

Machine Learning Methods for Neural Data Analysis

Lecture 6: Markerless Pose Tracking

Scott Linderman

STATS 220/320 (NBIO220, CS339N). Winter 2021.

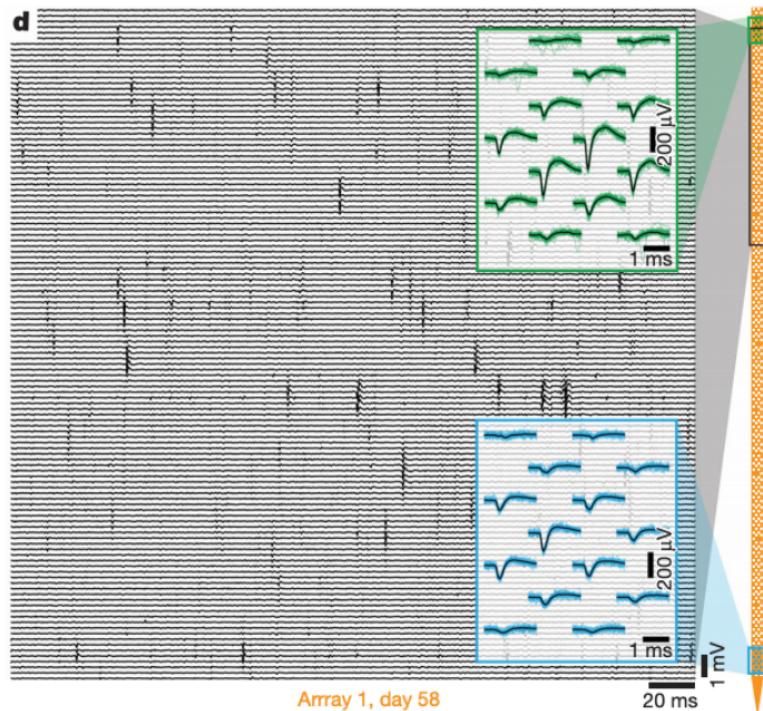
Announcements

- Plowing ahead today. The rest of the story on deconvolution under a point process story is in the course notes.
- Lab 3 Errata:
 - Background footprints should be normalized so that $\|u_0\|_2 = 1$, not maximum value one as it said in a code comment.
 - Lower diagonal of G should have $-e^{-1/\tau}$.

Recap

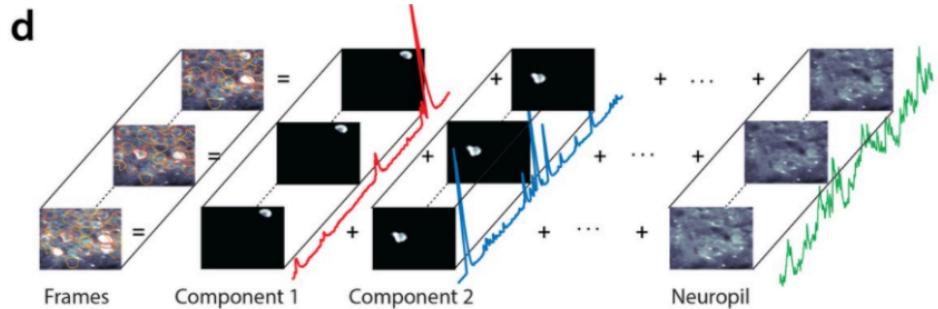
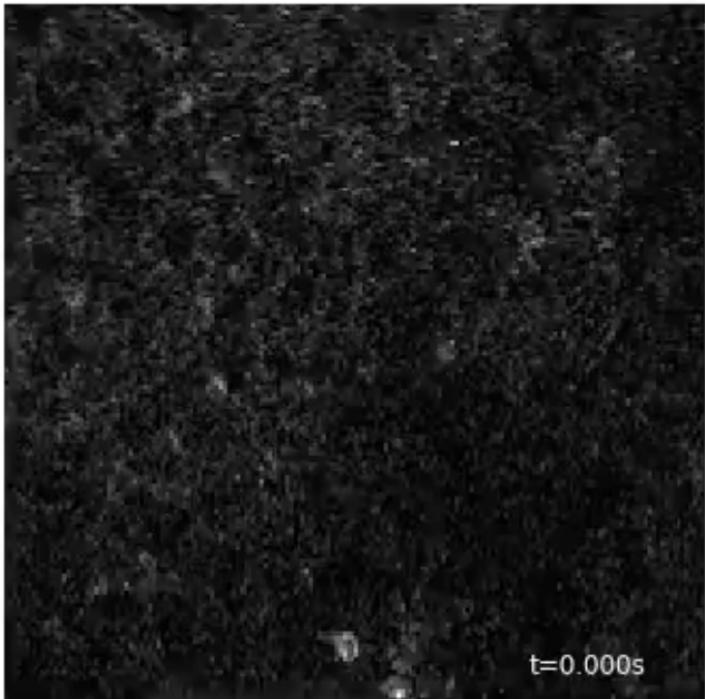
High-density probes (e.g. Neuropixels)

- The raw data is a **multidimensional time series** of **voltage measurements**, one for each recording site on the probe.
- When neurons near the probe fire an **action potential**, it registers a **spike in the voltage** on nearby channels.
- Typical recordings detect spikes from **O(100) neurons**.



Recap

2 photon calcium imaging



$$Y = U^T C + u_0 c_0^\top + \epsilon$$

$$U \in \mathbb{R}^{N \times P} \quad C \in \mathbb{R}_+^{N \times T} \quad u_0 \in \mathbb{R}^P \quad c_0 \in \mathbb{R}^T$$

$$\epsilon_{pt} \sim \mathcal{N}(0, \sigma^2)$$

**“The brain is worthy of study because it is
in charge of behavior”**

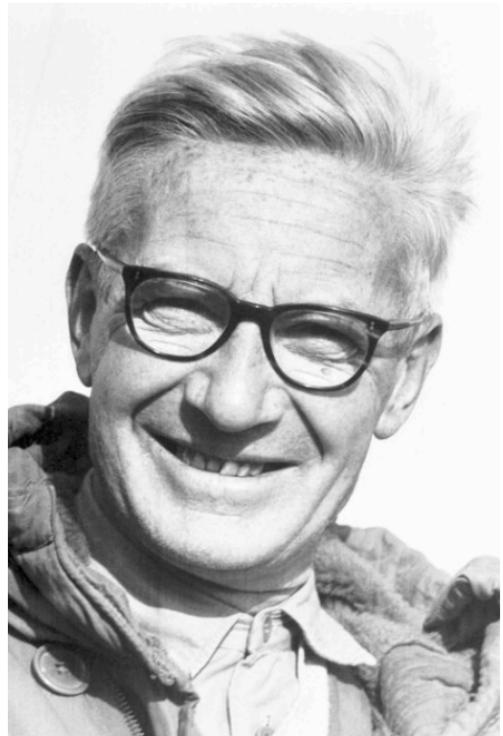
Datta, Anderson, Branson, Perona, and Leifer. Computational Neuroethology: A Call to Action. *Neuron* 2019.



Ethology

The study of (natural) behavior

- **Hypothesis:** “exposing the structure of behavior...will yield insights into how the brain creates behavior.” Datta et al.
- **Structure:** how behavior in the natural environment is built from components and organized over time in response to ecologically relevant stimuli.
- **Natural behavior:**
 - Exploring new environments
 - Foraging for food
 - Finding shelter
 - Identifying mates
 - ...



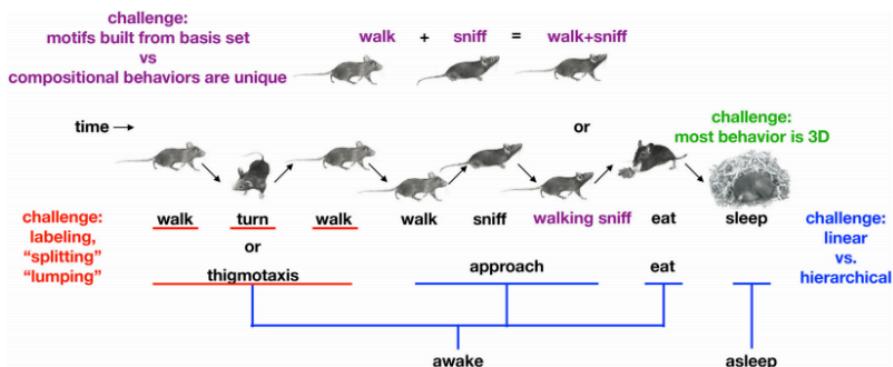
Nikolaas Tinbergen
Nobel Prize in Physiology or Medicine 1973



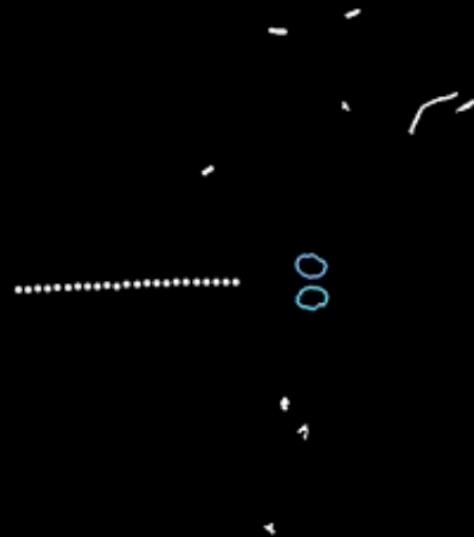
Computational (Neuro)Ethology

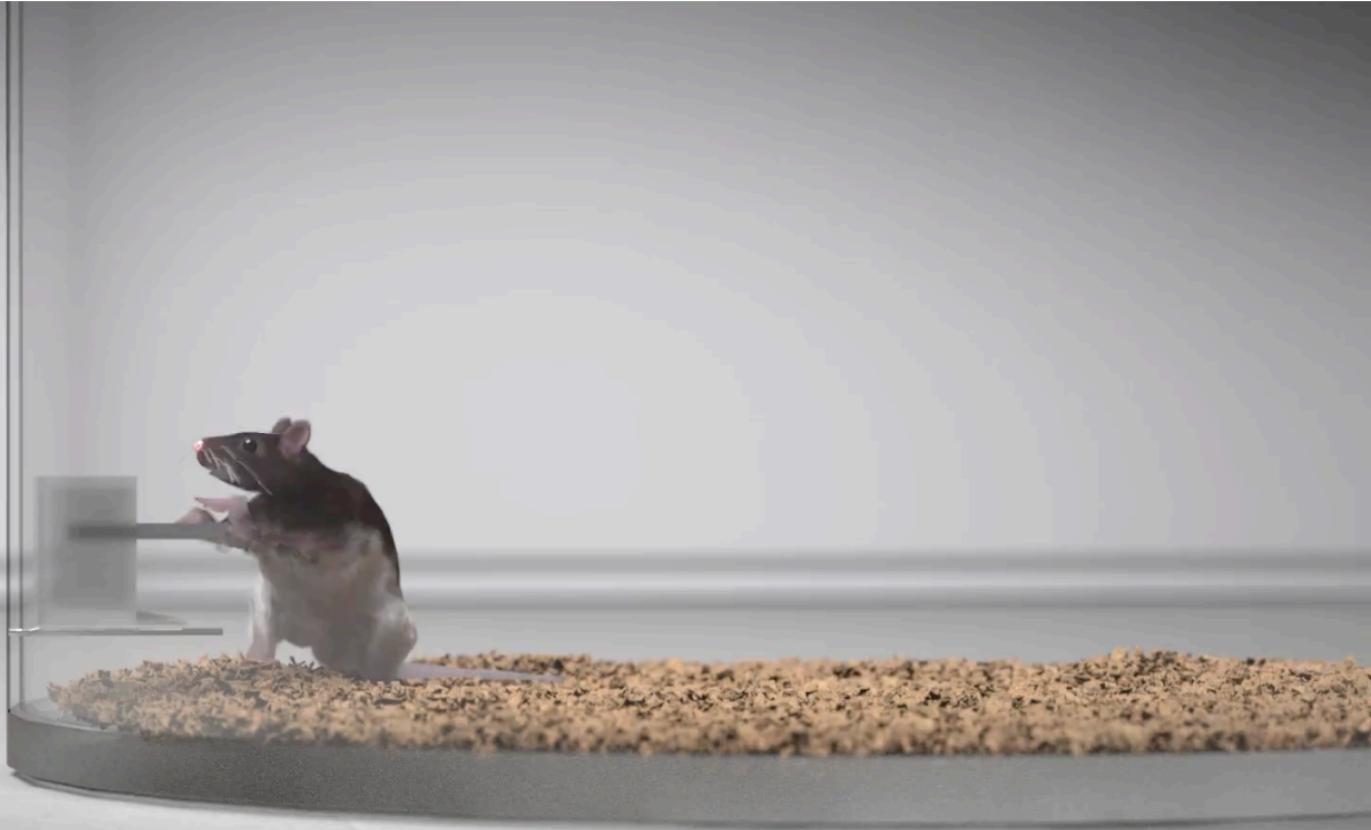
Quantifying natural behavior (and relating it to neural activity)

- Leveraging advances in **computer vision** and **machine learning** to extract behavioral features of interest from raw data.
- Modeling the dynamics of 3D pose as a function of sensory input and internal state.
- Decomposing behavior into stereotyped components and behavioral motifs.
- Correlating behavioral motifs with large scale neural recordings.
- Identifying causal relationships between neural activity and motor output.



