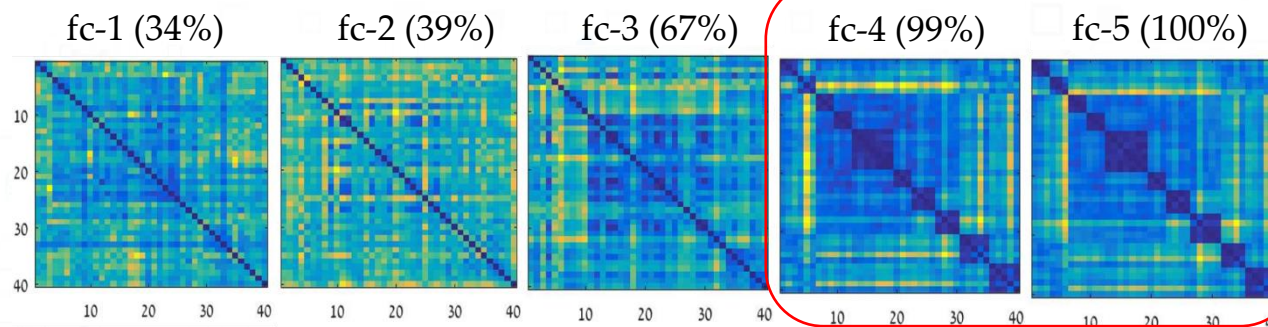Using location independence assumptions, a '$t-$explains$-s$' score is derived:

$$S^{l'}(s,t) = \sum_{t=1}^{K^{l'}} \left|\left\{q : h_q^{l'} = t, q \in R(p)\right\}\right| \log \frac{P\left(h_q^{l'} = t | h_p^l = s, q \in R(p)\right)}{P\left(h_q^{l'} = t\right)}$$

the number of times visual word $t$ appears in the receptive field of location $p$

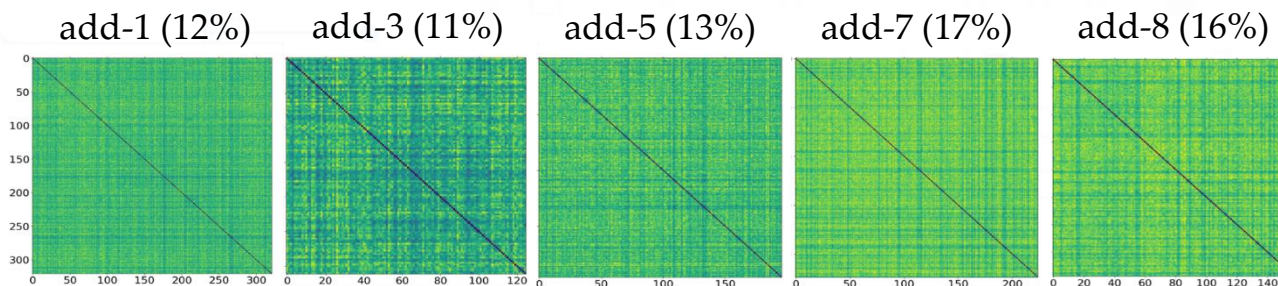how likely it is to see word $t$ in the receptive field of $s$ compared to seeing it in general

# Cluster Similarity Across Layers

fc-1 (34%)   fc-2 (39%)   fc-3 (67%)   fc-4 (99%)   fc-5 (100%)

MLP

increasing similarity between clusters representing the same class.

