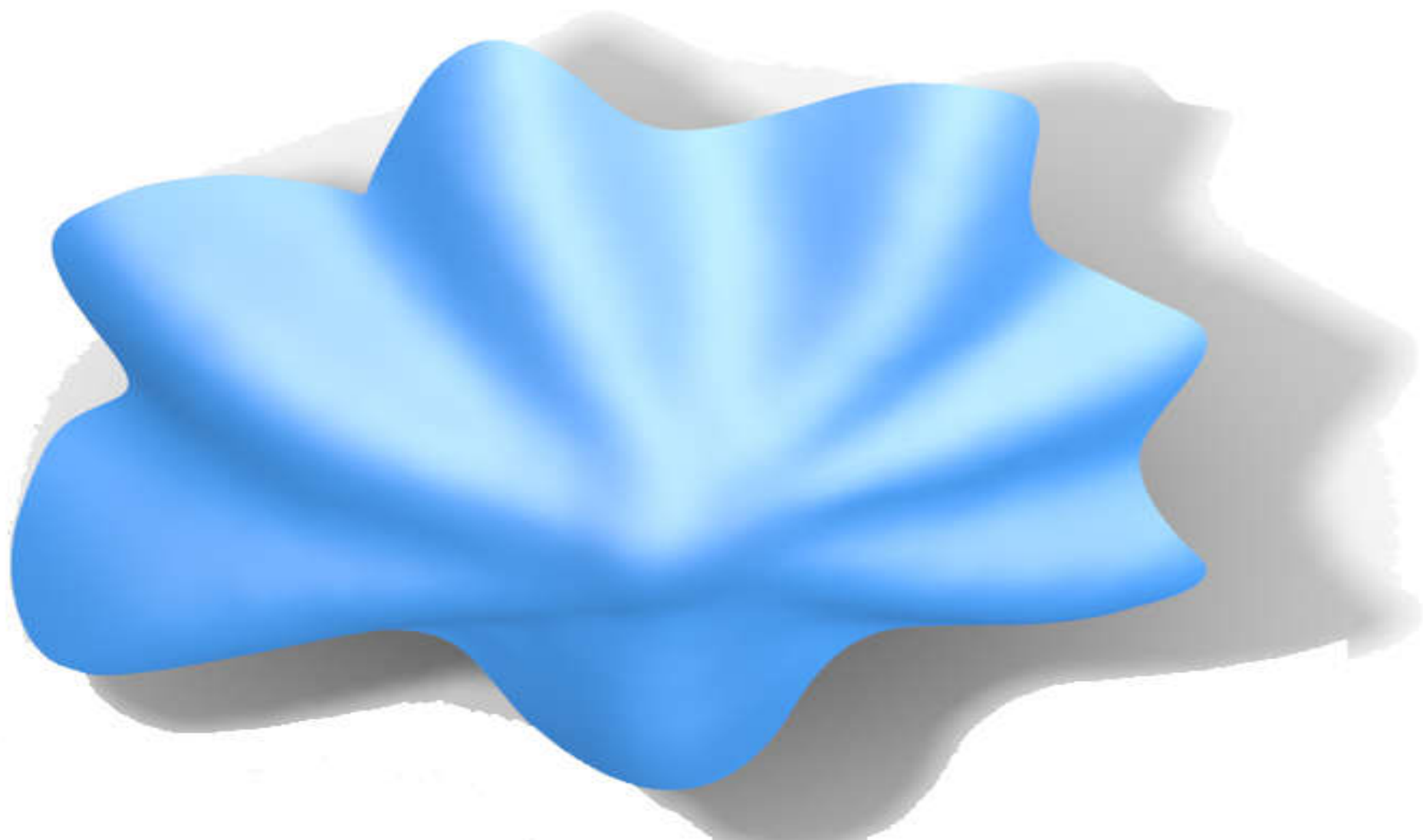
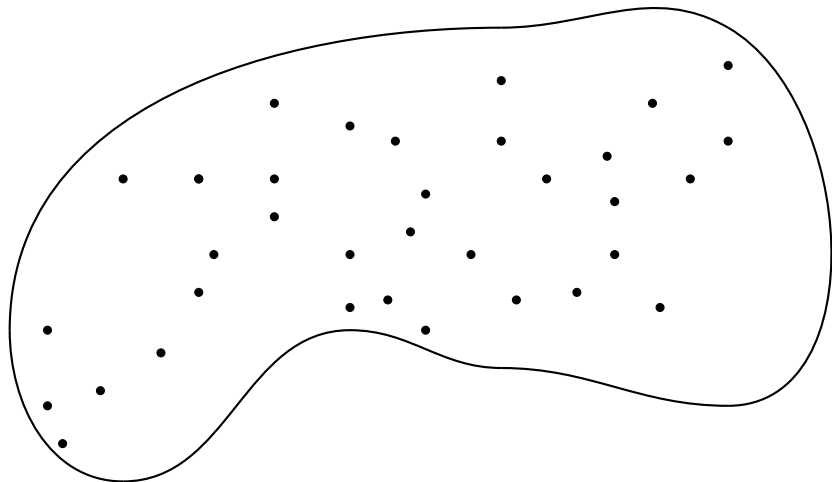


# Discovering low-dimensional manifolds in high-dimensional data sets

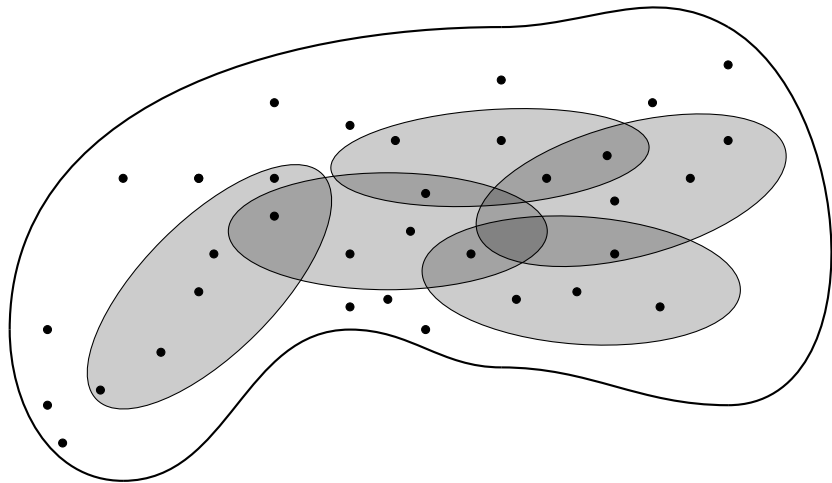


## Diffusion Maps: “Knit together” local geometry to get “better” distances



Small distances are much more reliable!

## Diffusion Maps: “knitting together” local geometry

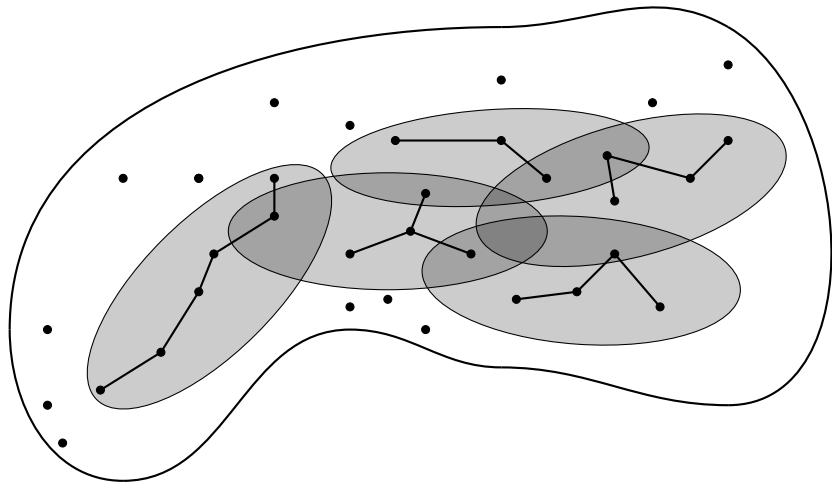




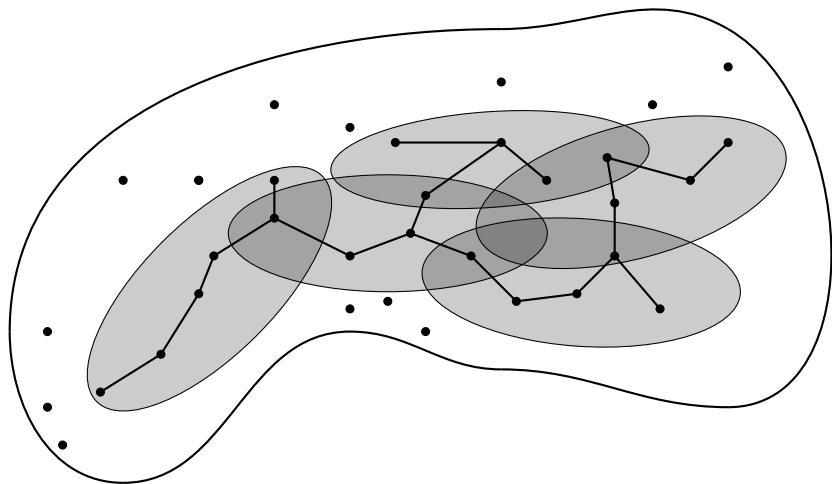
A collection of small tangent patches does give a good first approximation of a surface



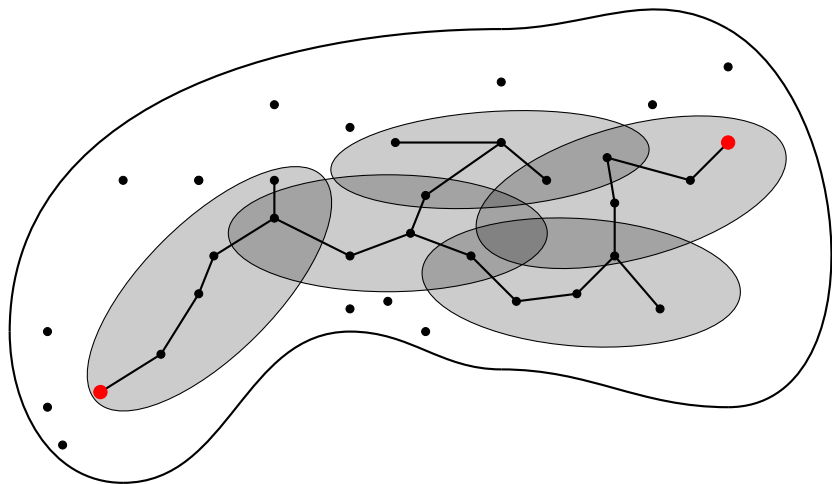
# Diffusion Maps: “knitting together” local geometry



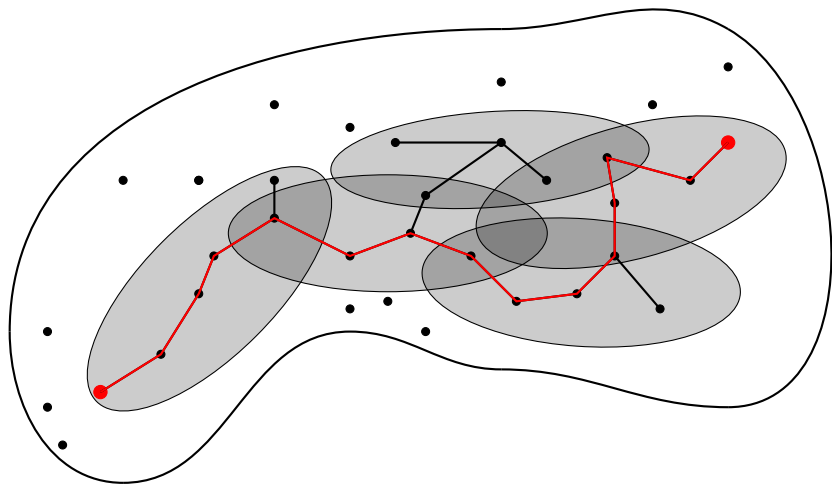
# Diffusion Maps: “knitting together” local geometry

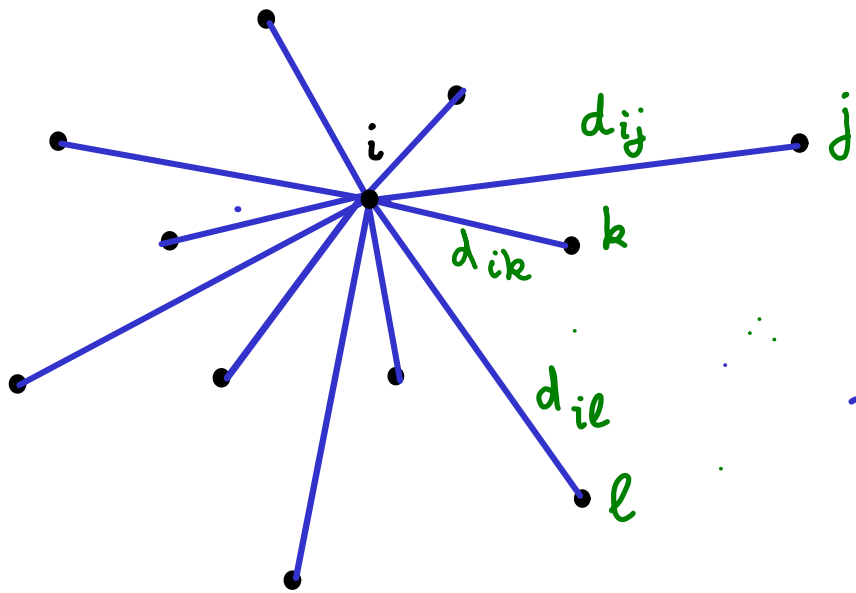


# Diffusion Maps: “knitting together” local geometry



## Diffusion Maps: “knitting together” local geometry





$$e^{-d_{ij}^2/2\tau} = W_{ij;\tau}$$

$$D_{i,i;\tau} = \sum_{j \neq i} W_{ij;\tau}$$

$D_{\tau}^{-1} W_{\tau}$  : defines a random walk on graph.