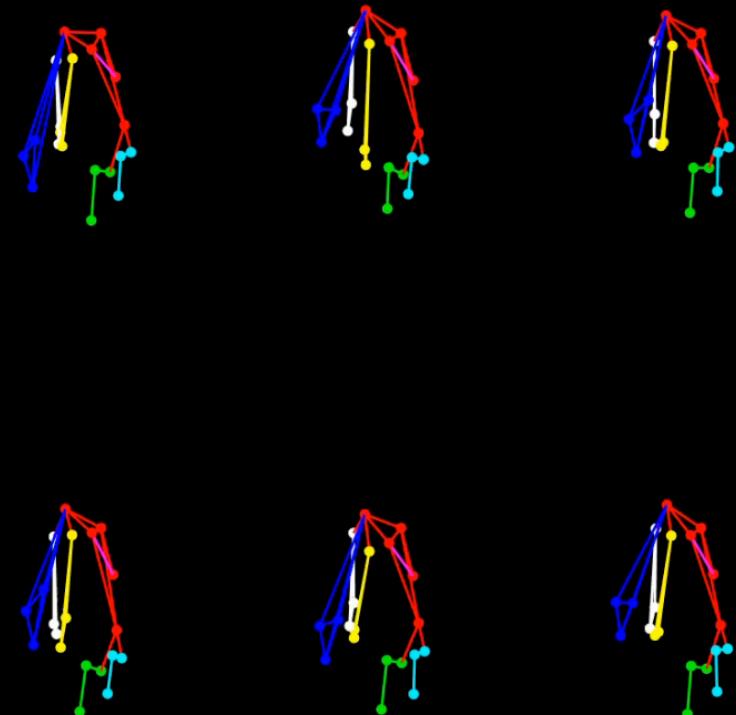
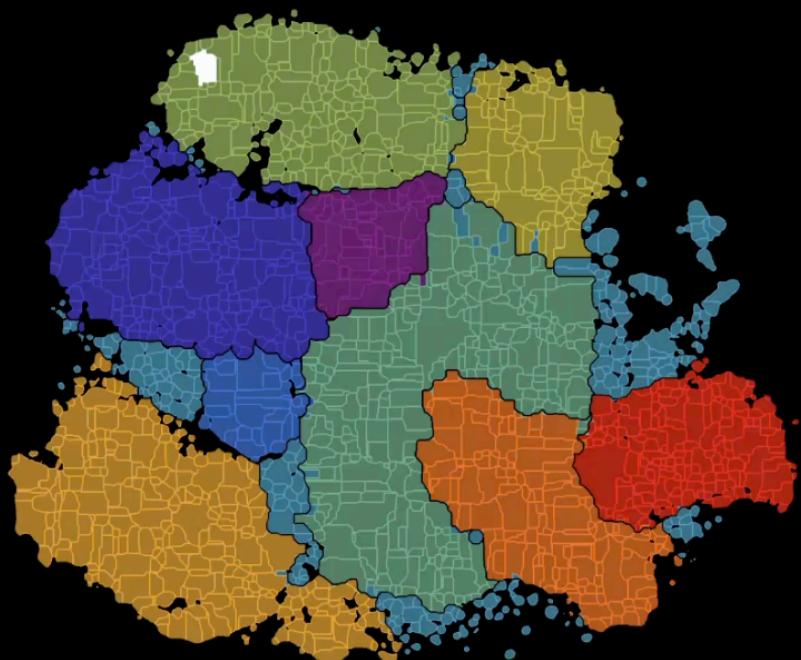
Datta et al (Neuron, 2019)





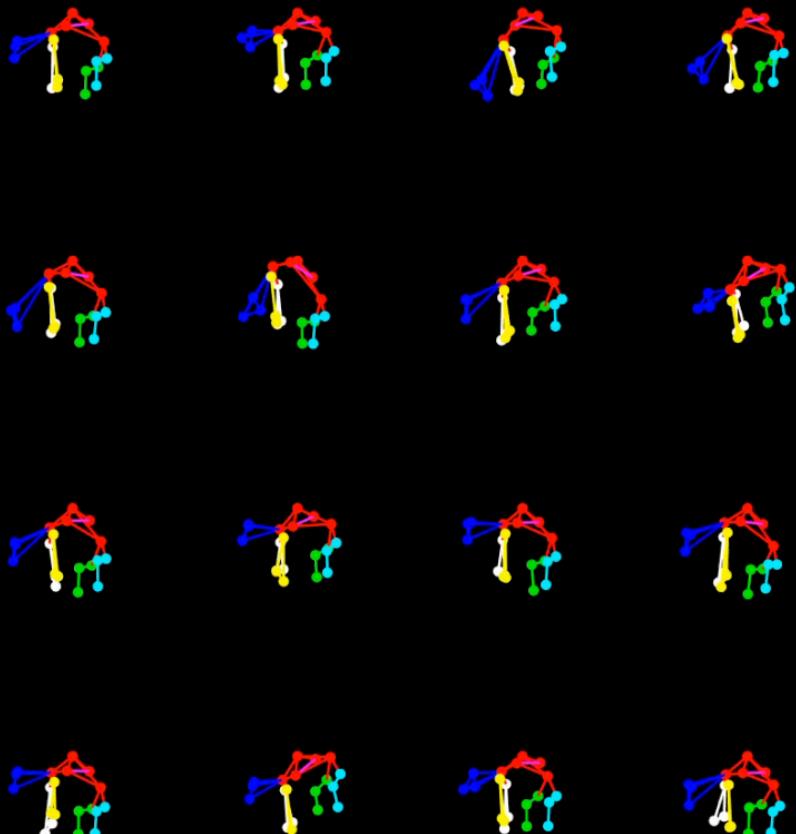
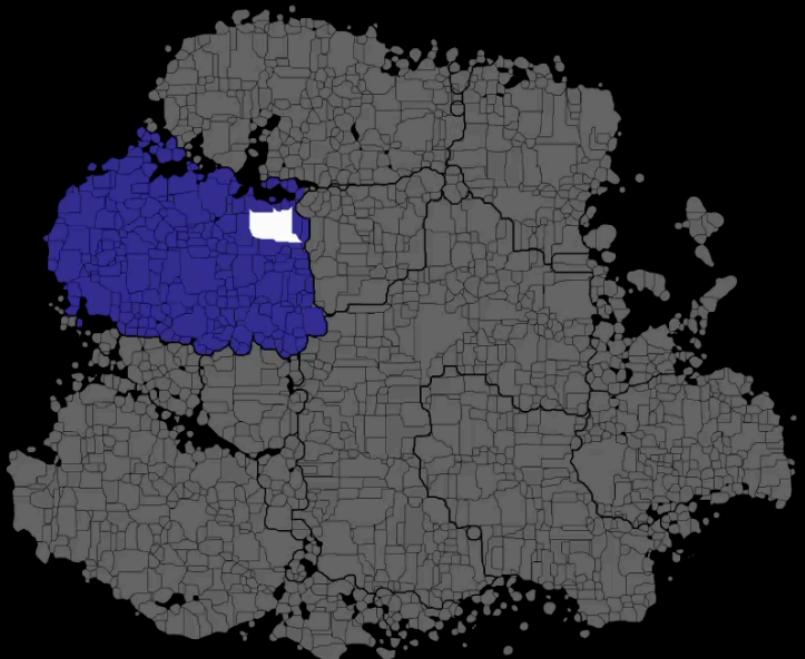
CAPTURE: Marshall et al (*Neuron*, 2020)

Left Groom

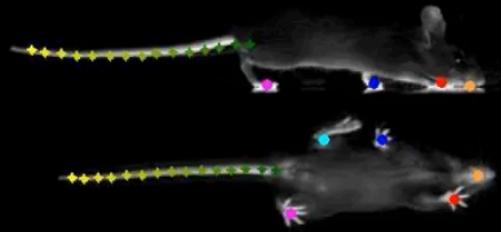
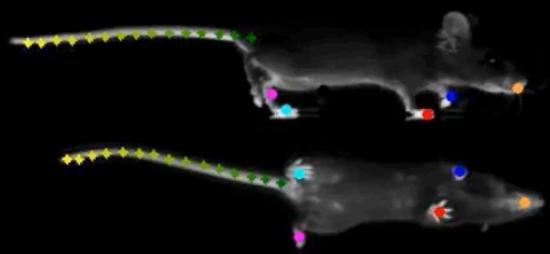


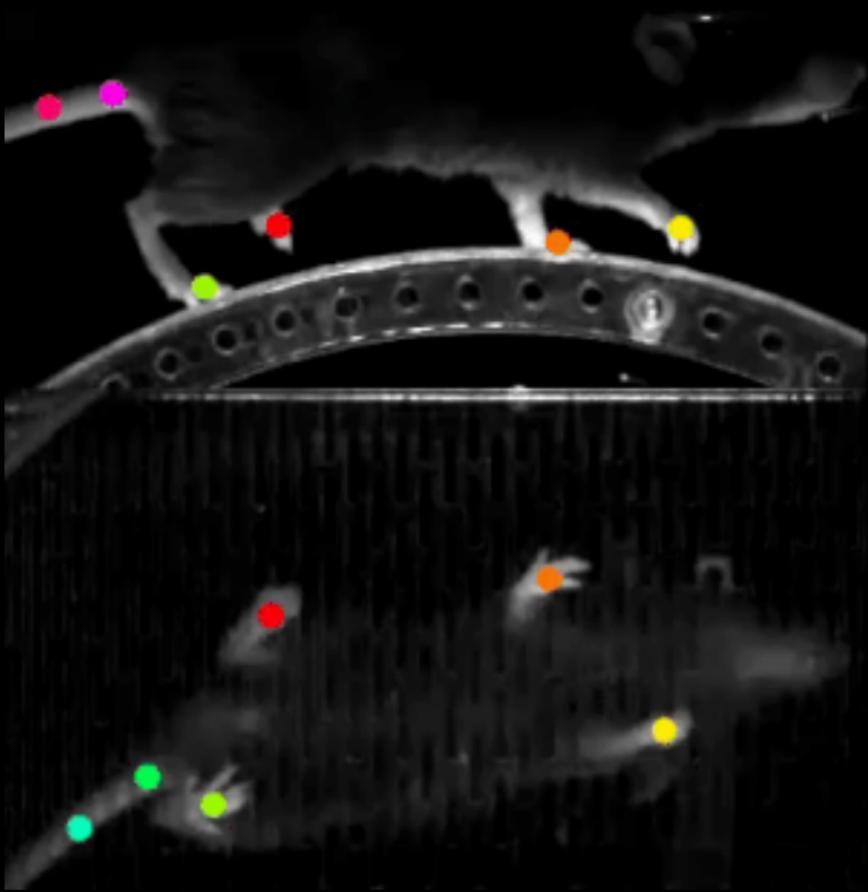
CAPTURE: Marshall et al (*Neuron*, 2020)

Right Groom - Low

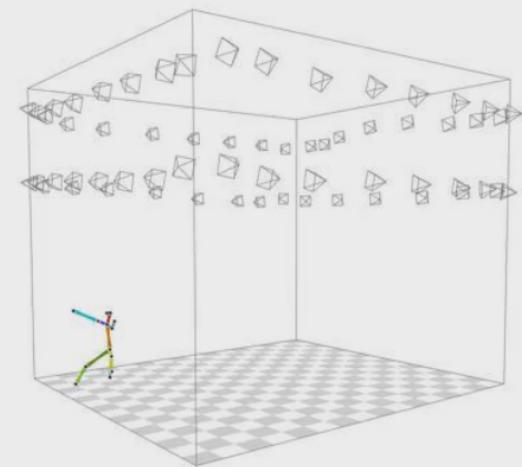
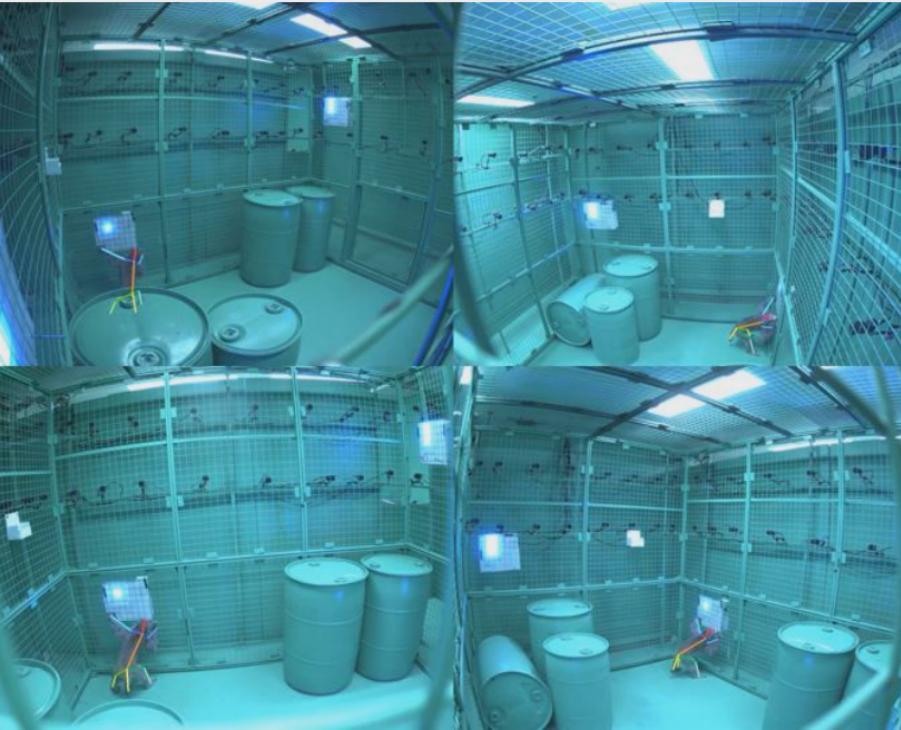


CAPTURE: Marshall et al (*Neuron*, 2020)

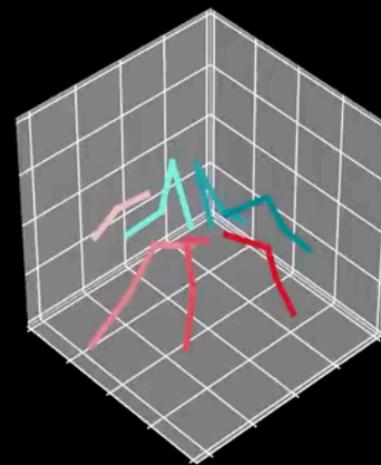
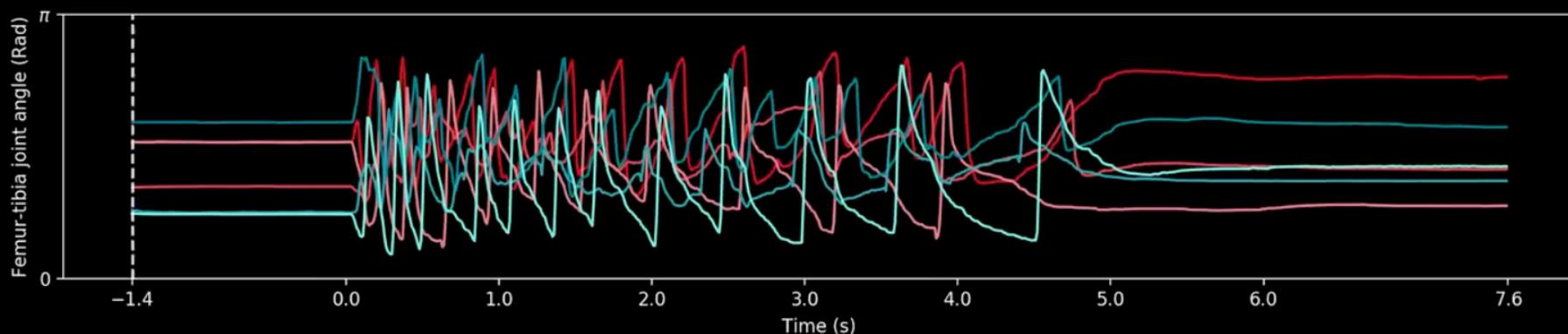


