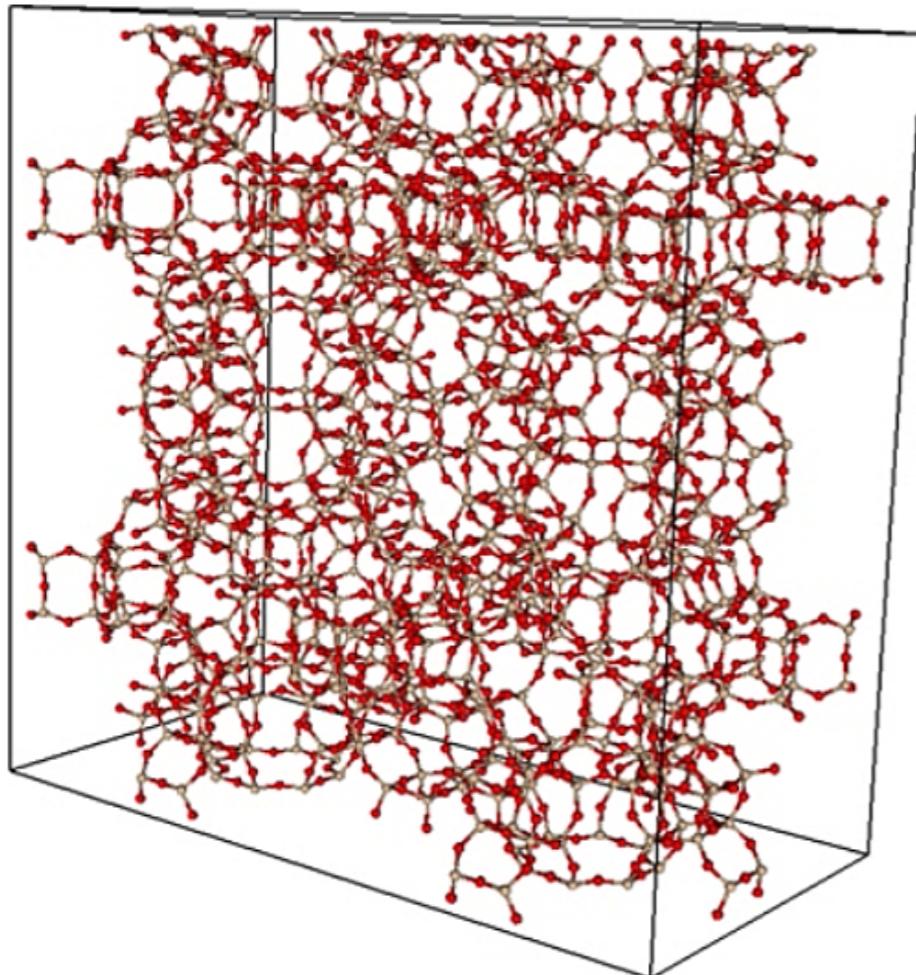


On metric and statistical properties of topological descriptors for geometric data

PhD defense of Mathieu Carrière

Topological Data Analysis

Modern data sets often have non-trivial underlying structures...
...and these structures (*shapes*) matter

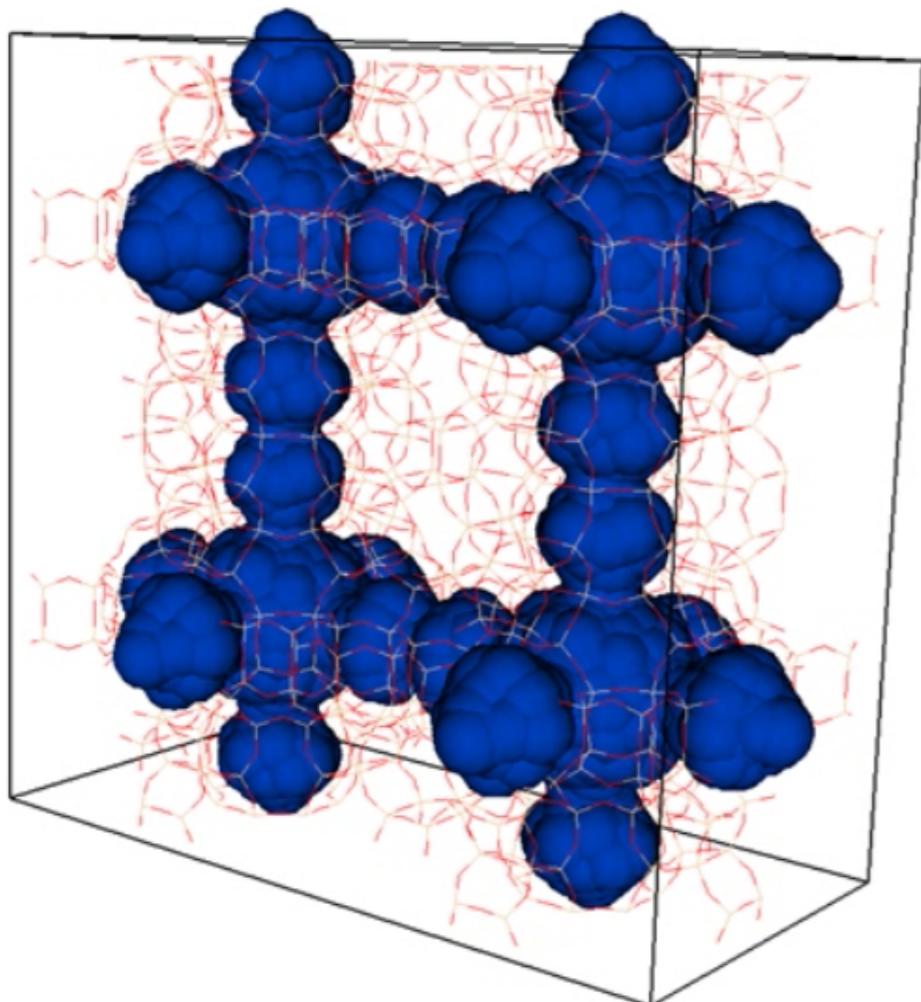


Example:

cavities in nanoporous materials
(e.g. zeolites for nanofilters)
determine their physical properties

Topological Data Analysis

Modern data sets often have non-trivial underlying structures...
...and these structures (*shapes*) matter

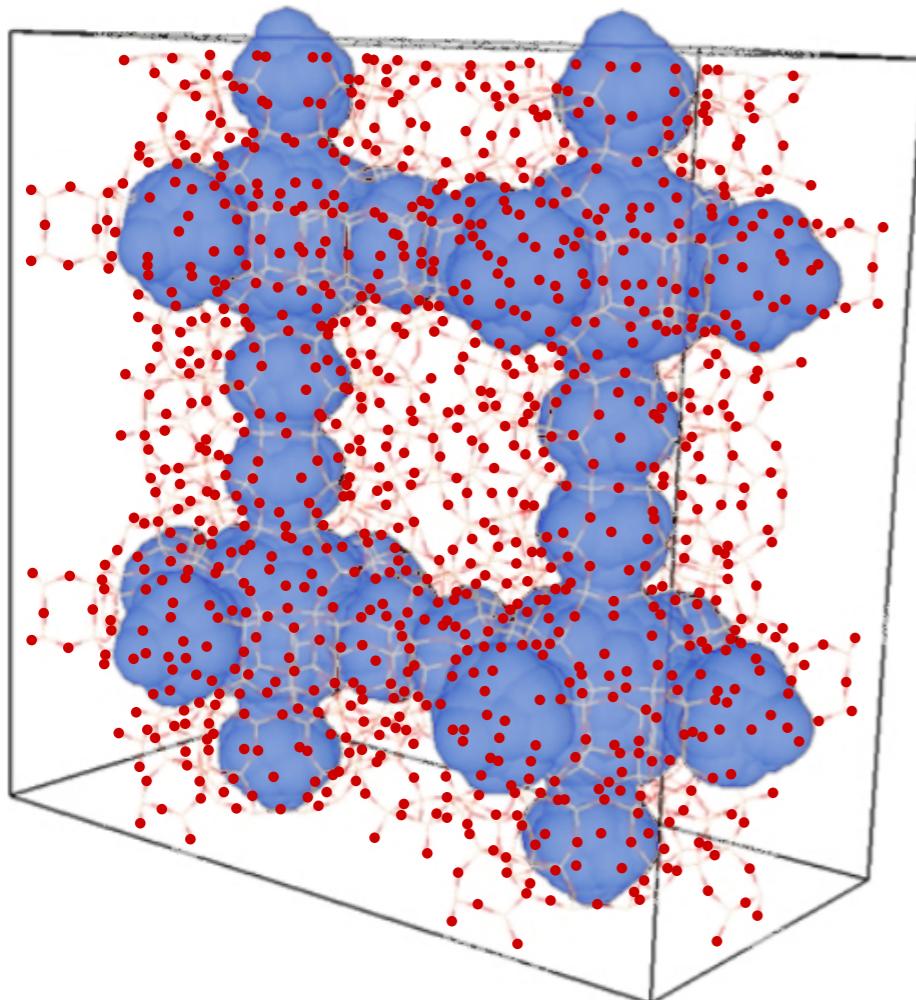


Example:

cavities in nanoporous materials
(e.g. zeolites for nanofilters)
determine their physical properties

Topological Data Analysis

Modern data sets often have non-trivial underlying structures...
...and these structures (*shapes*) matter



→ need **topological descriptors** that can:
extract the shapes hidden in data
encode these shapes compactly

Example:

cavities in nanoporous materials
(e.g. zeolites for nanofilters)
determine their physical properties

Topological Data Analysis

Why is topology interesting for data analysis?

- multiscale
- compact

Topological Data Analysis

Why is topology interesting for data analysis?

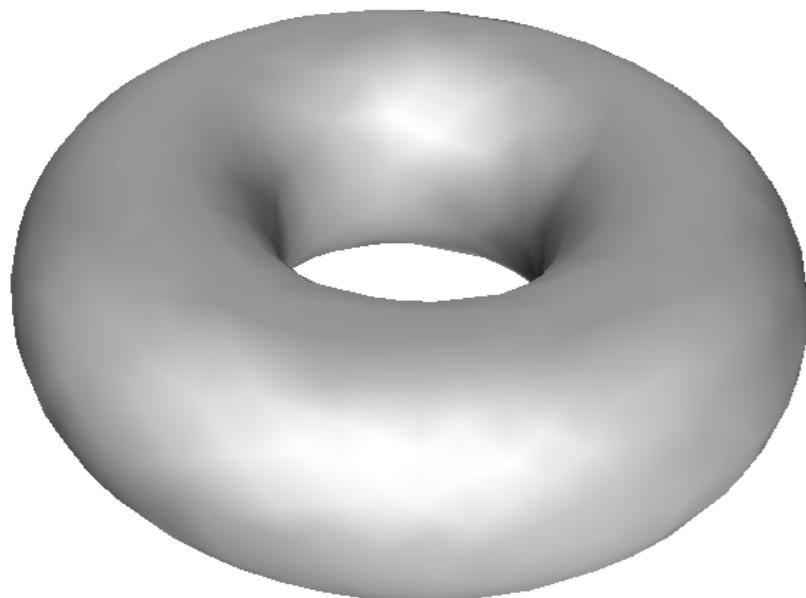
- multiscale
- compact
- complementary to other descriptors
- stable with respect to (small) perturbations
- invariant under coordinate changes/rigid transforms

Topological Data Analysis

Why is topology interesting for data analysis?

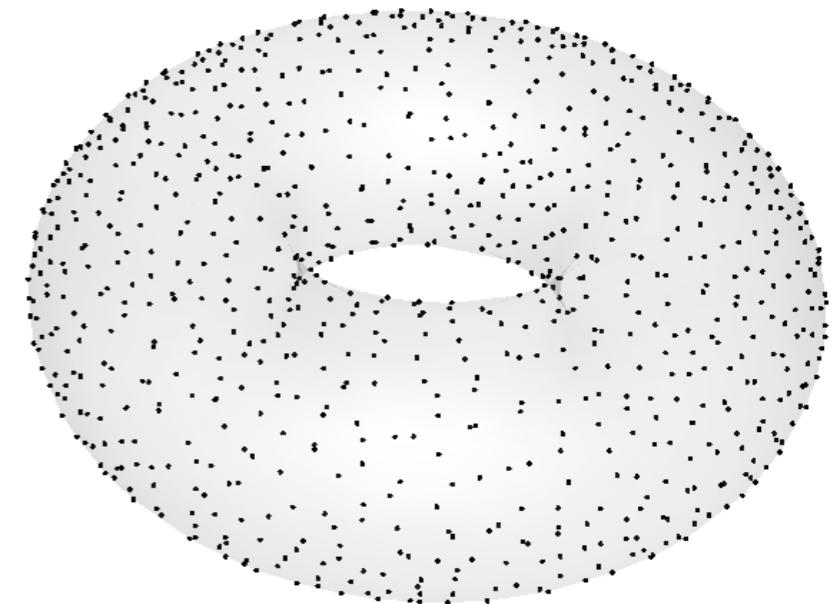
- multiscale
- compact

- complementary to other descriptors
- stable with respect to (small) perturbations
- invariant under coordinate changes/rigid transforms



topological space

← →
topological descriptors



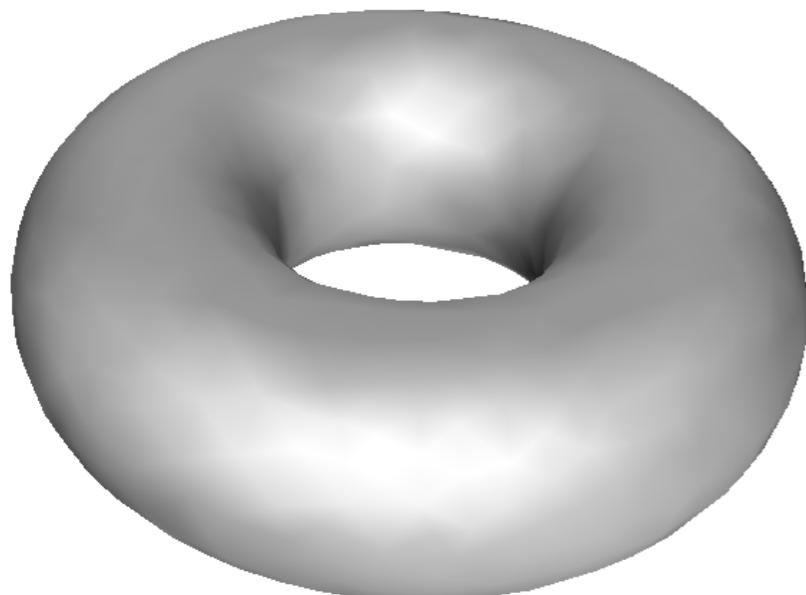
point cloud

Topological Data Analysis

Why is topology interesting for data analysis?

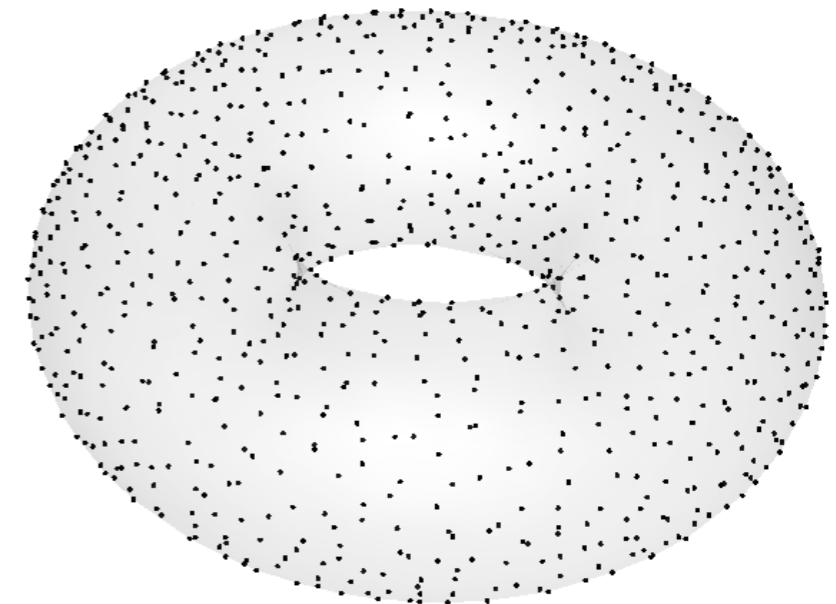
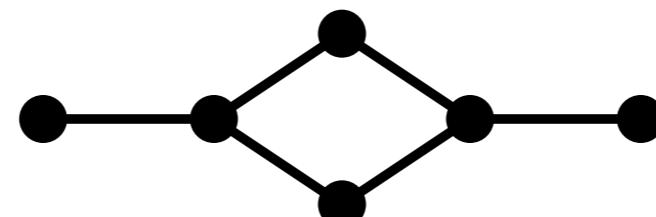
→ multiscale
→ compact

→ complementary to other descriptors
→ stable with respect to (small) perturbations
→ invariant under coordinate changes/rigid transforms

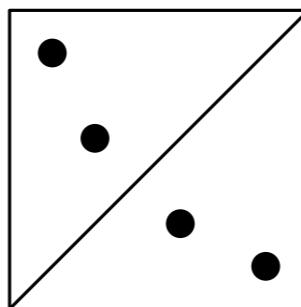


topological space

← →
topological descriptors



point cloud

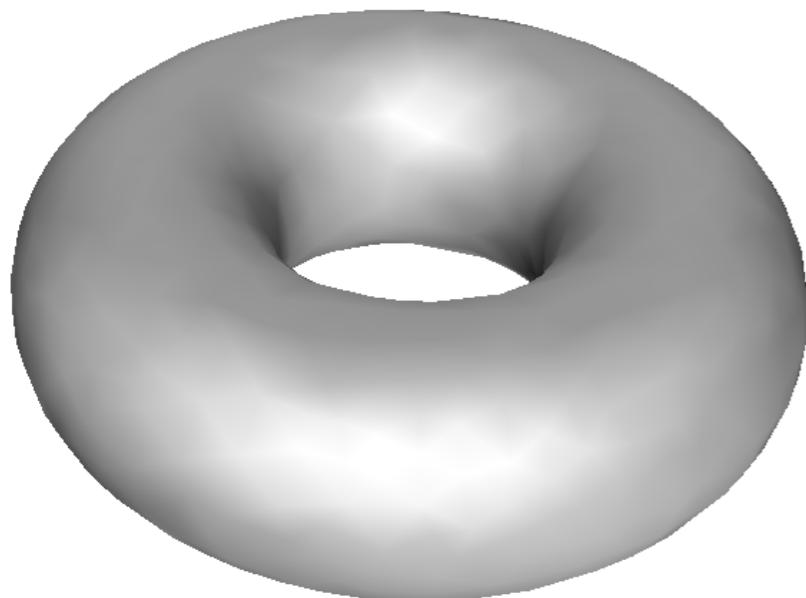


Topological Data Analysis

Why is topology interesting for data analysis?

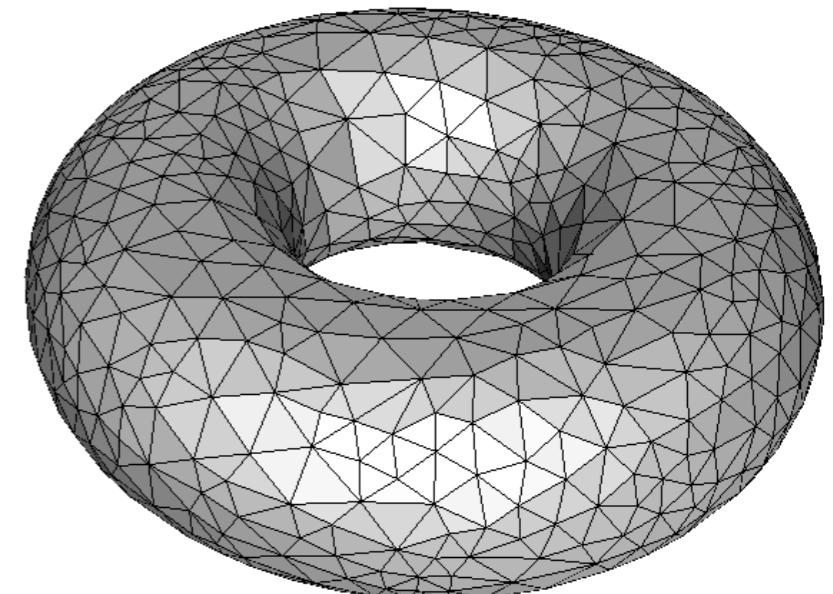
→ multiscale
→ compact

→ complementary to other descriptors
→ stable with respect to (small) perturbations
→ invariant under coordinate changes/rigid transforms

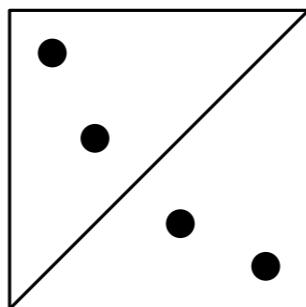
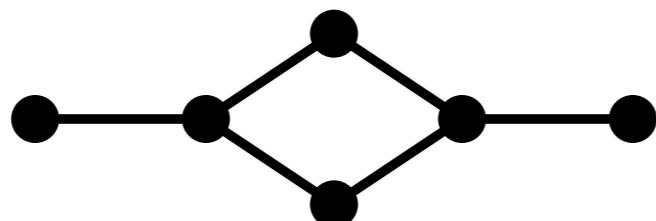


topological space

↔
topological descriptors



point cloud
triangulation

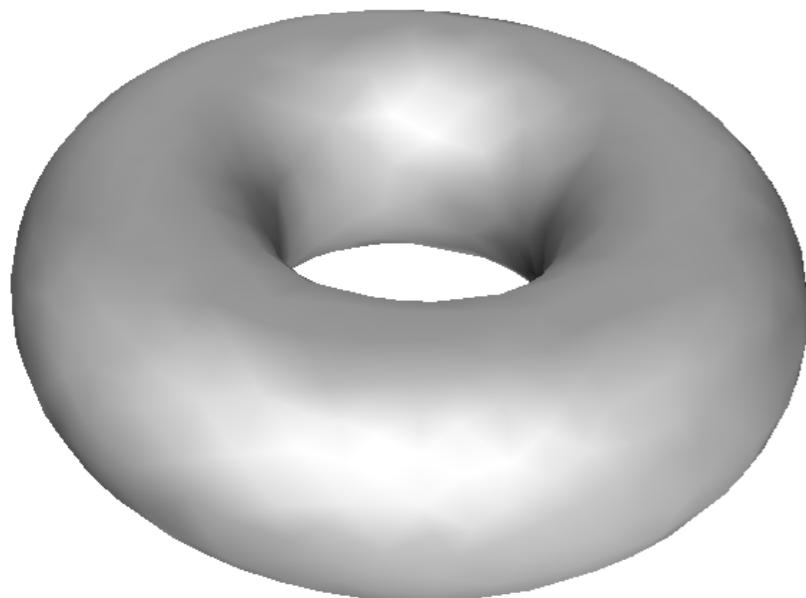


Topological Data Analysis

Why is topology interesting for data analysis?