How pick  $\tau$ ? want: approximation to diffusion on manifold

"true" diffusion: semi-group property

$$e^{-tL} e^{-sL} = e^{-(t+s)L}$$

Shan Shan

$\Rightarrow$  pick  $\tau$  so that  $(D_{\tau}^{-1} W_{\tau})^k \approx D_{k\tau}^{-1} W_{k\tau}$

spectral decomposition of  $D_{\tau}^{-1} W_{\tau}$

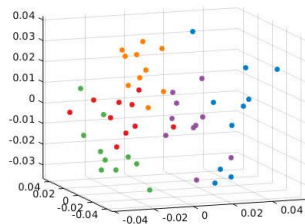
$\Rightarrow$  eigenvectors  $\psi_l$  with  
eigenvalues  $\lambda_{l, \tau}$

$$(e^{-\delta \tau L})_{ij} \simeq \sum_{l=1}^N (\lambda_{l, \tau})^{\delta} \psi_l(i) \psi_l(j)$$

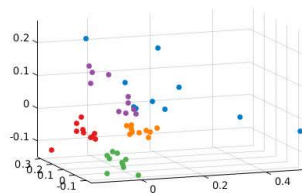
each data point  $j$  is mapped to the  
feature vector  $(\lambda_{l, \tau})^{\delta/2} \psi_l(j)_{l=1}^L$

$$D_{ij; t}^{\delta} = \sum_{l=1}^N (\lambda_{l, t})^t |\psi_l(i) - \psi_l(j)|^2$$

# MDS for cPD & DD



cPD



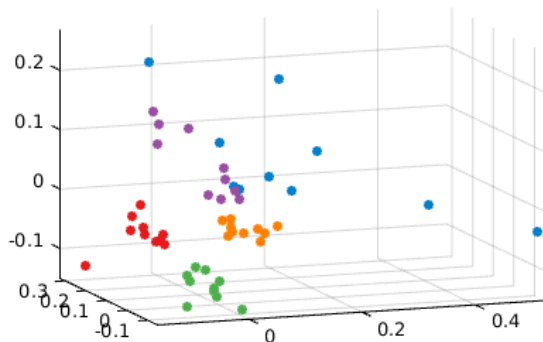
DD



# Diffusion Distance (DD)

Fix  $1 \leq m \leq N$ ,  $t \geq 0$ ,

$$D_m^t(S_i, S_j) = \left( \sum_{k=1}^m \lambda_k^t (u_k(i) - u_k(j))^2 \right)^{\frac{1}{2}}$$



It all started with a conversation with biologists....



Doug Boyer



Jukka Jernvall

More Precisely: biological morphologists



Study Teeth & Bones of  
extant & extinct animals

still live today

fossils

# Collaborators



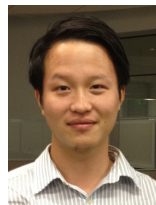
Rima Alaifari  
ETH Zürich



Doug Boyer  
Duke



Ingrid Daubechies  
Duke



Tingran Gao  
Duke



Yaron Lipman  
Weizmann



Roi Poranne  
ETH Zürich



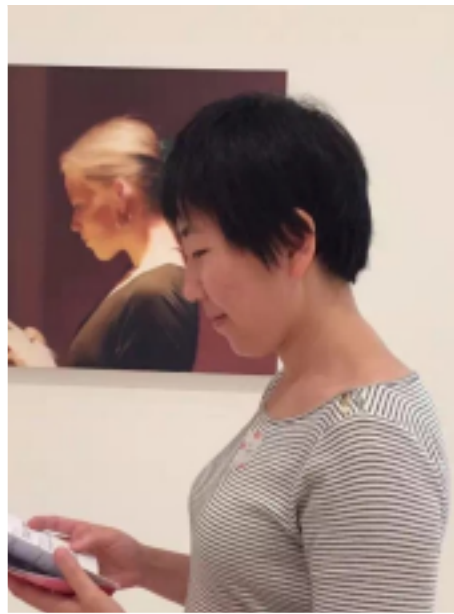
Jesús Puente  
J.P. Morgan



Robert Ravier  
Duke



Shahar Kovalsky



Shan Shan



Nadav Dym



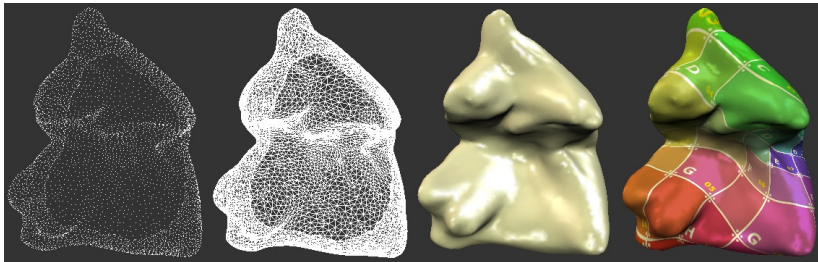
Chen-Yun Lin

First: project on “complexity” of teeth

First: project on “complexity” of teeth

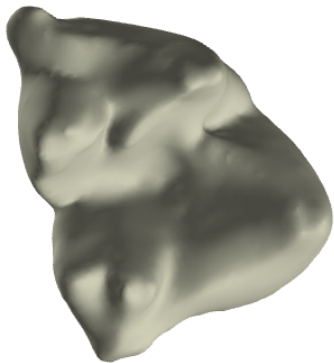
Then: find automatic way to compute Procrustes distances  
between surfaces — without landmarks

# Data Acquisition



Surface reconstructed from  $\mu$ CT-scanned voxel data

# Geometric Morphometrics

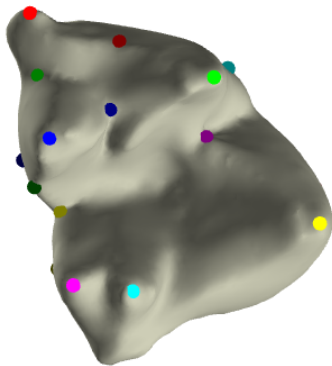


second mandibular molar of a Philippine flying lemur

- Manually put  $k$  landmarks



# Geometric Morphometrics

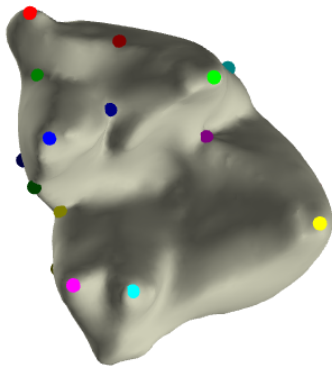


second mandibular molar of a Philippine flying lemur

- Manually put  $k$  landmarks

$$p_1, p_2, \dots, p_k$$

# Geometric Morphometrics



second mandibular molar of a Philippine flying lemur

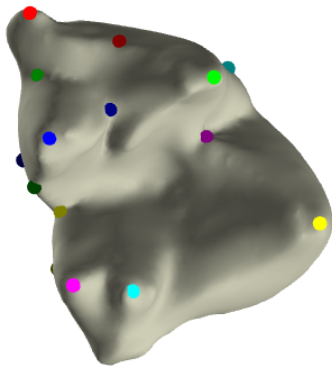
- Manually put  $k$  landmarks

$$p_1, p_2, \dots, p_k$$

- Use spatial coordinates of the landmarks as features

$$p_j = (x_j, y_j, z_j), j = 1, \dots, k$$

# Geometric Morphometrics



second mandibular molar of a Philippine flying lemur

- Manually put  $k$  **landmarks**

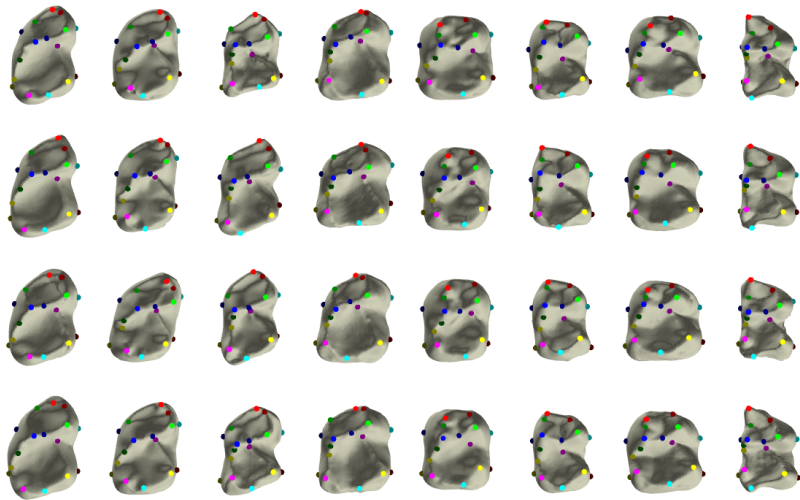
$$p_1, p_2, \dots, p_k$$

- Use **spatial** coordinates of the landmarks as features

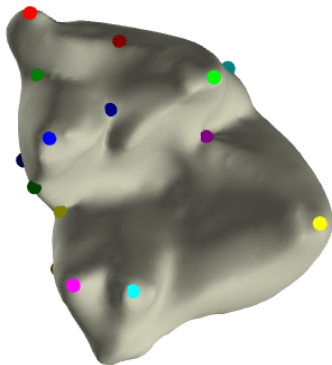
$$p_j = (x_j, y_j, z_j), j = 1, \dots, k$$

- Represent a shape in  $\mathbb{R}^{3 \times k}$

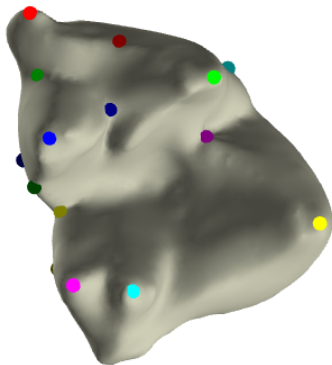
# The *Shape Space* of $k$ landmarks in $\mathbb{R}^3$



# Geometric Morphometrics: Limitation of Landmarks

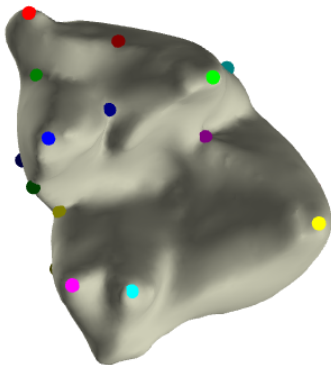


# Geometric Morphometrics: Limitation of Landmarks



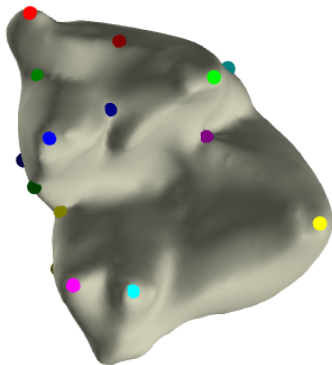
- **Landmark Placement:** tedious and time-consuming

# Geometric Morphometrics: Limitation of Landmarks



- **Landmark Placement:** tedious and time-consuming
- **Fixed Number of Landmarks:** lack of flexibility

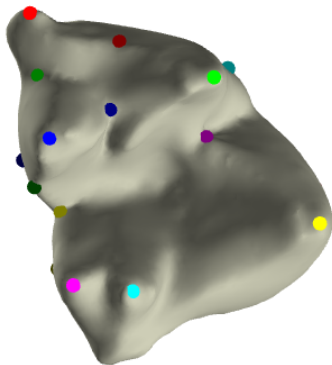
# Geometric Morphometrics: Limitation of Landmarks



- **Landmark Placement:** tedious and time-consuming
- **Fixed Number of Landmarks:** lack of flexibility
- **Domain Knowledge:** high degree of expertise needed, not easily accessible



# Geometric Morphometrics: Limitation of Landmarks



- **Landmark Placement:** tedious and time-consuming
- **Fixed Number of Landmarks:** lack of flexibility
- **Domain Knowledge:** high degree of expertise needed, not easily accessible
- **Subjectivity:** debates exist even among experts

First: project on “complexity” of teeth

Then: find automatic way to compute Procrustes distances  
between surfaces — without landmarks



Landmarked Teeth  $\longrightarrow$

$$d_{Procrustes}^2(S_1, S_2) = \min_{R \text{ rigid tr.}} \sum_{j=1}^J \|R(x_j) - y_j\|^2$$

First: project on “complexity” of teeth

Then: find automatic way to compute Procrustes distances  
between surfaces — without landmarks



Landmarked Teeth  $\longrightarrow$

$$d_{Procrustes}^2(S_1, S_2) = \min_{R \text{ rigid tr.}} \sum_{j=1}^J \|R(x_j) - y_j\|^2$$



Find way to compute a distance that does as well,  
for biological purposes, as Procrustes distance,  
based on expert-placed landmarks, automatically?