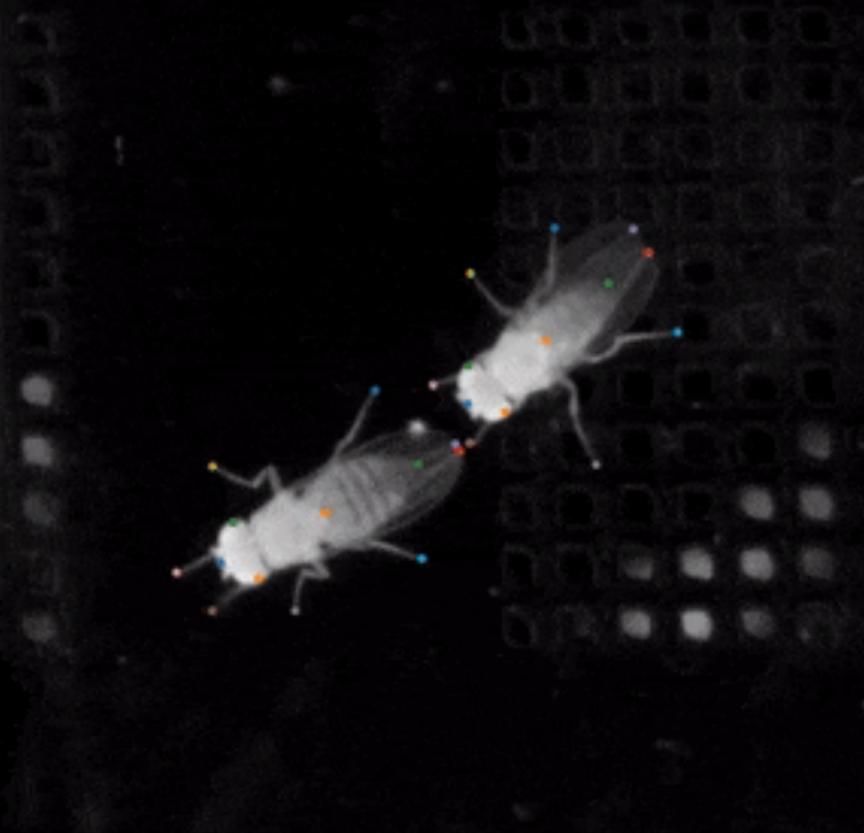


DeepLabCut: Mathis et al. (*Nat Neuro* 2018)



OpenMonkeyStudio: Bala et al (*Nature Comm.*, 2020)





SLEAP: Pereira et al (bioRxiv, 2020)

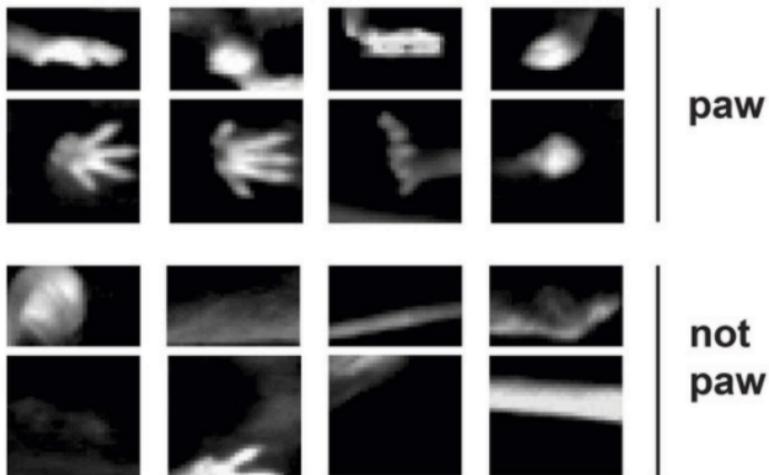
Agenda

1. Basics of markerless pose tracking
2. Transfer learning
3. Triangulating 3D pose from multiple 2D views

Basic pose tracking

Turn it into a supervised learning problem

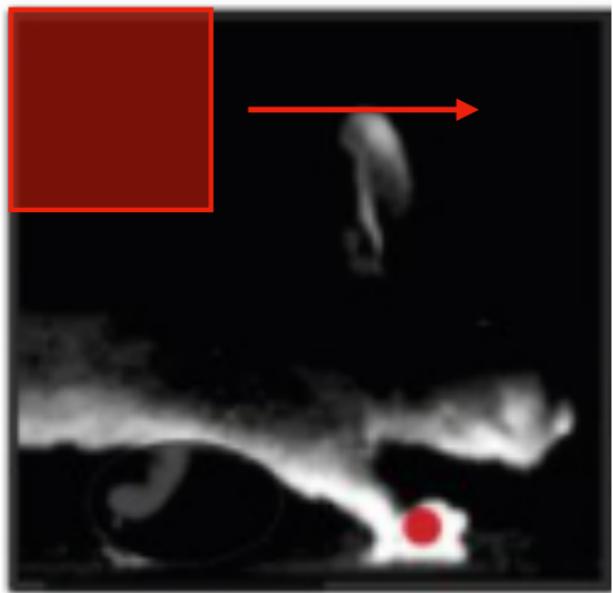
- Extract patches from the video frames and label them as positive or negative examples of a key point (e.g. paw).
- Train a binary classifier (logistic regression, SVM, neural network, etc.) to predict key point or not.



Basic pose tracking

Turn it into a supervised learning problem

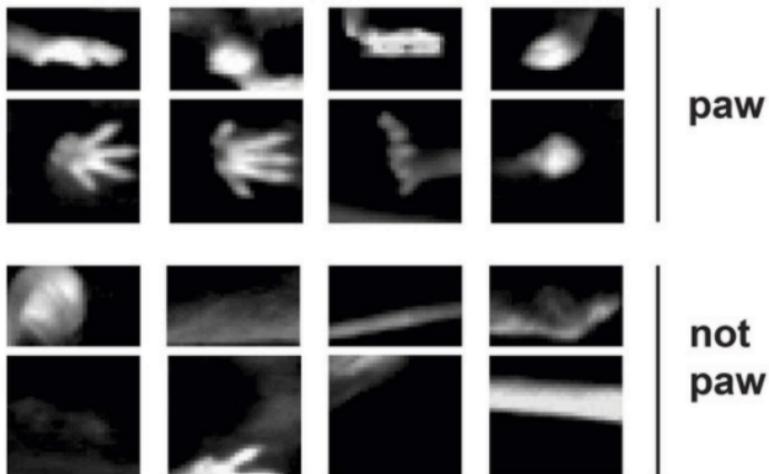
- Extract patches from the video frames and label them as positive or negative examples of a keypoint (e.g. paw).
- Train a binary classifier (logistic regression, SVM, neural network, etc.) to predict keypoint or not.
- At test time, classify each patch in the image and use a heuristic to pick the most likely keypoint location(s).



Basic pose tracking

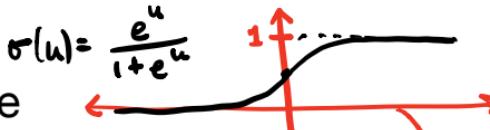
Mathematical formulation

- Let P_h and P_w be the height and width, respectively, of the patch (in pixels).
- N denote the number of patches
- $X_n \in \mathbb{R}^{P_h \times P_w}$ denote the n -th patch.
- $y_n \in \{0,1\}$ denote whether or not the patch is an instance of the keypoint.
- $W \in \mathbb{R}^{P_h \times P_w}$ denote the weights of our model.



Basic pose tracking

Via logistic regression

$$\sigma(u) = \frac{e^u}{1+e^u}$$

$$\langle W, X_n \rangle = \text{Tr}(W^T X_n) = \text{vec}(w)^T \text{vec}(X_n)$$

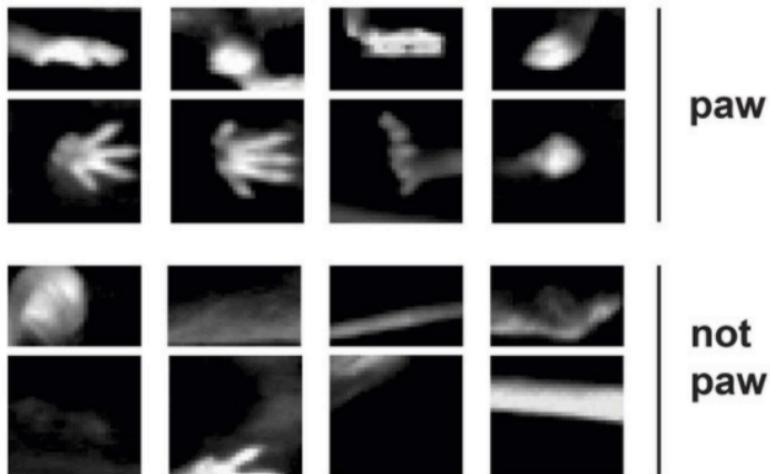
Assume

$$p(y_n | X_n, W) = \text{Bern}(y_n | \sigma(\langle W, X_n \rangle))$$

$$= \sigma(\langle W, X_n \rangle)^{y_n} (1 - \sigma(\langle W, X_n \rangle))^{1-y_n}$$

$$= \frac{e^{\langle W, X_n \rangle \cdot y_n}}{(1 + e^{\langle W, X_n \rangle})^{y_n}} \cdot \frac{1}{(1 + e^{\langle W, X_n \rangle})^{1-y_n}}$$

$$= \frac{e^{\langle W, X_n \rangle \cdot y_n}}{1 + e^{\langle W, X_n \rangle}}$$



Basic pose tracking

Maximum likelihood estimation

$$\mathcal{L}(W) = \sum_{n=1}^N \log p(y_n | x_n, W) = \sum_{n=1}^N \langle W, x_n \rangle \cdot y_n - \log(1 + e^{\langle W, x_n \rangle})$$

$$\nabla_w \mathcal{L}(w) = \sum_{n=1}^N x_n \cdot y_n - \frac{e^{\langle W, x_n \rangle}}{1 + e^{\langle W, x_n \rangle}} \cdot x_n$$

$$= \sum_{n=1}^N \underbrace{\left(y_n - \sigma(\langle W, x_n \rangle) \right)}_{\text{error}} \cdot x_n$$

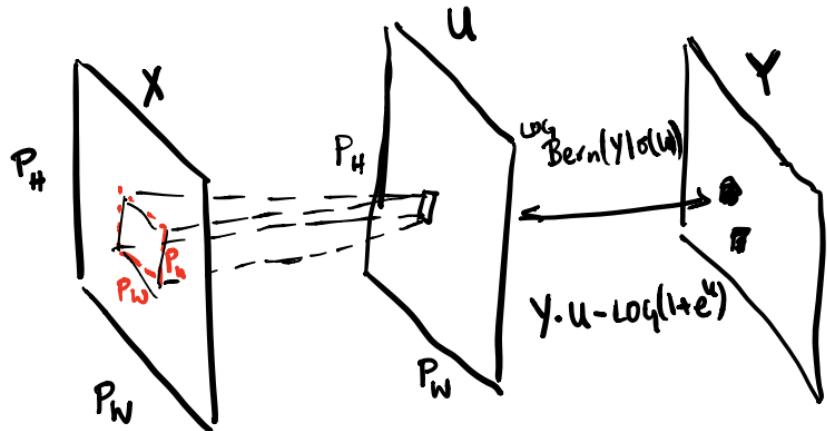
Gradient Ascent

$$\text{for } i=1 \dots \\ w^{(i+1)} = w^{(i)} + \alpha_i \nabla \mathcal{L}(w)$$

Basic pose tracking

With convolutional neural networks

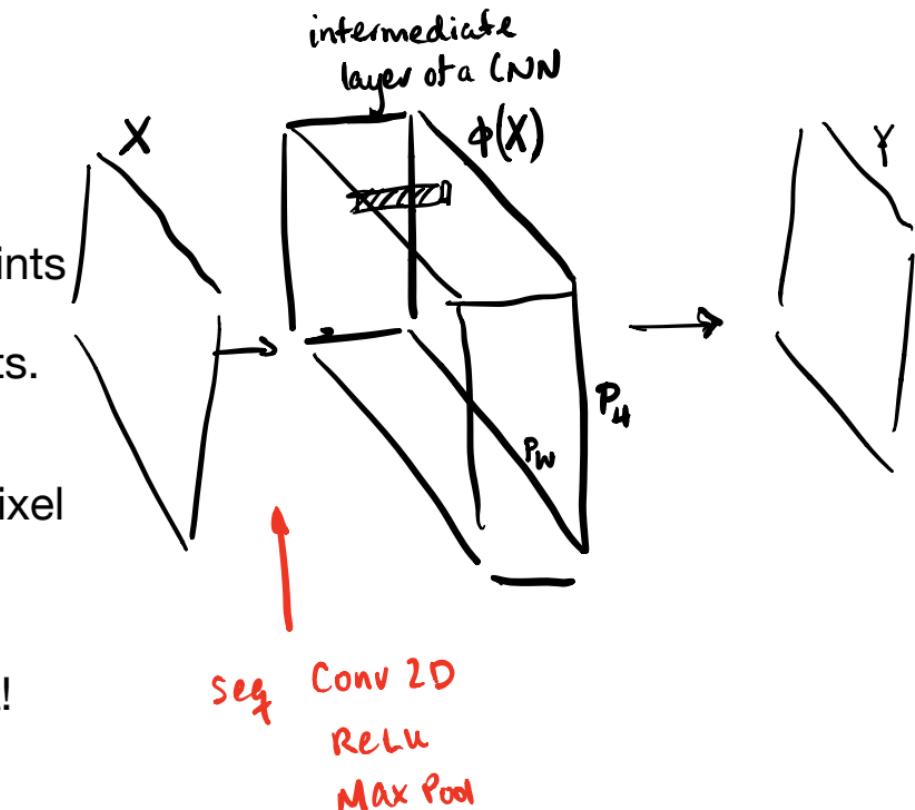
- Instead of working with patches, let's work with images directly.
- Let $X \in \mathbb{R}^{P_H \times P_W}$ denote an image (height P_H , width P_W)
- Let $Y \in \{0,1\}^{P_H \times P_W}$ indicate the location(s) of the keypoint.
- The 2D cross-correlation $X \star W$; is a sliding dot product of weights across all $P_h \times P_w$ patches in the image. It produces a $P_H \times P_W$ output.
- In PyTorch, it's implemented by the `F.conv2d` function and the Conv2D layer.