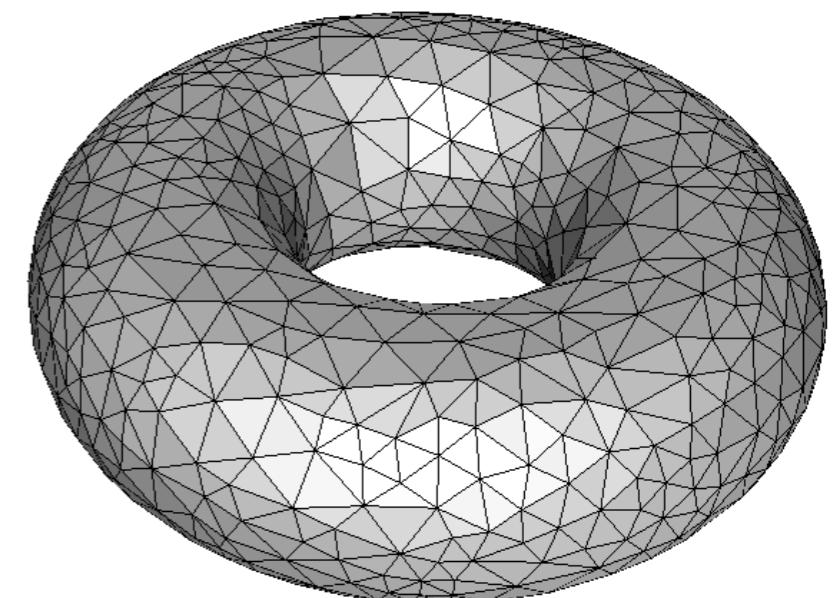
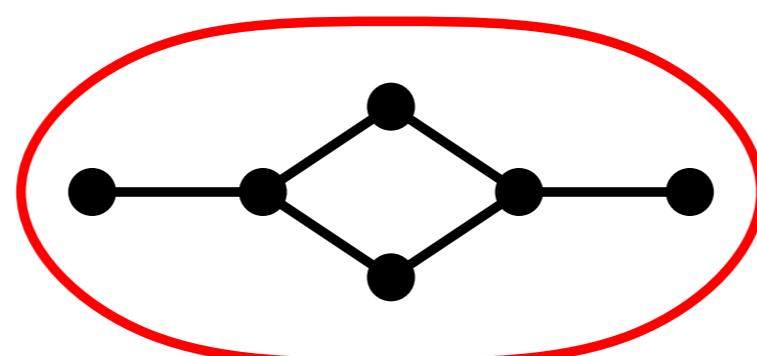
→ multiscale
→ compact

→ complementary to other descriptors
→ stable with respect to (small) perturbations
→ invariant under coordinate changes/rigid transforms

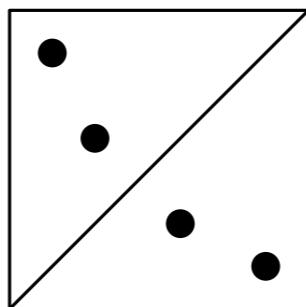


topological space

↔
topological descriptors



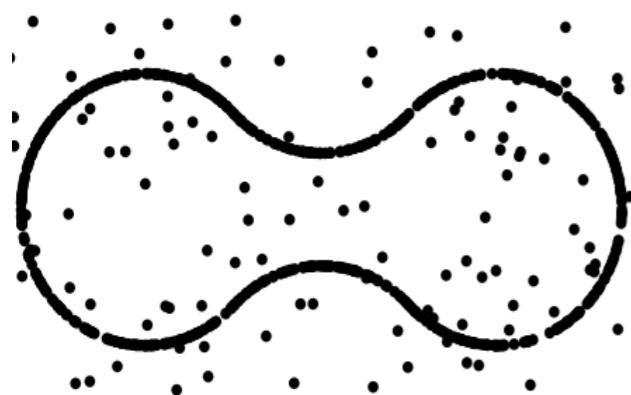
point cloud
triangulation



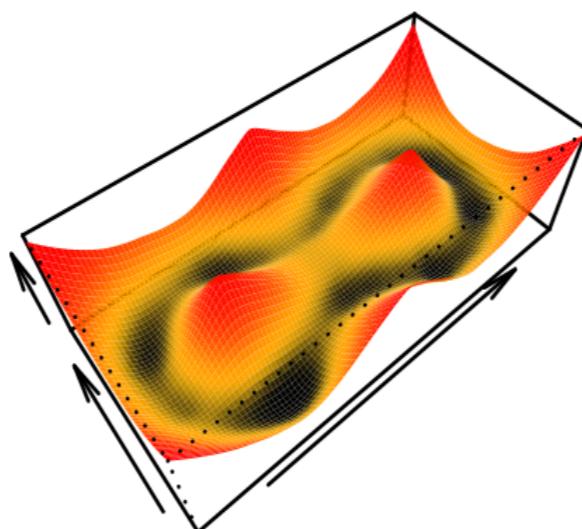
Topological Data Analysis

For **exploratory analysis, visualization**

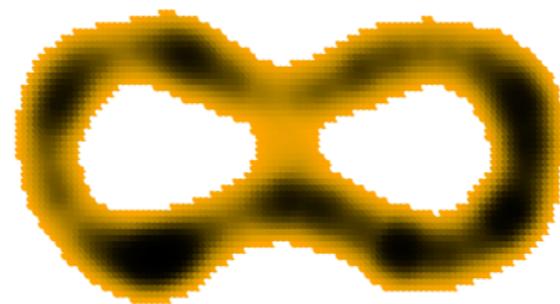
Cassini with Noise



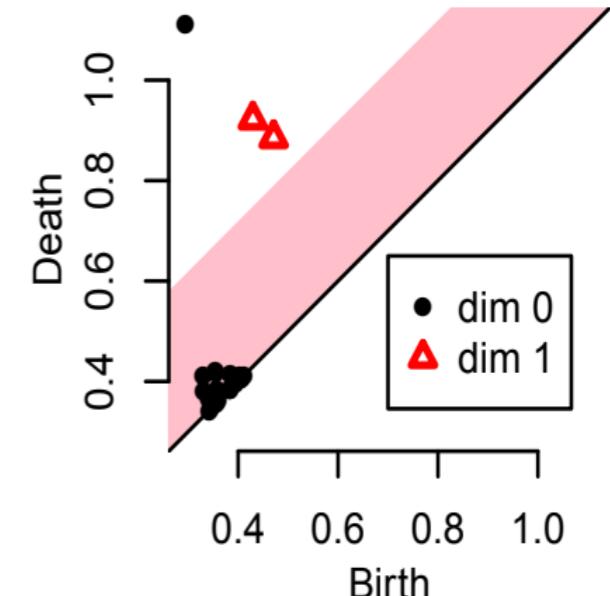
DTM



Sublevel Set, $t=0.5$

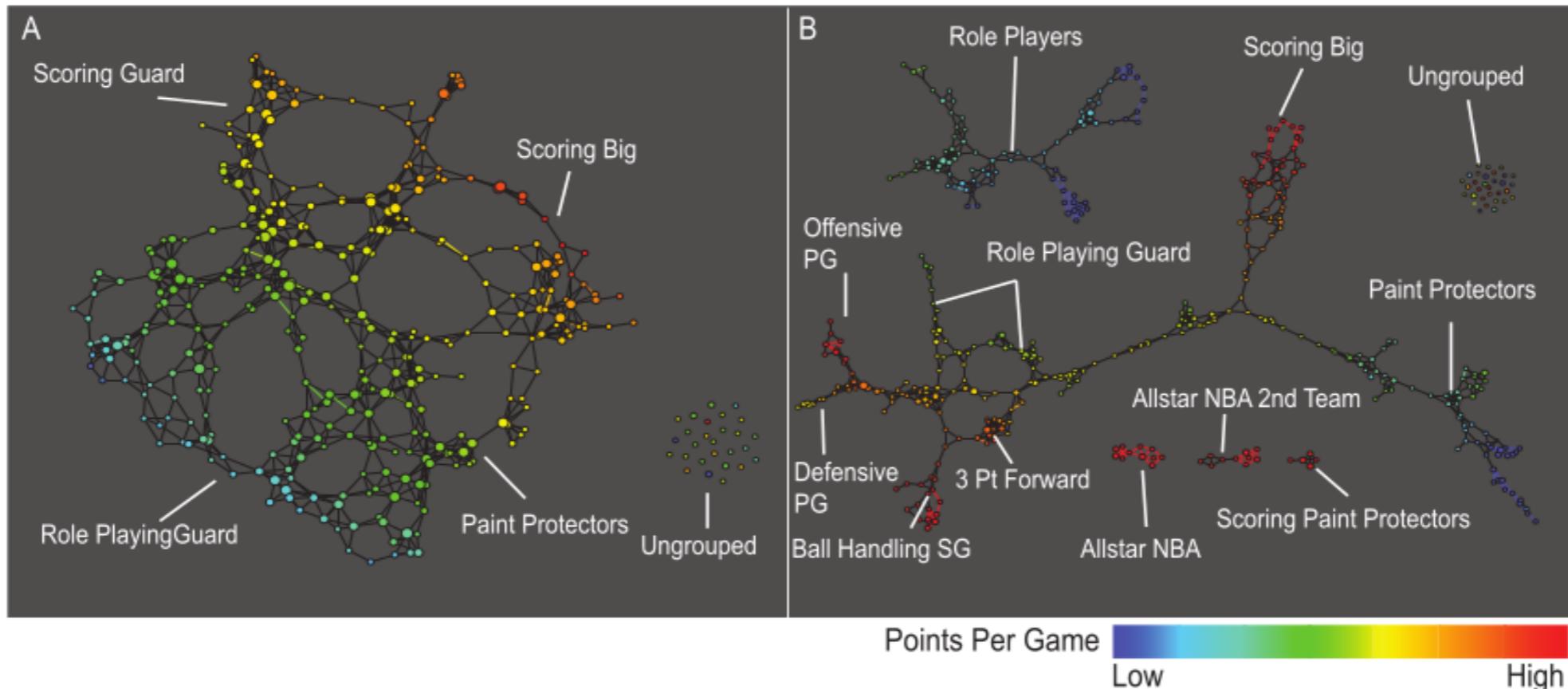


Persistence Diagram



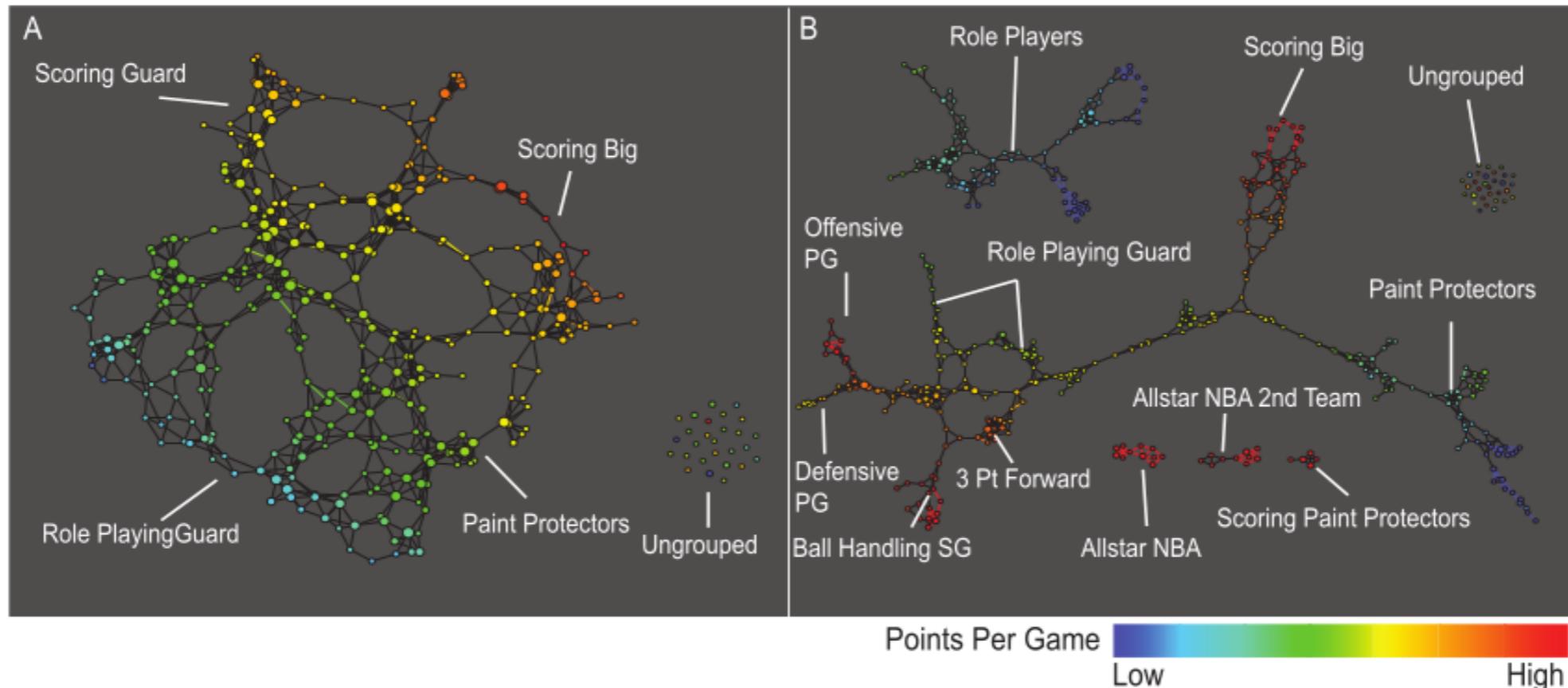
Topological Data Analysis

For exploratory analysis, visualization

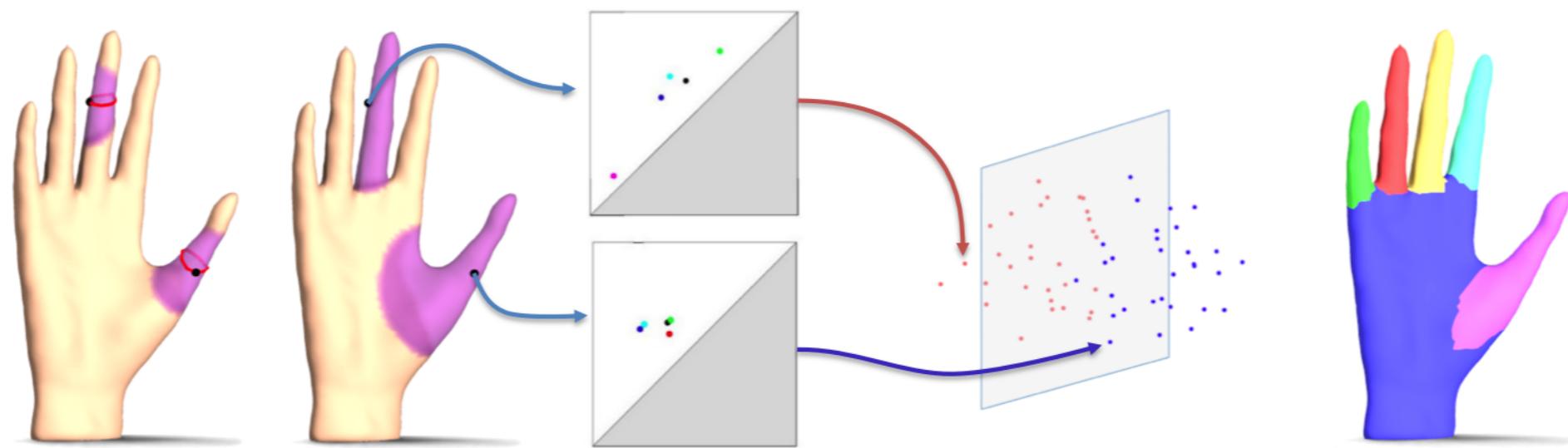


Topological Data Analysis

For exploratory analysis, visualization



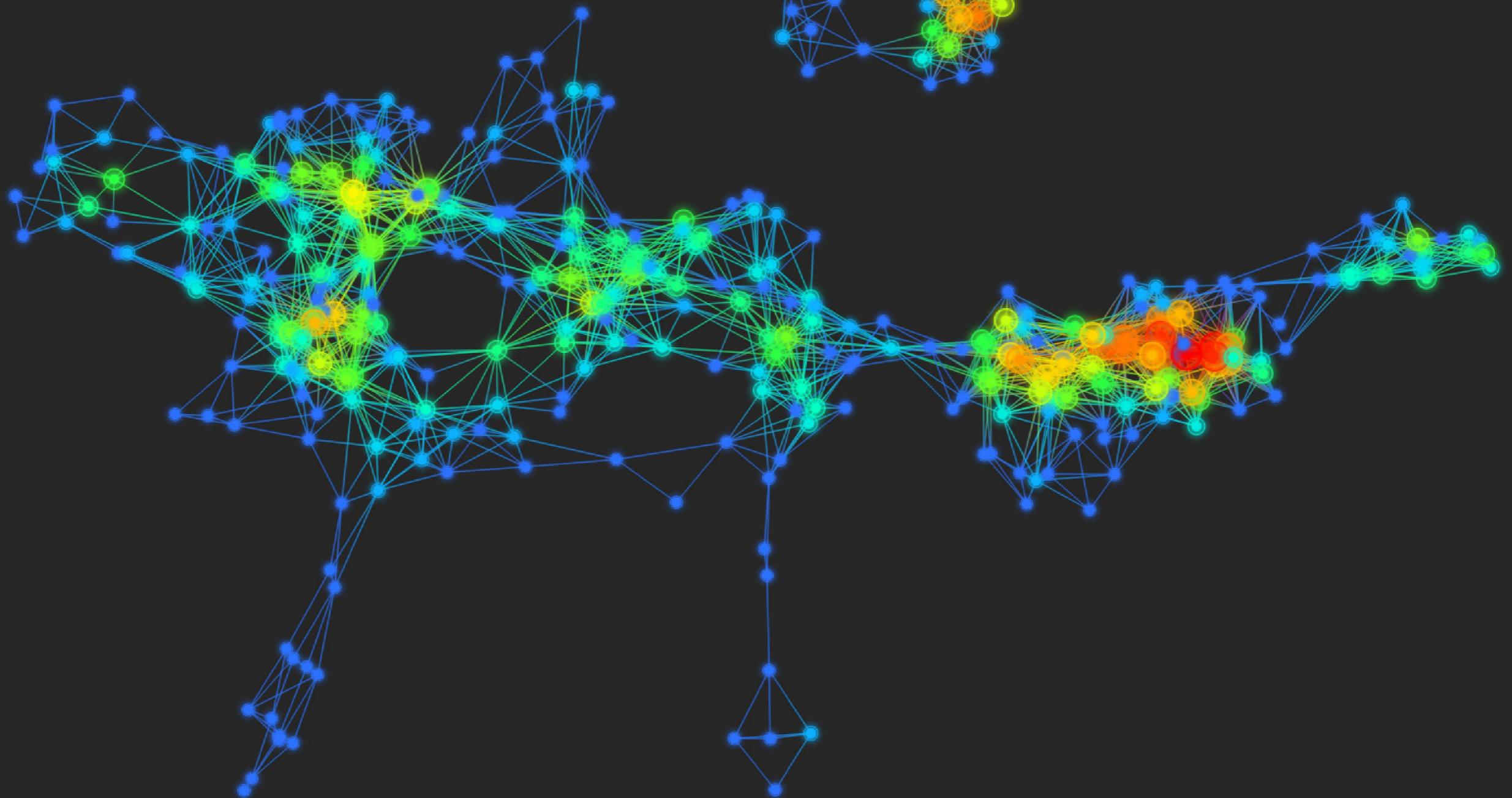
For feature extraction and statistical learning



Reeb Graphs and Mapper

[Reeb 1946]

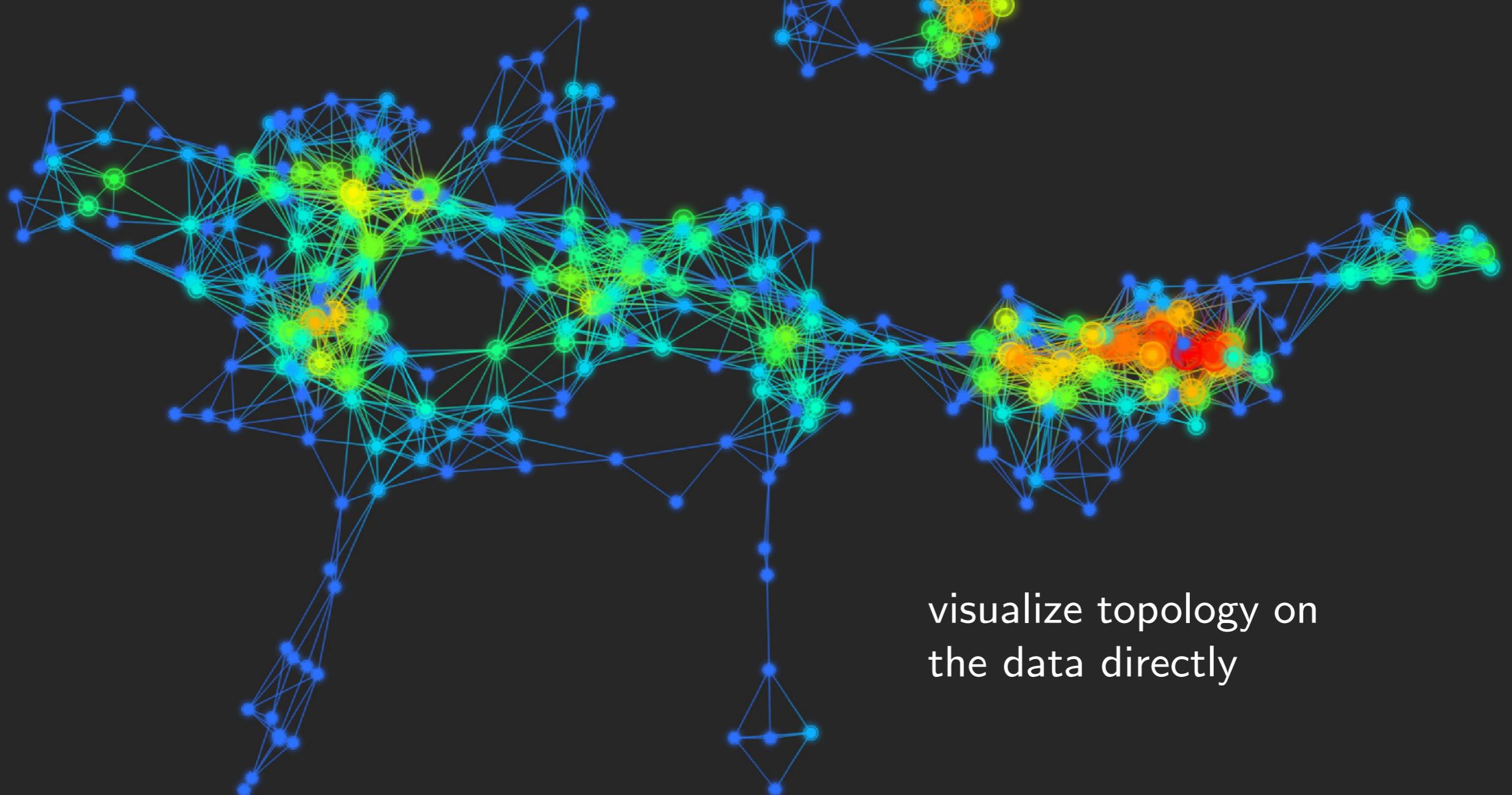
[Singh, Mémoli, Carlsson 2007]