First: project on “complexity” of teeth

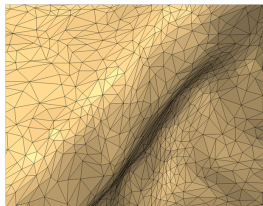
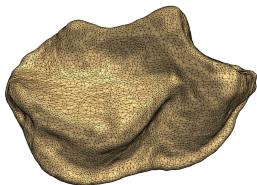
Then: find automatic way to compute Procrustes distances between surfaces — without landmarks

Landmarked Teeth  $\longrightarrow$

$$d_{Procrustes}^2(S_1, S_2) = \min_{R \text{ rigid tr.}} \sum_{j=1}^J \|R(x_j) - y_j\|^2$$

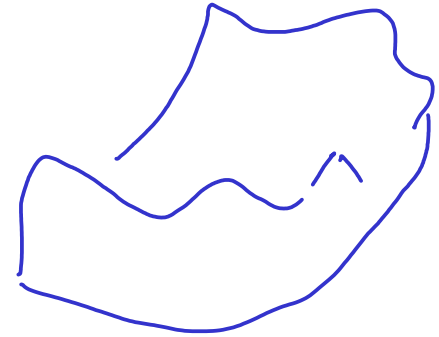
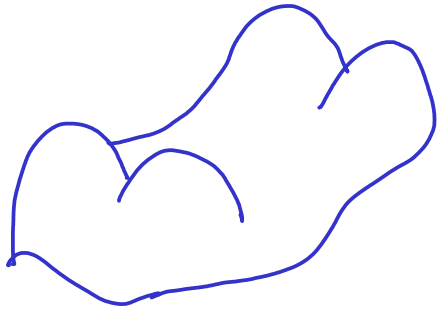
Find way to compute a distance that does as well, for biological purposes, as Procrustes distance, based on expert-placed landmarks, **automatically**?

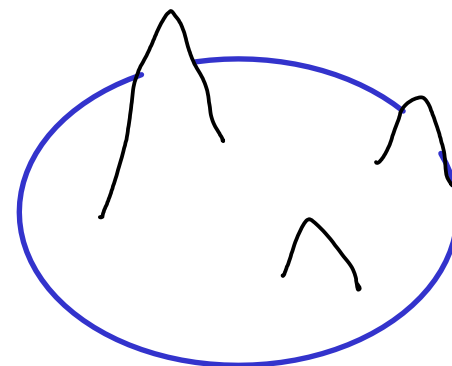
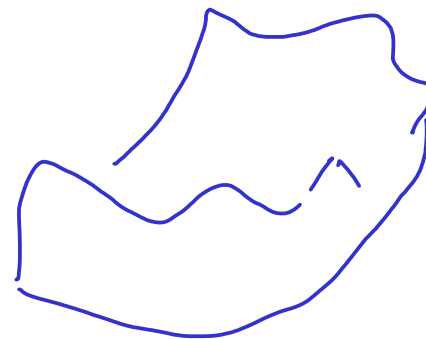
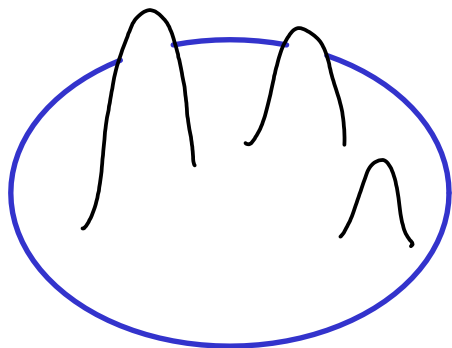
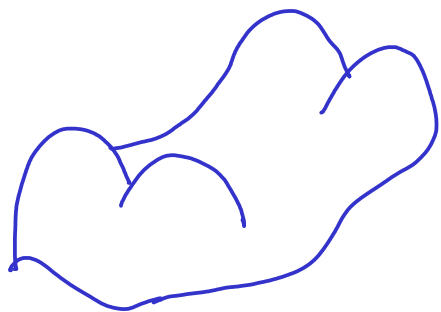
examples: finely discretized triangulated surfaces

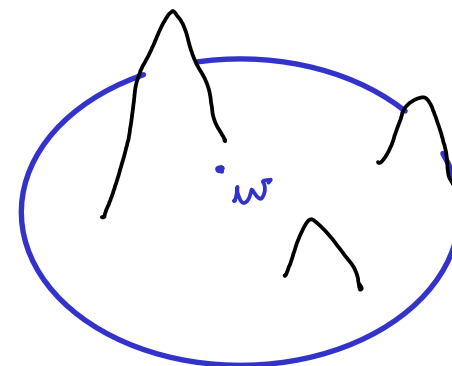
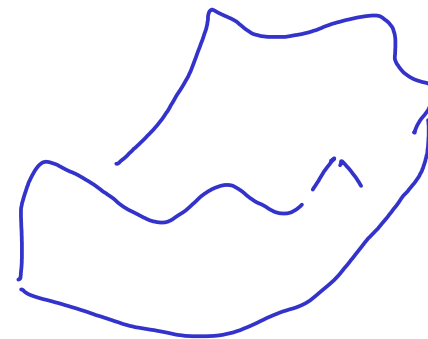
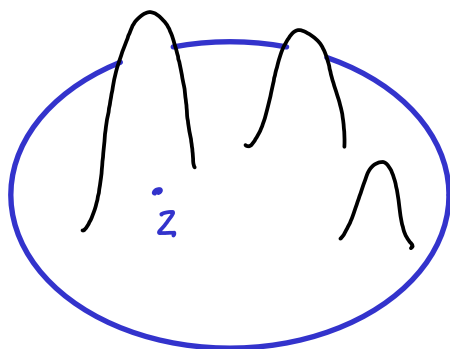
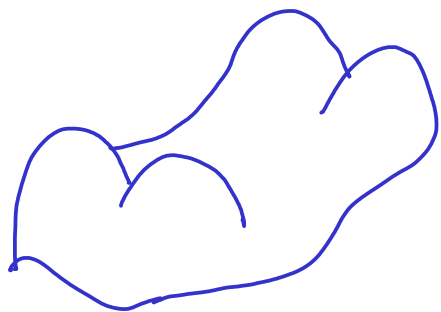


## We defined 2 different distances

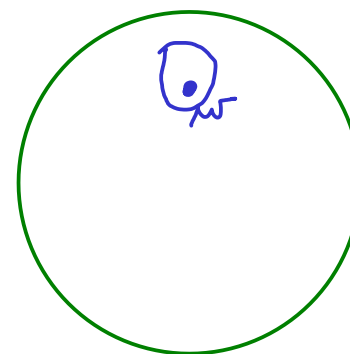
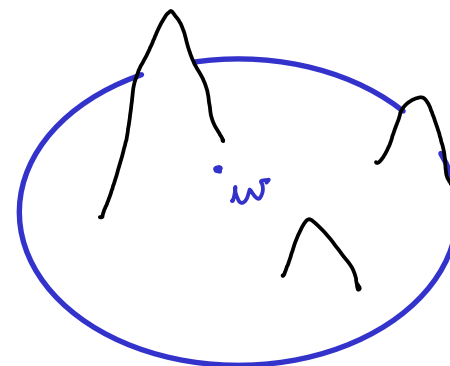
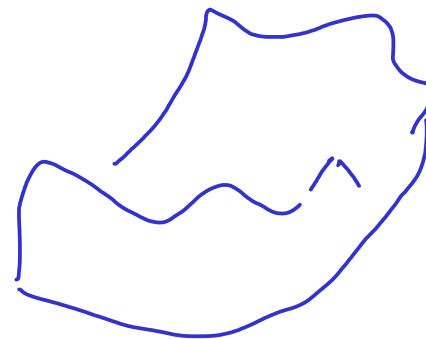
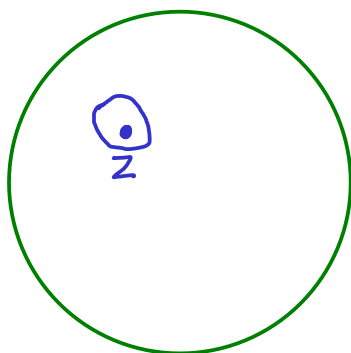
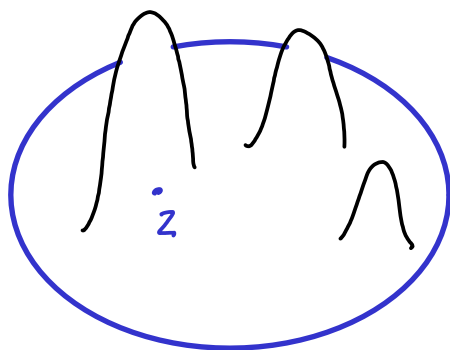
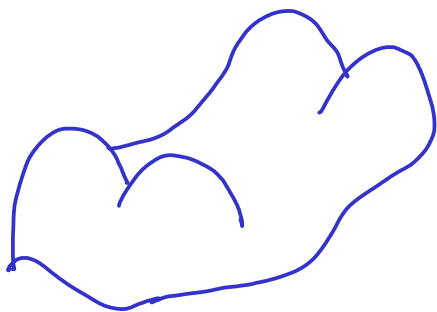
- $d_{cWn}(S_1, S_2)$ :
  - conformal flattening
  - comparison of neighborhood geometry
  - optimal mass transport
- $d_{cP}(S_1, S_2)$ : continuous Procrustes distance

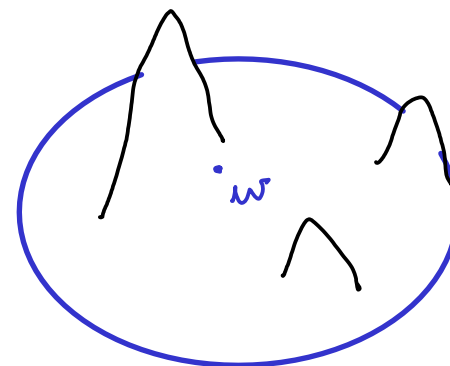
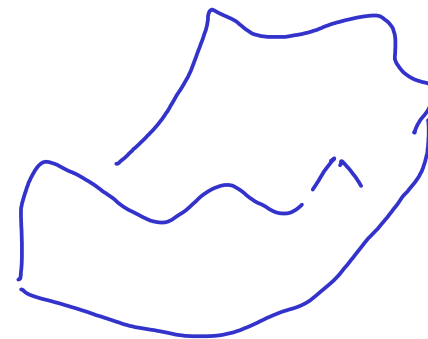
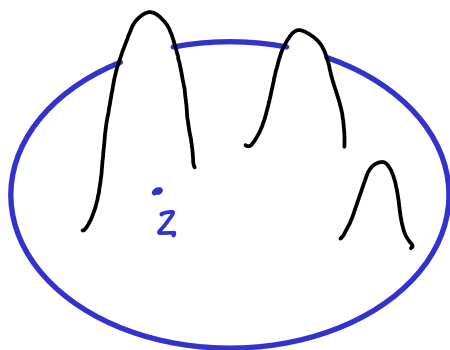
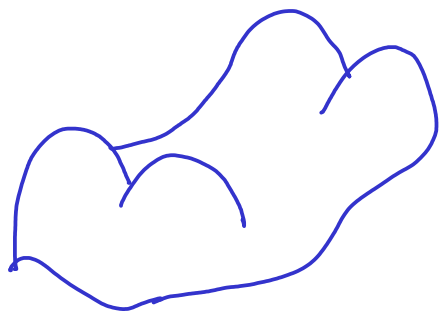