Basic pose tracking

Feature learning in CNNs

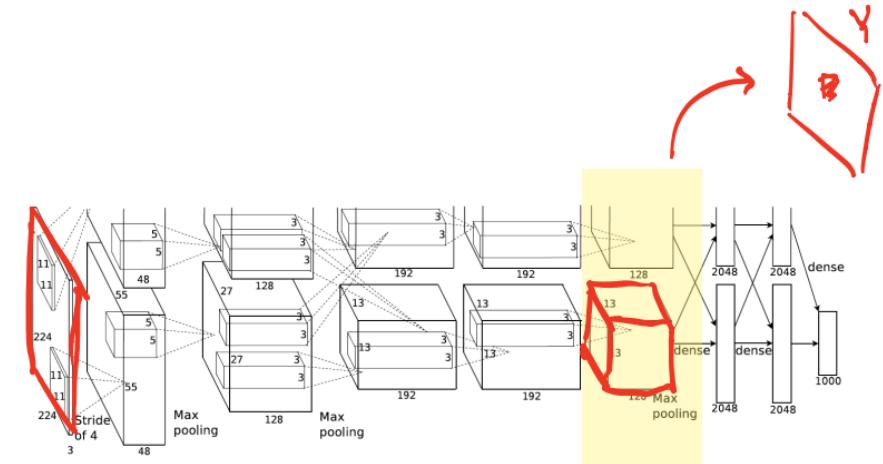
- This simple model assumes keypoints can be detected with a **linear classifier** using raw pixels as inputs.
- We can perform **nonlinear classification** by encoding each pixel with a vector of features.
- Rather than handcrafting these features, **learn them** from the data!



Transfer Learning

Transfer Learning

- **Idea:** rather than handcrafting features or learning them from scratch, **use a pre-trained network** for a related task.
- **Example:** use the features of a deep neural network for image classification.
- **Reroute** the output of an intermediate layer to a **new loss function**.
- Optionally, **fine tune** the weights in the early layers via stochastic gradient descent on the new loss.
- With good starting features, you **only need a few training examples** to perform animal pose estimation.

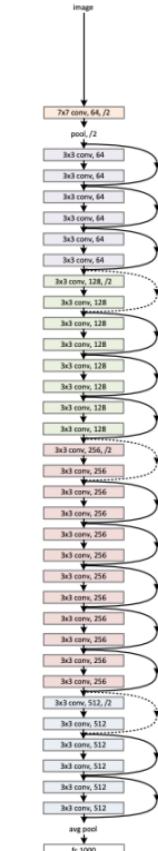


Transfer Learning In DeepLabCut, SLEAP, etc.

- DLC and SLEAP repurpose state-of-the-art deep networks for human pose detection.
- DLC starts with a residual network (resnet-50) and adds “deconvolutional” layers, as in DeeperCut for human pose estimation.
- SLEAP starts with “stacked hourglass networks” for human pose estimation.

Deep Residual Networks
(resnet-50)

34-layer residual



Stacked Hourglass Networks

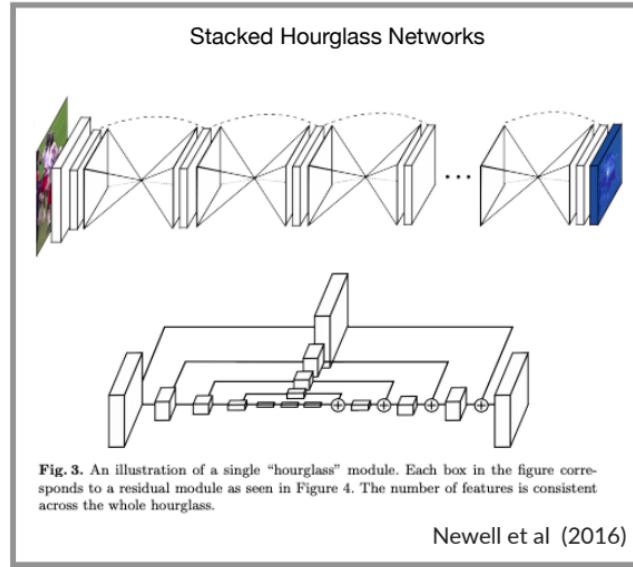
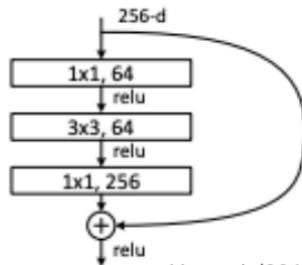


Fig. 3. An illustration of a single “hourglass” module. Each box in the figure corresponds to a residual module as seen in Figure 4. The number of features is consistent across the whole hourglass.

Newell et al (2016)

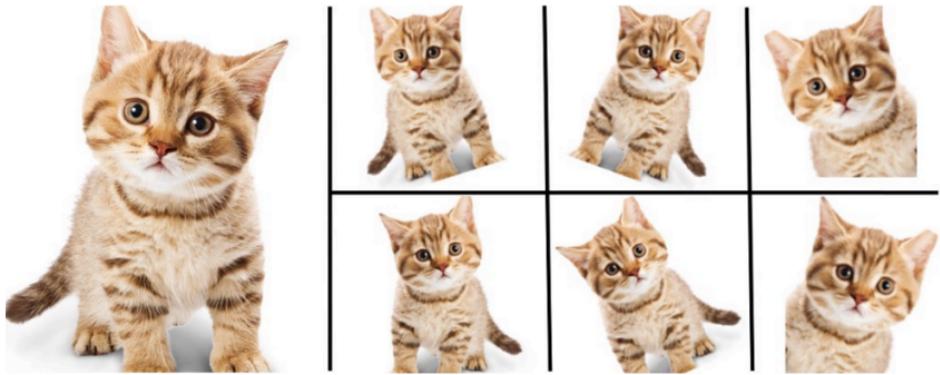


He et al (2015)

Transfer Learning

Data augmentation

- Labeling data is tedious.
- **Idea:** Make the most of each training example by making alterations your classifier should be robust to.
- Eg a cropped, rotated, and scaled paw is still a paw. A partially occluded paw is still a paw.



3D Triangulation

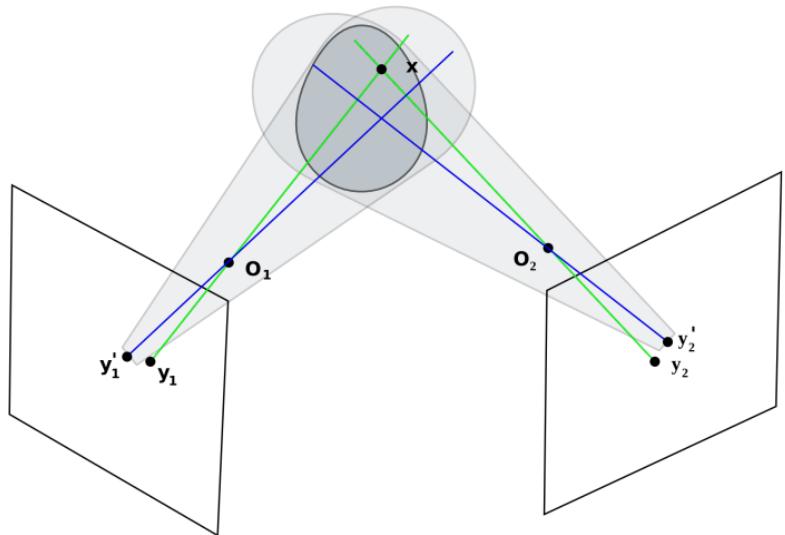
Basic triangulation

Projective geometry

- Projective geometry makes far away objects appear smaller.

$$\vec{y}_c \approx f_c(\vec{x})$$

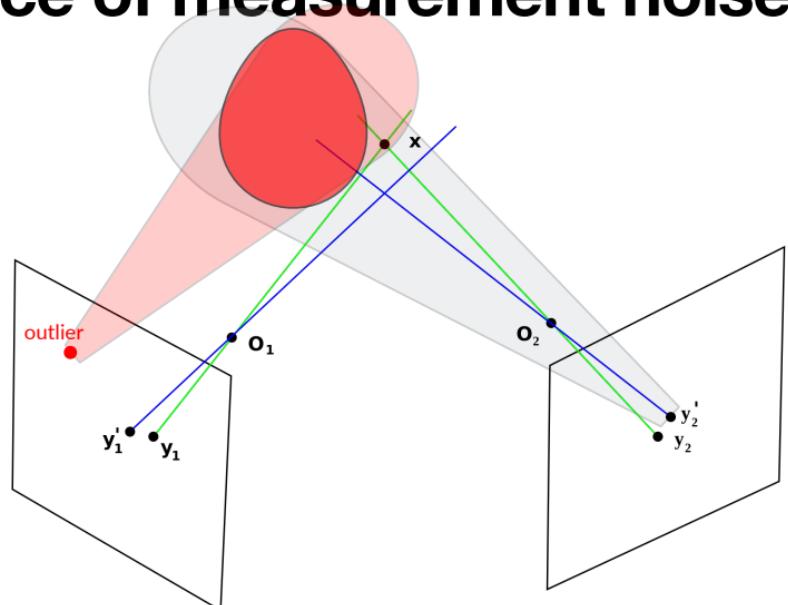
$$f_c(\vec{x}) = \frac{1}{w}(u, v)^\top \text{ where } (u, v, w)^\top = A_c \vec{x} + b_c,$$



Modified from wikipedia.org

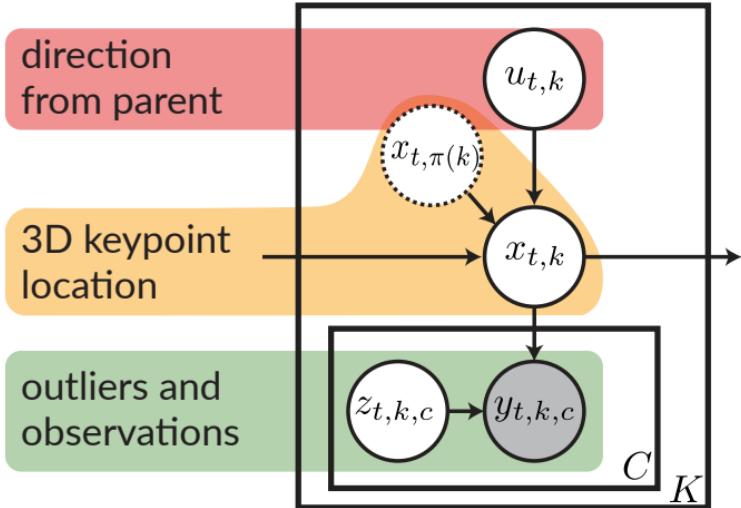
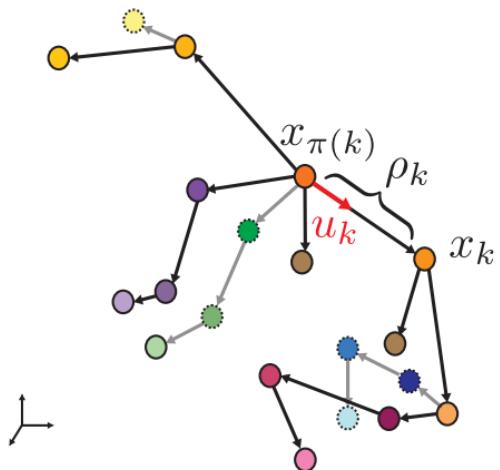
Triangulation in the presence of measurement noise

- Projective geometry makes far away objects appear smaller.
- Outliers in 2D estimates can severely affect 3D triangulation.
- Typical approaches:
 - More data
 - Temporal constraints
 - Median filtering (DLC-3D) / RANSAC
 - Robust noise models