Reeb Graphs and Mapper

[Reeb 1946]

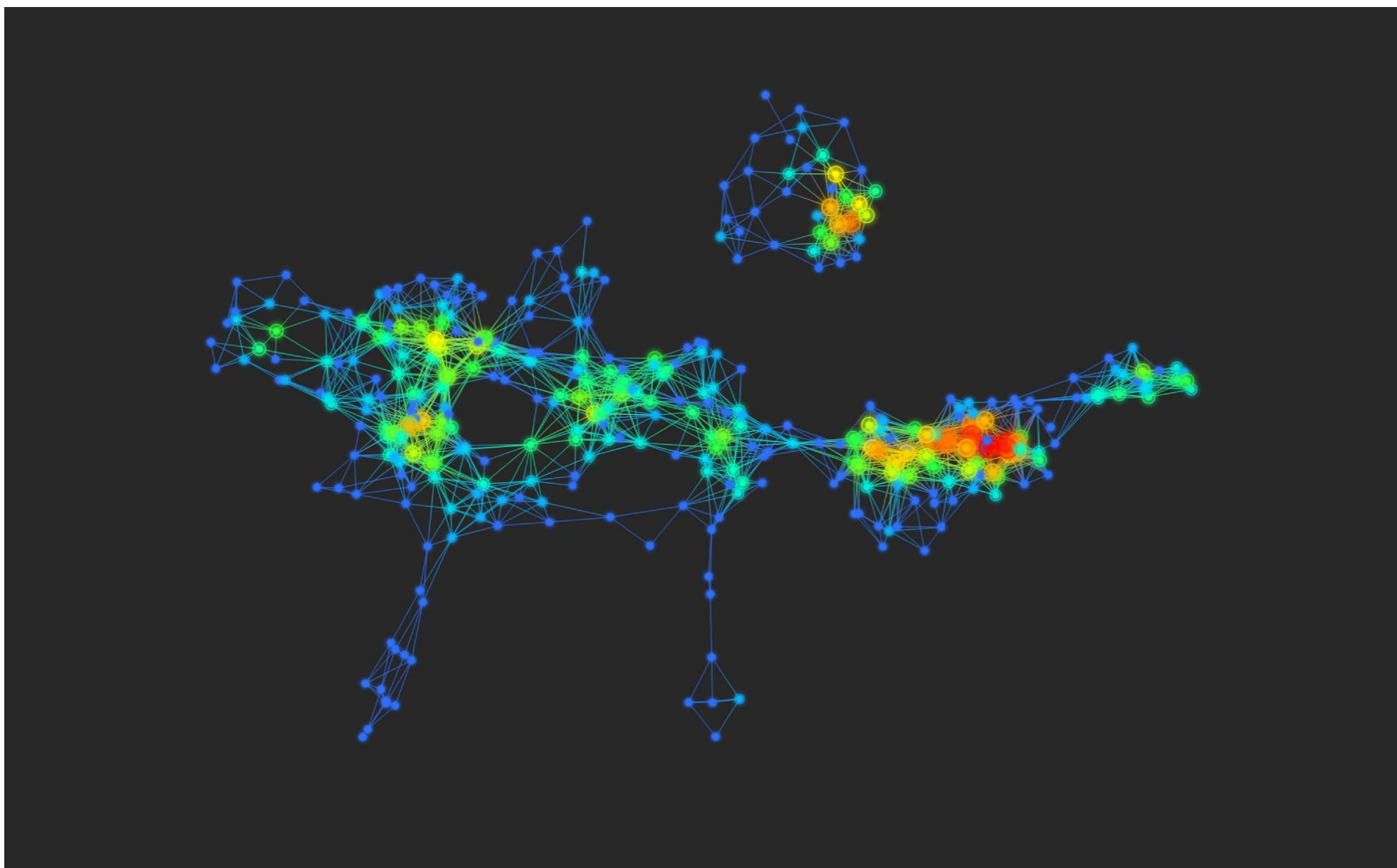
[Singh, Mémoli, Carlsson 2007]



visualize topology
on the data directly

Applications

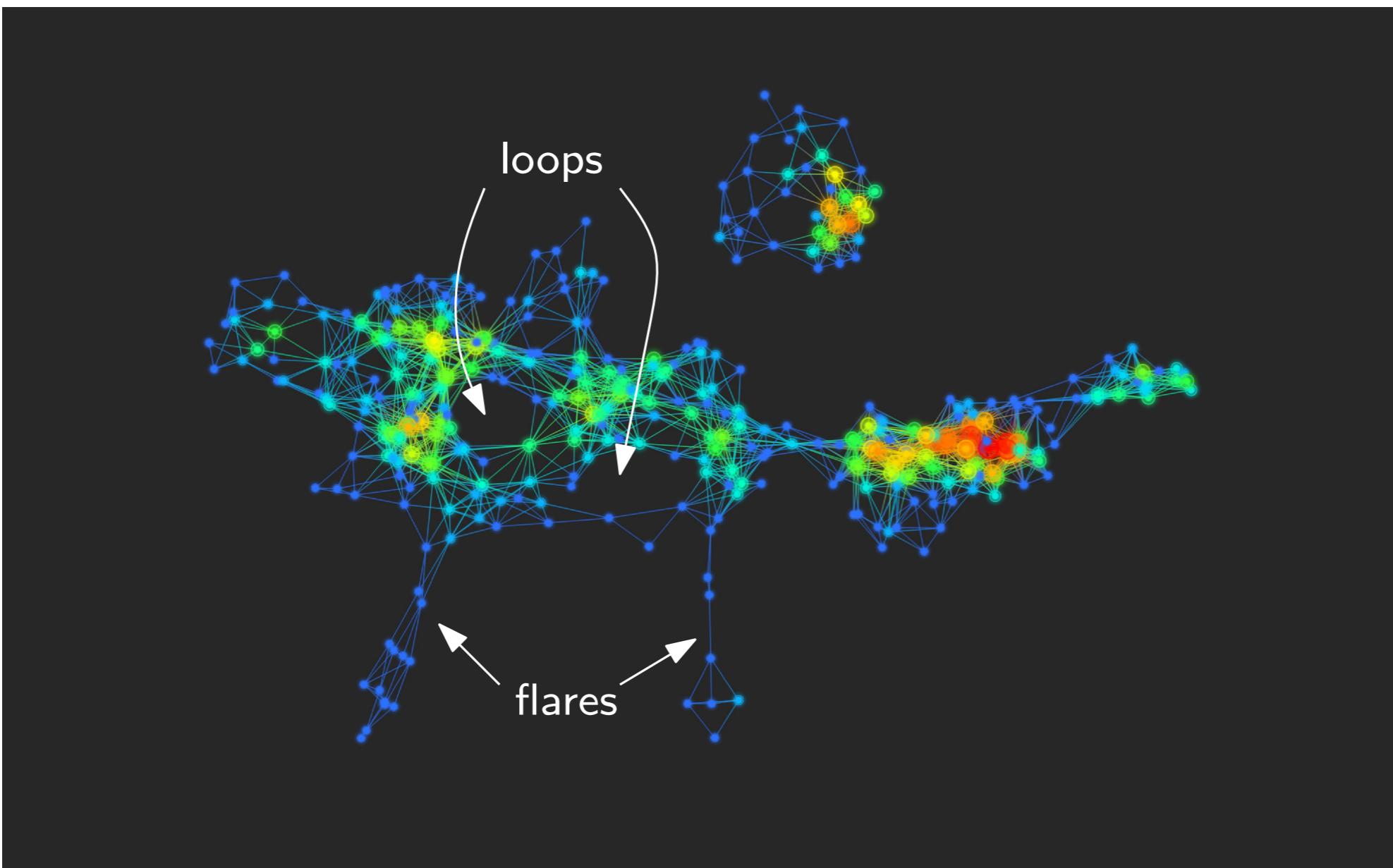
- visualization
- clustering
- feature selection



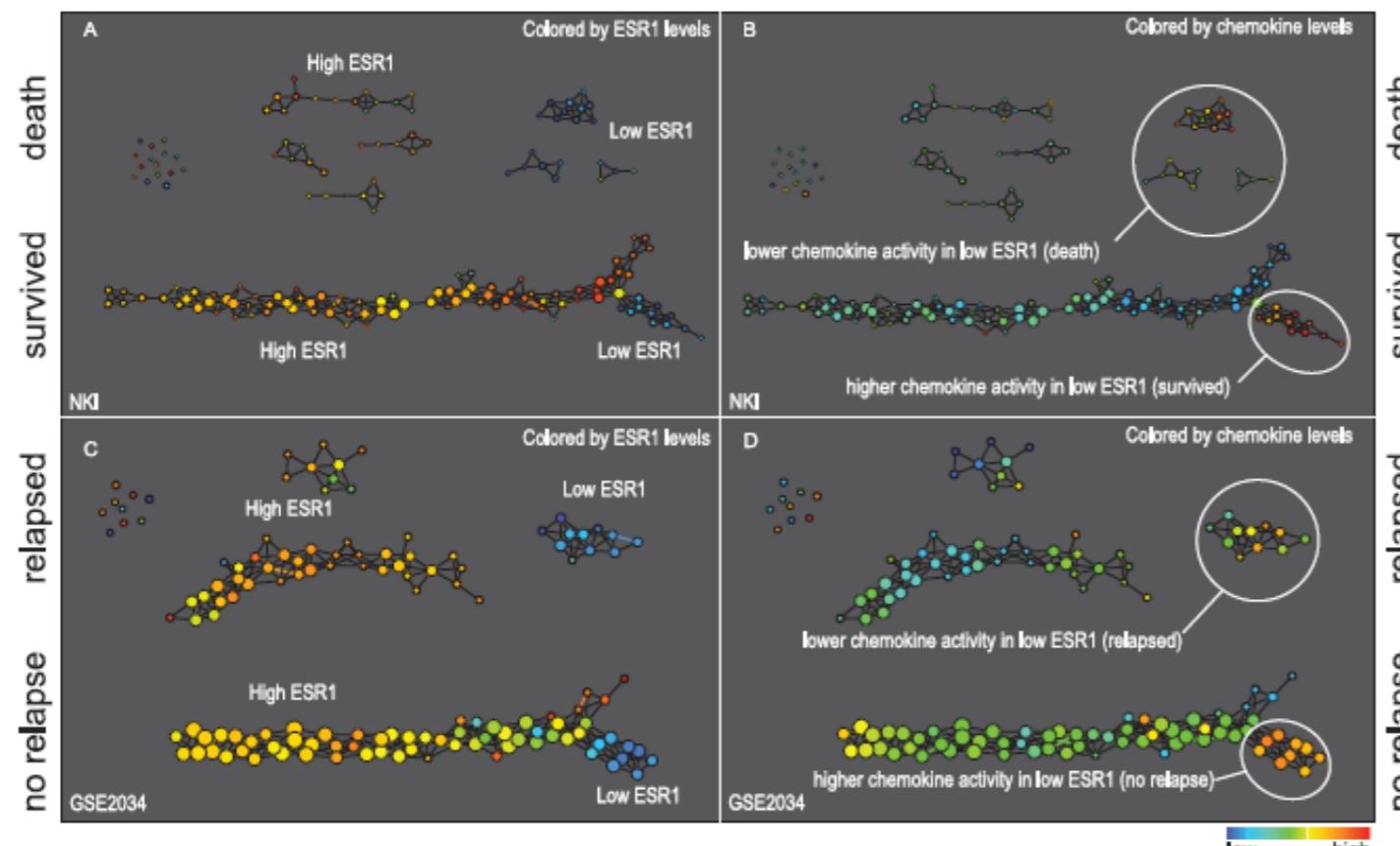
Applications

- visualization
- clustering
- feature selection

Principle: identify statistically relevant sub-populations through **patterns** (flares, loops)

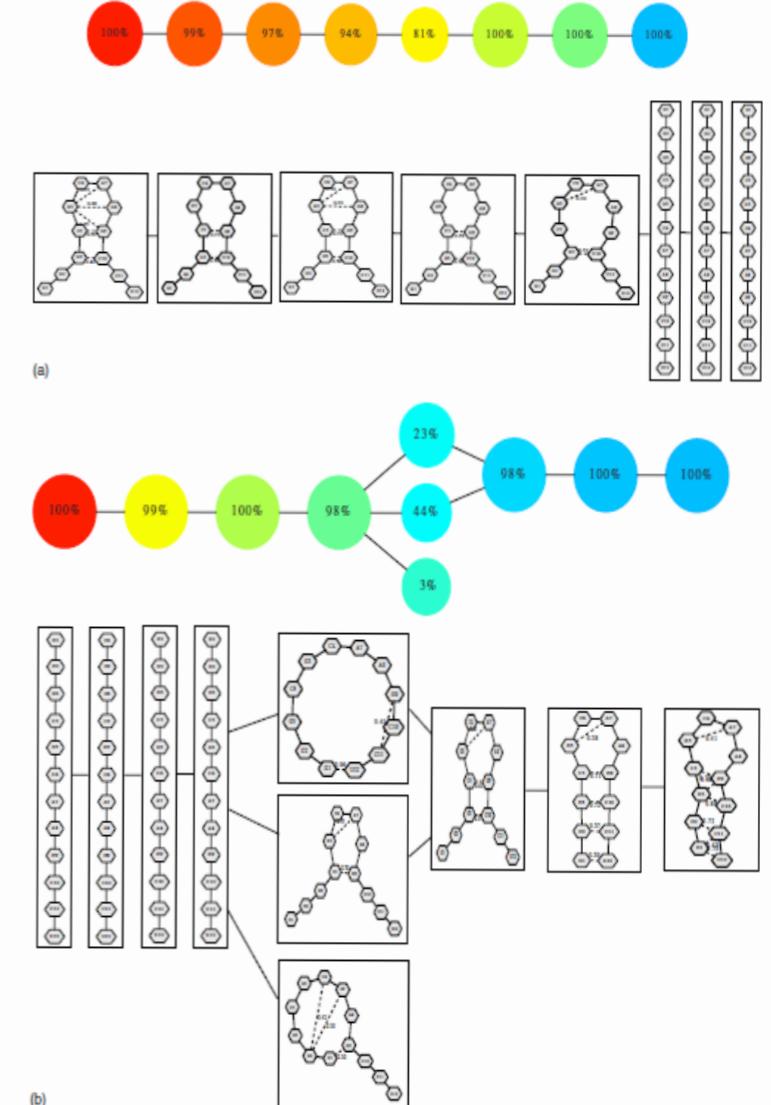


Applications



breast cancer subtype identification

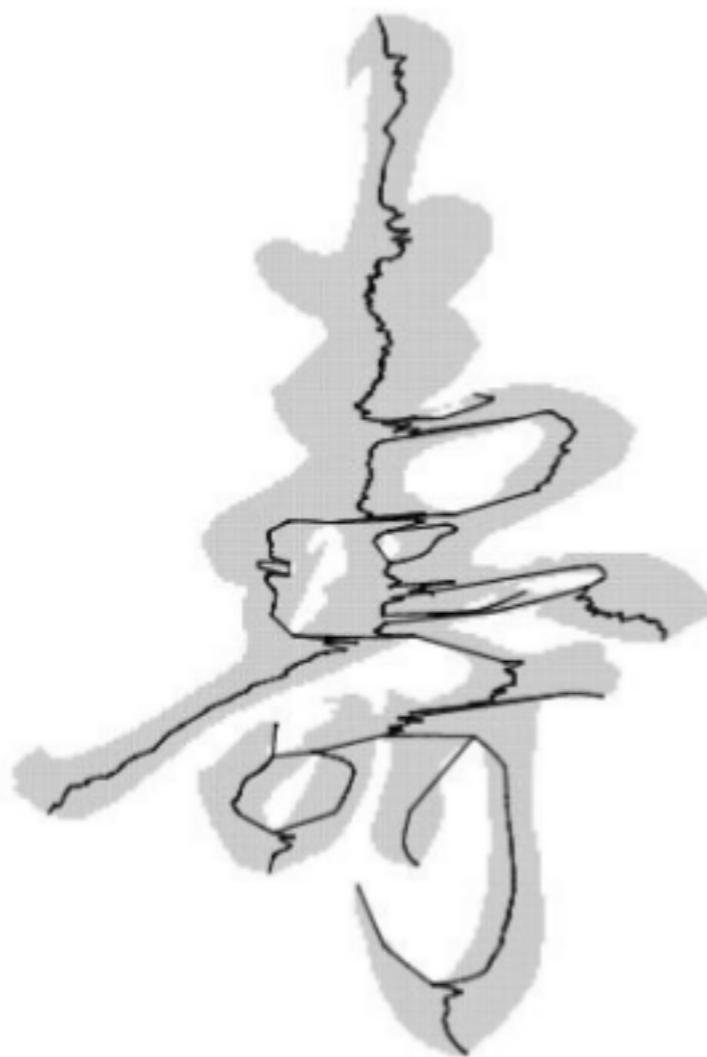
[Nicolau et al. 2011]



protein folding pathways

[Yao et al. 2009]

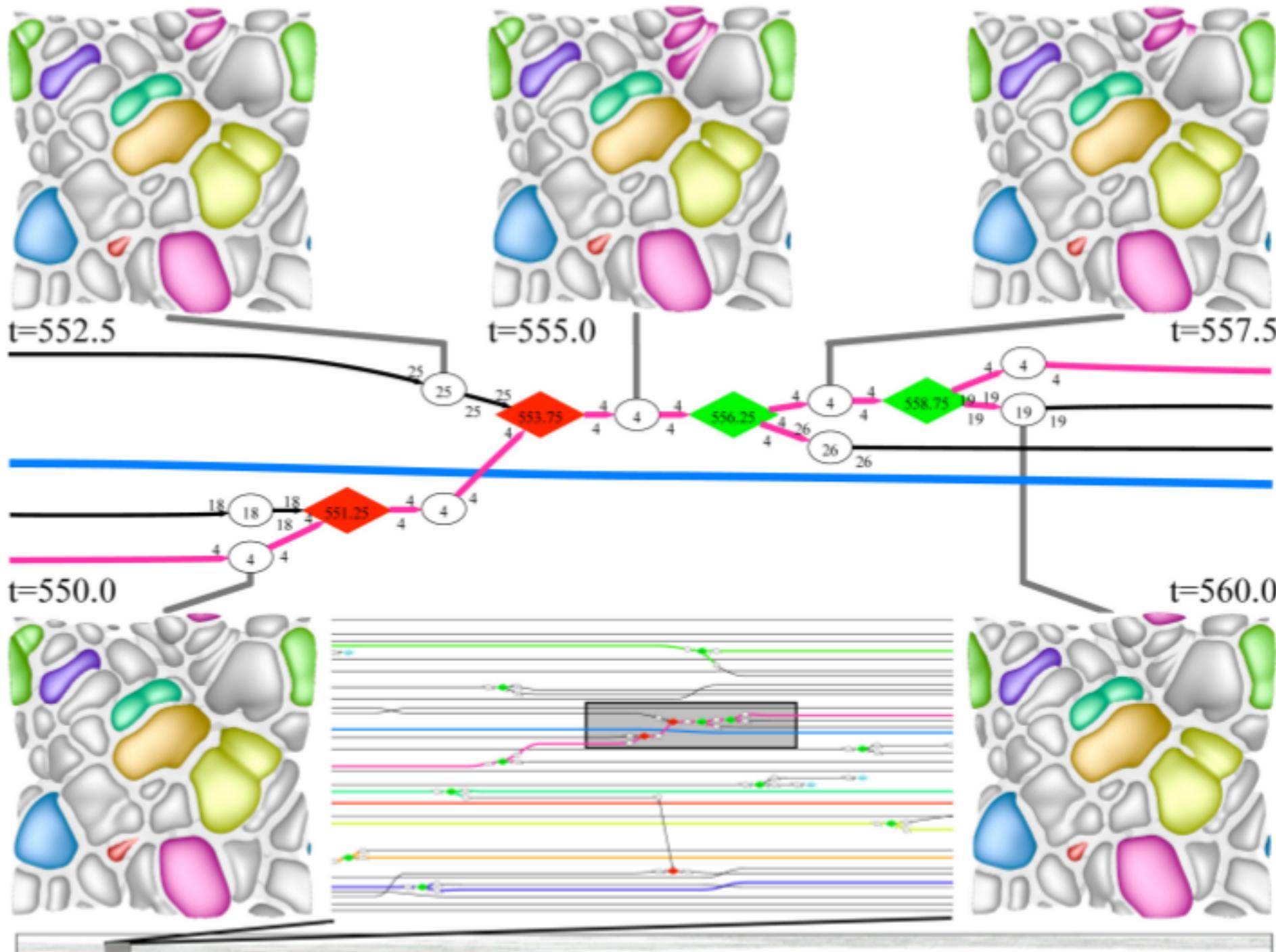
Applications



Data Skeletonization

[Ge et al. 2011]