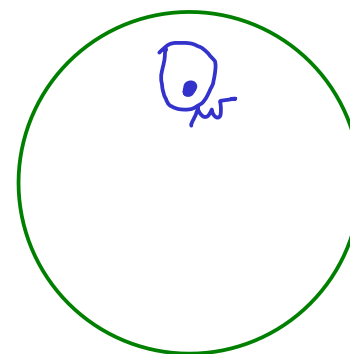
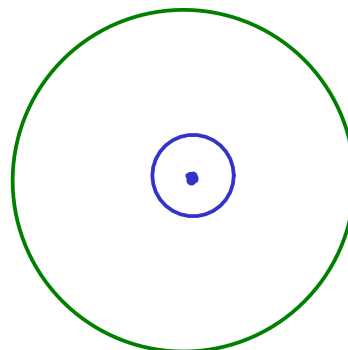
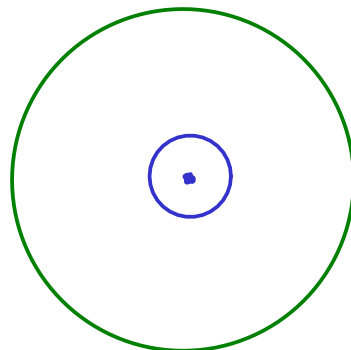
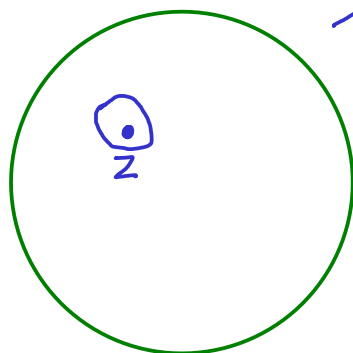


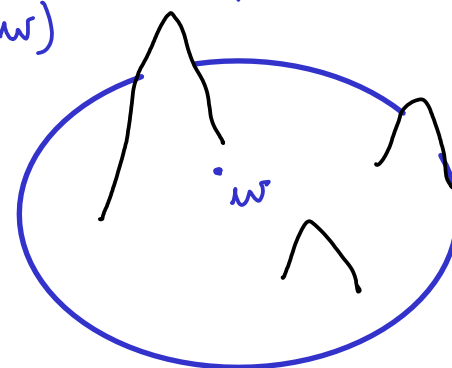
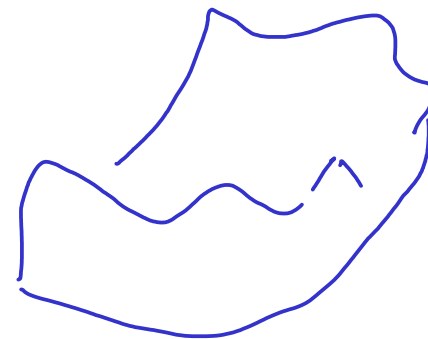
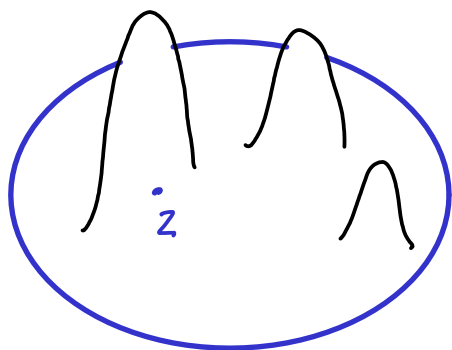
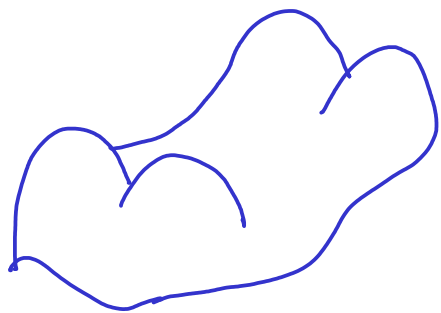






$$d_{R}^{\mu, \nu}(z, w)$$

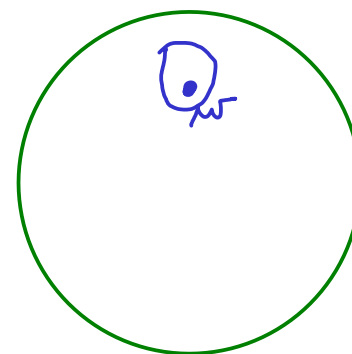
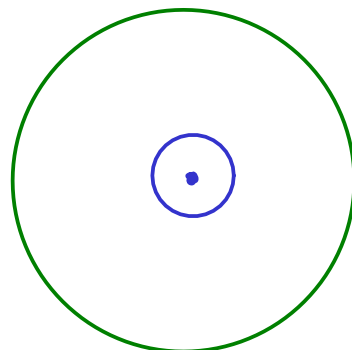
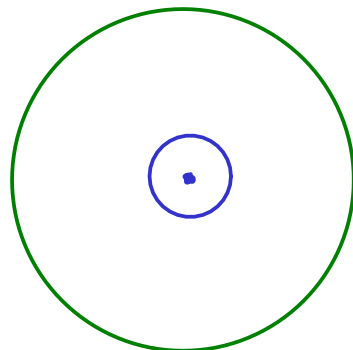
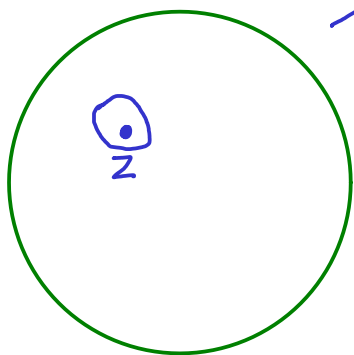




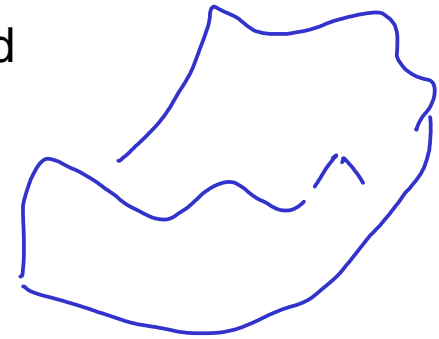
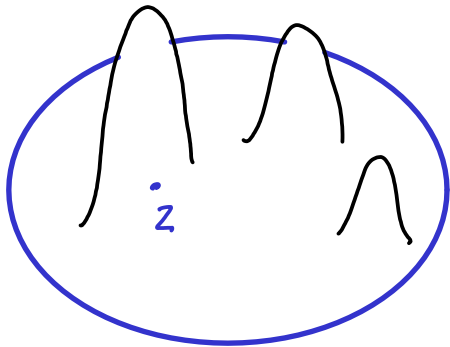
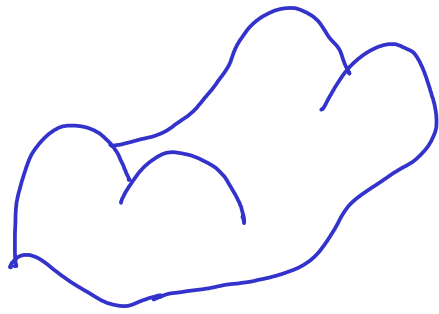
$$\mathbb{D}(S_1, S_2) = \inf_{\pi \in \Pi(\mu, \nu)} \int d_R^{\mu, \nu}(z, w) d\pi(z, w)$$



$$d_R^{\mu, \nu}(z, w)$$



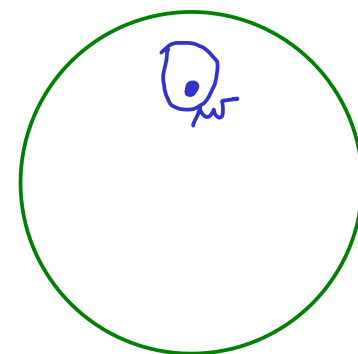
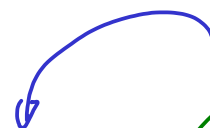
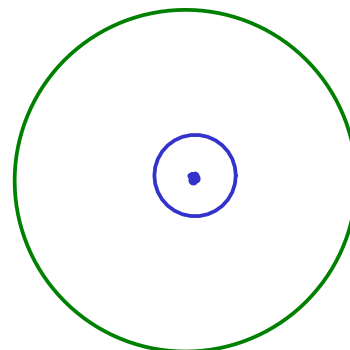
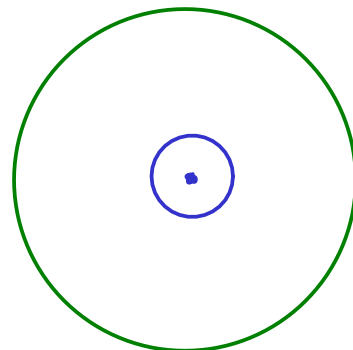
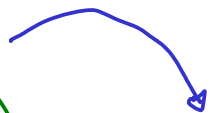
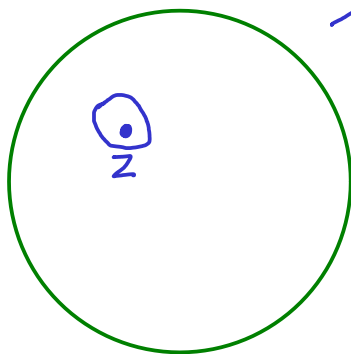
conformal Wasserstein neighborhood distance



$$\mathcal{D}(S_1, S_2) = \inf_{\pi \in \Pi(\mu, \nu)} \int d_R^{\mu, \nu}(z, w) d\pi(z, w)$$



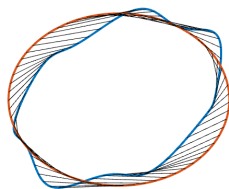
$$d_R^{\mu, \nu}(z, w)$$



# Continuous Procrustes Distance (cPD)

$$D_{\text{cP}}(S_1, S_2) = \left( \int_{S_1} \|x - \mathcal{C}(x)\|^2 d\text{vol}_{S_1}(x) \right)^{\frac{1}{2}},$$

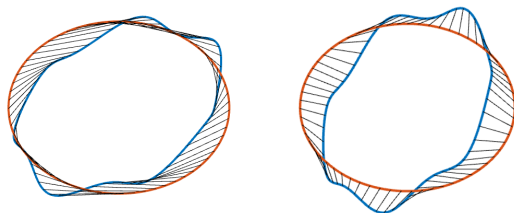
where  $\mathcal{C} : S_1 \rightarrow S_2$  is an area-preserving diffeomorphism.



# Continuous Procrustes Distance (cPD)

$$D_{\text{cP}}(S_1, S_2) = \left( \inf_{R \in \mathbb{E}(3)} \int_{S_1} \|R(x) - \mathcal{C}(x)\|^2 d\text{vol}_{S_1}(x) \right)^{\frac{1}{2}},$$

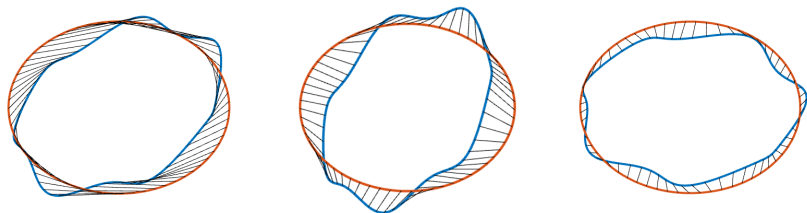
where  $\mathcal{C} : S_1 \rightarrow S_2$  is an **area-preserving diffeomorphism**, and  $\mathbb{E}_3$  is the Euclidean group on  $\mathbb{R}^3$ .



# Continuous Procrustes Distance (cPD)

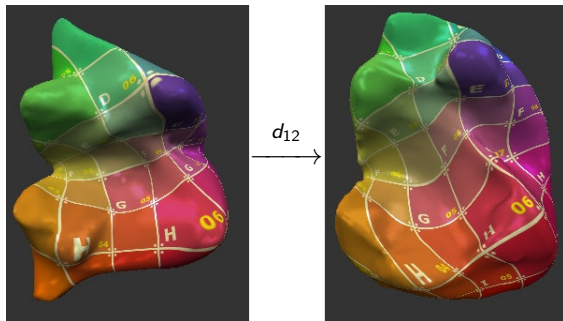
$$D_{\text{cP}}(S_1, S_2) = \left( \inf_{\mathcal{C} \in \mathcal{A}(S_1, S_2)} \inf_{R \in \mathbb{E}(3)} \int_{S_1} \|R(x) - \mathcal{C}(x)\|^2 d\text{vol}_{S_1}(x) \right)^{\frac{1}{2}},$$

where  $\mathcal{A}(S_1, S_2)$  is the set of **area-preserving diffeomorphisms** between  $S_1$  and  $S_2$ , and  $\mathbb{E}_3$  is the Euclidean group on  $\mathbb{R}^3$ .