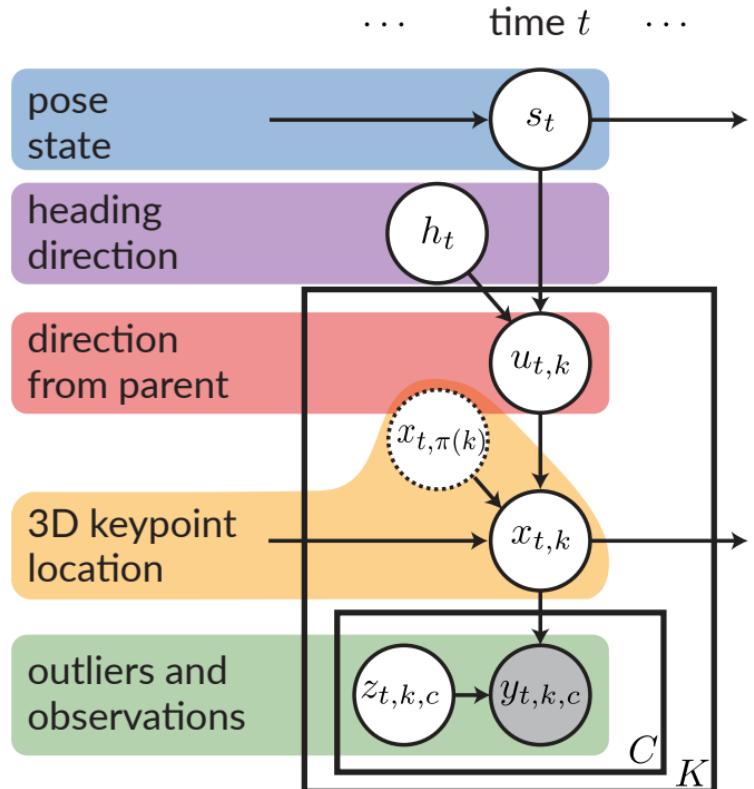
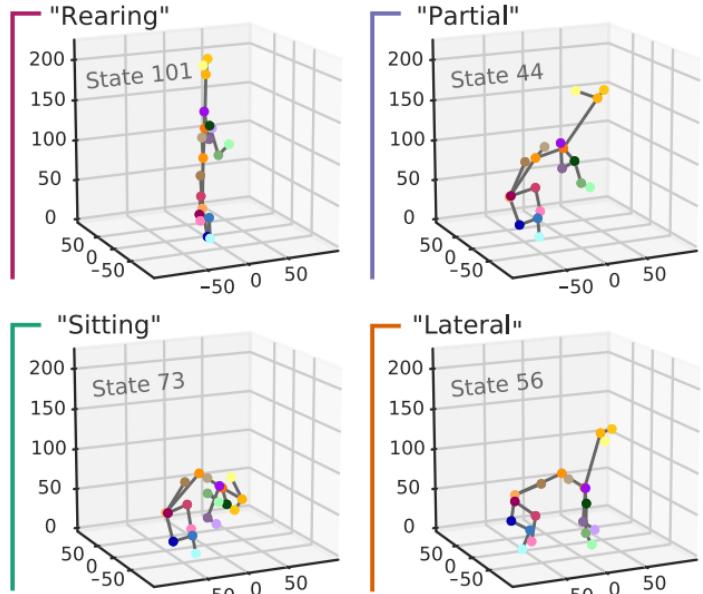
Modified from wikipedia.org

Bayesian triangulation



Zhang et al (AISTATS, 2021)

GIMBAL: Capturing correlations in direction vectors with pose states



Conclusion

- **Precise behavior quantifications** are critical for understanding how neural activity relates to behavioral output.
- **Markerless pose tracking** methods have made it much easier to obtain such quantifications.
- **Convolutional neural networks** are naturally suited to this task.
- With **transfer learning**, we can leverage state-of-the-art deep networks for image classification to warm-start pose tracking.
- We can **triangulate 3D pose** from 2D images using projecting geometry and spatiotemporal priors.

Further reading

- Datta, Sandeep Robert, et al. "Computational neuroethology: a call to action." *Neuron* 104.1 (2019): 11-24.
- Mathis, Alexander, et al. "DeepLabCut: markerless pose estimation of user-defined body parts with deep learning." *Nature neuroscience* 21.9 (2018): 1281-1289.
- Pereira, Talmo D., et al. "Fast animal pose estimation using deep neural networks." *Nature methods* 16.1 (2019): 117-125.
- He, Kaiming, et al. "Deep residual learning for image recognition." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.

Machine Learning Methods for Neural Data Analysis

Lecture 6: Markerless Pose Tracking

Scott Linderman

STATS 220/320 (NBIO220, CS339N). Winter 2021.

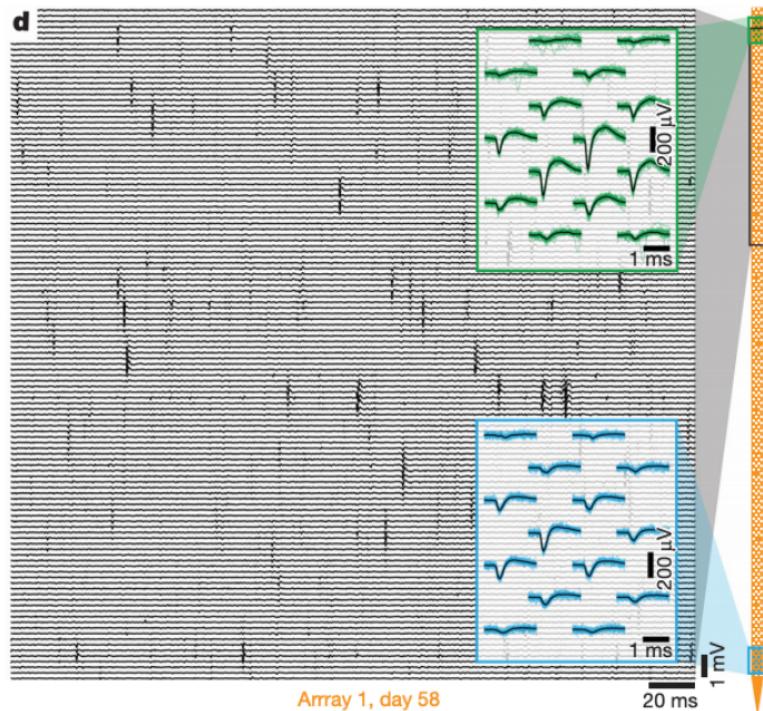
Announcements

- Plowing ahead today. The rest of the story on deconvolution under a point process story is in the course notes.
- Lab 3 Errata:
 - Background footprints should be normalized so that $\|u_0\|_2 = 1$, not maximum value one as it said in a code comment.
 - Lower diagonal of G should have $-e^{-1/\tau}$.

Recap

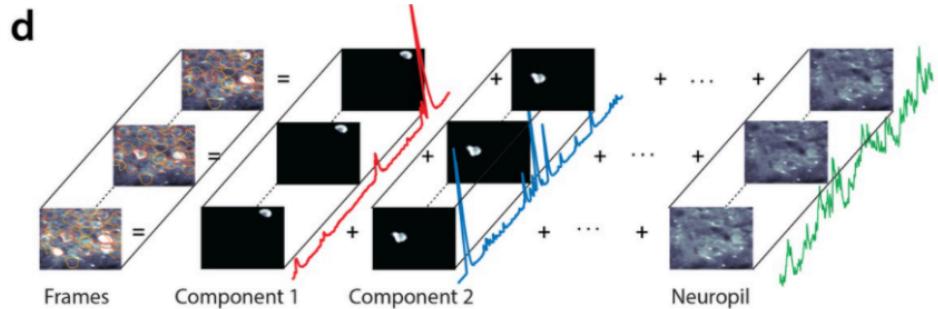
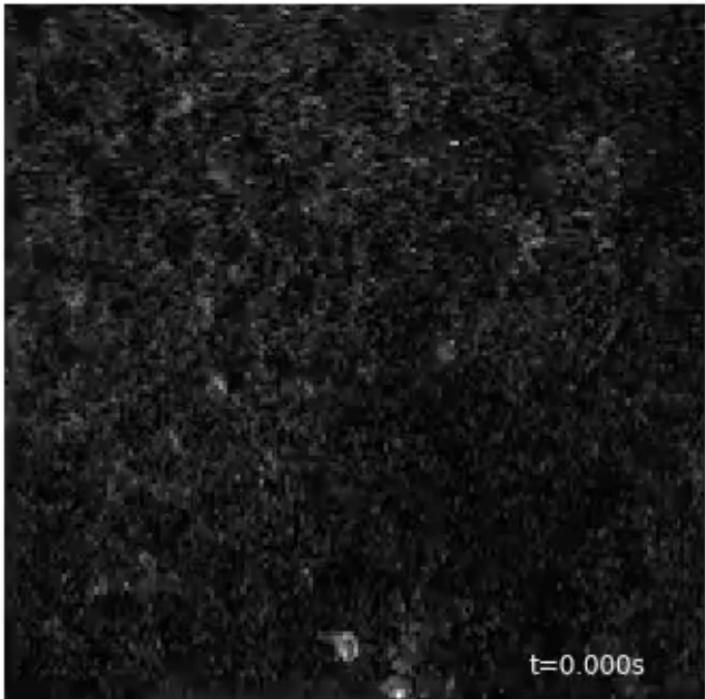
High-density probes (e.g. Neuropixels)

- The raw data is a **multidimensional time series** of **voltage measurements**, one for each recording site on the probe.
- When neurons near the probe fire an **action potential**, it registers a **spike in the voltage** on nearby channels.
- Typical recordings detect spikes from **O(100) neurons**.



Recap

2 photon calcium imaging



$$Y = U^T C + u_0 c_0^\top + \epsilon$$

$$U \in \mathbb{R}^{N \times P} \quad C \in \mathbb{R}_+^{N \times T} \quad u_0 \in \mathbb{R}^P \quad c_0 \in \mathbb{R}^T$$

$$\epsilon_{pt} \sim \mathcal{N}(0, \sigma^2)$$

**“The brain is worthy of study because it is
in charge of behavior”**

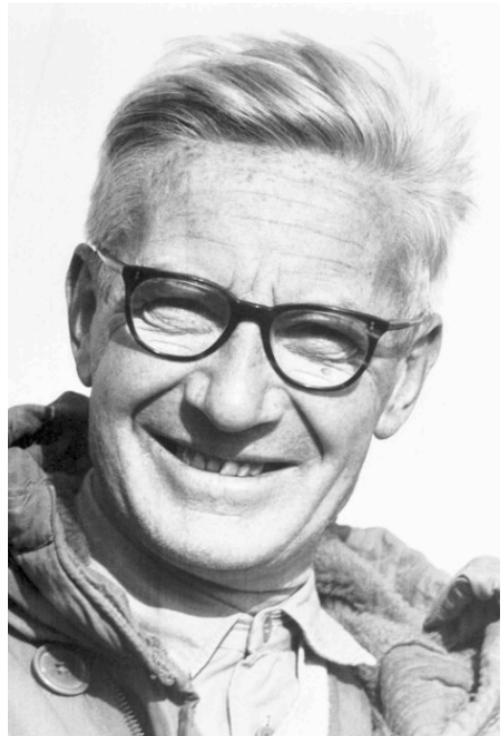
Datta, Anderson, Branson, Perona, and Leifer. Computational Neuroethology: A Call to Action. *Neuron* 2019.



Ethology

The study of (natural) behavior

- **Hypothesis:** “exposing the structure of behavior...will yield insights into how the brain creates behavior.” Datta et al.
- **Structure:** how behavior in the natural environment is built from components and organized over time in response to ecologically relevant stimuli.
- **Natural behavior:**
 - Exploring new environments
 - Foraging for food
 - Finding shelter
 - Identifying mates
 - ...



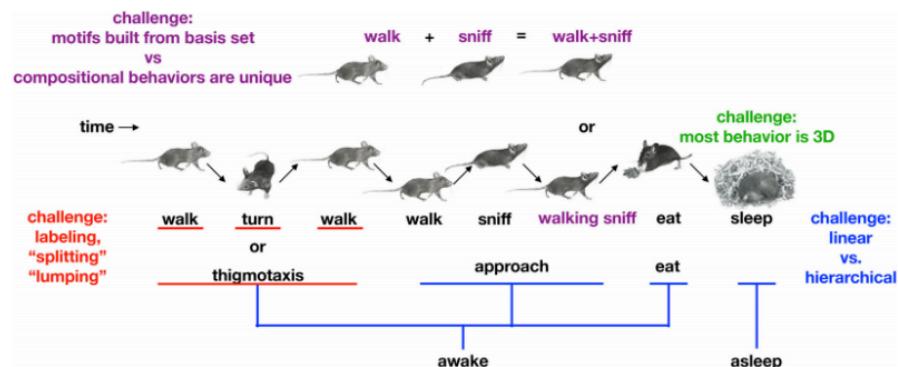
Nikolaas Tinbergen
Nobel Prize in Physiology or Medicine 1973



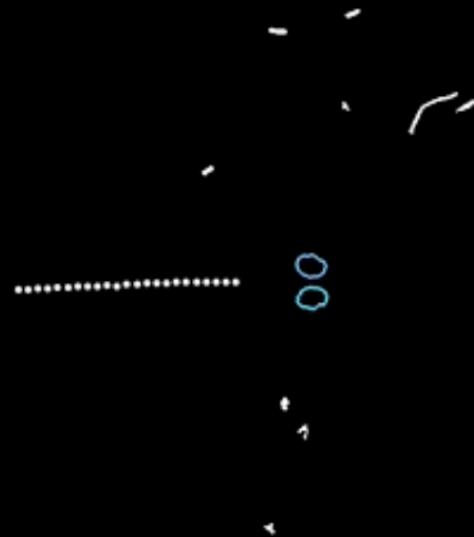
Computational (Neuro)Ethology

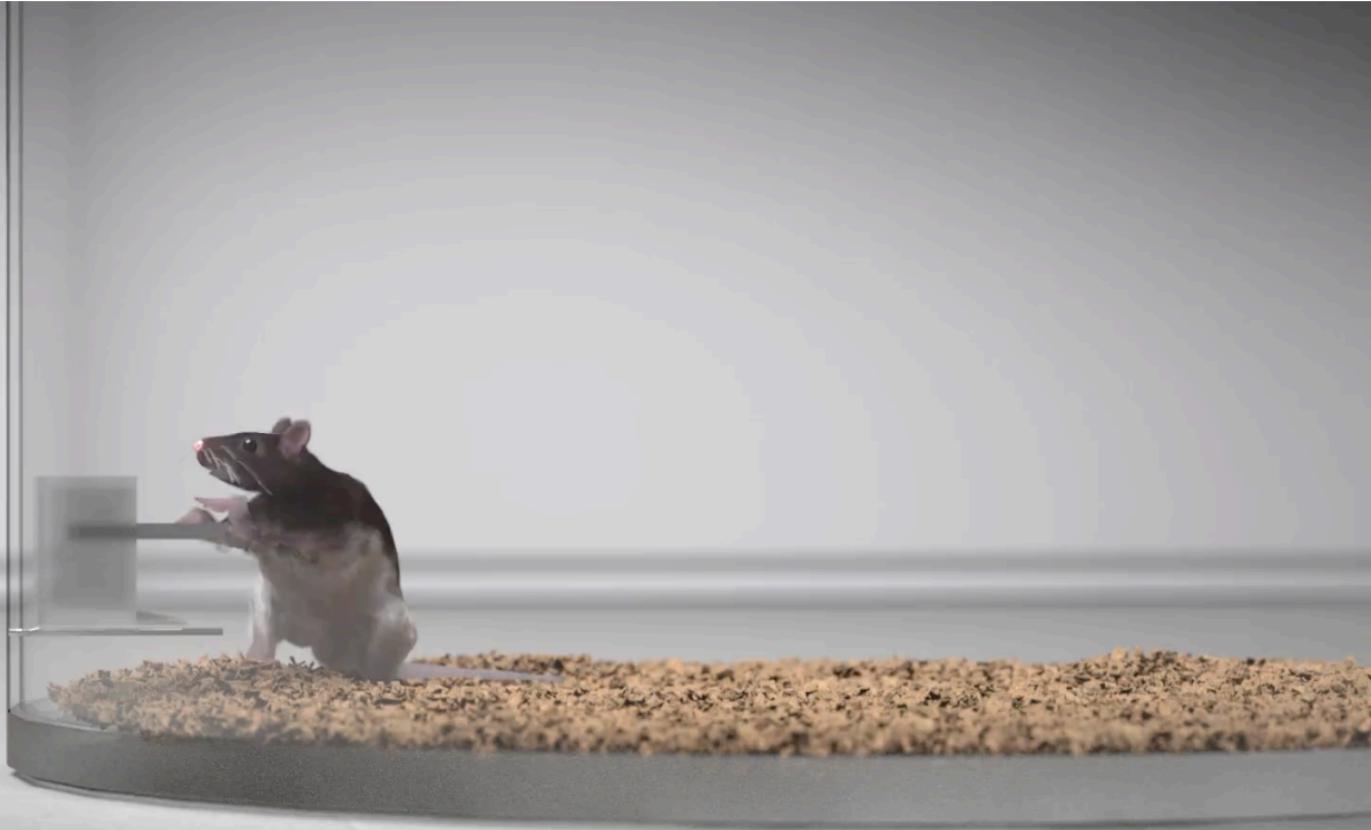
Quantifying natural behavior (and relating it to neural activity)

- Leveraging advances in **computer vision** and **machine learning** to extract behavioral features of interest from raw data.
- Modeling the dynamics of 3D pose as a function of sensory input and internal state.
- Decomposing behavior into stereotyped components and behavioral motifs.
- Correlating behavioral motifs with large scale neural recordings.
- Identifying causal relationships between neural activity and motor output.