Applications

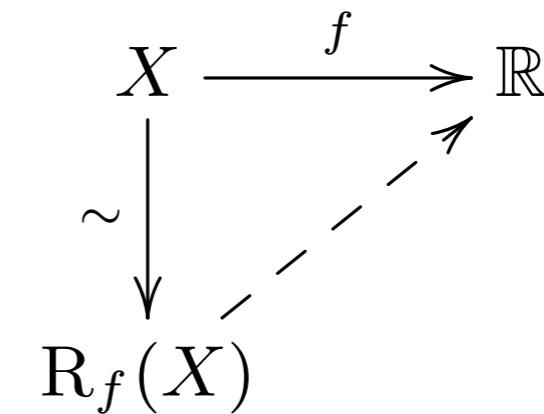
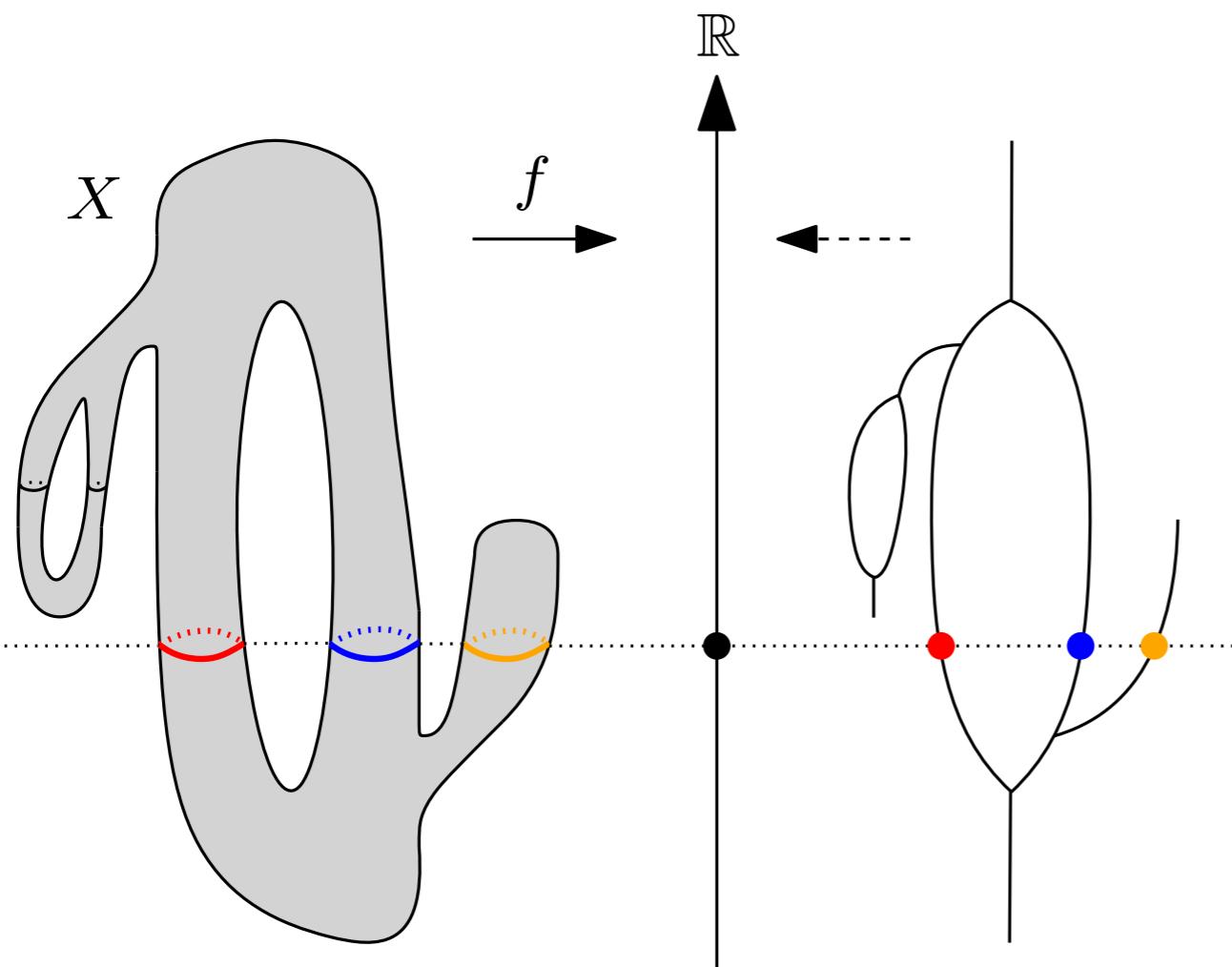


Burning regions of lean hydrogen flame over time
[Weber et al. 2011]

Reeb Graphs

$x \sim y \iff [f(x) = f(y) \text{ and } x, y \text{ belong to same cc of } f^{-1}(f(x))]$

Def: $R_f(X) = X / \sim$



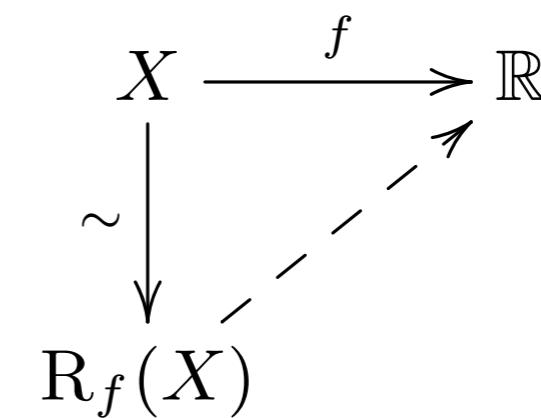
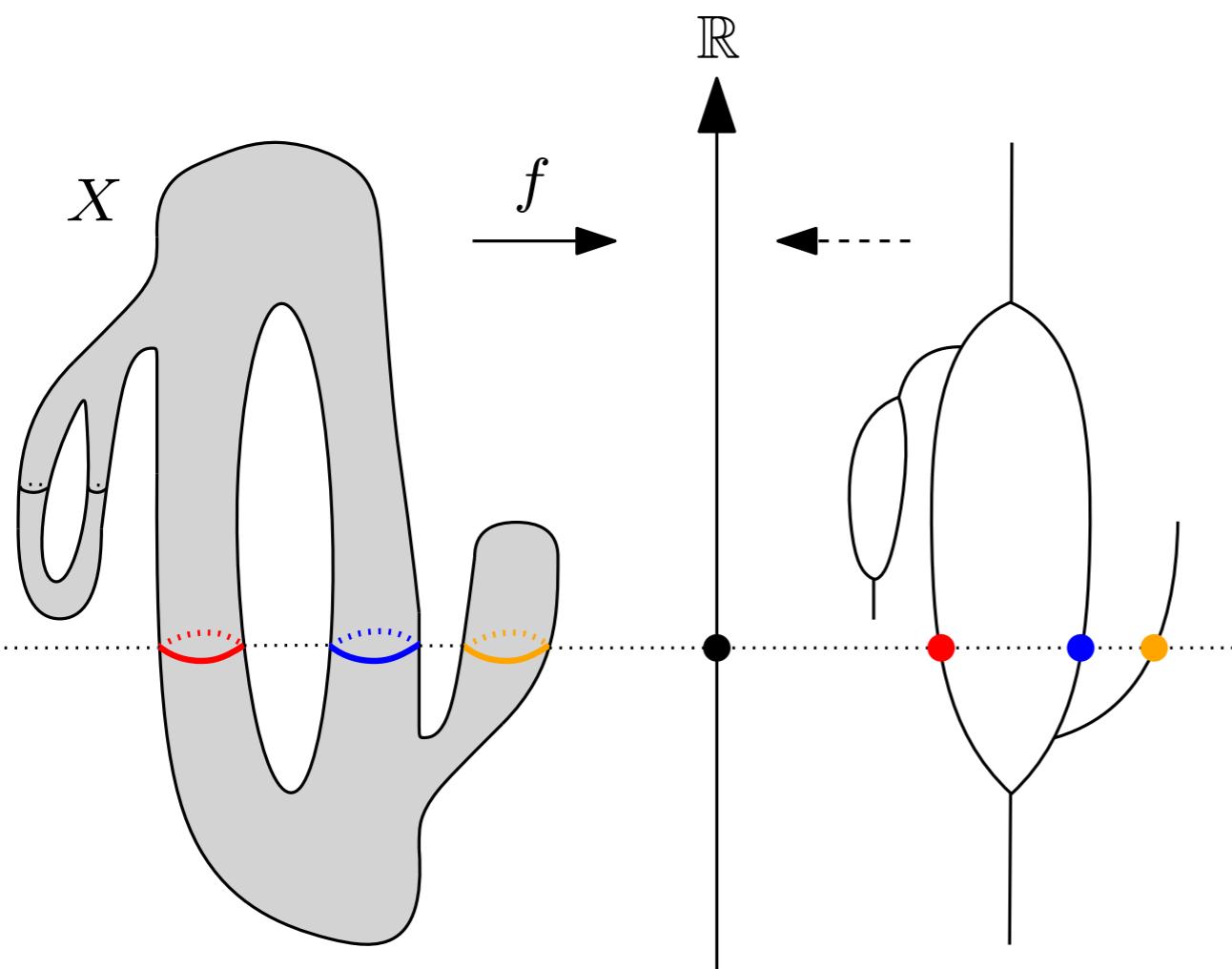
Prop: $R_f(X)$ is a graph if (X, f) is of **Morse type**

Reeb Graphs

$x \sim y \iff [f(x) = f(y) \text{ and } x, y \text{ belong to same cc of } f^{-1}(f(x))]$

Def: $R_f(X) = X / \sim$

Caveat: computation from point cloud is difficult



Prop: $R_f(X)$ is a graph if (X, f) is of **Morse type**

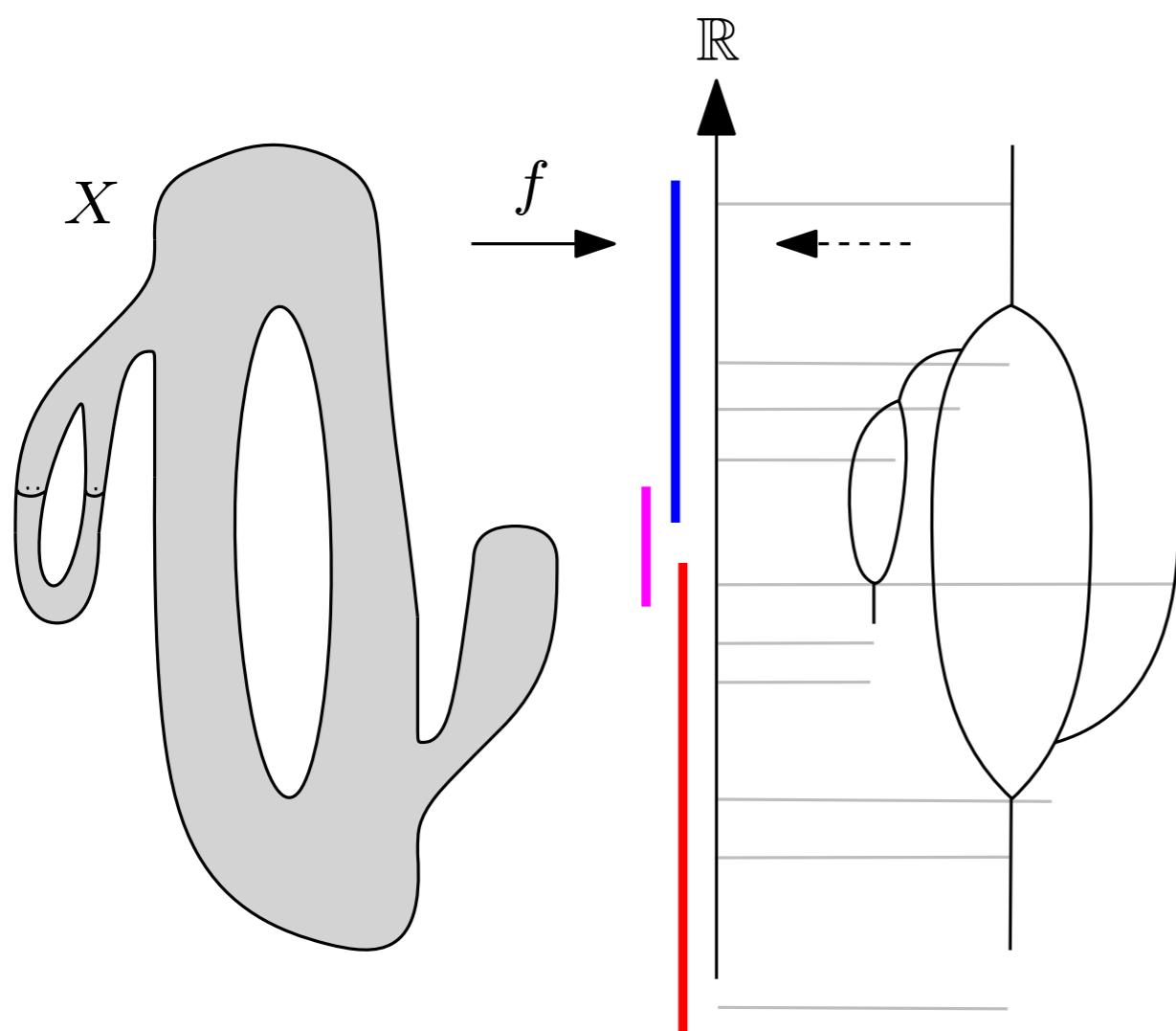
Reeb Graphs

$x \sim y \iff [f(x) = f(y) \text{ and } x, y \text{ belong to same cc of } f^{-1}(f(x))]$

Def: $R_f(X) = X / \sim$

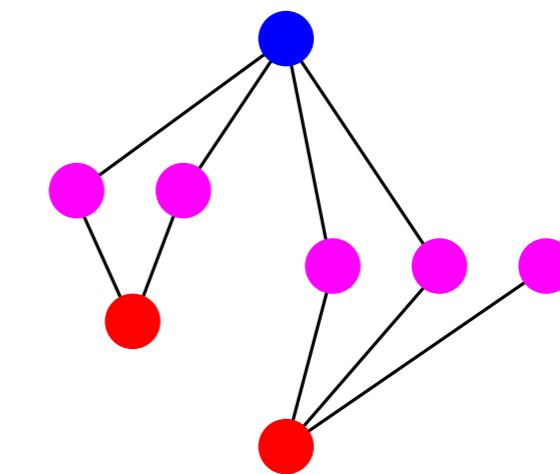
Caveat: computation from point cloud is difficult

→ Proxy: **Mapper**

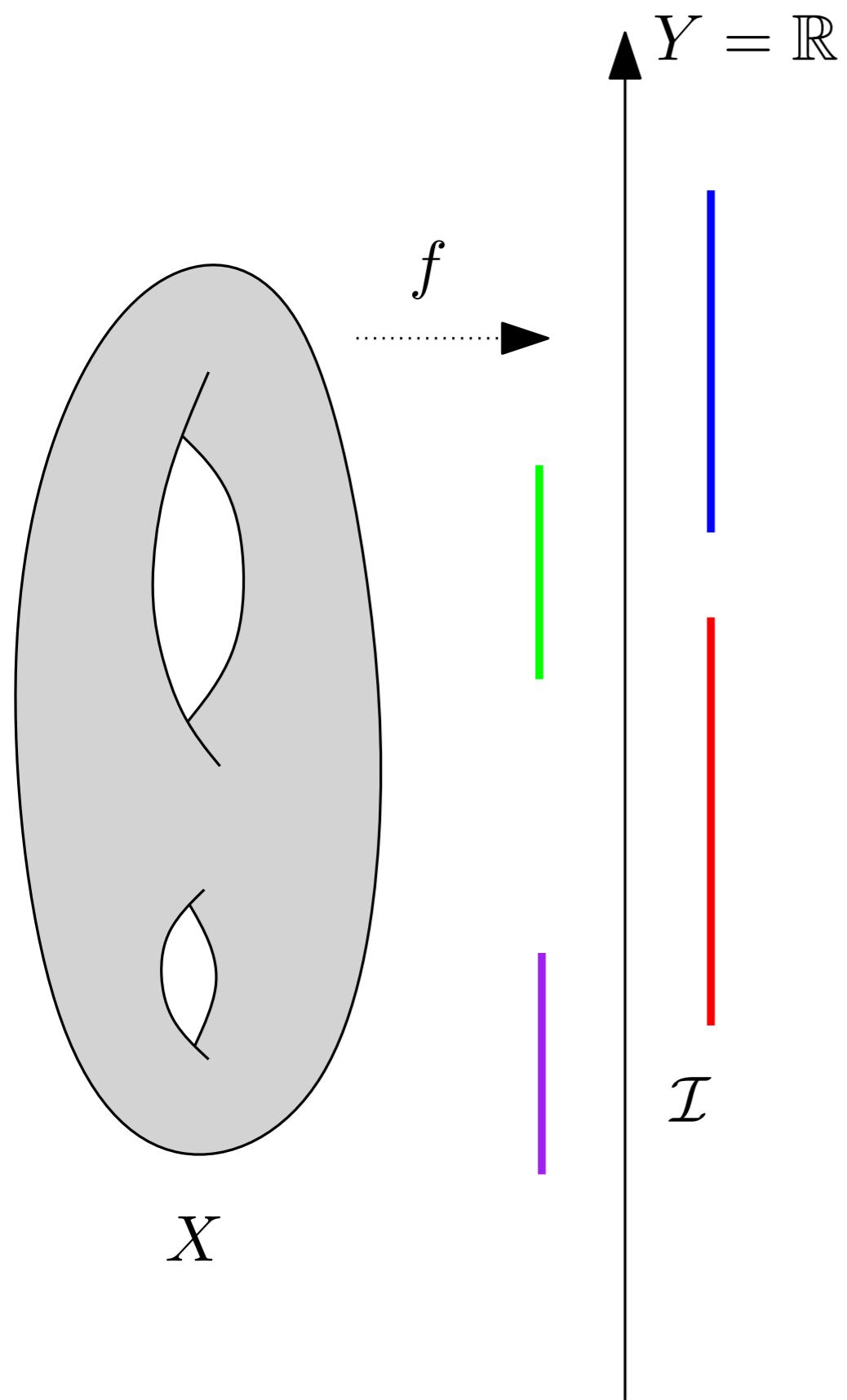


Pixelized Reeb graph

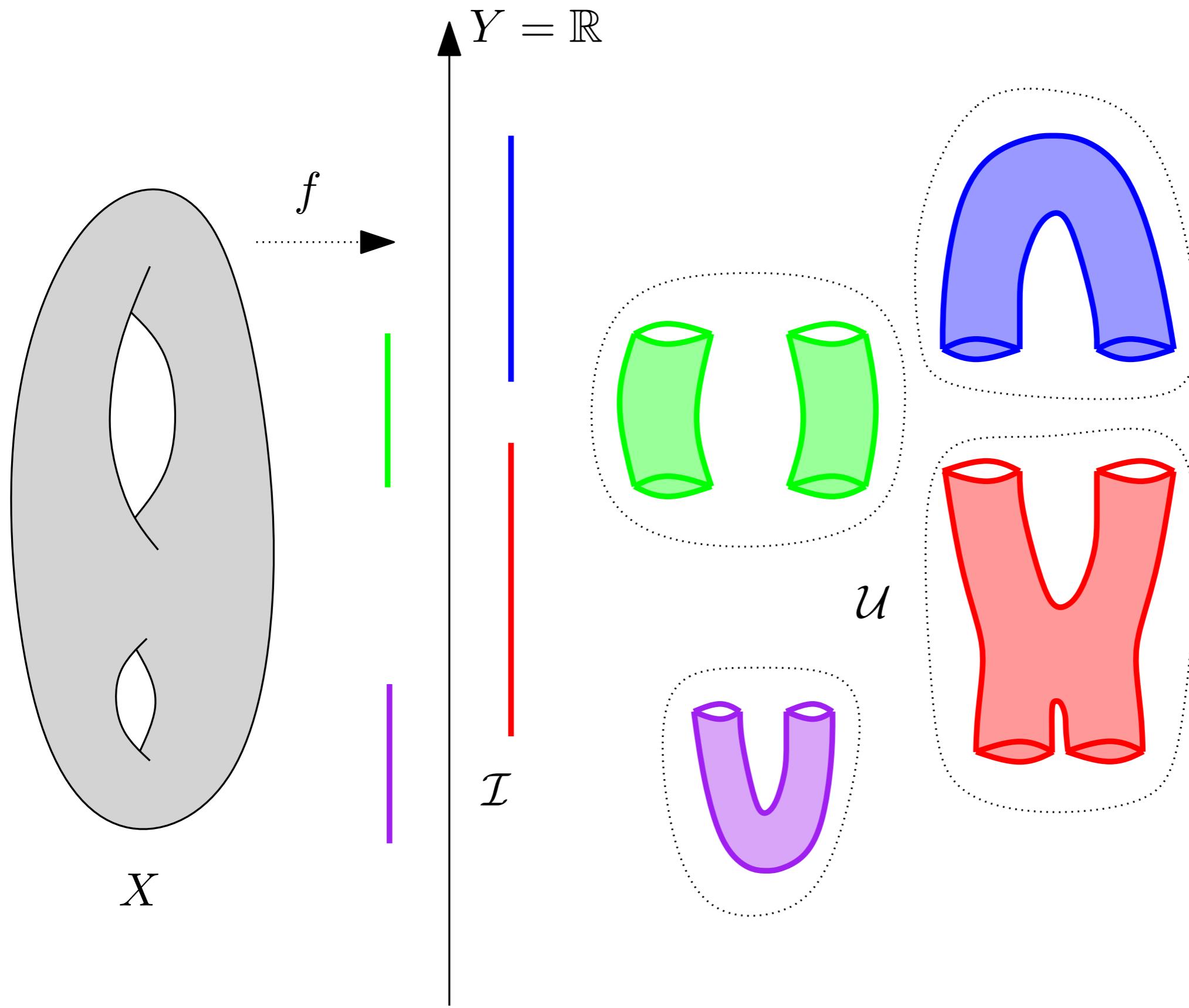
Singleton → interval: $f^{-1}(I)$



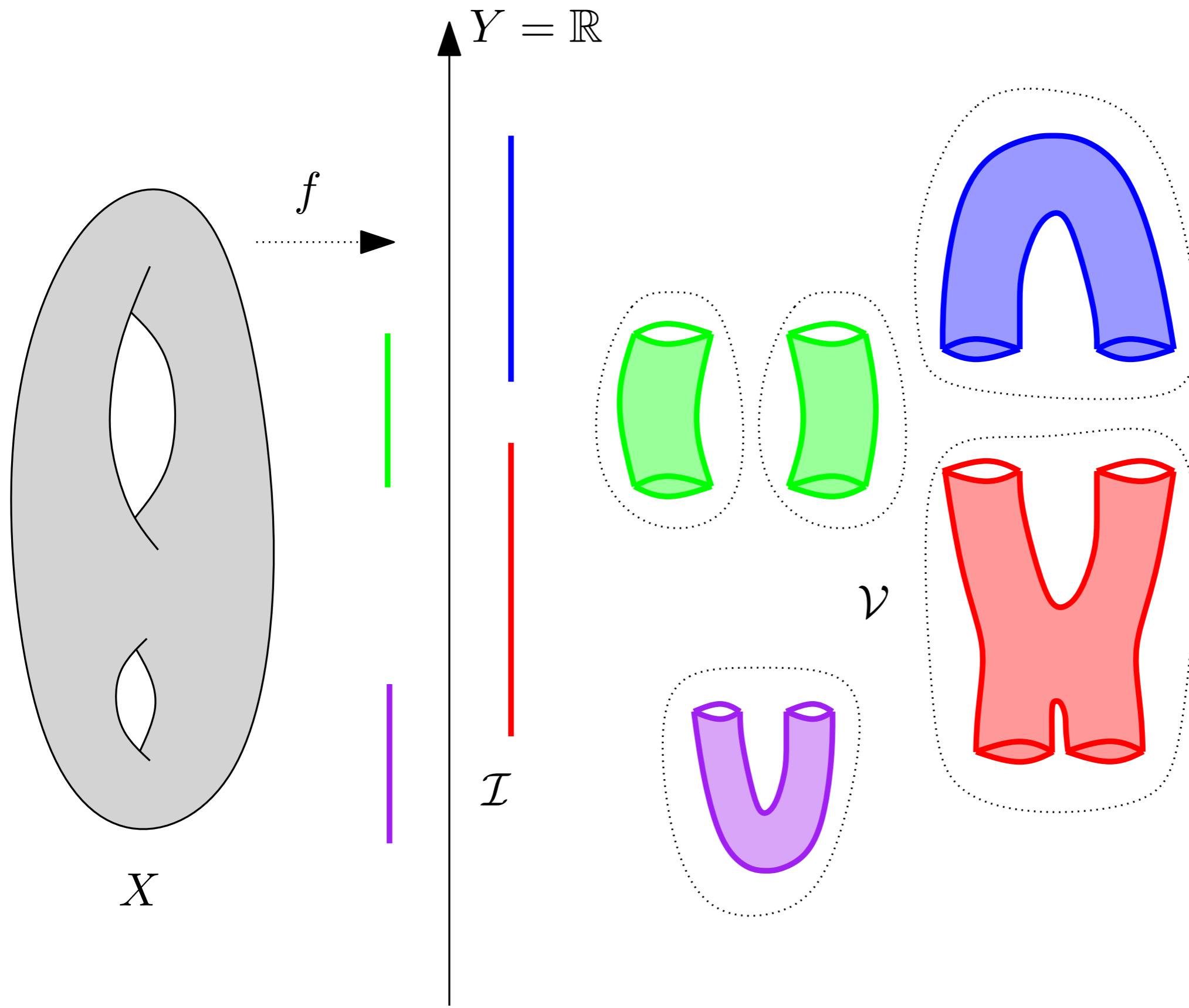
Mapper (continuous setting)