add-1 (12%)   add-3 (11%)   add-5 (13%)   add-7 (17%)   add-8 (16%)
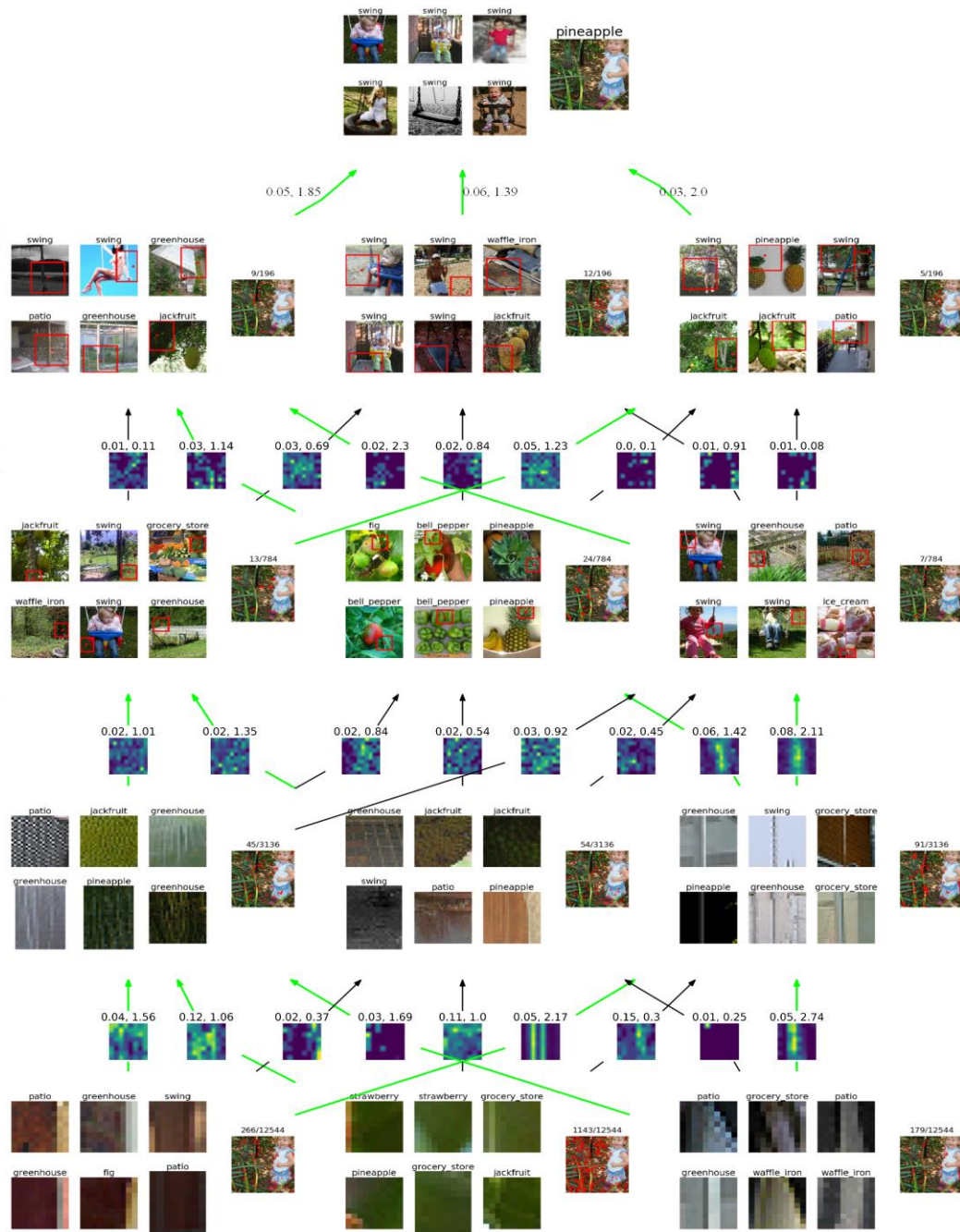
CNN

the clusters stay local and diverse, even at the uppermost layers.