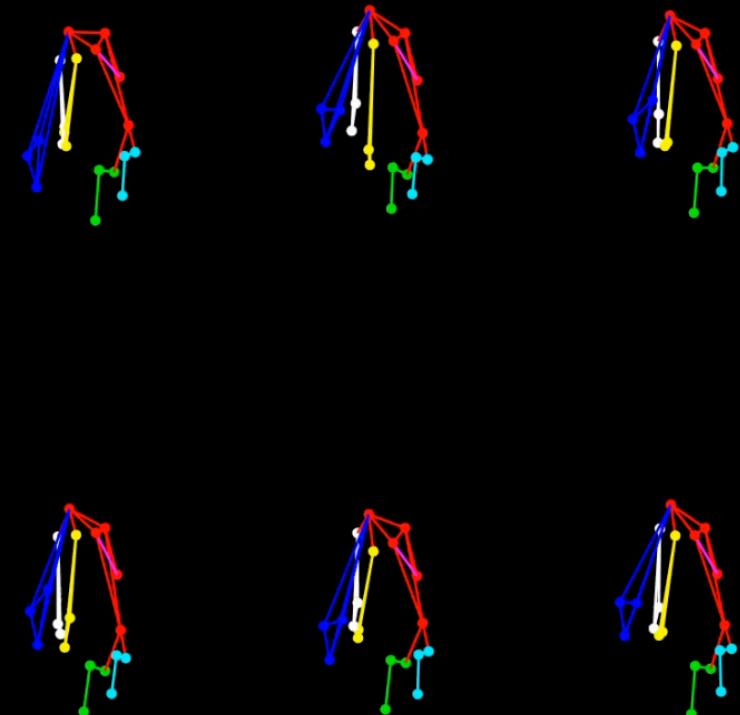
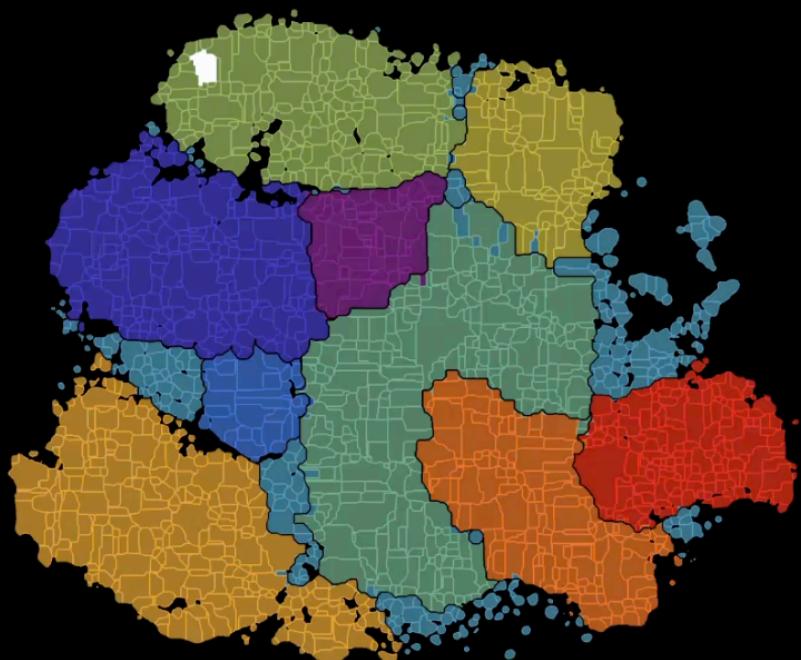
Datta et al (Neuron, 2019)





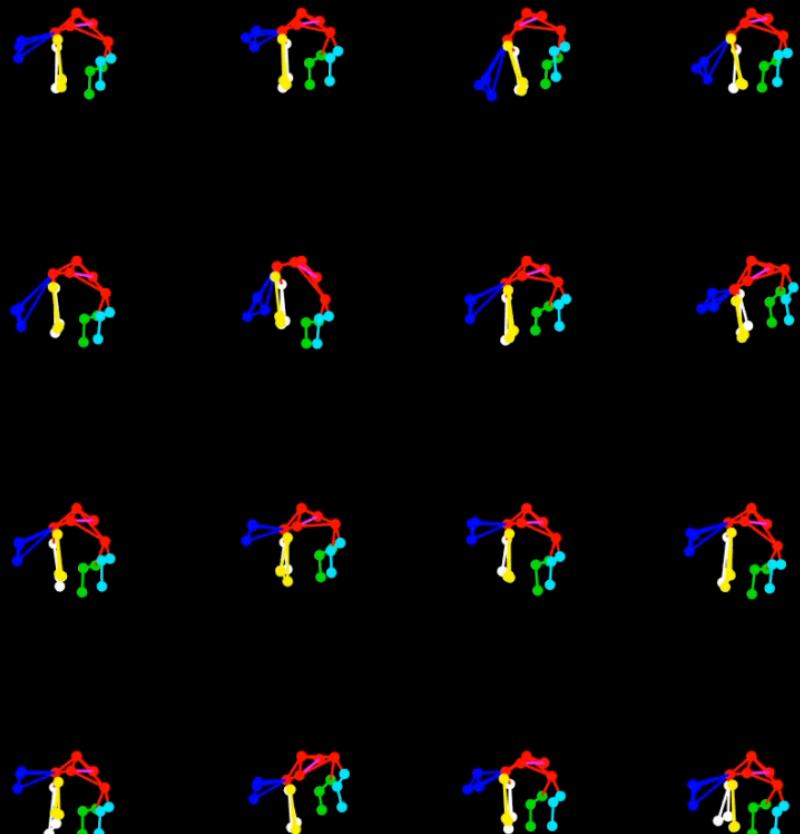
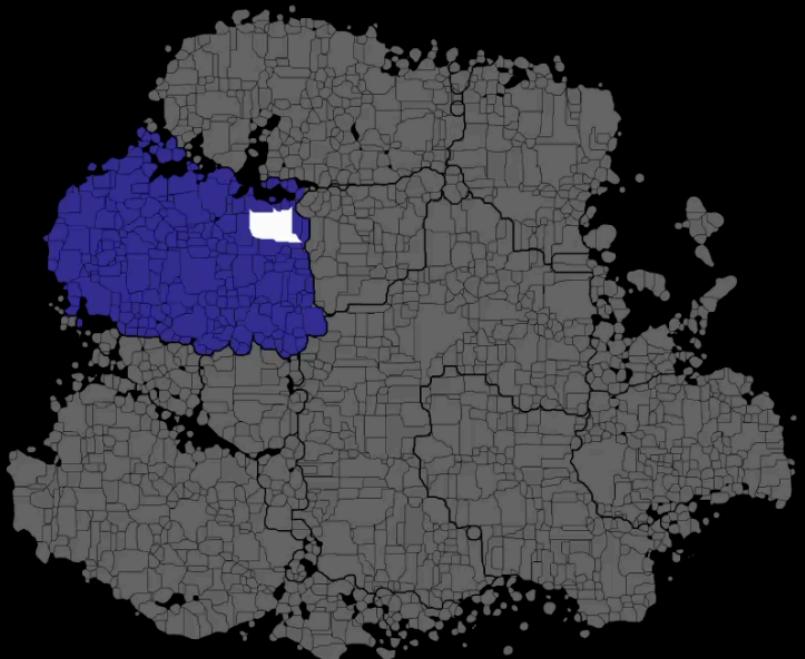
CAPTURE: Marshall et al (*Neuron*, 2020)

Left Groom

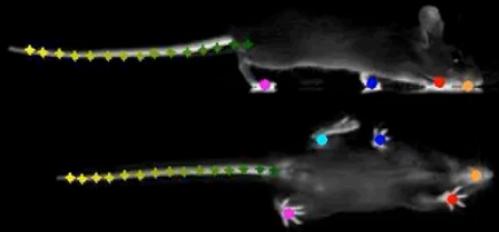
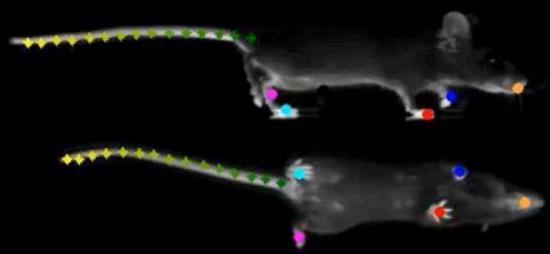


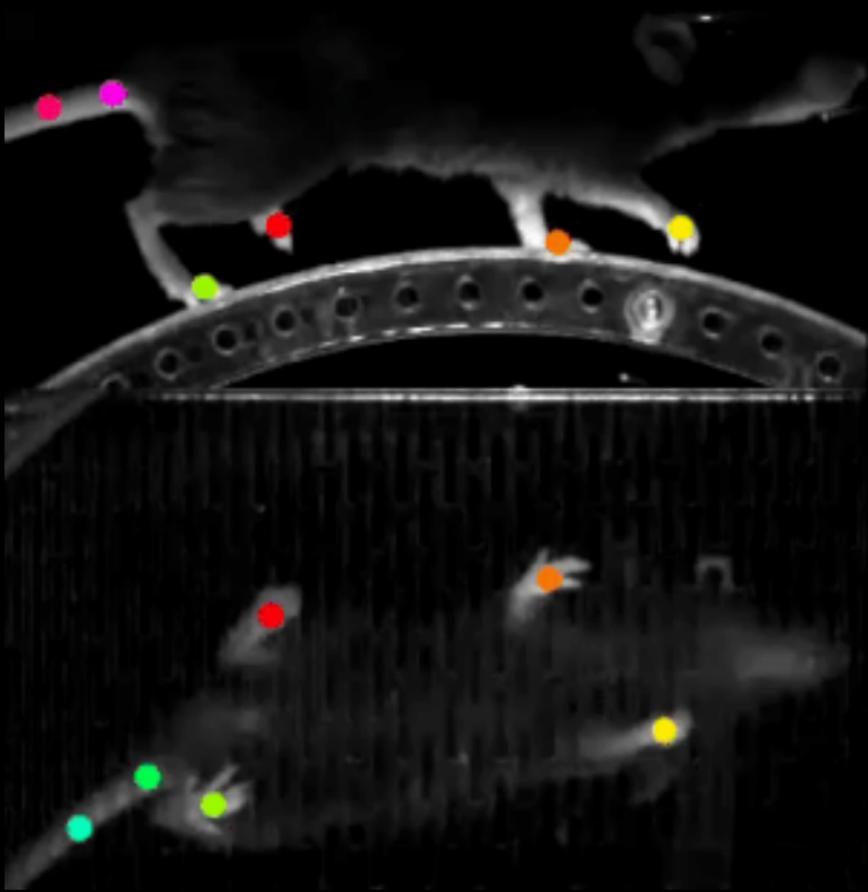
CAPTURE: Marshall et al (*Neuron*, 2020)

Right Groom - Low

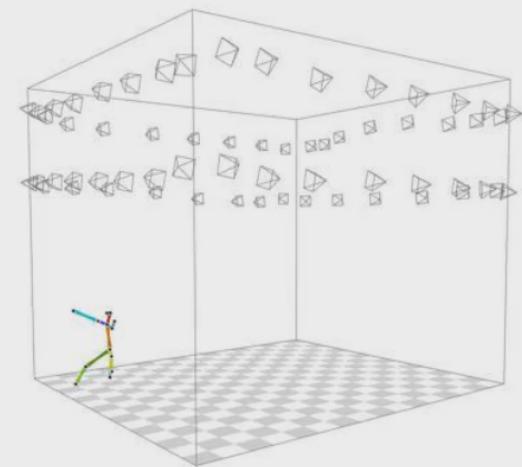
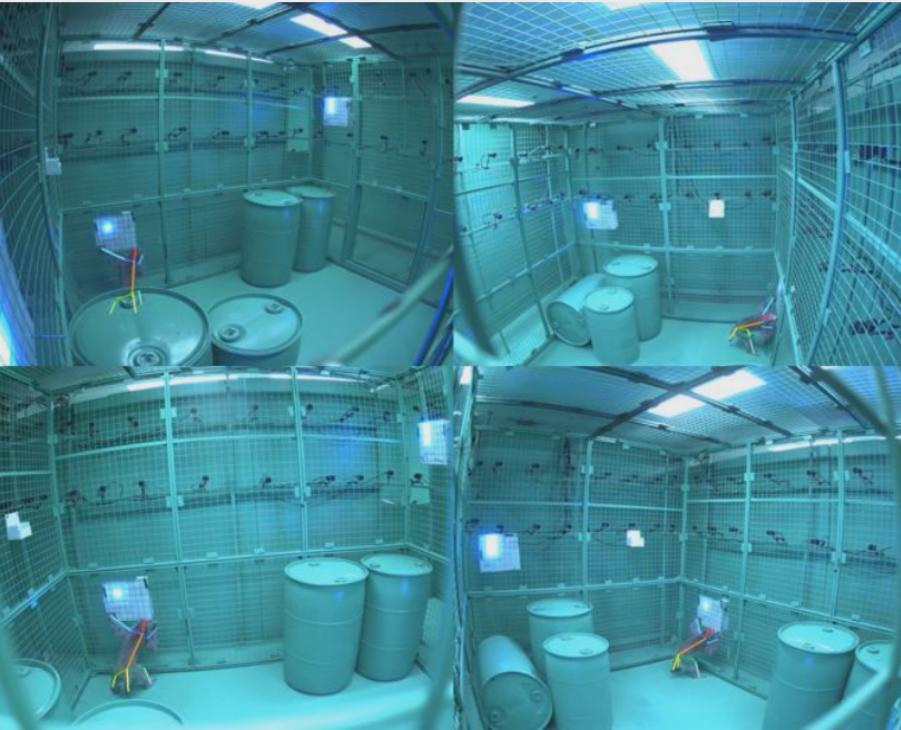