# Continuous Procrustes Distance (cPD)

$$d_{cP}(S_1, S_2) = \inf_{C \in \mathcal{A}} \inf_{R \in \mathbb{E}_3} \left( \int_{S_1} \|R(x) - C(x)\|^2 d\text{vol}_{S_1}(x) \right)^{1/2}$$

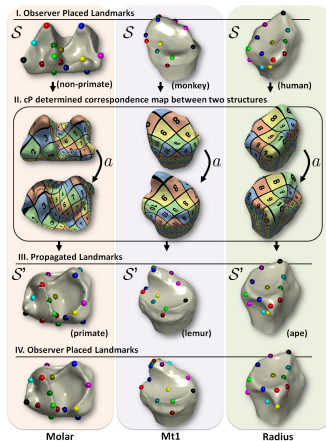
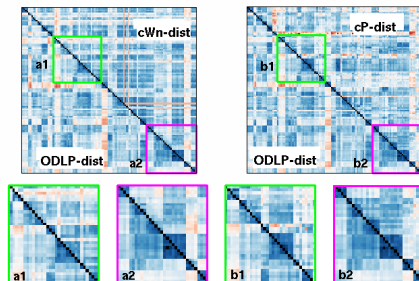




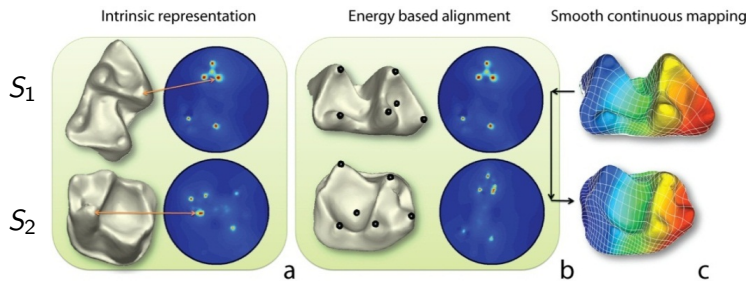
# We defined 2 different distances

$d_{cWn}(S_1, S_2)$ : conformal flattening  
comparison of neighborhood geometry  
optimal mass transport

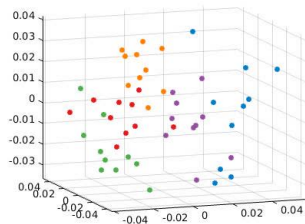
$d_{cP}(S_1, S_2)$ : continuous Procrustes distance



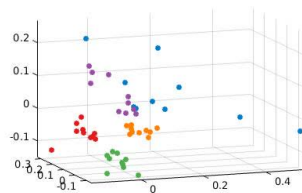
# Bypass Explicit Feature Extraction



# MDS for cPD & DD

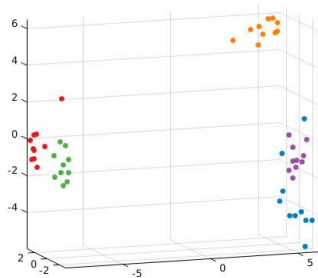


cPD

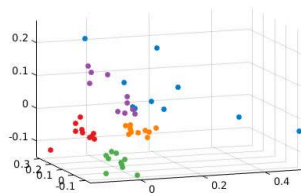


DD

Even better can be obtained!

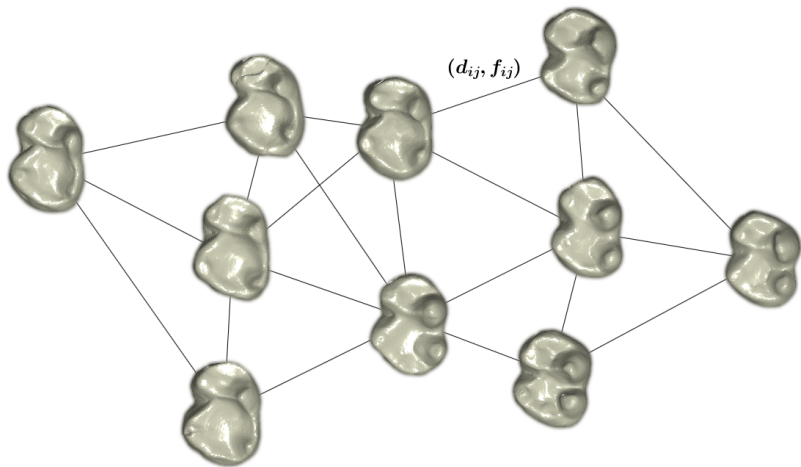


HBDD



DD

to get Diffusion Distance : used local distances  
knitted together  
→ spectral parametrization  
→ distance.



to get Diffusion Distance : used local distances  
knitted together  
→ spectral parametrization  
→ distance.

mappings were used only to obtain numerical  
values for local distances.

to get Diffusion Distance : used local distances  
knitted together  
→ spectral parametrization  
→ distance.

mappings were used only to obtain numerical  
values for local distances.

but they can do much more for us!

in fact: we have a fiber bundle.  
( because of the mappings )

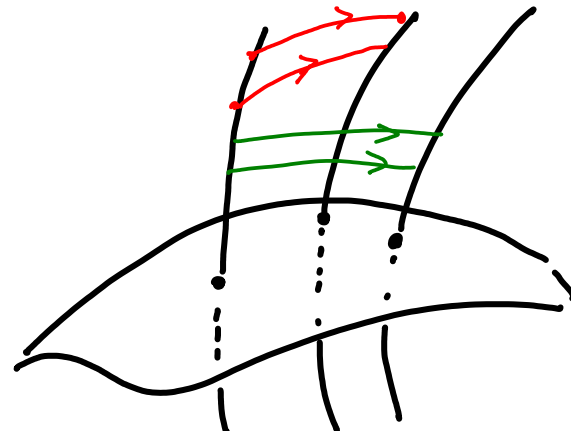
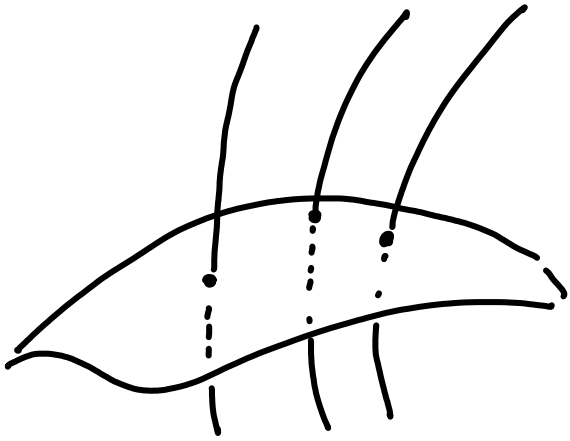


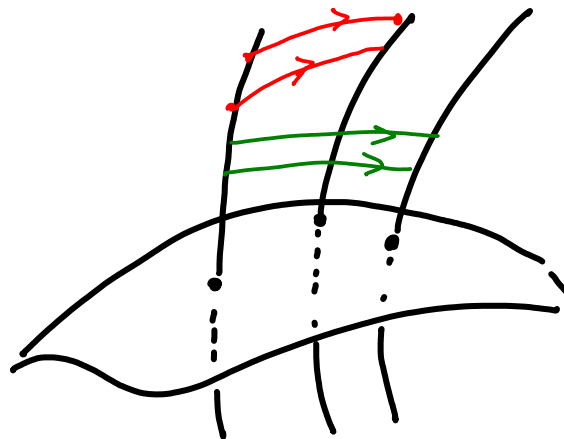
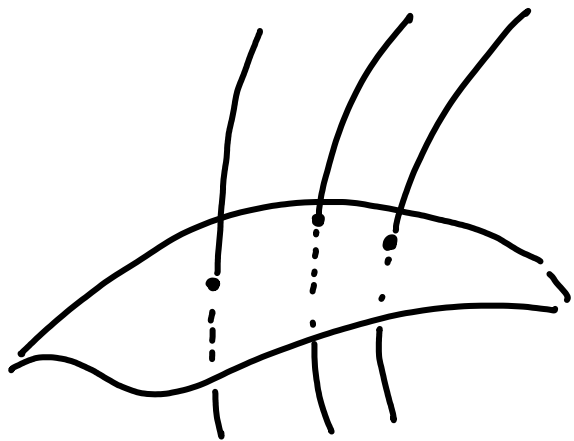
to get Diffusion Distance : used local distances  
knitted together  
→ spectral parametrization  
→ distance.

mappings were used only to obtain numerical  
values for local distances.

but they can do much more for us!

in fact: we have a fiber bundle.  
(because of the mappings)

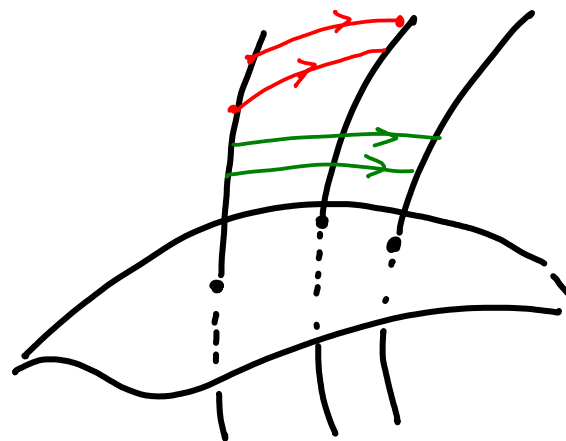
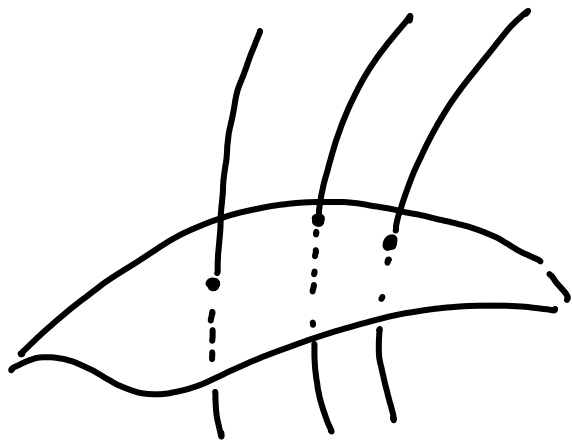




Connection.



family of mappings between fibers



Connection.



family of mappings between fibers

Tigran Gao: use these to define a much more detailed diffusion structure on the higher-dimensional object  
 → "project" at a later stage to obtain "horizontal" part of diffusion.

# Horizontal Random Walk on a Fibre Bundle

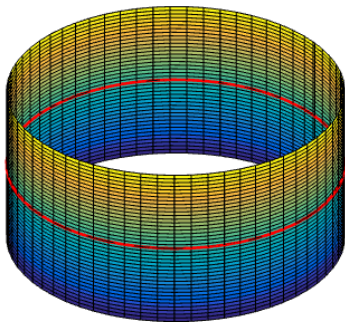
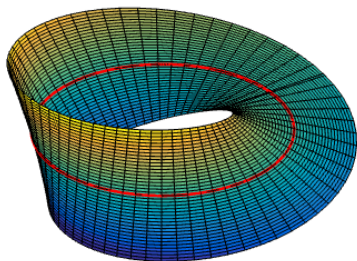
**Fibre Bundle**  $\mathcal{E} = (E, M, F, \pi)$

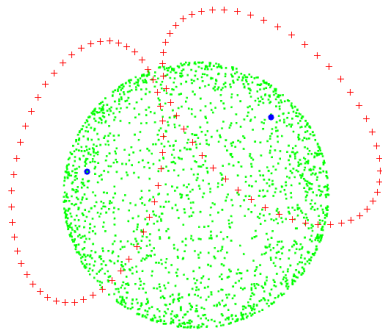
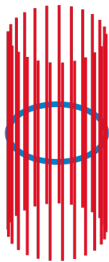
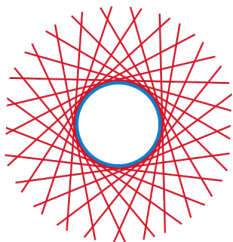
- ▶  $E$ : *total* manifold
- ▶  $M$ : *base* manifold
- ▶  $\pi : E \rightarrow M$ : smooth surjective map (*bundle projection*)
- ▶  $F$ : *fibre* manifold

# Horizontal Random Walk on a Fibre Bundle

**Fibre Bundle**  $\mathcal{E} = (E, M, F, \pi)$

- ▶  $E$ : total manifold
- ▶  $M$ : base manifold
- ▶  $\pi : E \rightarrow M$ : smooth surjective map (*bundle projection*)
- ▶  $F$ : fibre manifold
- ▶ *local triviality*: for “small” open set  $U \subset M$ ,  $\pi^{-1}(U)$  is diffeomorphic to  $U \times F$

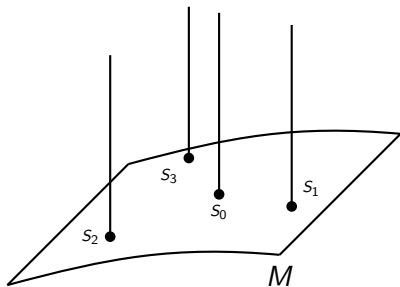
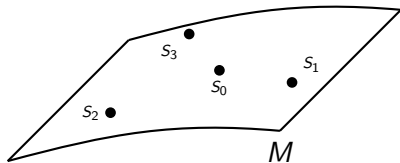