# VGG-16 Image Inference Graph

**Ground truth**: Pineapple

**Network prediction**: Swing

# VGG-16 Image Inference Graph

"Rope"

swing    swing    greenhouse    9/196

patio    greenhouse    jackfruit

"Grass"

swing    pineapple    swing    5/196

jackfruit    jackfruit    patio

0.02, 2.3    0.03, 1.14    0.05, 1.23

"Vertical-stripes"

swing    greenhouse    patio    7/784

swing    swing    ice_cream

"Green-rounded-line"

fig    bell_pepper    pineapple    24/784

bell_pepper    bell_pepper    pineapple

"Vegetation"

jackfruit    swing    grocery_store    13/784

waffle_iron    swing    greenhouse

0.08, 2.11    0.06, 1.42    0.02, 1.01    0.02, 1.35

"Isolated-vertical-line"

greenhouse    swing    grocery_store    91/3136

pineapple    greenhouse    grocery_store

"Ground-structures"

greenhouse    jackfruit    jackfruit    54/3136

swing    patio    pineapple

"Vertical-stripe-structure"

patio    jackfruit    greenhouse    45/3136
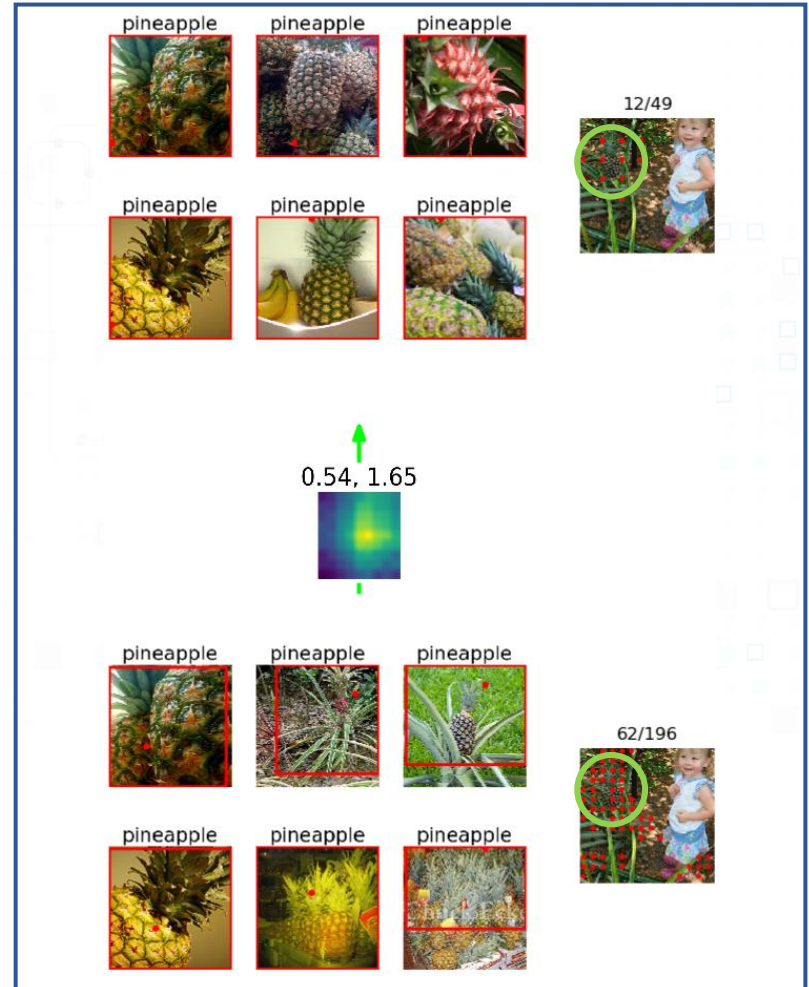
greenhouse    pineapple    greenhouse

# ResNet50 Image Inference Graph

**Ground truth**: Pineapple

**Network prediction**: Pineapple

# Inference Graphs for CNN Interpretation

Yael Konforti, Alon Shpigler, Boaz Lerner, Aharon Bar Hillel

Ben-Gurion University of the Negev, Beer Sheva, Israel

{yaelkonf, alonshp}@post.bgu.ac.il; {boaz, barhille}@bgu.ac.il

Code: https://github.com/yaelkon/GMM-CNN