CAPTURE: Marshall et al (Neuron, 2020)

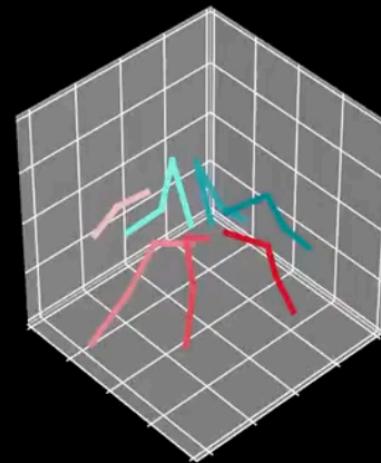
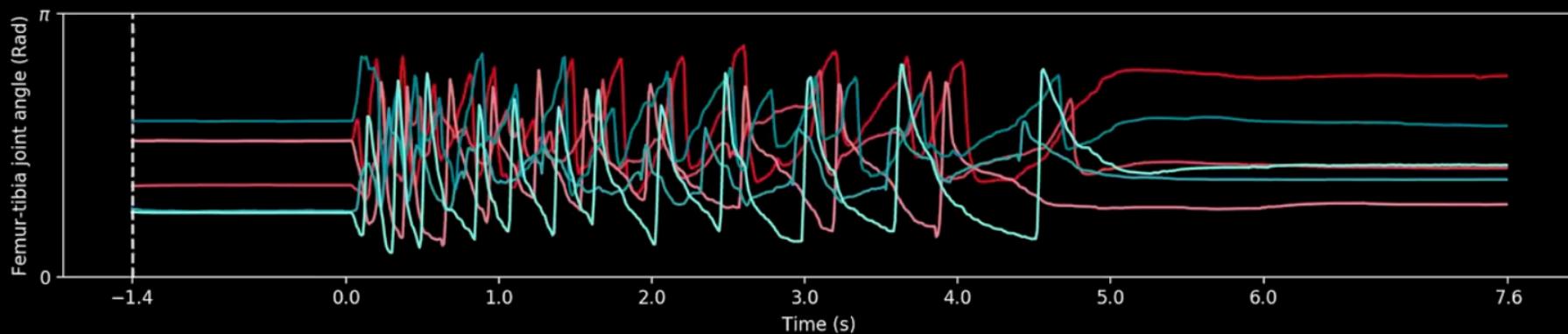


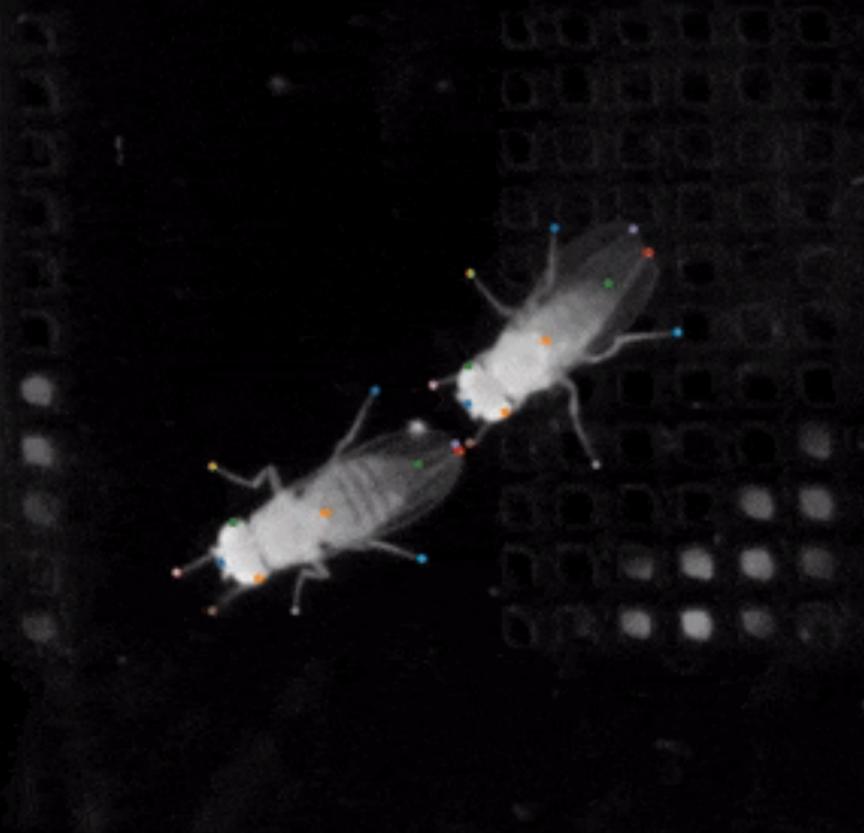


DeepLabCut: Mathis et al. (*Nat Neuro* 2018)



OpenMonkeyStudio: Bala et al (*Nature Comm.*, 2020)





SLEAP: Pereira et al (bioRxiv, 2020)

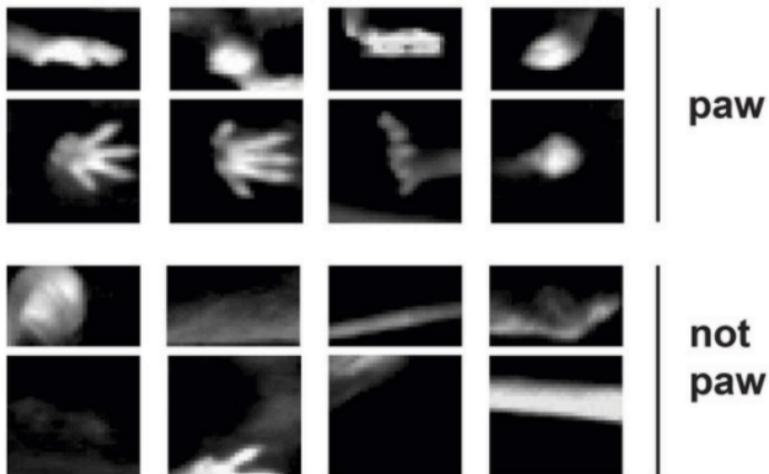
Agenda

1. Basics of markerless pose tracking
2. Transfer learning
3. Triangulating 3D pose from multiple 2D views

Basic pose tracking

Turn it into a supervised learning problem

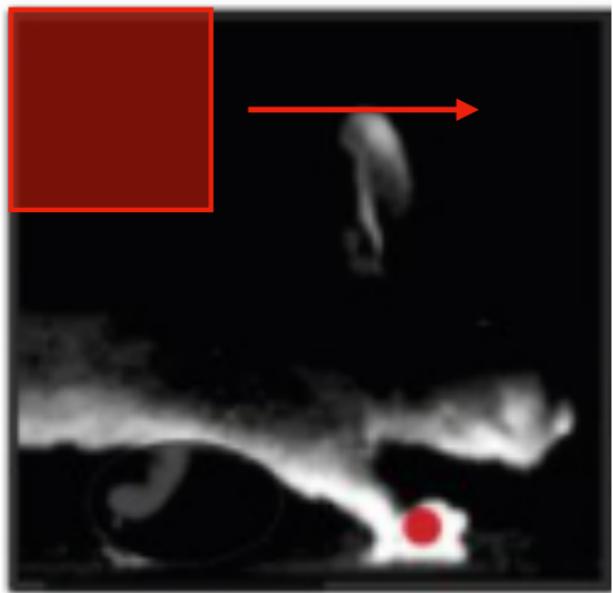
- Extract patches from the video frames and label them as positive or negative examples of a key point (e.g. paw).
- Train a binary classifier (logistic regression, SVM, neural network, etc.) to predict key point or not.



Basic pose tracking

Turn it into a supervised learning problem

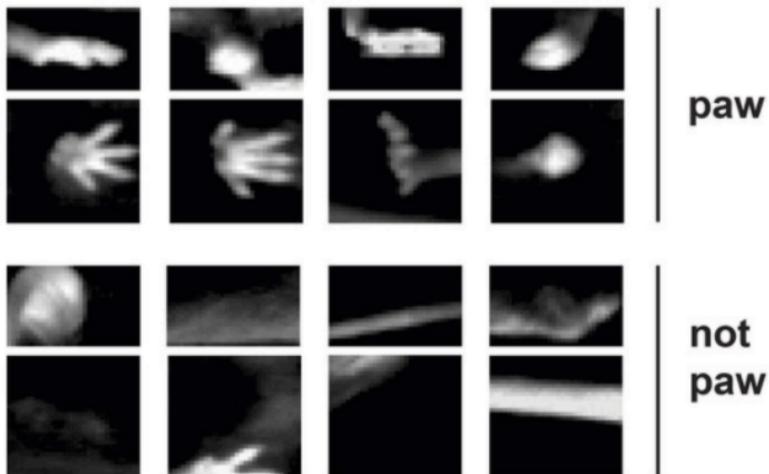
- Extract patches from the video frames and label them as positive or negative examples of a keypoint (e.g. paw).
- Train a binary classifier (logistic regression, SVM, neural network, etc.) to predict keypoint or not.
- At test time, classify each patch in the image and use a heuristic to pick the most likely keypoint location(s).



Basic pose tracking

Mathematical formulation

- Let P_h and P_w be the height and width, respectively, of the patch (in pixels).
- N denote the number of patches
- $X_n \in \mathbb{R}^{P_h \times P_w}$ denote the n -th patch.
- $y_n \in \{0,1\}$ denote whether or not the patch is an instance of the keypoint.
- $W \in \mathbb{R}^{P_h \times P_w}$ denote the weights of our model.

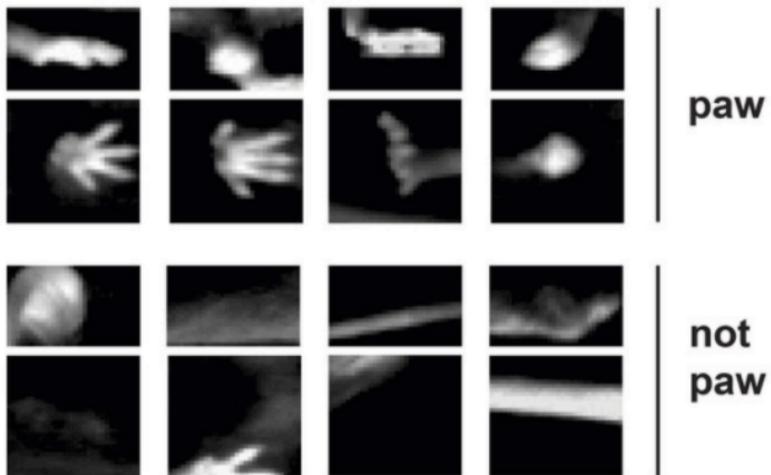


Basic pose tracking

Via logistic regression

Assume

$$p(y_n \mid x_n, w) = \text{Bern}\left(y_n \mid \sigma(\langle W, X_n \rangle)\right)$$



Basic pose tracking

Maximum likelihood estimation

$$\mathcal{L}(W) =$$

$$\nabla \mathcal{L}(w) =$$

Basic pose tracking

With convolutional neural networks

- Instead of working with patches, let's work with images directly.
- Let $X \in \mathbb{R}^{P_H \times P_W}$ denote an image (height P_H , width P_W)
- Let $Y \in \{0,1\}^{P_H \times P_W}$ indicate the location(s) of the keypoint.
- The 2D cross-correlation $X \star W$; is a sliding dot product of weights across all $P_h \times P_w$ patches in the image. It produces a $P_H \times P_W$ output.
- In PyTorch, it's implemented by the `F.conv2d` function and the `Conv2D` layer.

Basic pose tracking

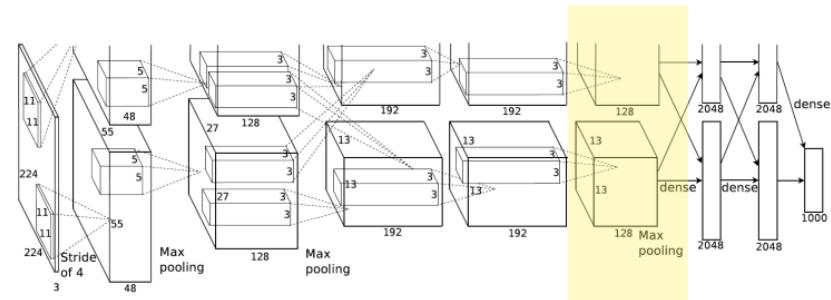
Feature learning in CNNs

- This simple model assumes keypoints can be detected with a **linear classifier** using raw pixels as inputs.
- We can perform **nonlinear classification** by encoding each pixel with a vector of features.
- Rather than handcrafting these features, **learn them** from the data!