# Horizontal Random Walk on a Fibre Bundle

**Fibre Bundle**  $\mathcal{C} = (E, M, F, \pi)$

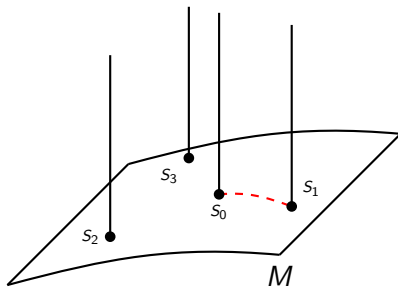
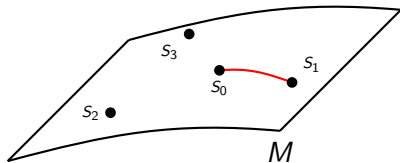
- ▶  $E$ : total manifold
- ▶  $M$ : base manifold
- ▶  $\pi : E \rightarrow M$ : smooth surjective map (*bundle projection*)
- ▶  $F$ : fibre manifold
- ▶ *local triviality*: for “small” open set  $U \subset M$ ,  $\pi^{-1}(U)$  is diffeomorphic to  $U \times F$



# Horizontal Random Walk on a Fibre Bundle

Fibre Bundle  $\mathcal{C} = (E, M, F, \pi)$

- ▶  $E$ : total manifold
- ▶  $M$ : base manifold
- ▶  $\pi : E \rightarrow M$ : smooth surjective map (*bundle projection*)
- ▶  $F$ : fibre manifold
- ▶ *local triviality*: for “small” open set  $U \subset M$ ,  $\pi^{-1}(U)$  is diffeomorphic to  $U \times F$

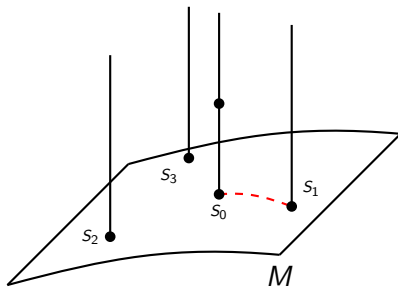
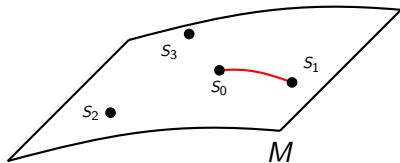




# Horizontal Random Walk on a Fibre Bundle

Fibre Bundle  $\mathcal{C} = (E, M, F, \pi)$

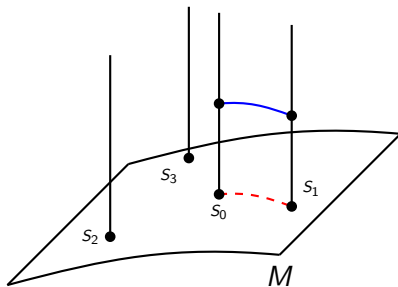
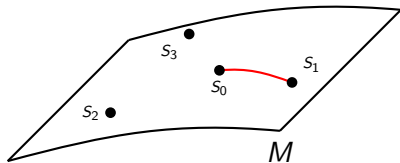
- ▶  $E$ : total manifold
- ▶  $M$ : base manifold
- ▶  $\pi : E \rightarrow M$ : smooth surjective map (*bundle projection*)
- ▶  $F$ : fibre manifold
- ▶ *local triviality*: for “small” open set  $U \subset M$ ,  $\pi^{-1}(U)$  is diffeomorphic to  $U \times F$



# Horizontal Random Walk on a Fibre Bundle

Fibre Bundle  $\mathcal{C} = (E, M, F, \pi)$

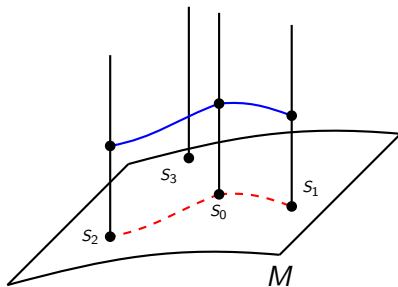
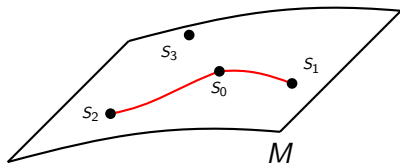
- ▶  $E$ : total manifold
- ▶  $M$ : base manifold
- ▶  $\pi : E \rightarrow M$ : smooth surjective map (*bundle projection*)
- ▶  $F$ : fibre manifold
- ▶ *local triviality*: for “small” open set  $U \subset M$ ,  $\pi^{-1}(U)$  is diffeomorphic to  $U \times F$



# Horizontal Random Walk on a Fibre Bundle

Fibre Bundle  $\mathcal{C} = (E, M, F, \pi)$

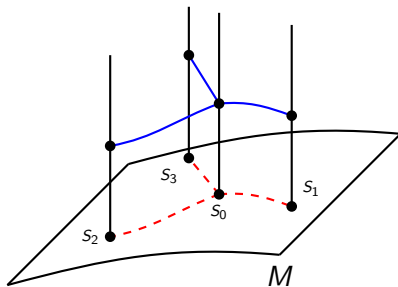
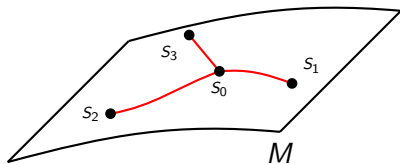
- ▶  $E$ : total manifold
- ▶  $M$ : base manifold
- ▶  $\pi : E \rightarrow M$ : smooth surjective map (*bundle projection*)
- ▶  $F$ : fibre manifold
- ▶ *local triviality*: for “small” open set  $U \subset M$ ,  $\pi^{-1}(U)$  is diffeomorphic to  $U \times F$



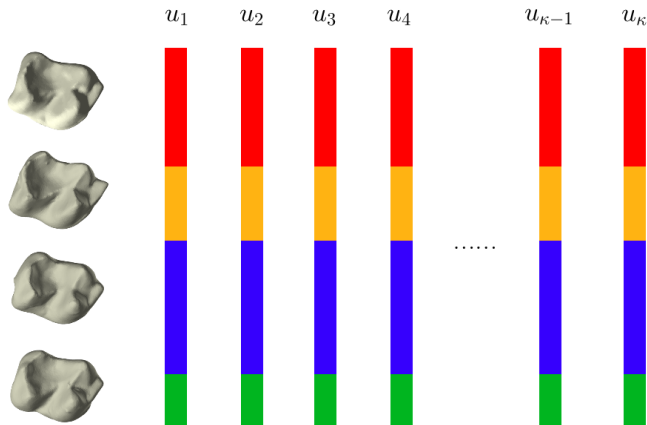
# Horizontal Random Walk on a Fibre Bundle

Fibre Bundle  $\mathcal{C} = (E, M, F, \pi)$

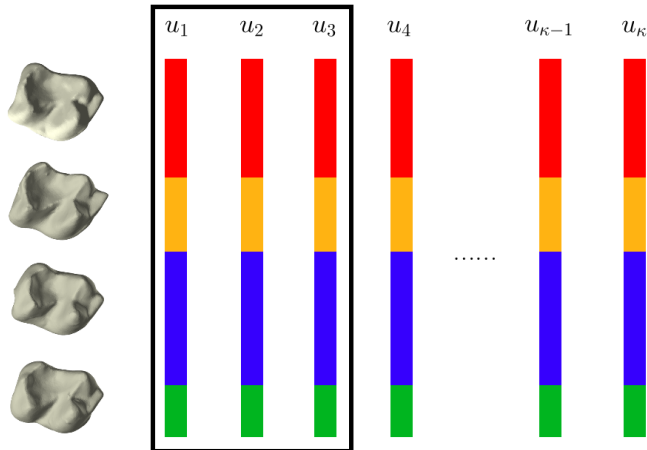
- ▶  $E$ : total manifold
- ▶  $M$ : base manifold
- ▶  $\pi : E \rightarrow M$ : smooth surjective map (*bundle projection*)
- ▶  $F$ : fibre manifold
- ▶ *local triviality*: for “small” open set  $U \subset M$ ,  $\pi^{-1}(U)$  is diffeomorphic to  $U \times F$



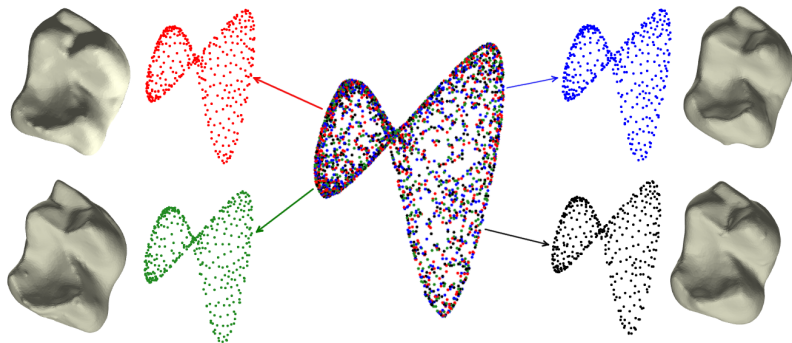
# Horizontal Diffusion Maps: Embedding the Entire Bundle



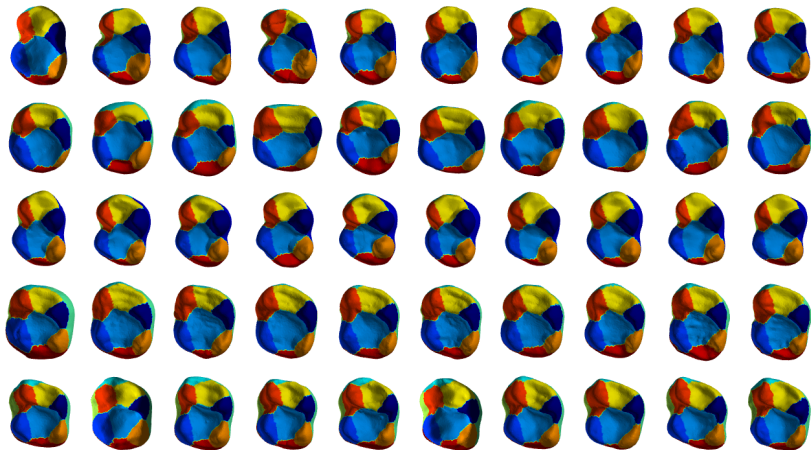
# Horizontal Diffusion Maps: Embedding the Entire Bundle



# Horizontal Diffusion Maps

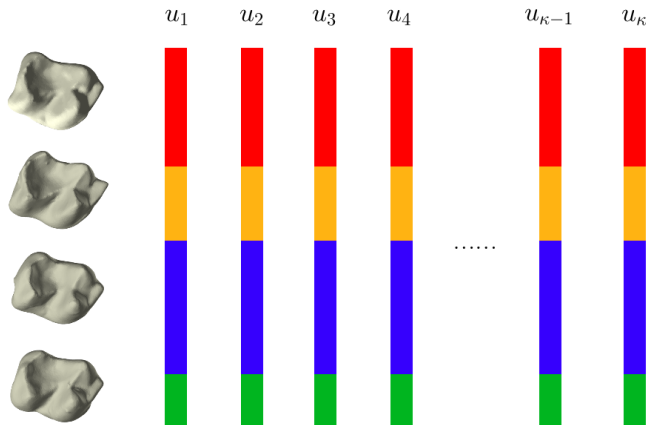


## Automatic Landmarking — Interpretability

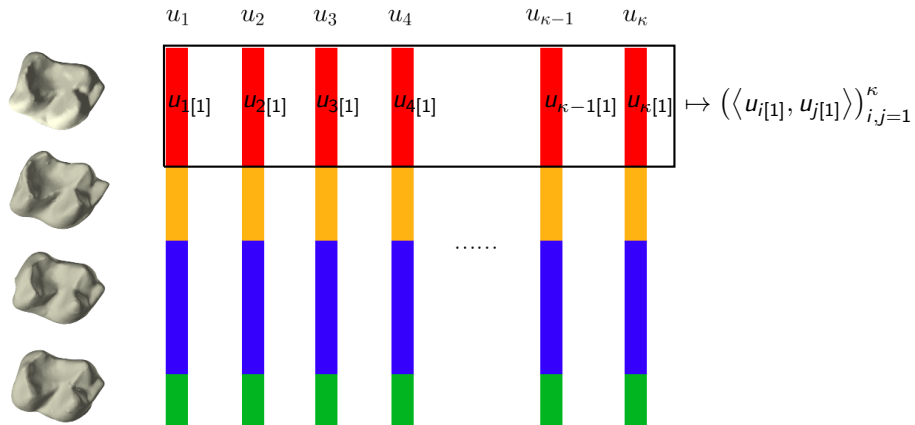




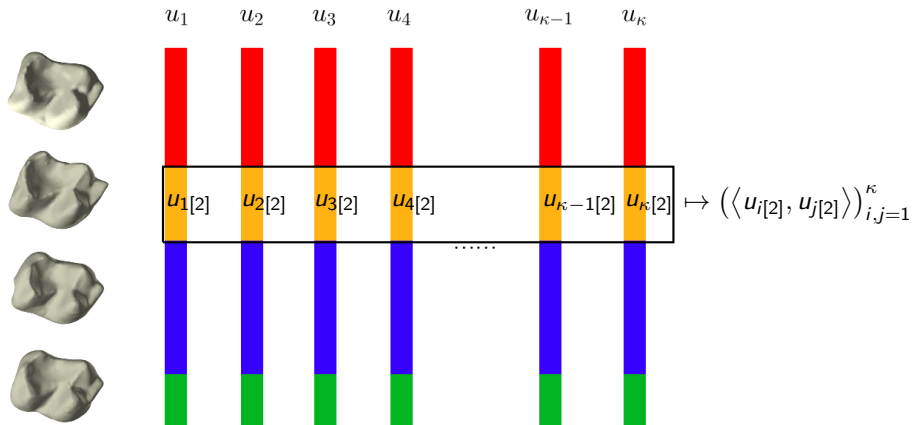
# Horizontal Diffusion Maps: Embedding the Base Manifold