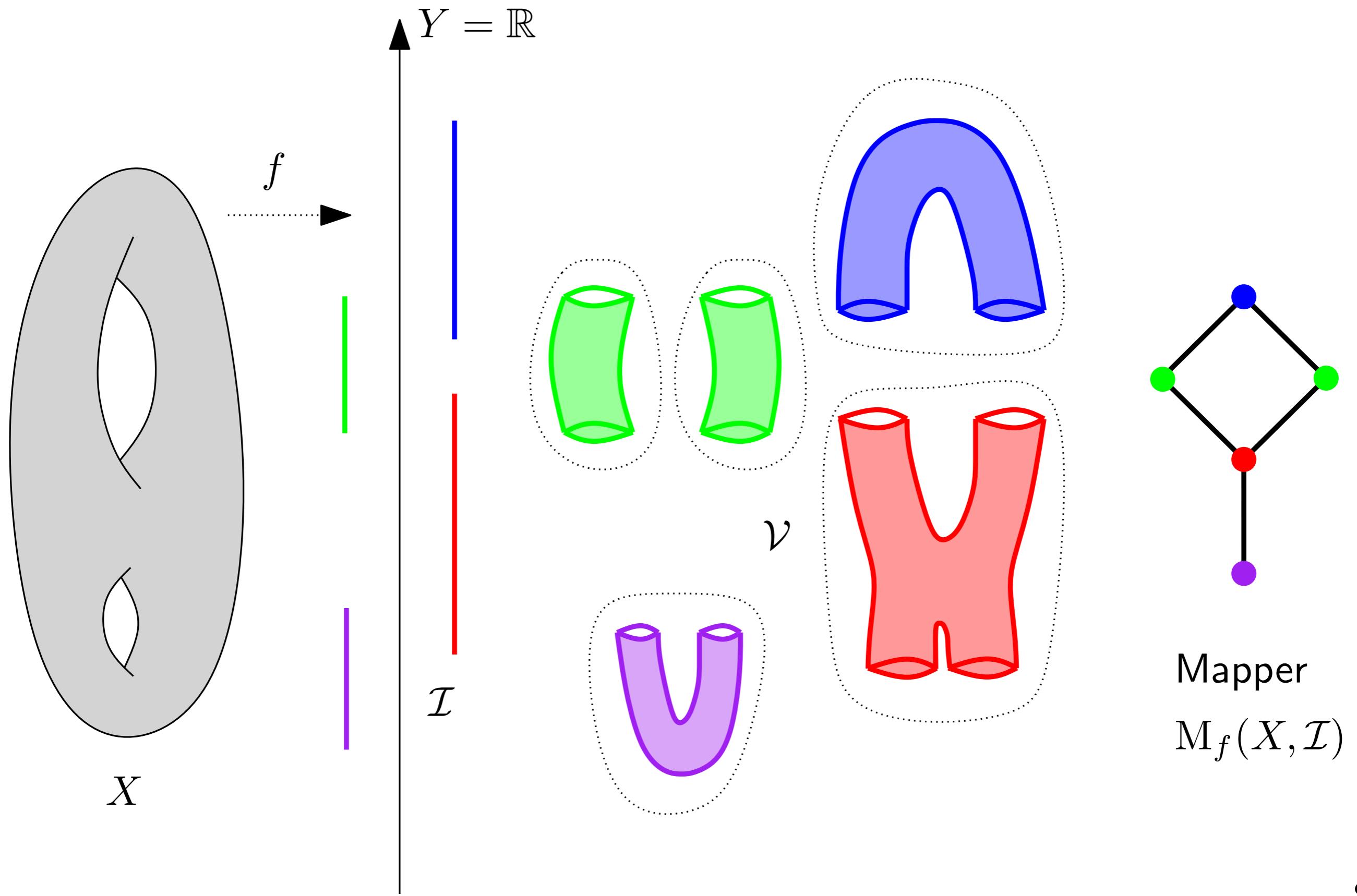
Mapper (continuous setting)



Mapper (continuous setting)



Mapper (continuous setting)



Mapper (continuous setting)

Input:

- topological space X
- continuous function $f : X \rightarrow \mathbb{R}$
- cover \mathcal{I} of $\text{im}(f)$ by open intervals: $\text{im}(f) \subseteq \bigcup_{I \in \mathcal{I}} I$

Method:

- Compute *pullback cover* \mathcal{U} of X : $\mathcal{U} = \{f^{-1}(I)\}_{I \in \mathcal{I}}$
- Refine \mathcal{U} by separating the connected components
- The Mapper is the *nerve* of \mathcal{V} :
 - 1 vertex per element $V \in \mathcal{V}$
 - 1 edge per intersection $V \cap V' \neq \emptyset, V, V' \in \mathcal{V}$
 - 1 k -simplex per $(k + 1)$ -fold intersection $\bigcap_{i=0}^k V_i \neq \emptyset, V_0, \dots, V_k \in \mathcal{V}$

Mapper (discrete setting)

Input:

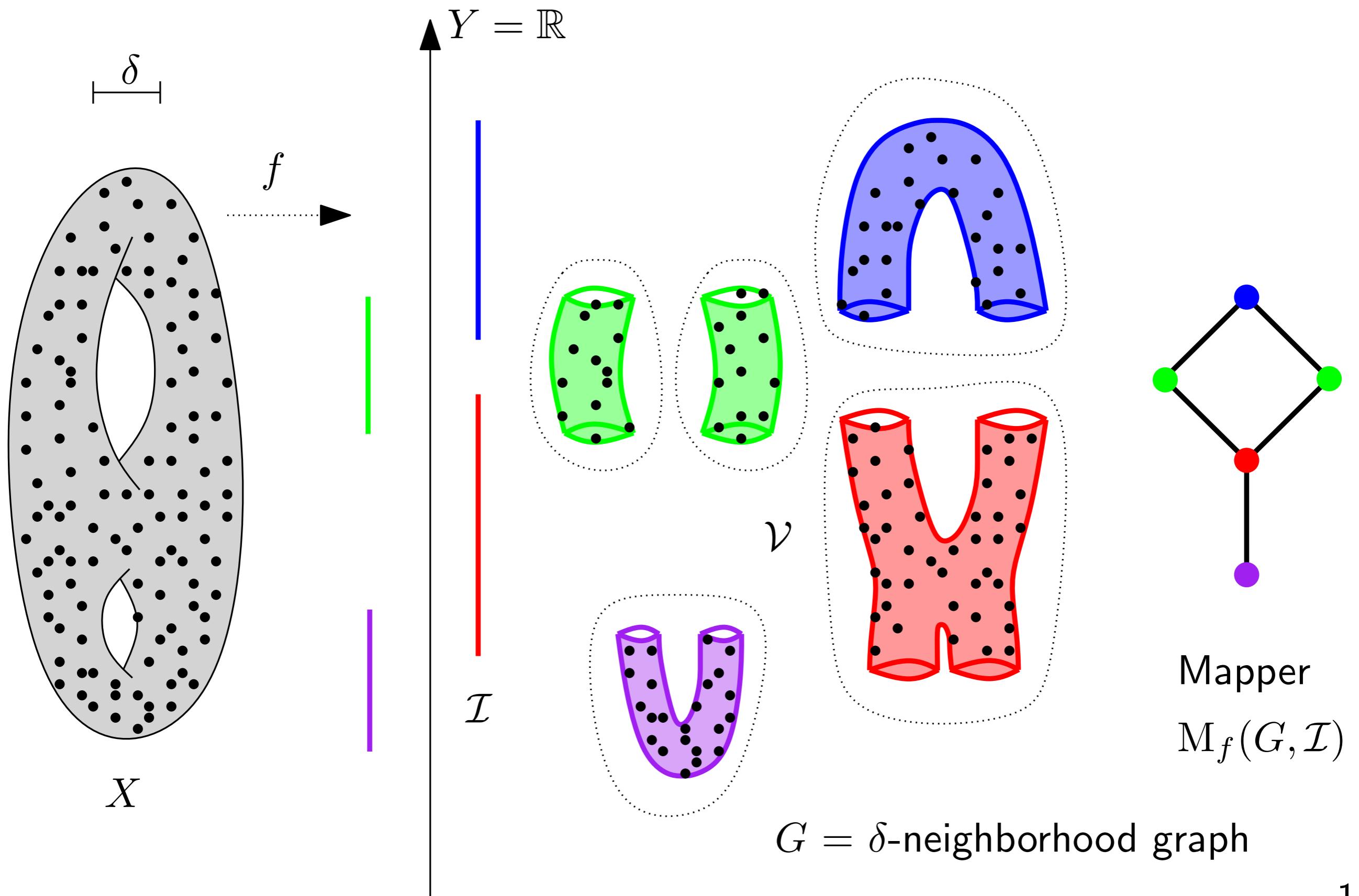
- point cloud $P \subseteq X$ with metric d_P
- continuous function $f : P \rightarrow \mathbb{R}$
- cover \mathcal{I} of $\text{im}(f)$ by open intervals: $\text{im} f \subseteq \bigcup_{I \in \mathcal{I}} I$

Method:

- Compute *pullback cover* \mathcal{U} of P : $\mathcal{U} = \{f^{-1}(I)\}_{I \in \mathcal{I}}$
- Refine \mathcal{U} by separating each of its elements into its various connected components in G → connected cover \mathcal{V}
- The Mapper is the *nerve* of \mathcal{V} :
 - 1 vertex per element $V \in \mathcal{V}$
 - 1 edge per intersection $V \cap V' \neq \emptyset, V, V' \in \mathcal{V}$
 - 1 k -simplex per $(k + 1)$ -fold intersection $\bigcap_{i=0}^k V_i \neq \emptyset, V_0, \dots, V_k \in \mathcal{V}$

(intersections materialized
by data points)

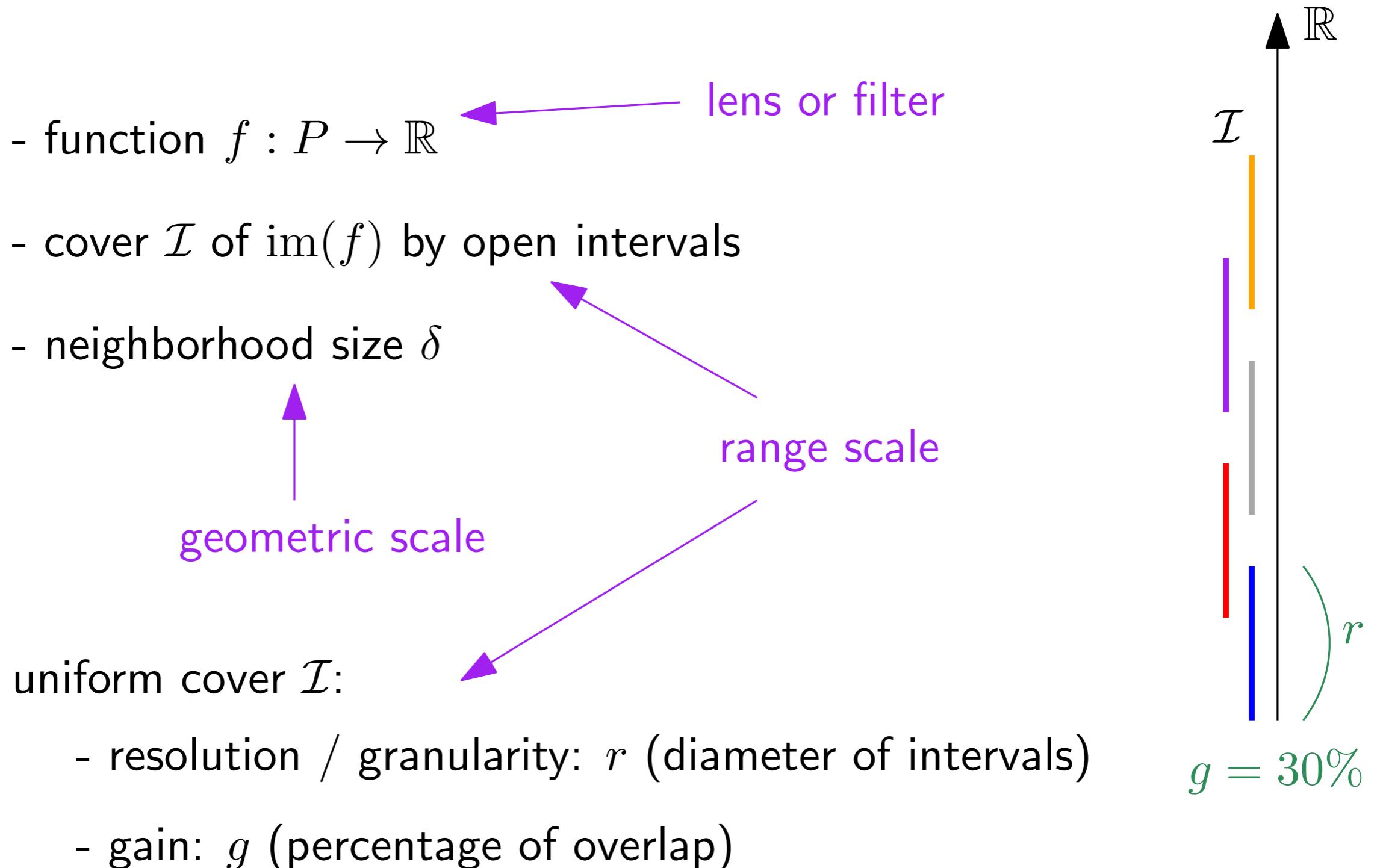
Mapper (discrete setting)



Parameters

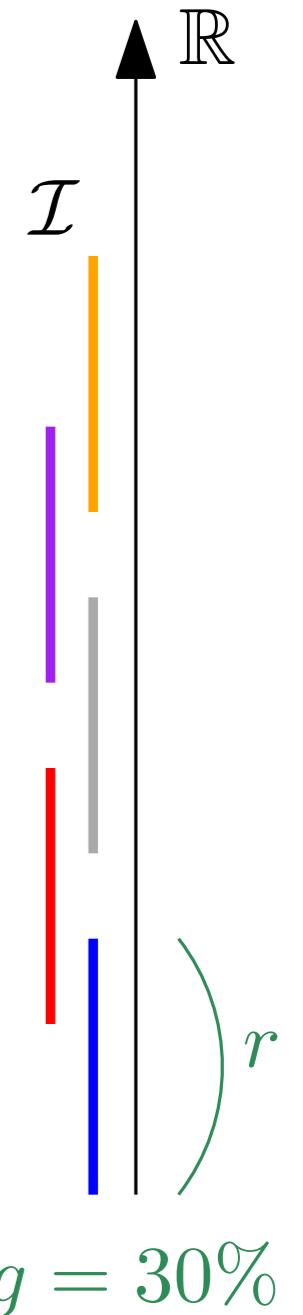
- function $f : P \rightarrow \mathbb{R}$
 - cover \mathcal{I} of $\text{im}(f)$ by open intervals
 - neighborhood size δ
- geometric scale
- range scale
- lens or filter
-
- ```
graph TD; A["function f : P → ℝ"] --> B["cover I of im(f) by open intervals"]; A --> C["neighborhood size δ"]; D["geometric scale"] --> C; E["range scale"] --> C; F["lens or filter"] --> A;
```

# Parameters



# Parameters

- function  $f : P \rightarrow \mathbb{R}$
  - cover  $\mathcal{I}$  of  $\text{im}(f)$  by open intervals
  - neighborhood size  $\delta$
- uniform cover  $\mathcal{I}$ :
- resolution / granularity:  $r$  (diameter of intervals)
  - gain:  $g$  (percentage of overlap)



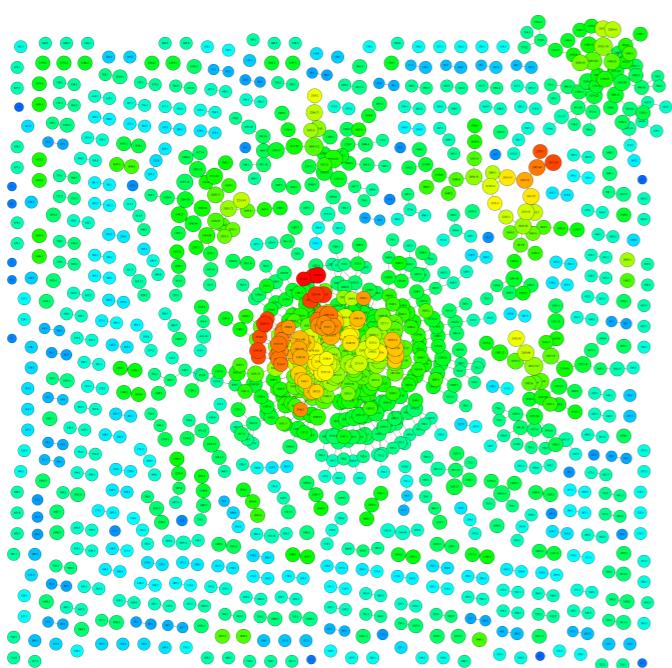
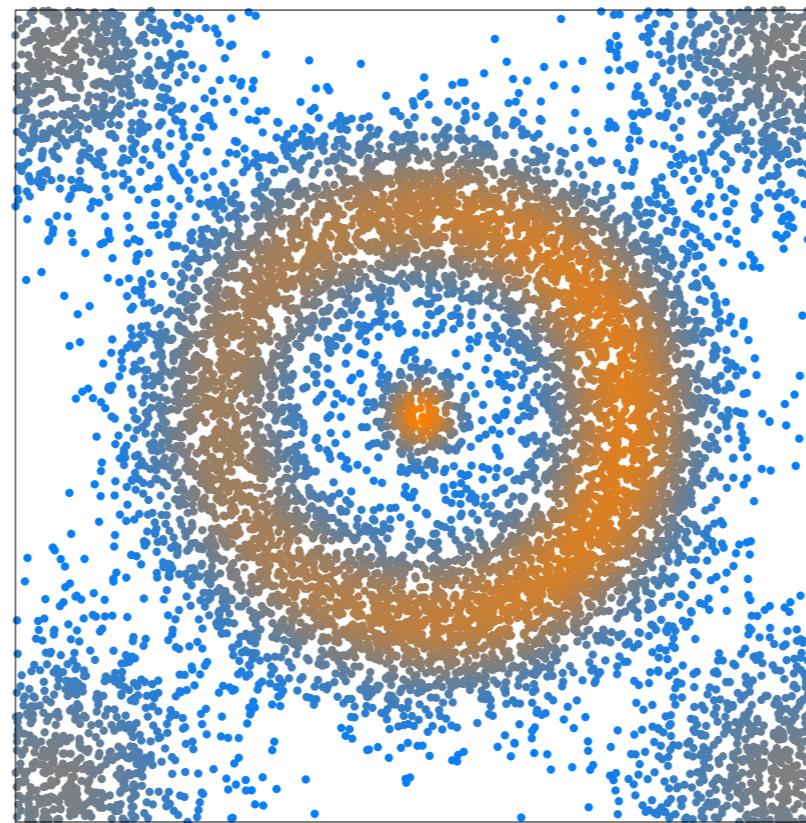
**In practice:** trial-and-error