Transfer Learning

Transfer Learning

- **Idea:** rather than handcrafting features or learning them from scratch, **use a pre-trained network** for a related task.
- **Example:** use the features of a deep neural network for image classification.
- **Reroute** the output of an intermediate layer to a **new loss function**.
- Optionally, **fine tune** the weights in the early layers via stochastic gradient descent on the new loss.
- With good starting features, you **only need a few training examples** to perform animal pose estimation.

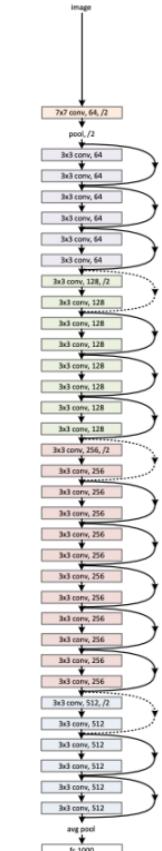


Transfer Learning In DeepLabCut, SLEAP, etc.

- DLC and SLEAP repurpose state-of-the-art deep networks for human pose detection.
- DLC starts with a residual network (resnet-50) and adds “deconvolutional” layers, as in DeeperCut for human pose estimation.
- SLEAP starts with “stacked hourglass networks” for human pose estimation.

Deep Residual Networks
(resnet-50)

34-layer residual



Stacked Hourglass Networks

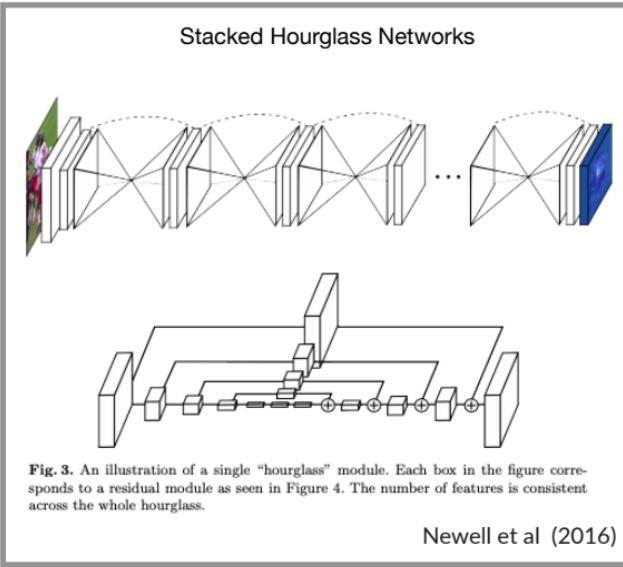
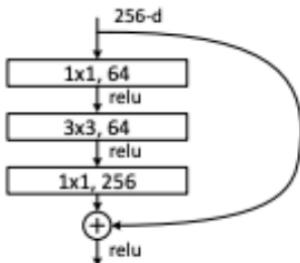


Fig. 3. An illustration of a single “hourglass” module. Each box in the figure corresponds to a residual module as seen in Figure 4. The number of features is consistent across the whole hourglass.

Newell et al (2016)

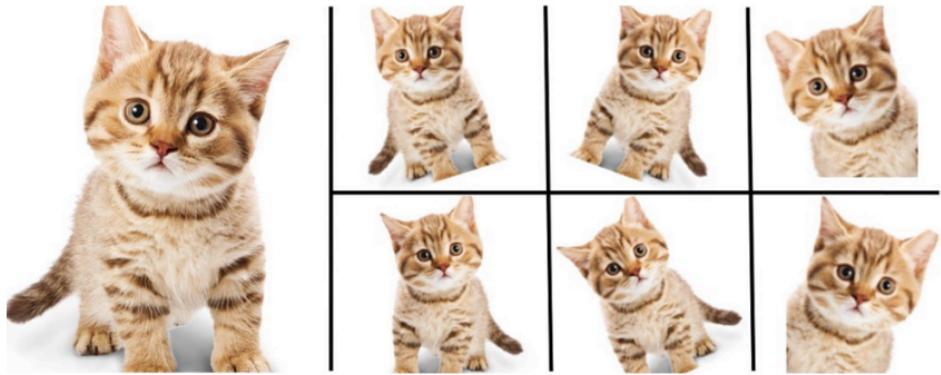


He et al (2015)

Transfer Learning

Data augmentation

- Labeling data is tedious.
- **Idea:** Make the most of each training example by making alterations your classifier should be robust to.
- Eg a cropped, rotated, and scaled paw is still a paw. A partially occluded paw is still a paw.



3D Triangulation

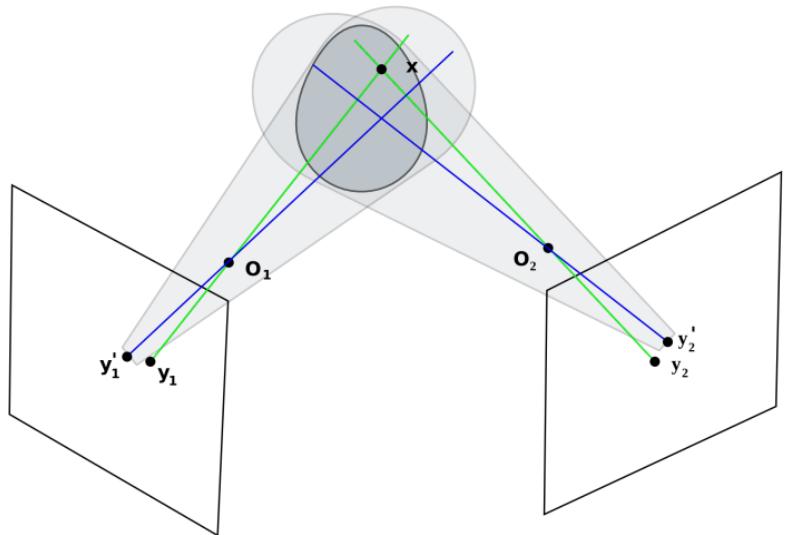
Basic triangulation

Projective geometry

- Projective geometry makes far away objects appear smaller.

$$\vec{y}_c \approx f_c(\vec{x})$$

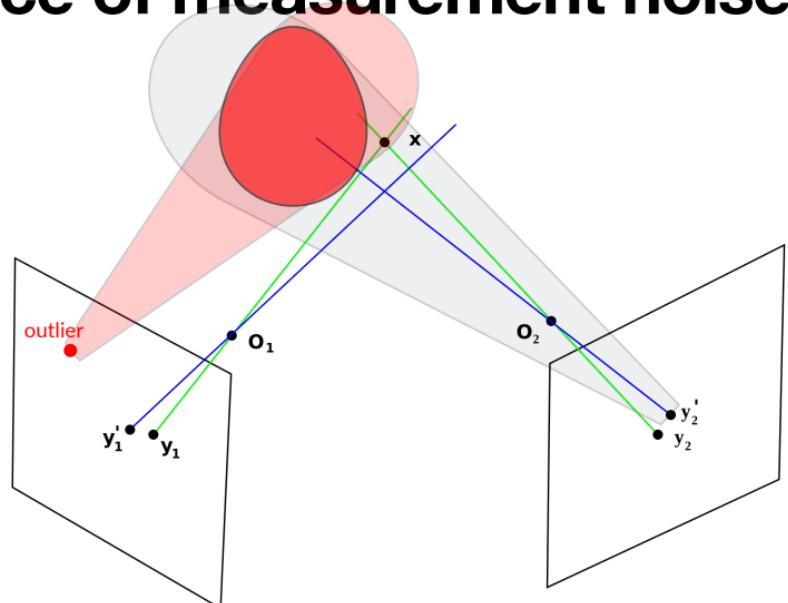
$$f_c(\vec{x}) = \frac{1}{w}(u, v)^\top \text{ where } (u, v, w)^\top = A_c \vec{x} + b_c,$$



Modified from wikipedia.org

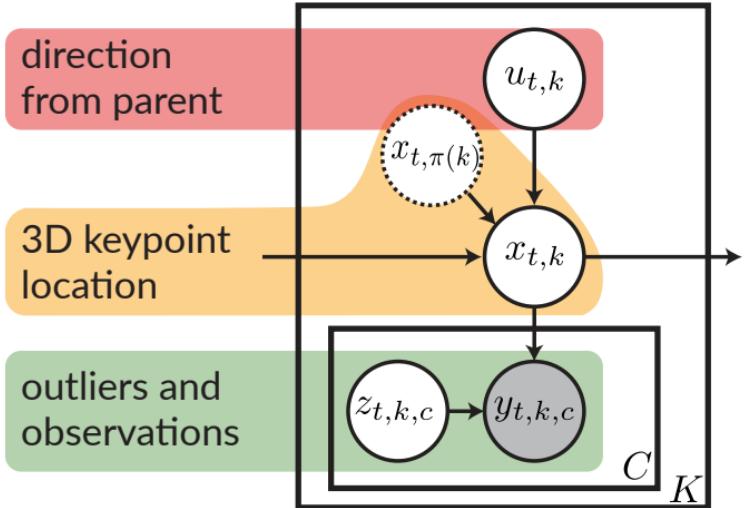
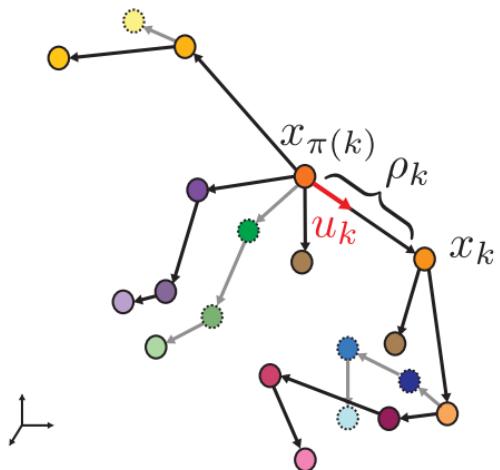
Triangulation in the presence of measurement noise

- Projective geometry makes far away objects appear smaller.
- Outliers in 2D estimates can severely affect 3D triangulation.
- Typical approaches:
 - More data
 - Temporal constraints
 - Median filtering (DLC-3D) / RANSAC
 - Robust noise models



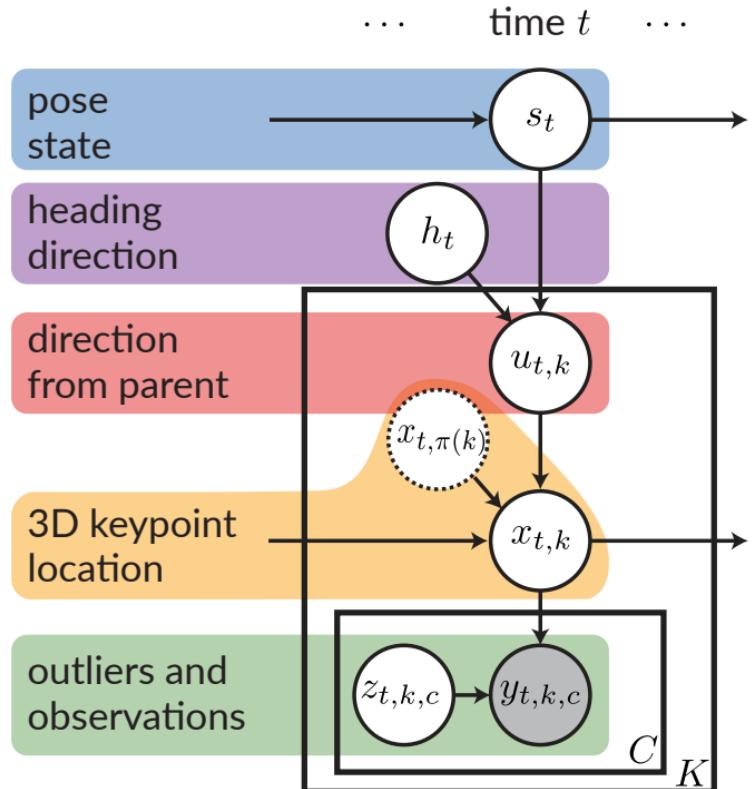
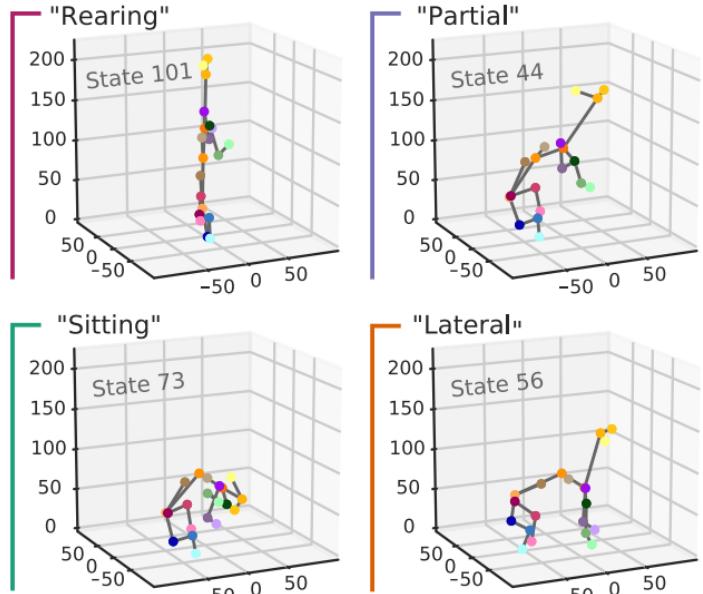
Modified from wikipedia.org

Bayesian triangulation



Zhang et al (AISTATS, 2021)

GIMBAL: Capturing correlations in direction vectors with pose states



Conclusion

- **Precise behavior quantifications** are critical for understanding how neural activity relates to behavioral output.
- **Markerless pose tracking** methods have made it much easier to obtain such quantifications.
- **Convolutional neural networks** are naturally suited to this task.
- With **transfer learning**, we can leverage state-of-the-art deep networks for image classification to warm-start pose tracking.
- We can **triangulate 3D pose** from 2D images using projecting geometry and spatiotemporal priors.

Further reading

- Datta, Sandeep Robert, et al. "Computational neuroethology: a call to action." *Neuron* 104.1 (2019): 11-24.
- Mathis, Alexander, et al. "DeepLabCut: markerless pose estimation of user-defined body parts with deep learning." *Nature neuroscience* 21.9 (2018): 1281-1289.
- Pereira, Talmo D., et al. "Fast animal pose estimation using deep neural networks." *Nature methods* 16.1 (2019): 117-125.
- He, Kaiming, et al. "Deep residual learning for image recognition." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.