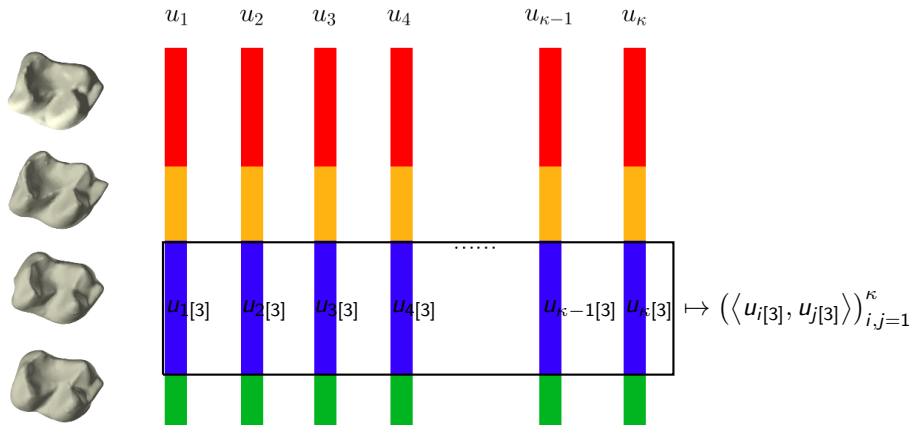
# Horizontal Diffusion Maps: Embedding the Base Manifold



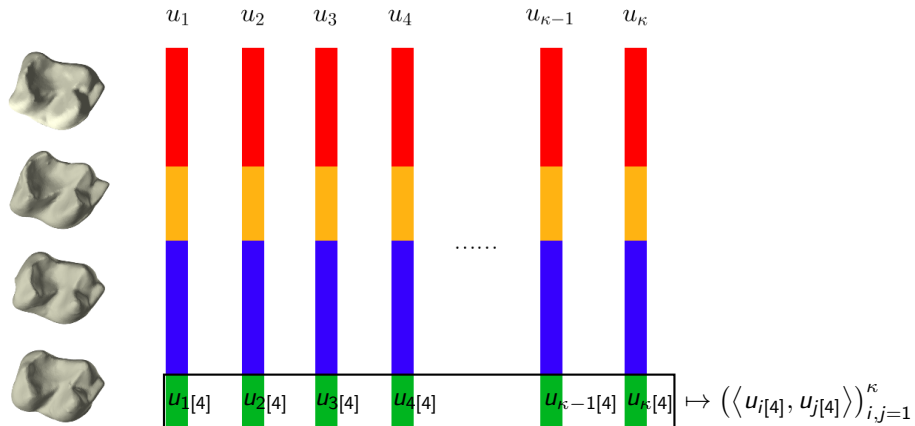
# Horizontal Diffusion Maps: Embedding the Base Manifold



# Horizontal Diffusion Maps: Embedding the Base Manifold



# Horizontal Diffusion Maps: Embedding the Base Manifold



spectral coordinates for points in fiber bundle:

$$(j, p) \longrightarrow \left( u_k(j, p) \right)_{k=1, \dots, K}$$

$\swarrow$   
 $S_j$

$\nwarrow$   
pt  $p$   
on  $S_j$

spectral coordinates for points in fiber bundle:

$$(j, p) \longrightarrow \left( u_k(j, p) \right)_{k=1, \dots, K}$$

$\swarrow$   
 $S_j$

$\nwarrow$   
pt  $p$   
on  $S_j$

$\downarrow$  "project" to geometry  
on base manifold

spectral coordinates for points in fiber bundle:

$$\begin{array}{ccc}
 (j, p) & \longrightarrow & (u_k(j, p))_{k=1, \dots, K} \\
 \swarrow \text{ } \nearrow & & \\
 S_i & \text{pt } p \text{ on } S_j & \\
 & \downarrow \text{"project" to geometry} & \\
 & \text{on base manifold} & 
 \end{array}$$

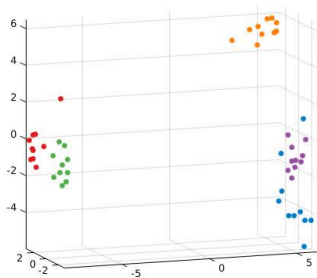
hor. dist  $(S_i, S_j)$

= dist. between corresponding  
point clouds in  $K$ -dim space.

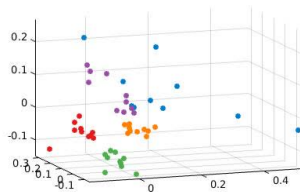
$$= \left[ \sum_{p, q} \lambda_k^{\varepsilon, \delta}(p, q) |u_k(i, p) - u_k(j, q)|^2 \right]^{1/2}$$



# Species Clustering

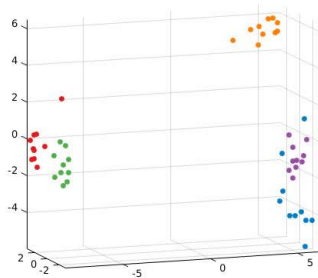


Horizontal Base Diffusion Distance (**with** Maps)

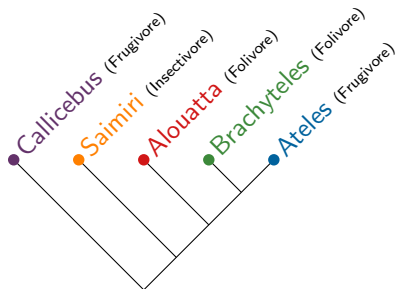


Diffusion Distance (**without** Maps)

# Species Clustering



Horizontal Base Diffusion Distance (with Maps)



# Learning the base manifold in a fibre bundle

This fibre bundle learning method applies much more generally!

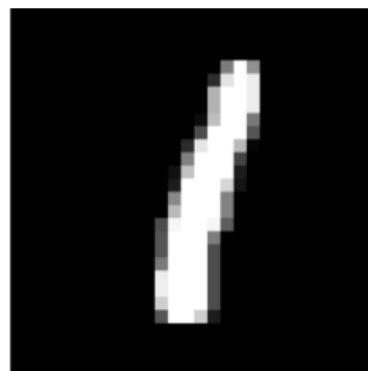
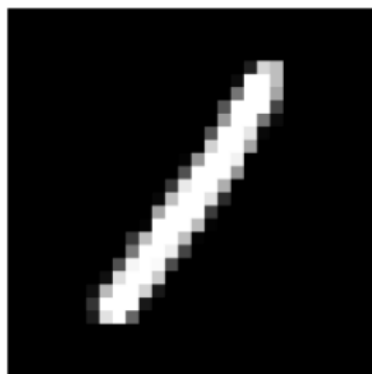
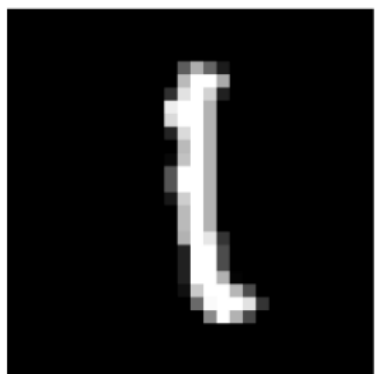
# Learning the base manifold in a fibre bundle

This fibre bundle learning method applies much more generally!



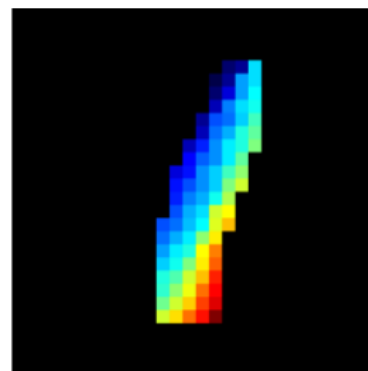
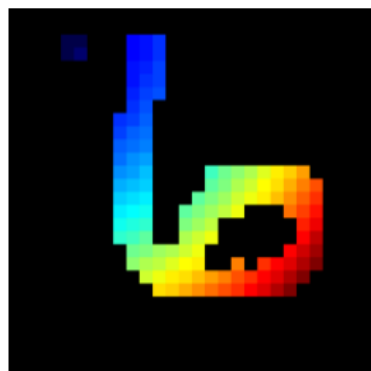
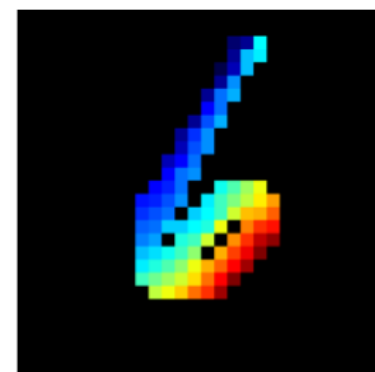
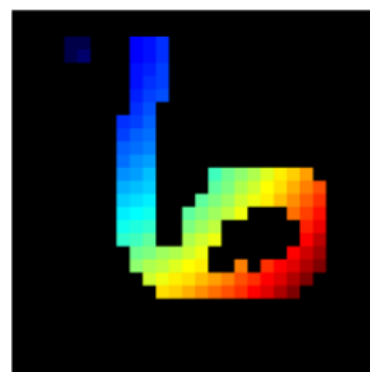
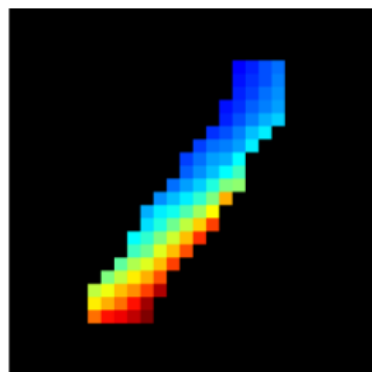
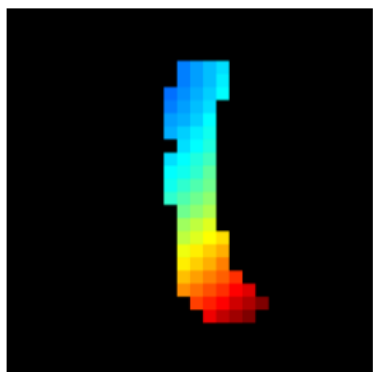
# Learning the base manifold in a fibre bundle

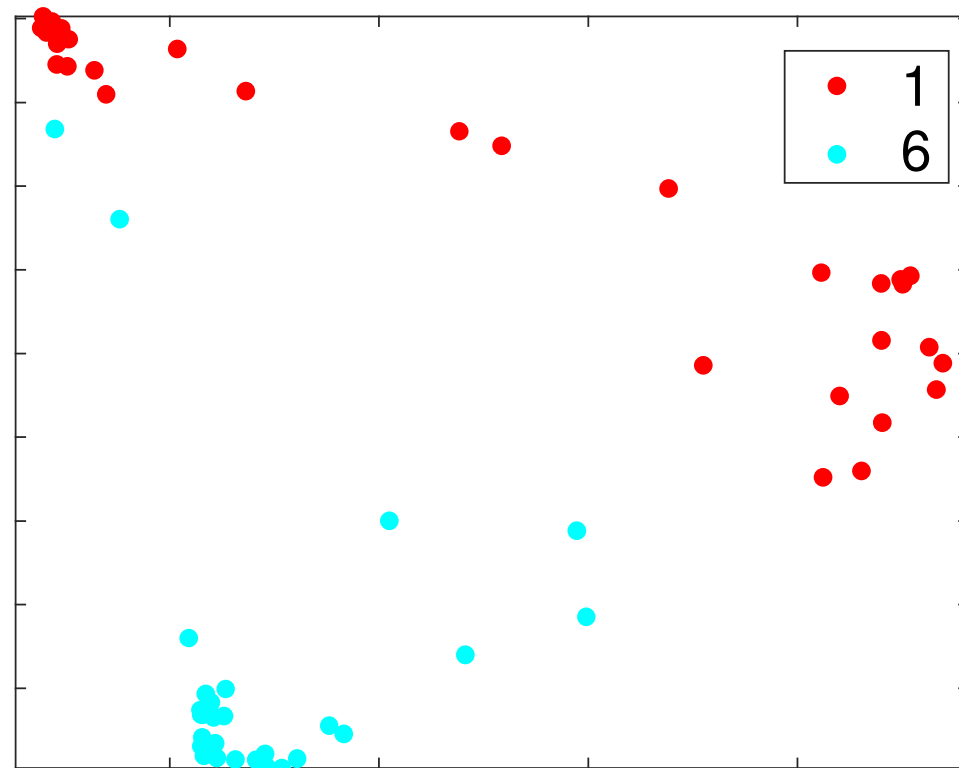
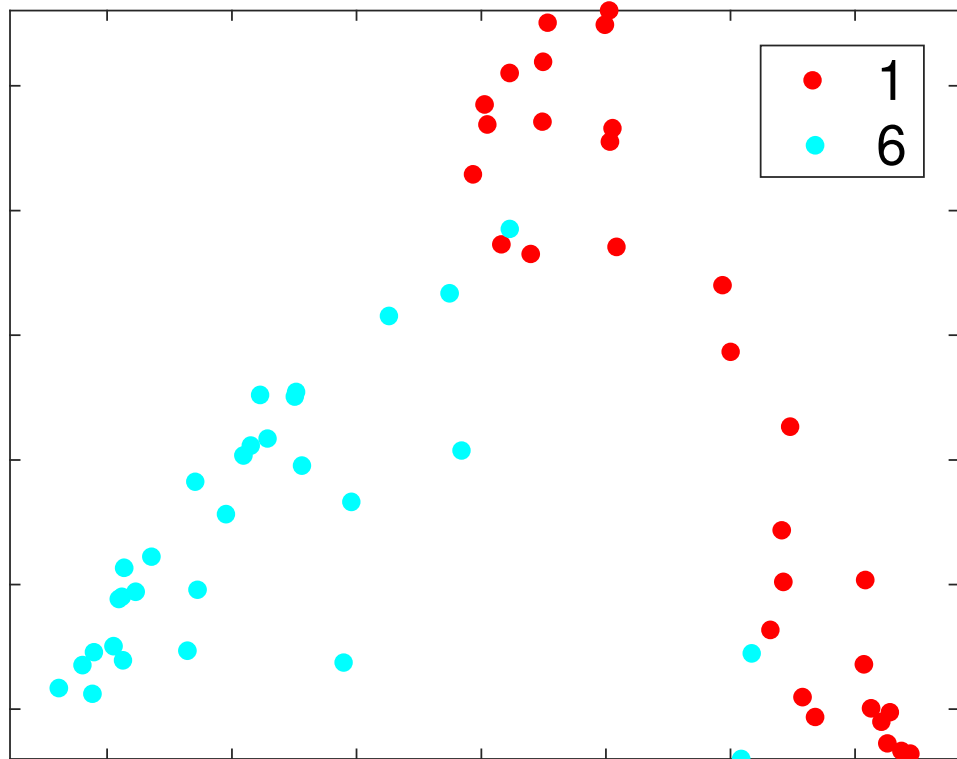
This fibre bundle learning method applies much more generally!

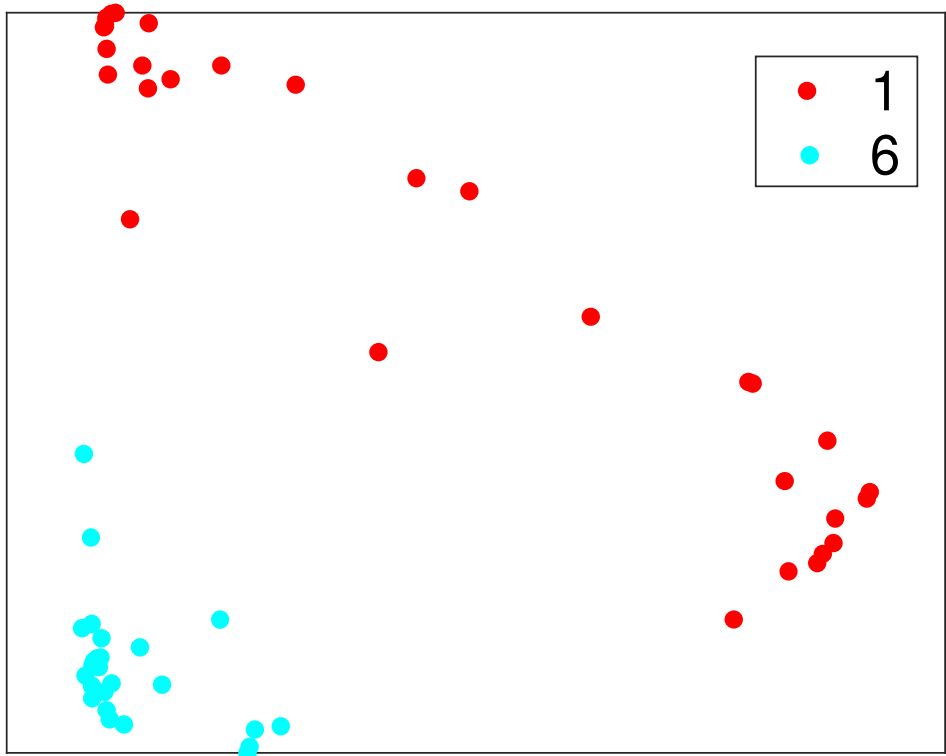
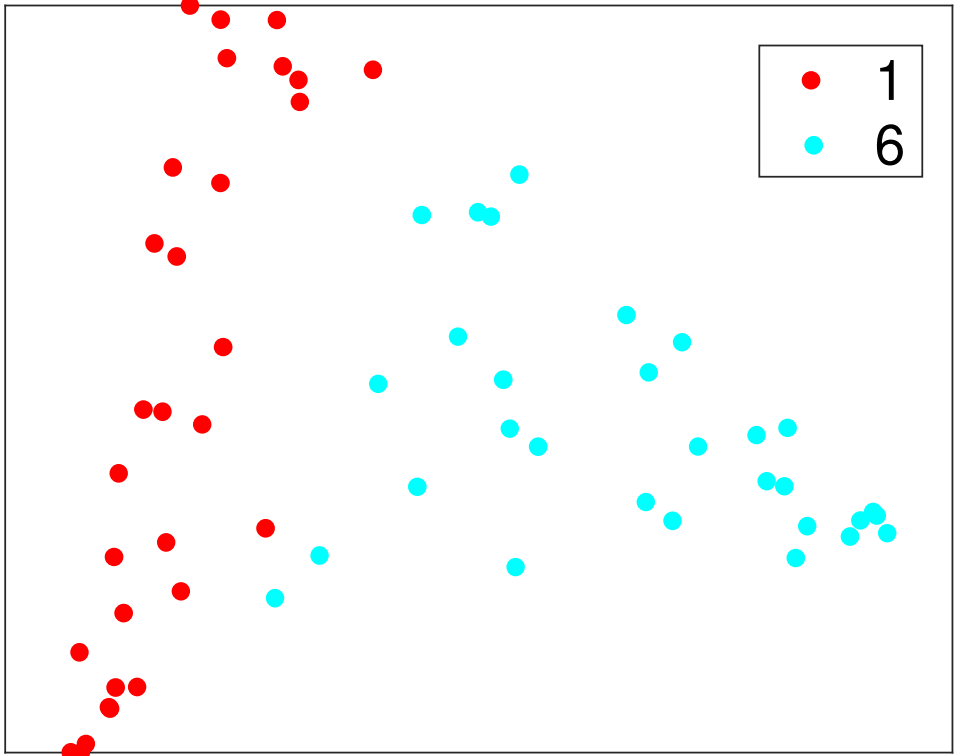


# Learning the base manifold in a fibre bundle

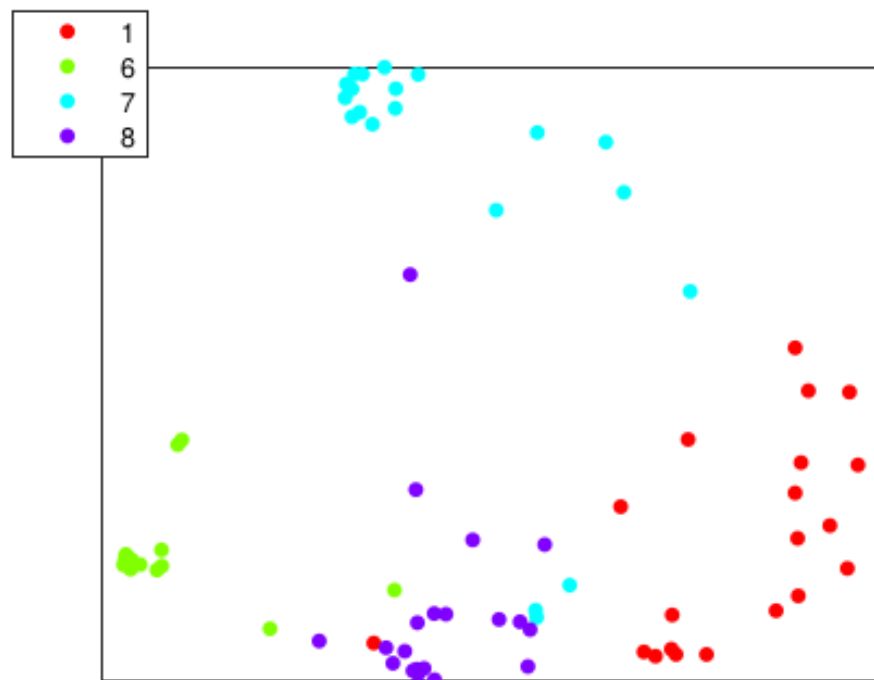
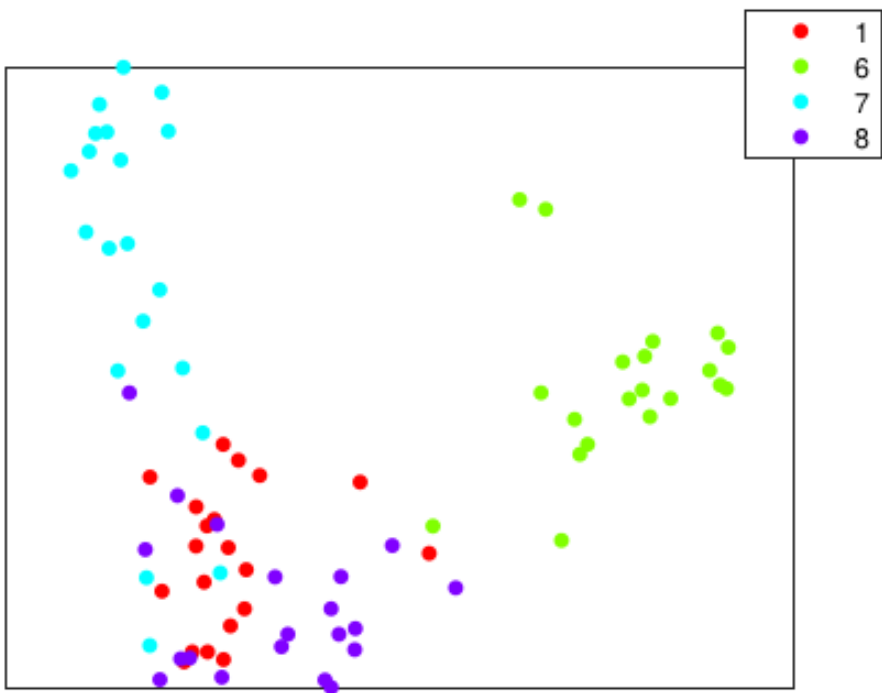
This fibre bundle learning method applies much more generally!





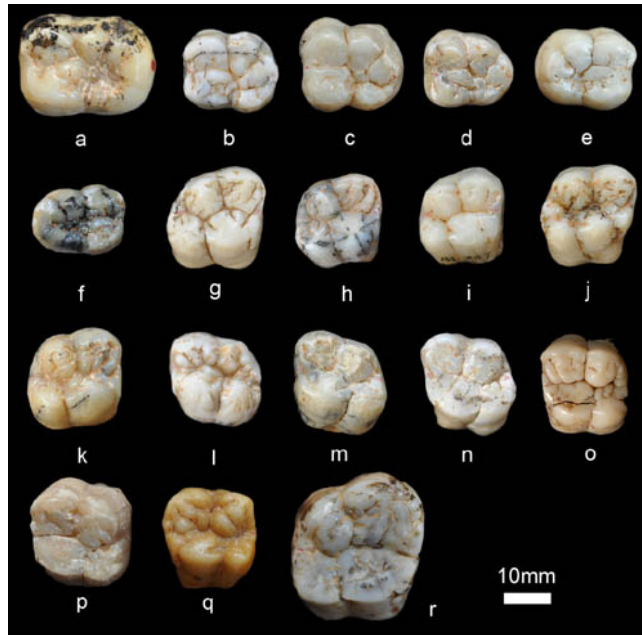






- multi-resolution ; coarse- & fine-graining.

connection is reasonable for bones/teeth of closely related species.



primate molars



crabeater seal molars