# Examples

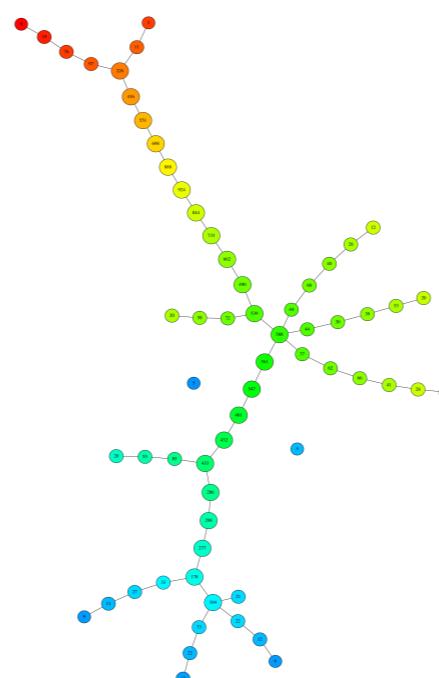
---

$\hat{f}$  = density estimator

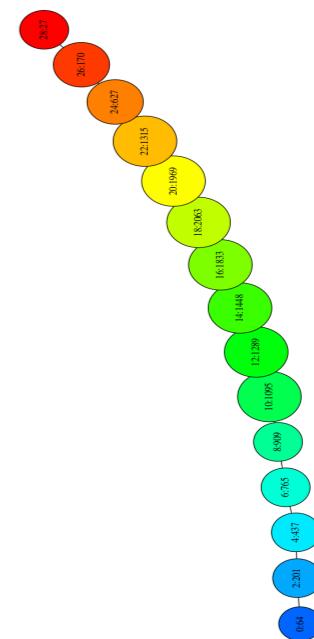
$r = 0.3, g = 20\%$



$\delta = 0.1\%$

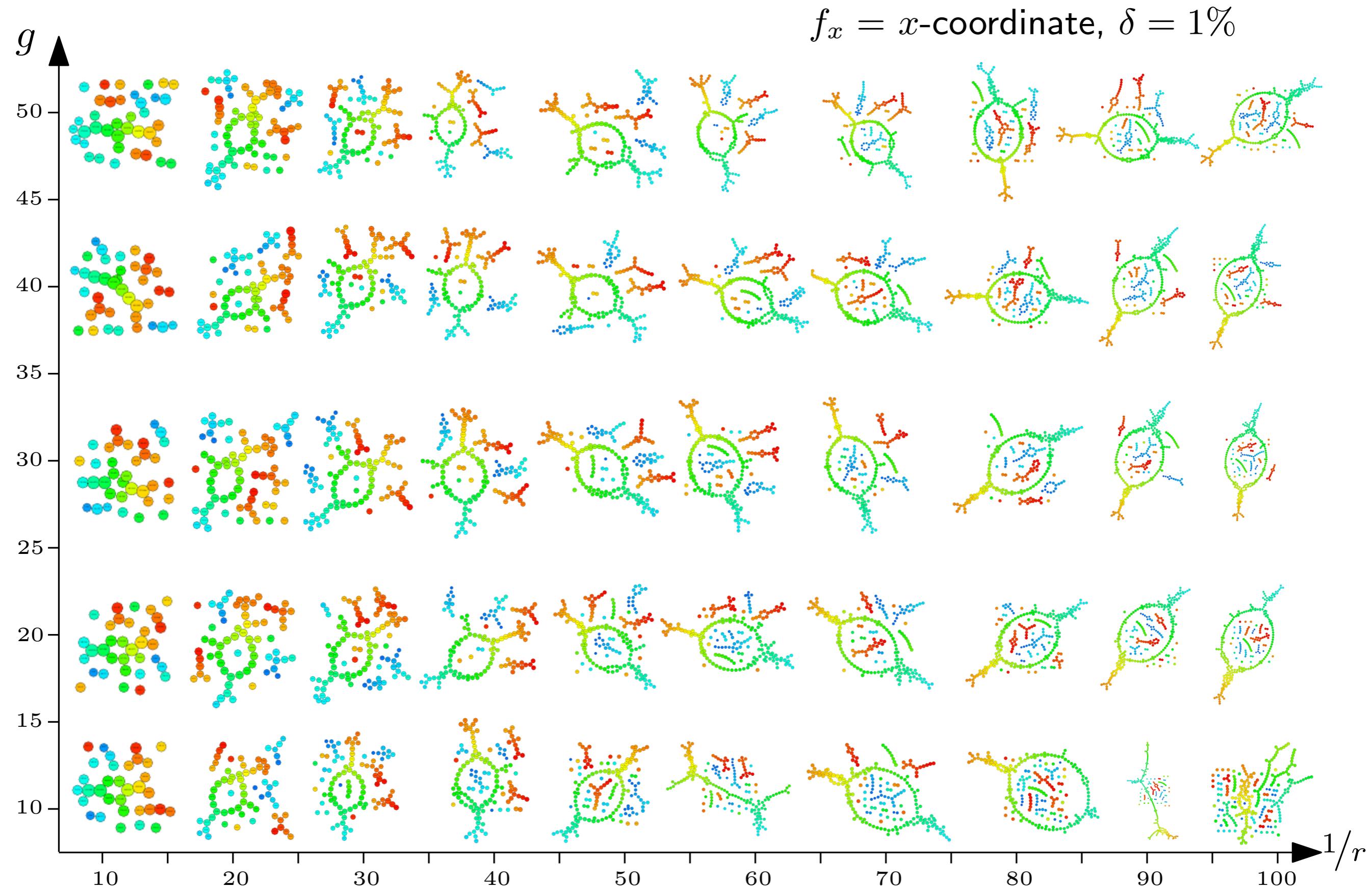


$\delta = 1\%$

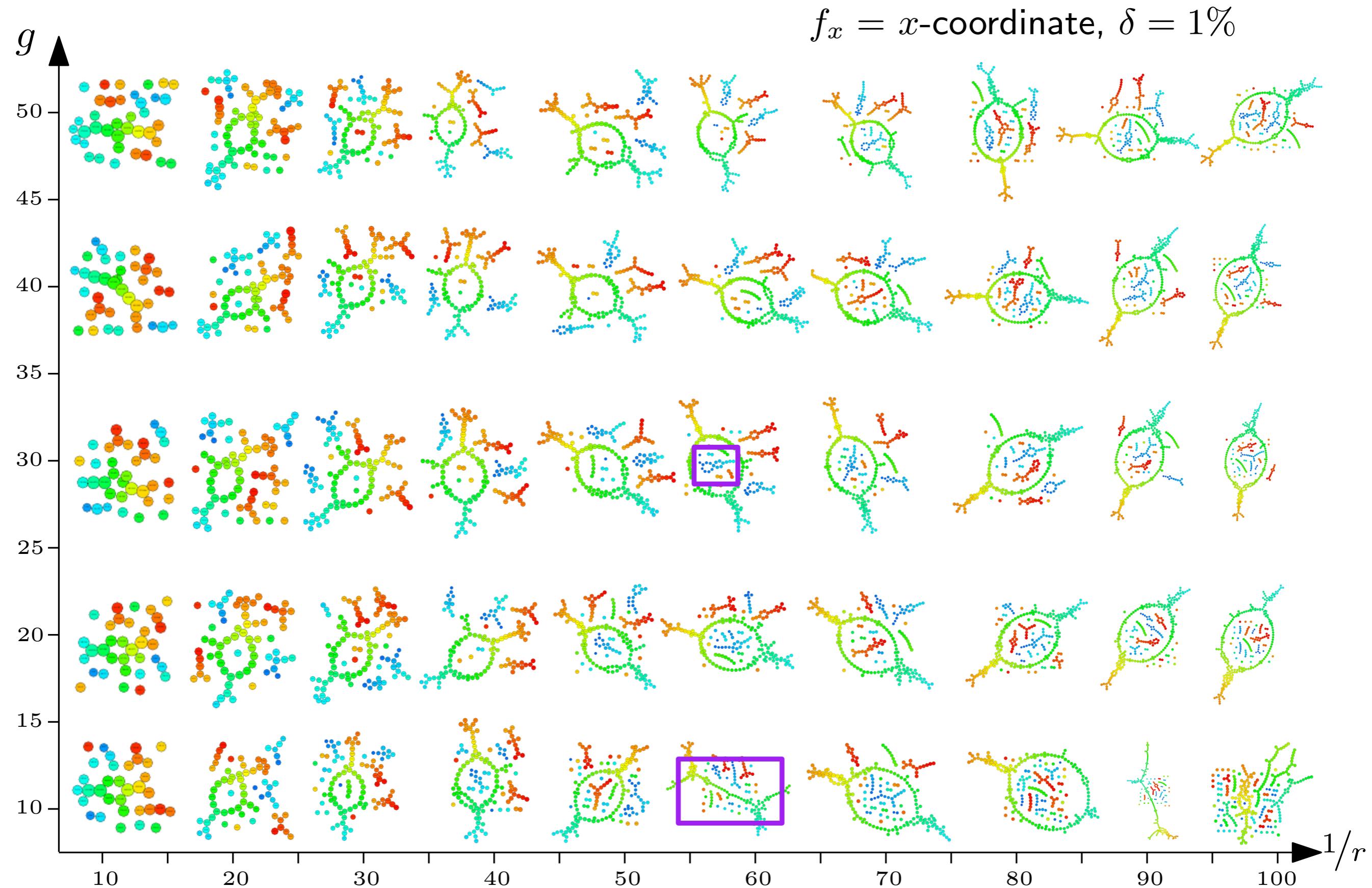


$\delta = 10\%$

# Examples



# Examples



# Contributions

---

Structure of Mapper depending on  $r, g, \delta$  [C., Oudot 2016-17]

Heuristics and Statistics for Mapper [C., Michel, Oudot 2017]

Local Equivalence of Mapper and PDs [C., Oudot 2017]

Machine Learning for Mapper (through PDs):

Finite dimensional kernel [C., Ovsjanikov, Oudot 2015]

Gaussian kernel [C., Cuturi, Oudot 2017]

# Contributions

---

Structure of Mapper depending on  $r, g, \delta$  [C., Oudot 2016-17]

Heuristics and Statistics for Mapper [C., Michel, Oudot 2017]

Local Equivalence of Mapper and PDs [C., Oudot 2017]

Machine Learning for Mapper (through PDs):

Finite dimensional kernel [C., Ovsjanikov, Oudot 2015]

Gaussian kernel [C., Cuturi, Oudot 2017]

# Extended Persistence Diagram

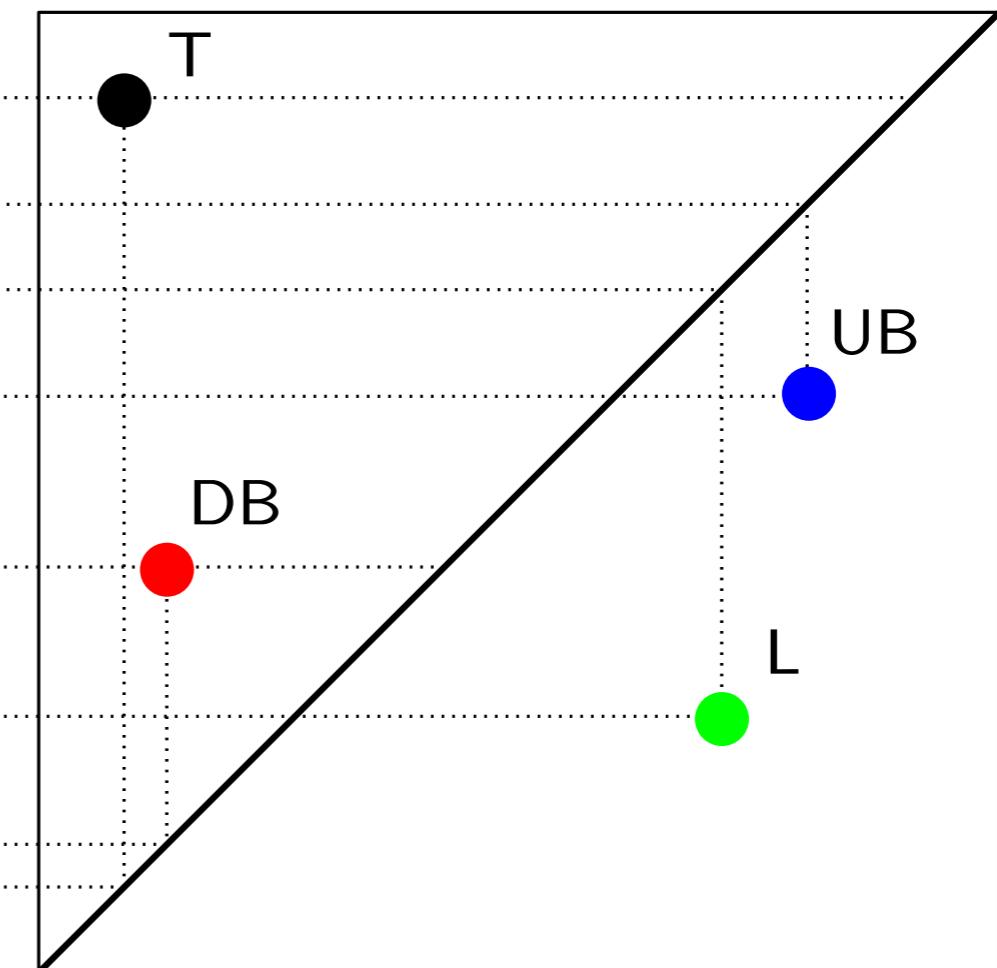
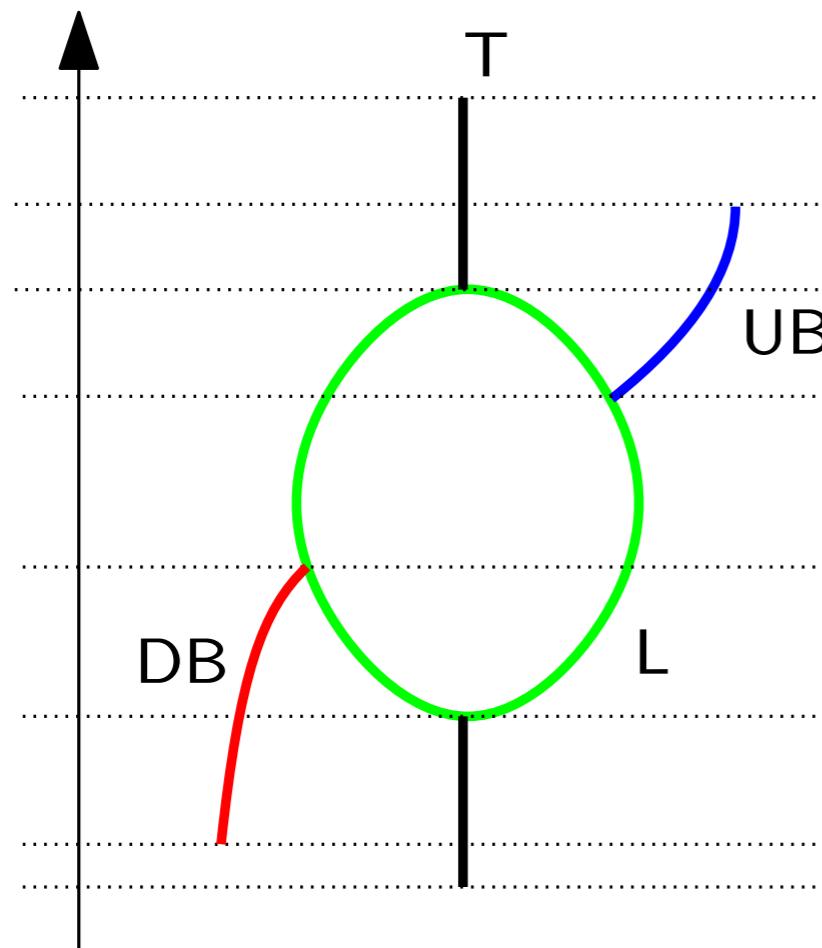
**Def:**  $Dg(R_f) = \text{bag-of-features}$  descriptor for  $R_f(X)$ :

$DB(R_f) \longleftrightarrow \text{downward branches}$

$UB(R_f) \longleftrightarrow \text{upward branches}$

$T(R_f) \longleftrightarrow \text{trunks (cc)}$

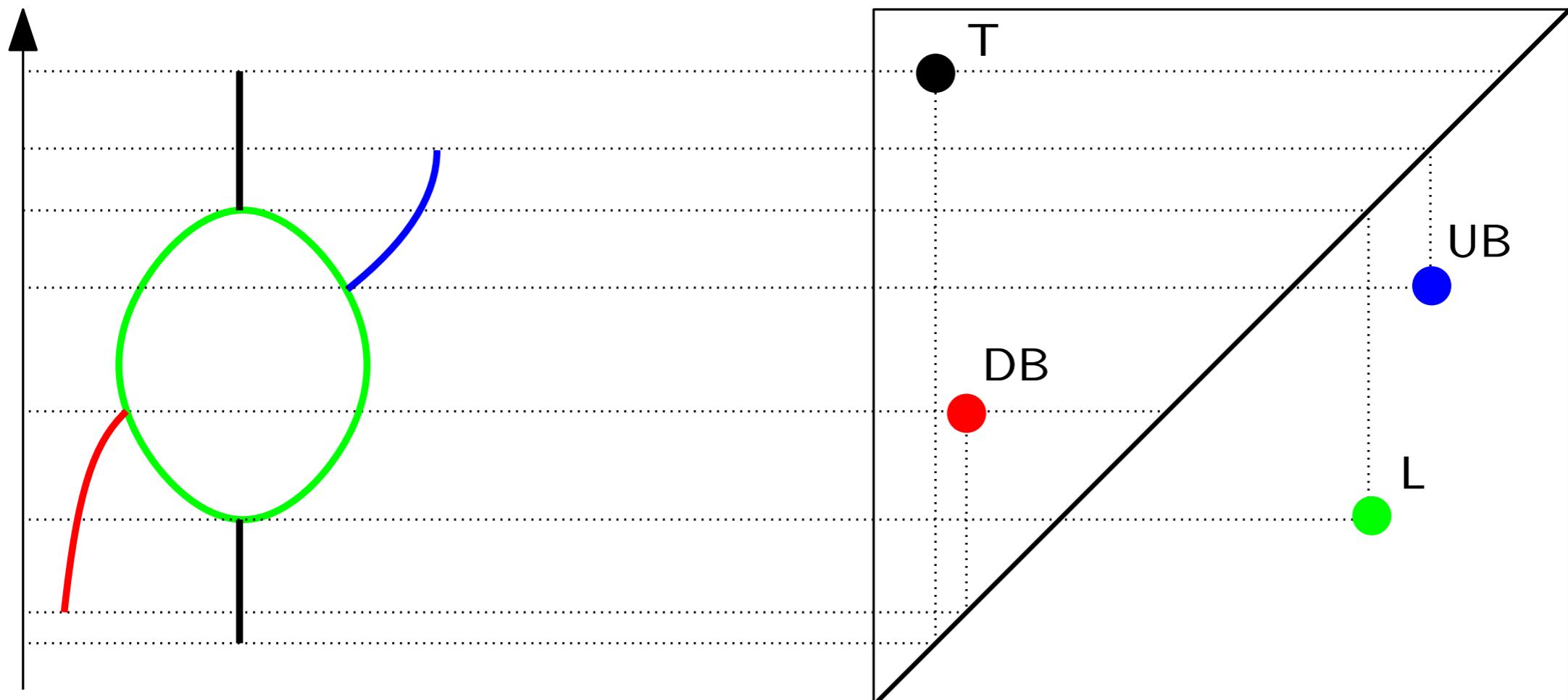
$L(R_f) \longleftrightarrow \text{loops}$



# Extended Persistence Diagram

Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

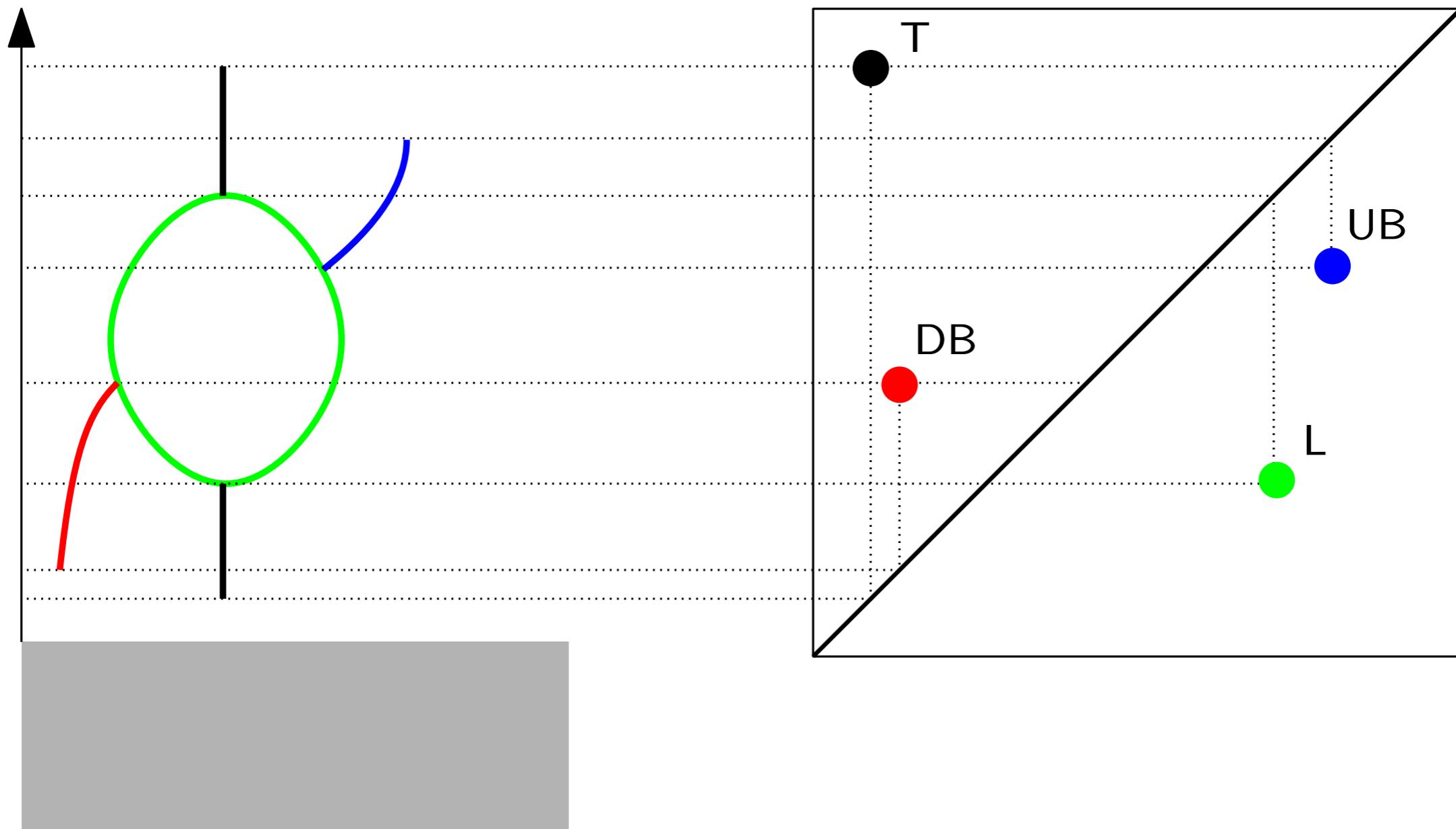
- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Extended Persistence Diagram

Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

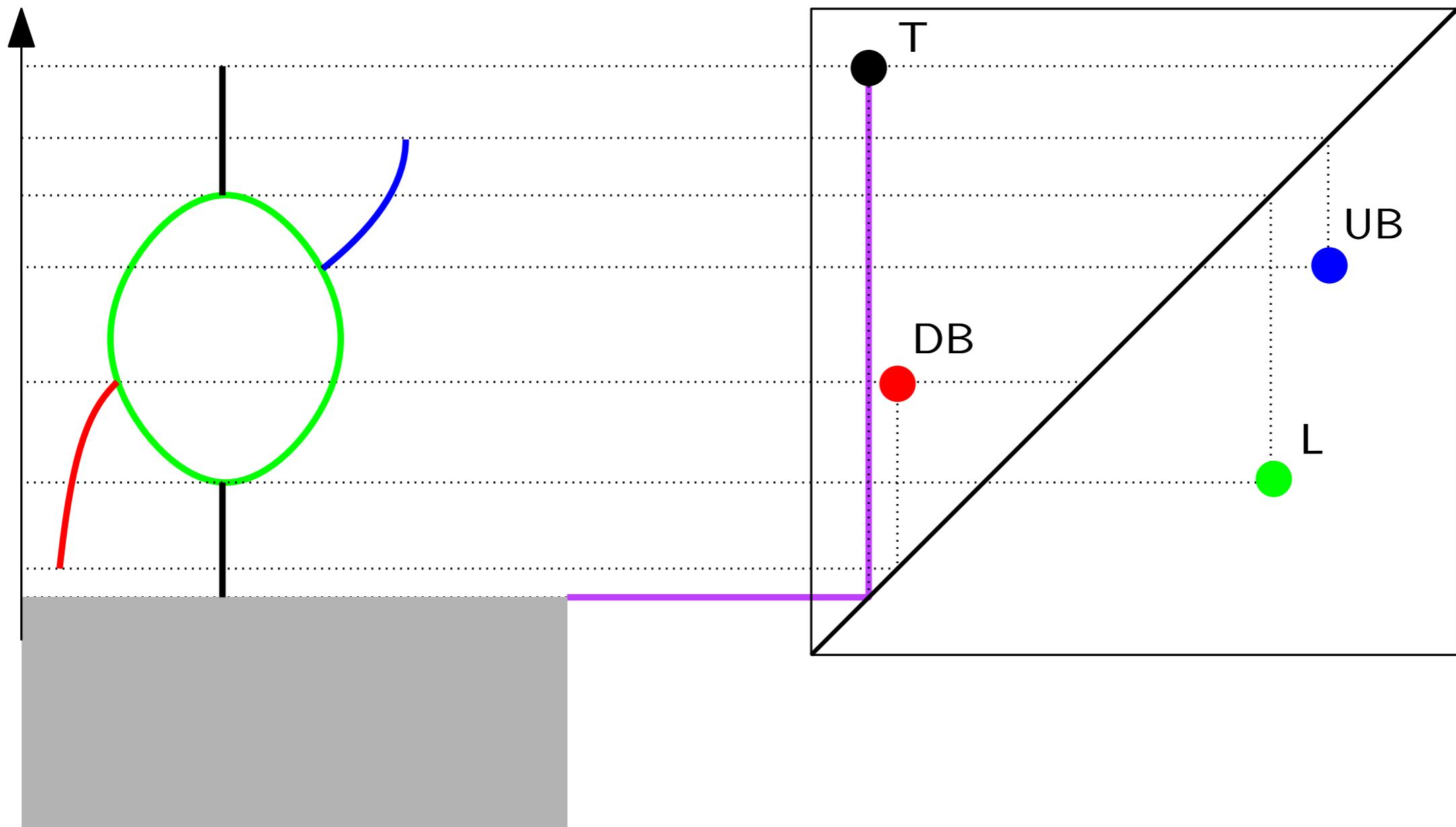
- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Extended Persistence Diagram

Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

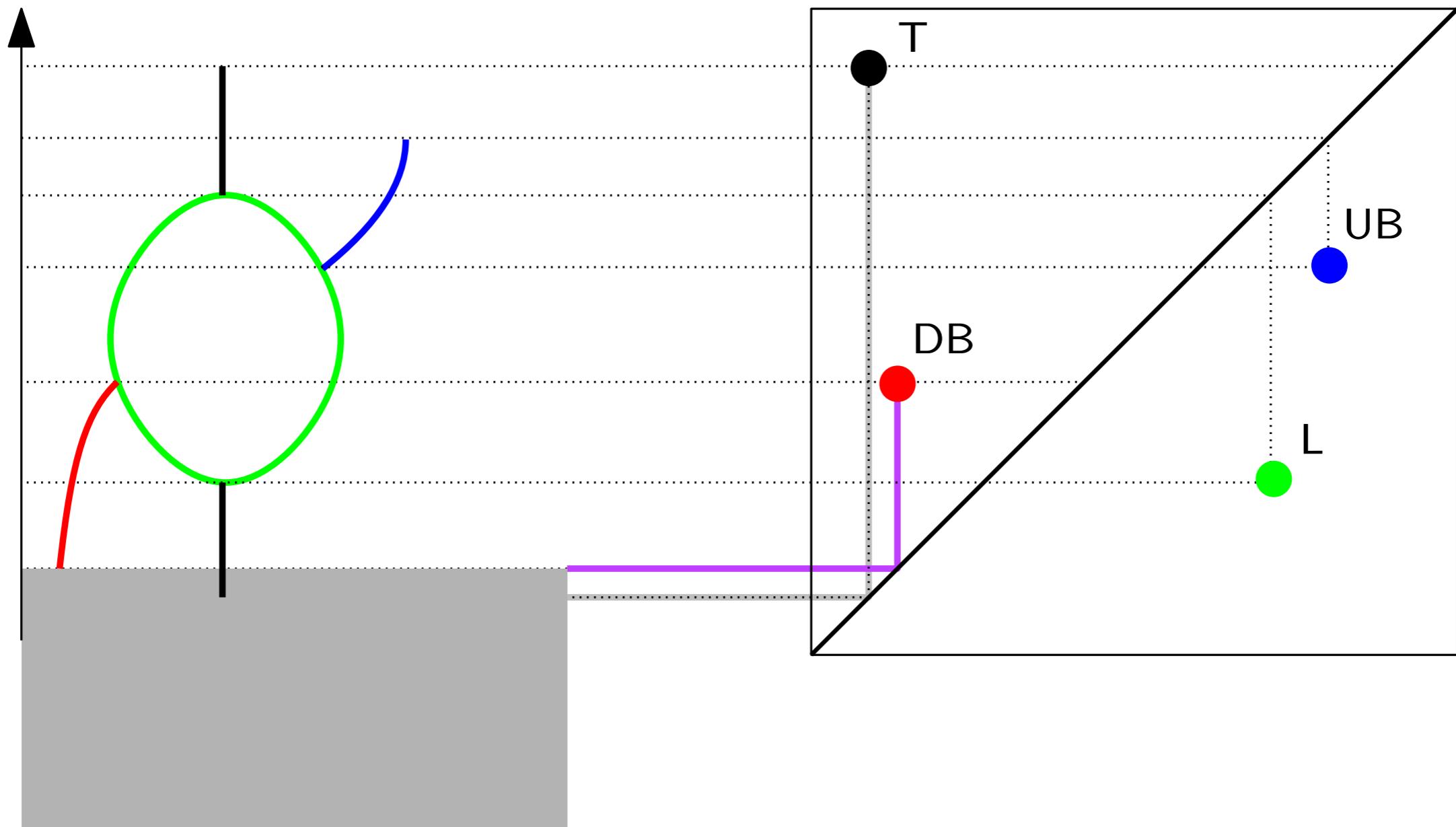
- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Extended Persistence Diagram

Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

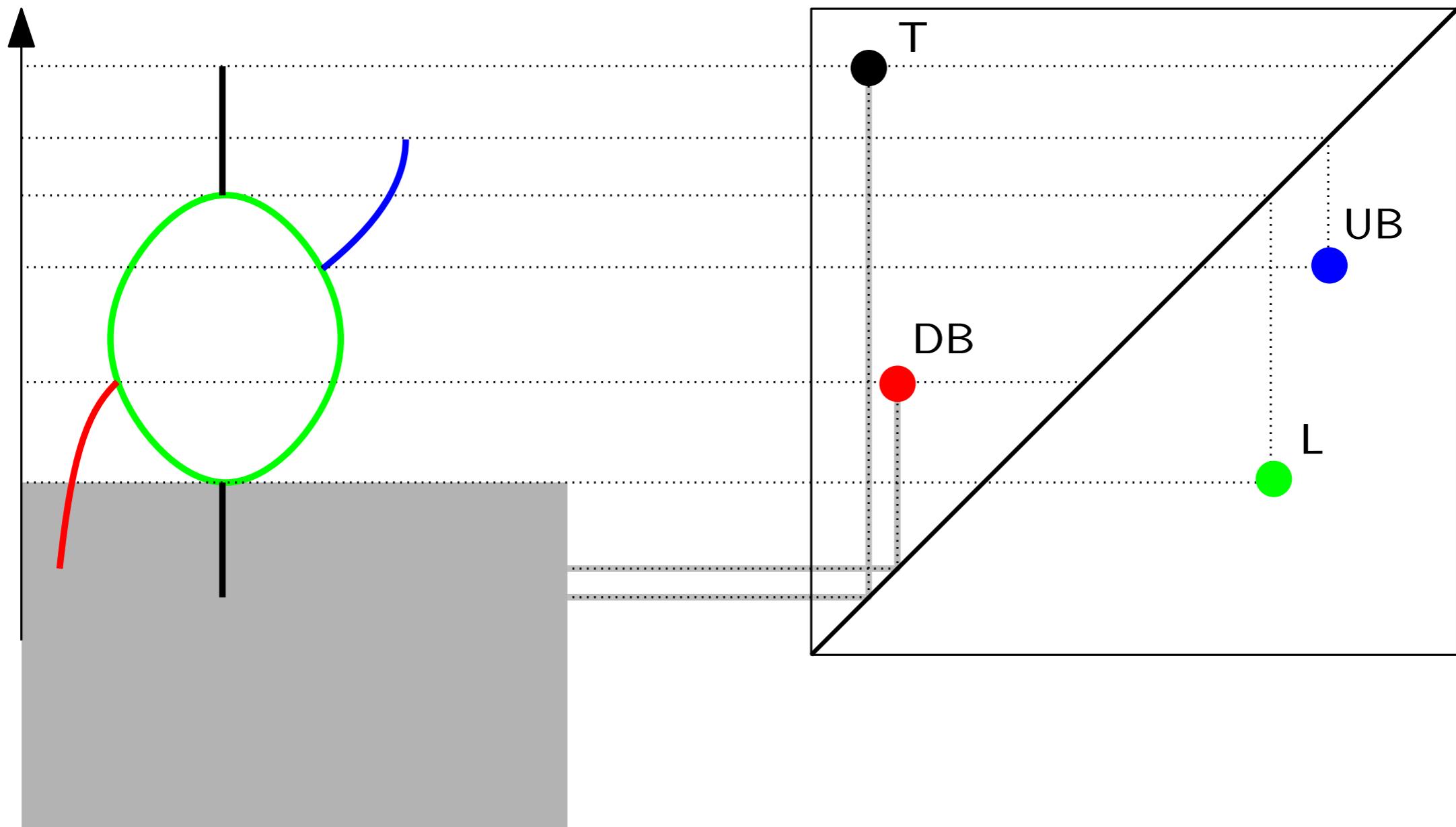
- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Extended Persistence Diagram

Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

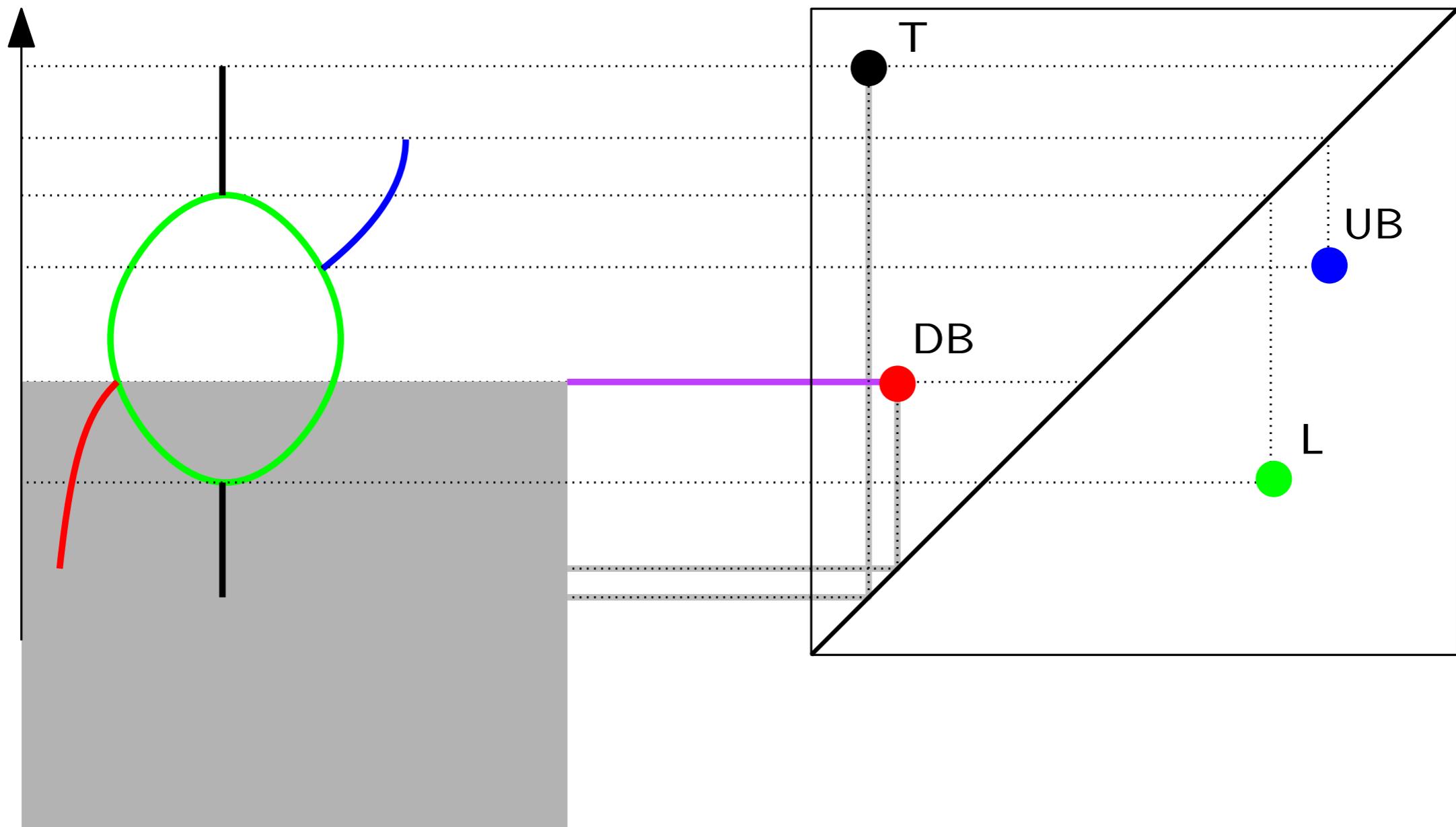
- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Extended Persistence Diagram

Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

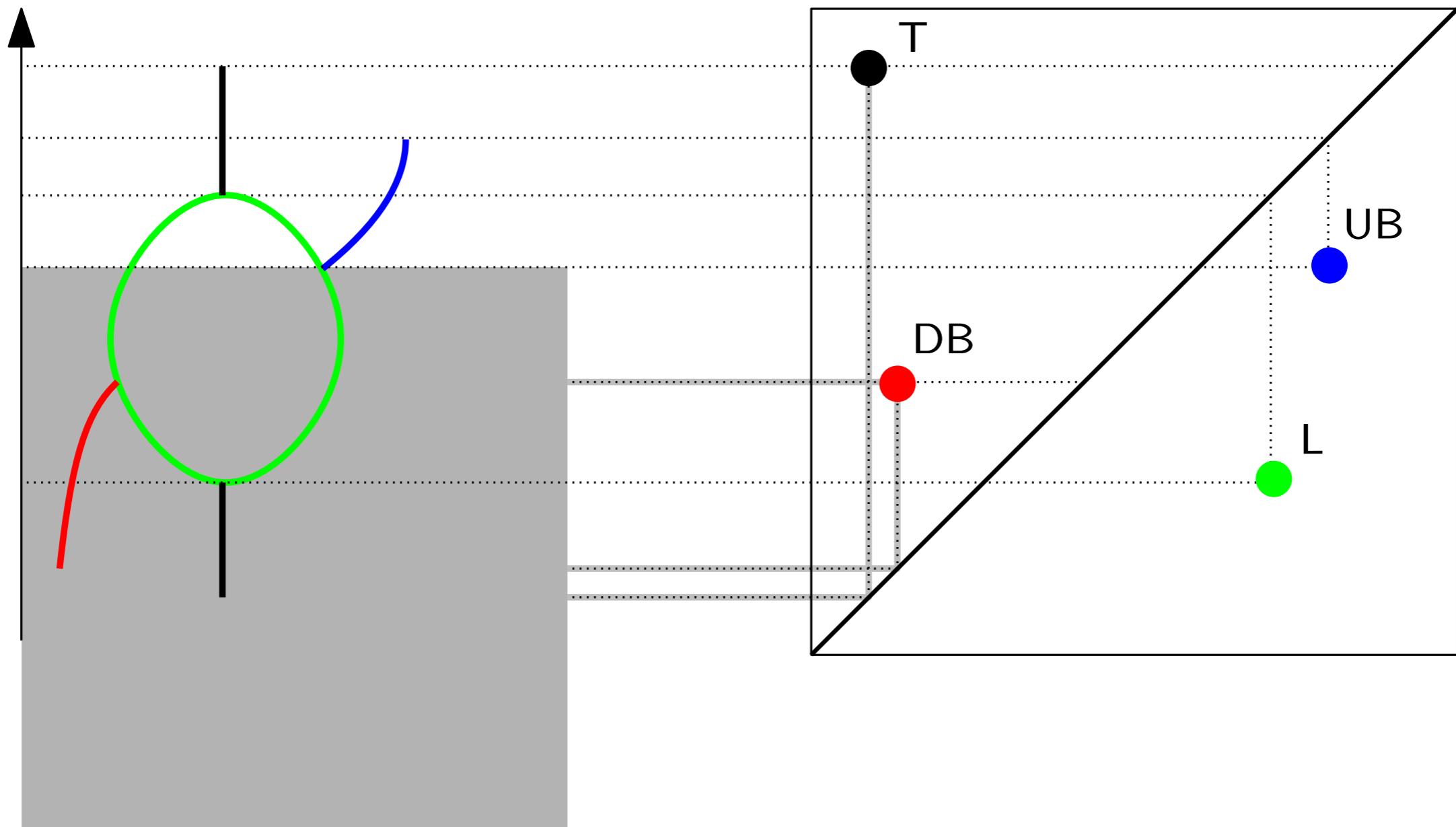
- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Extended Persistence Diagram

Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

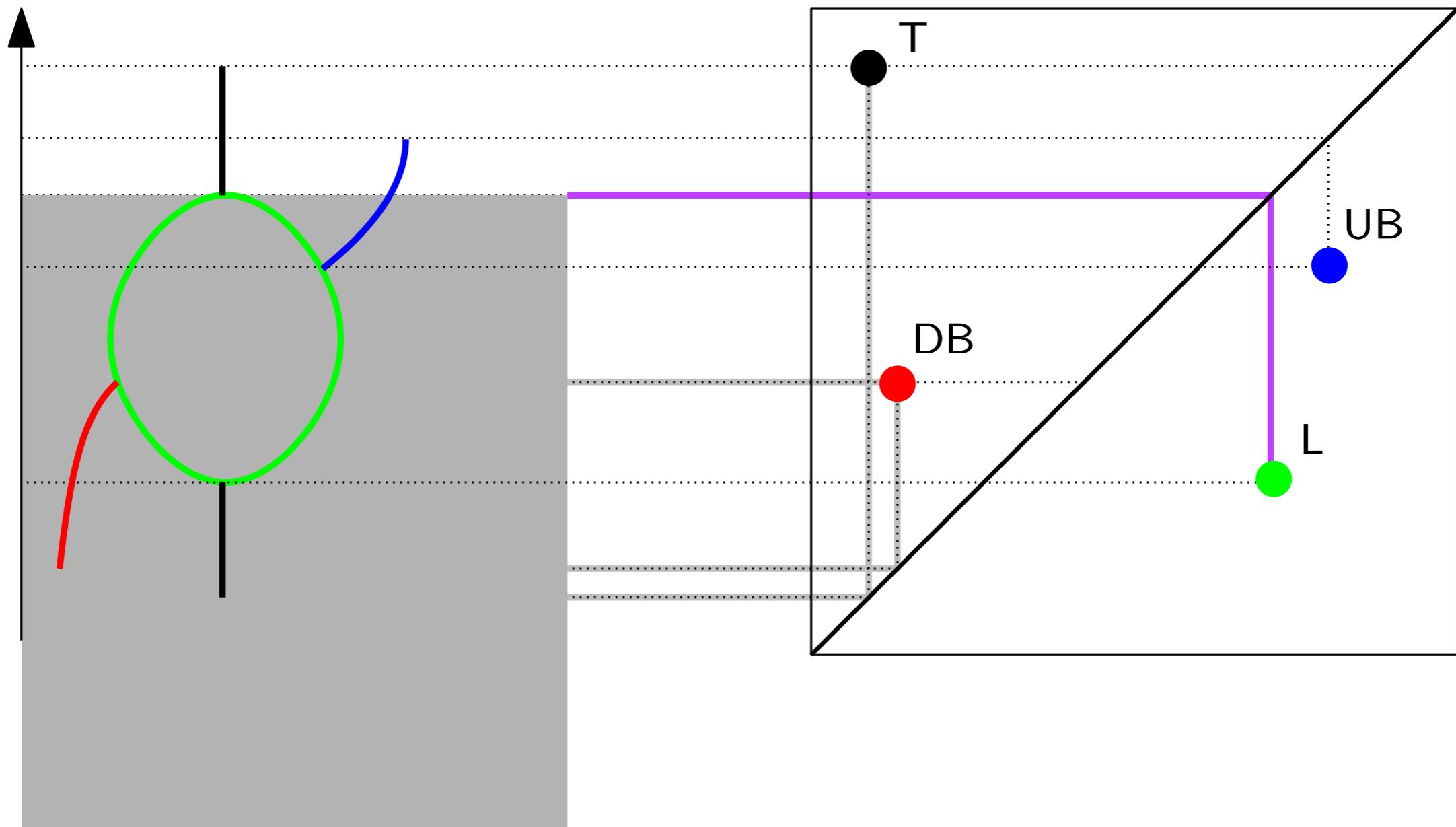
- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Extended Persistence Diagram

Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

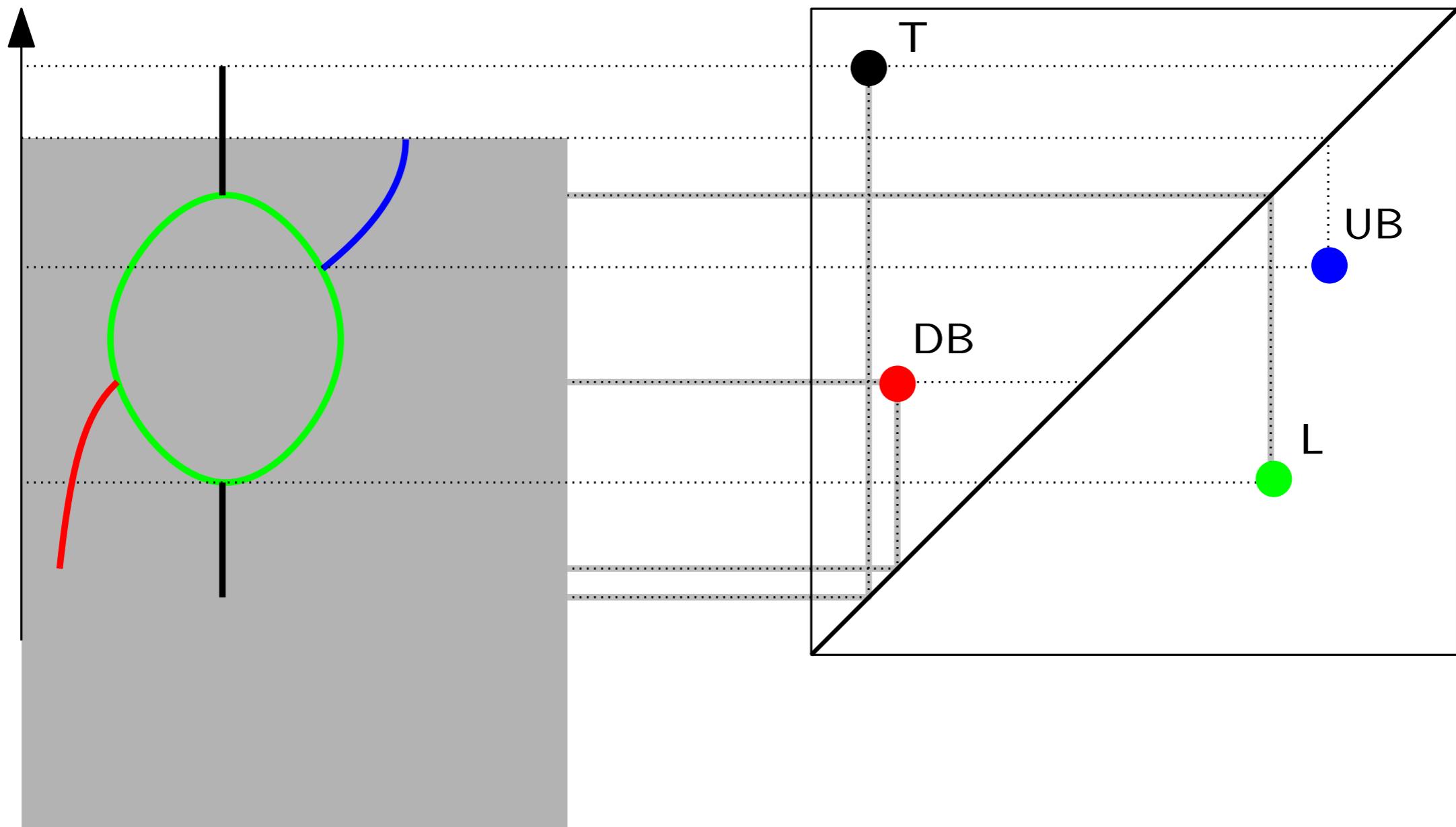
- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Extended Persistence Diagram

Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

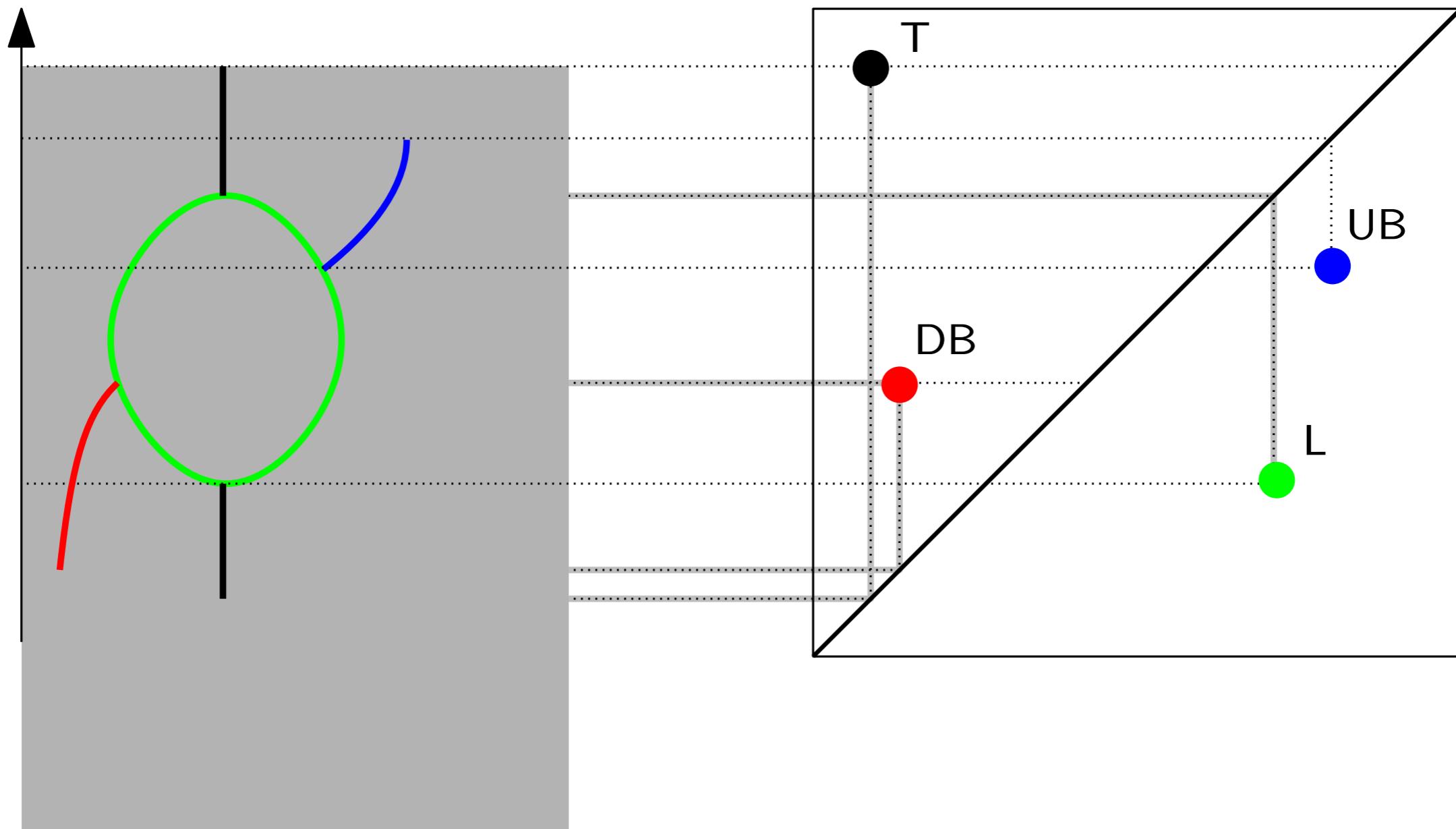
- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Extended Persistence Diagram

Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

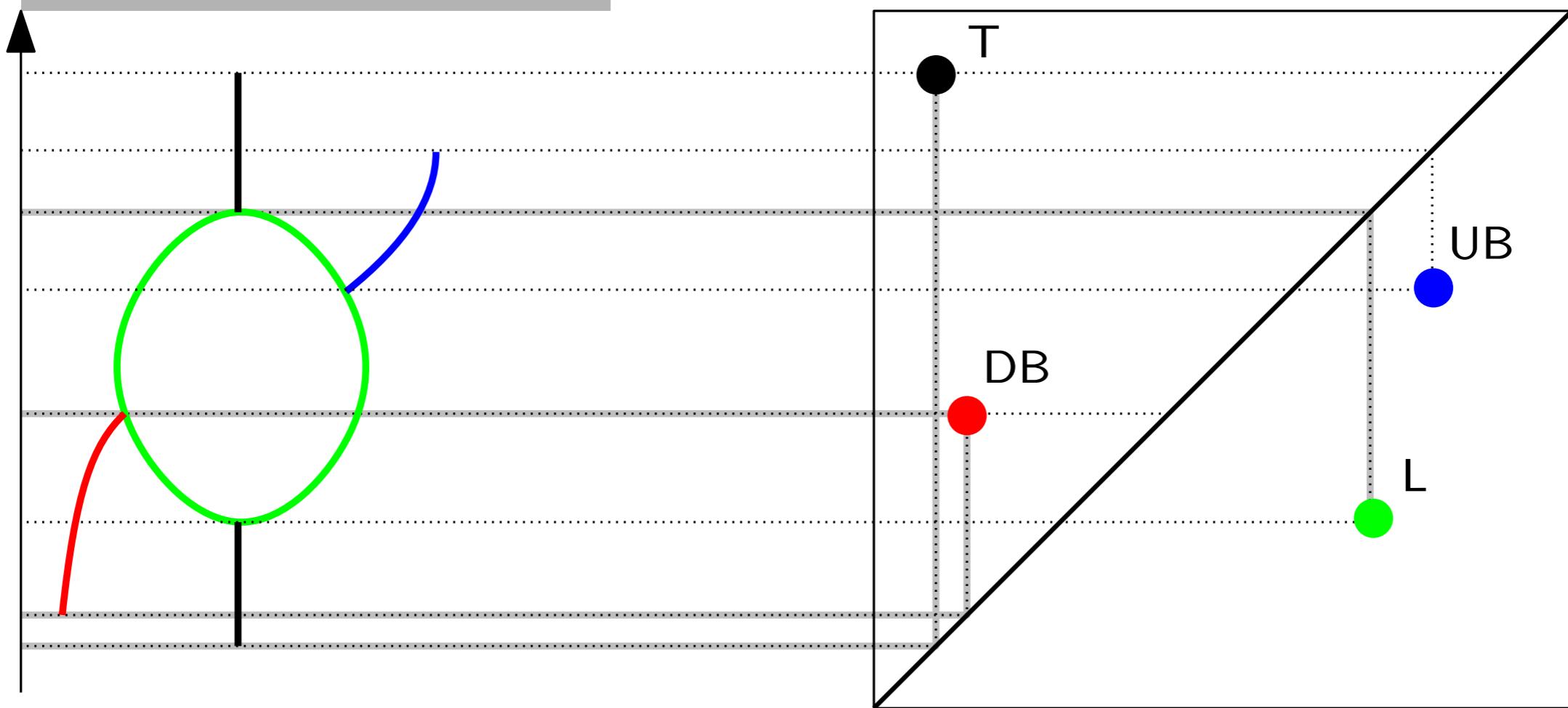
- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Extended Persistence Diagram

Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

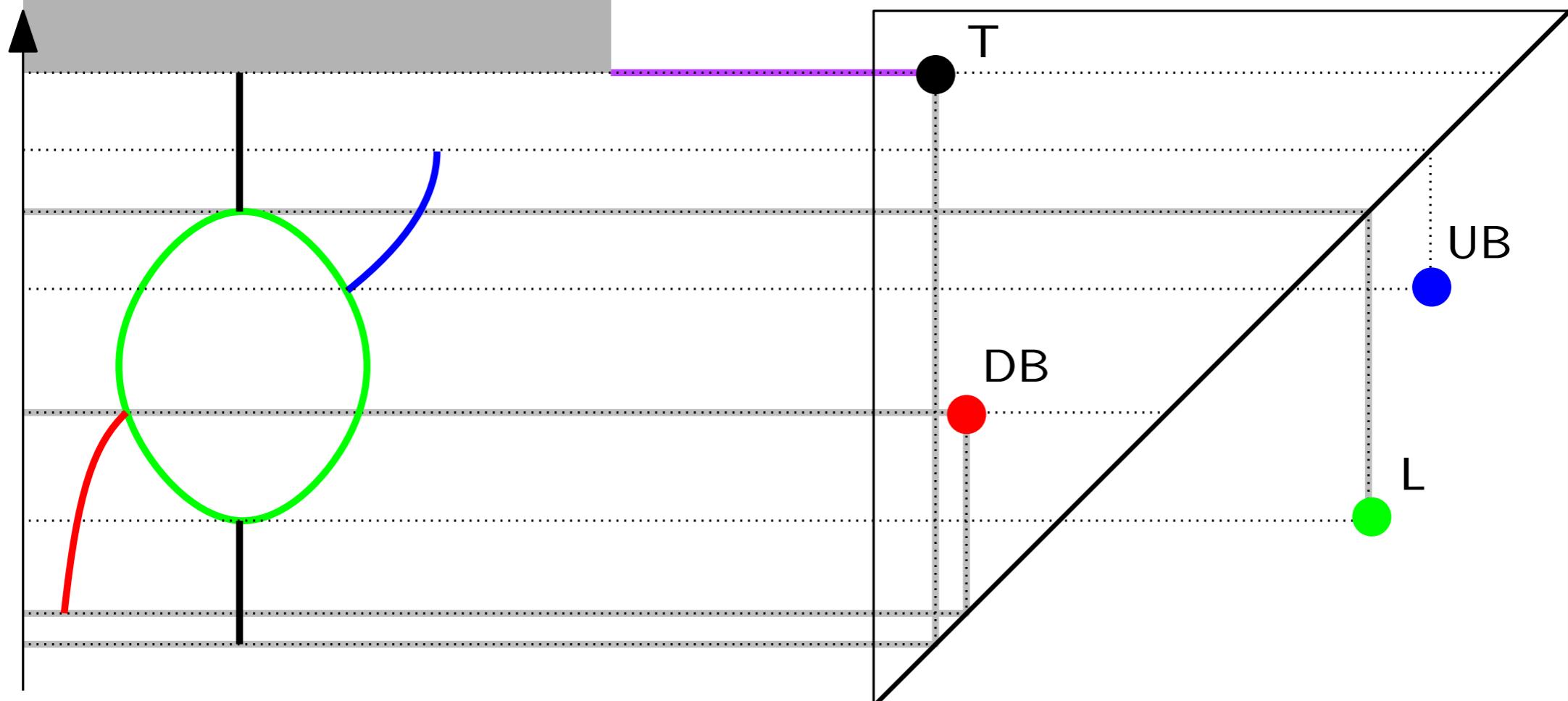
- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Extended Persistence Diagram

Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

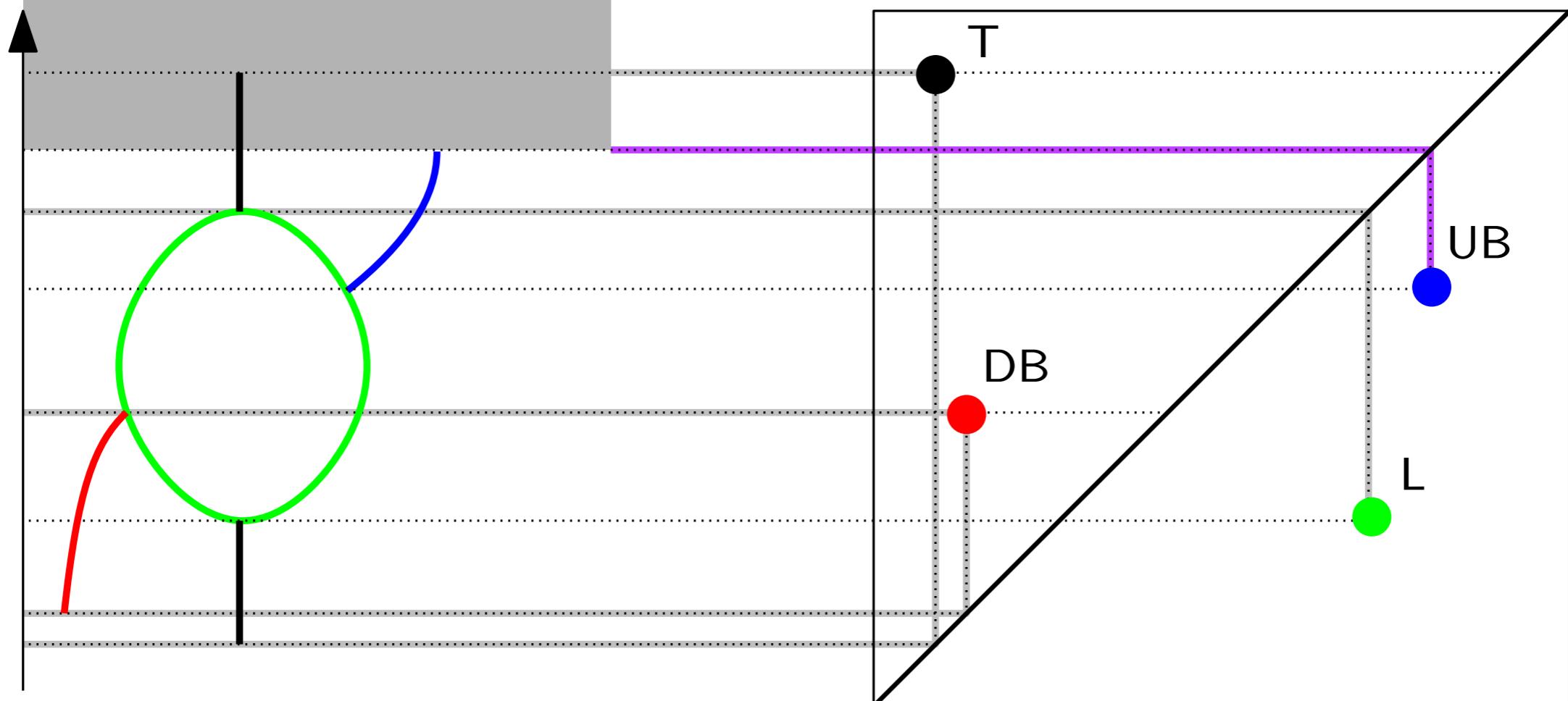
- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Extended Persistence Diagram

Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

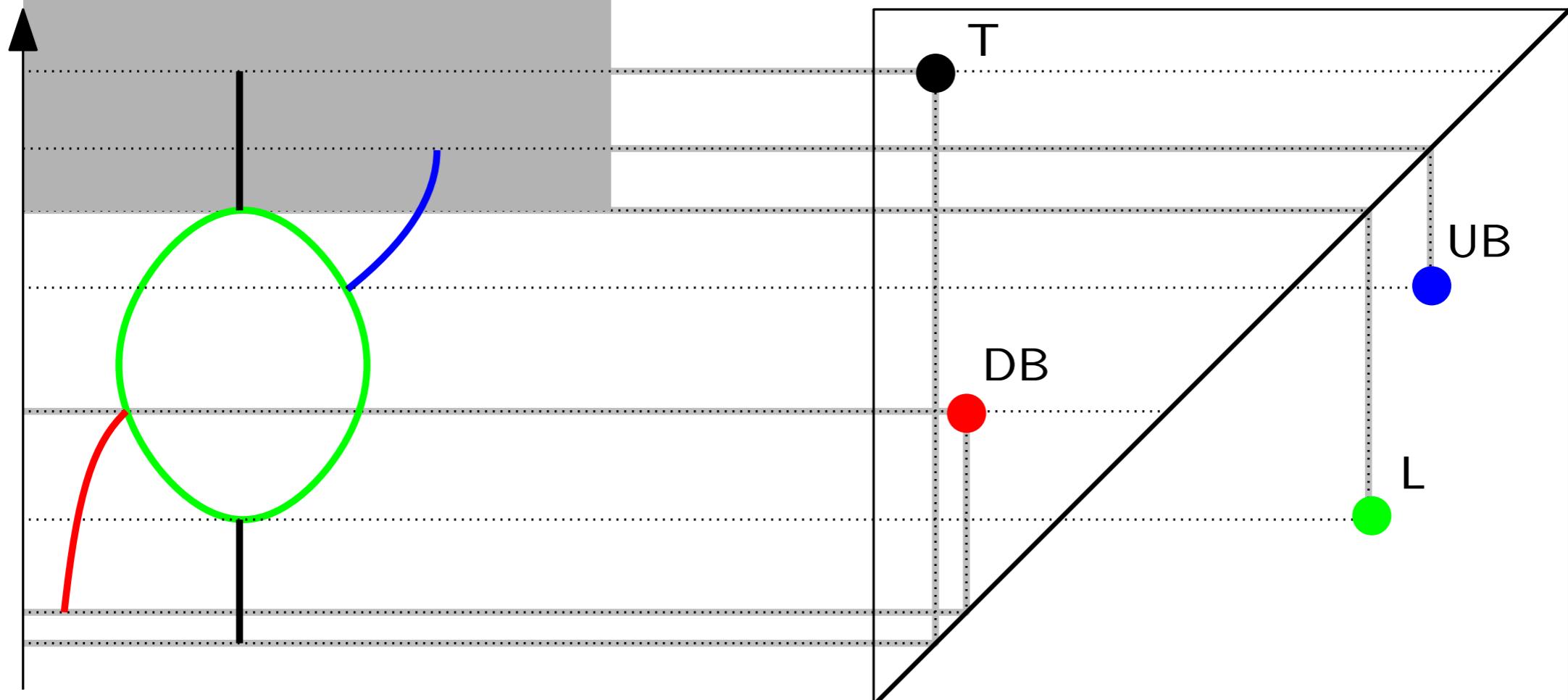
- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Extended Persistence Diagram

Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

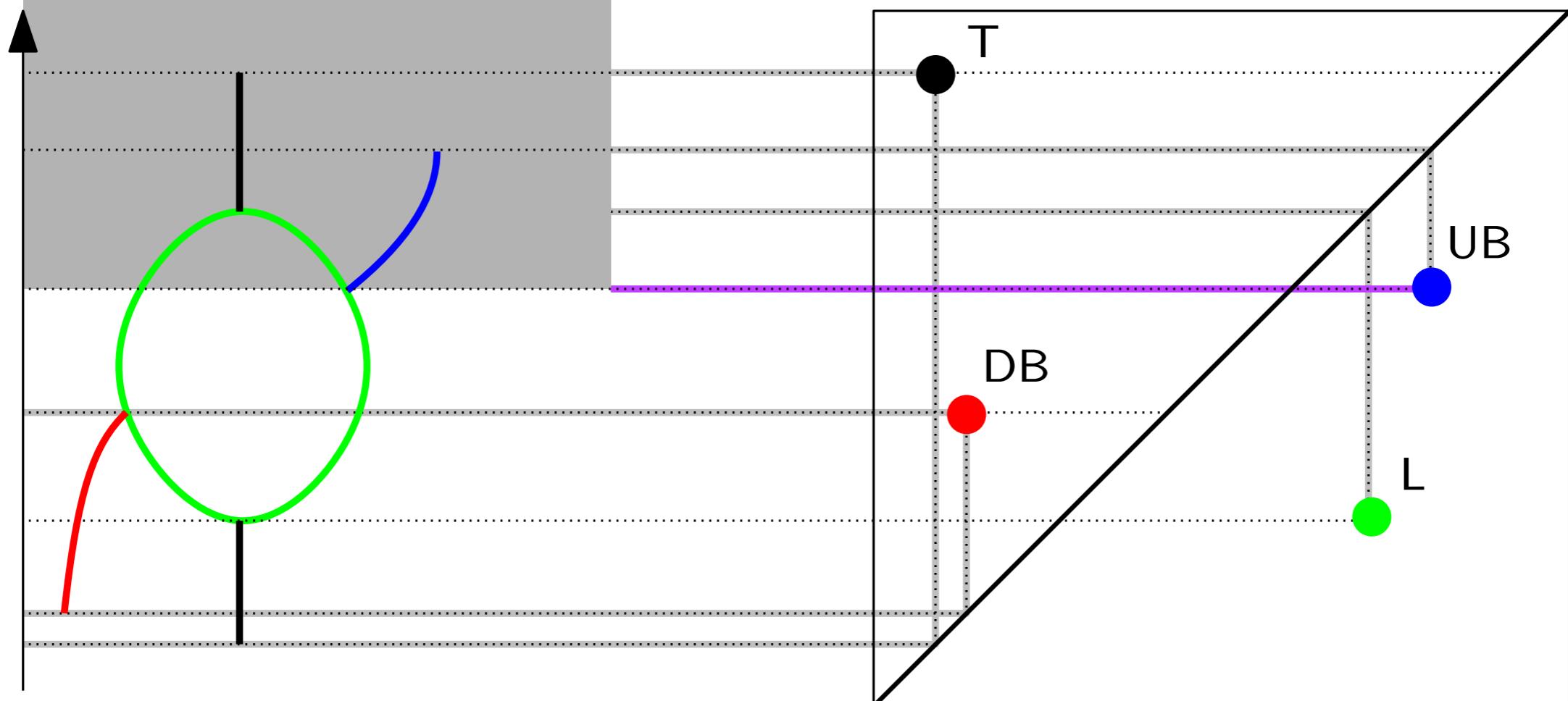
- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Extended Persistence Diagram

Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

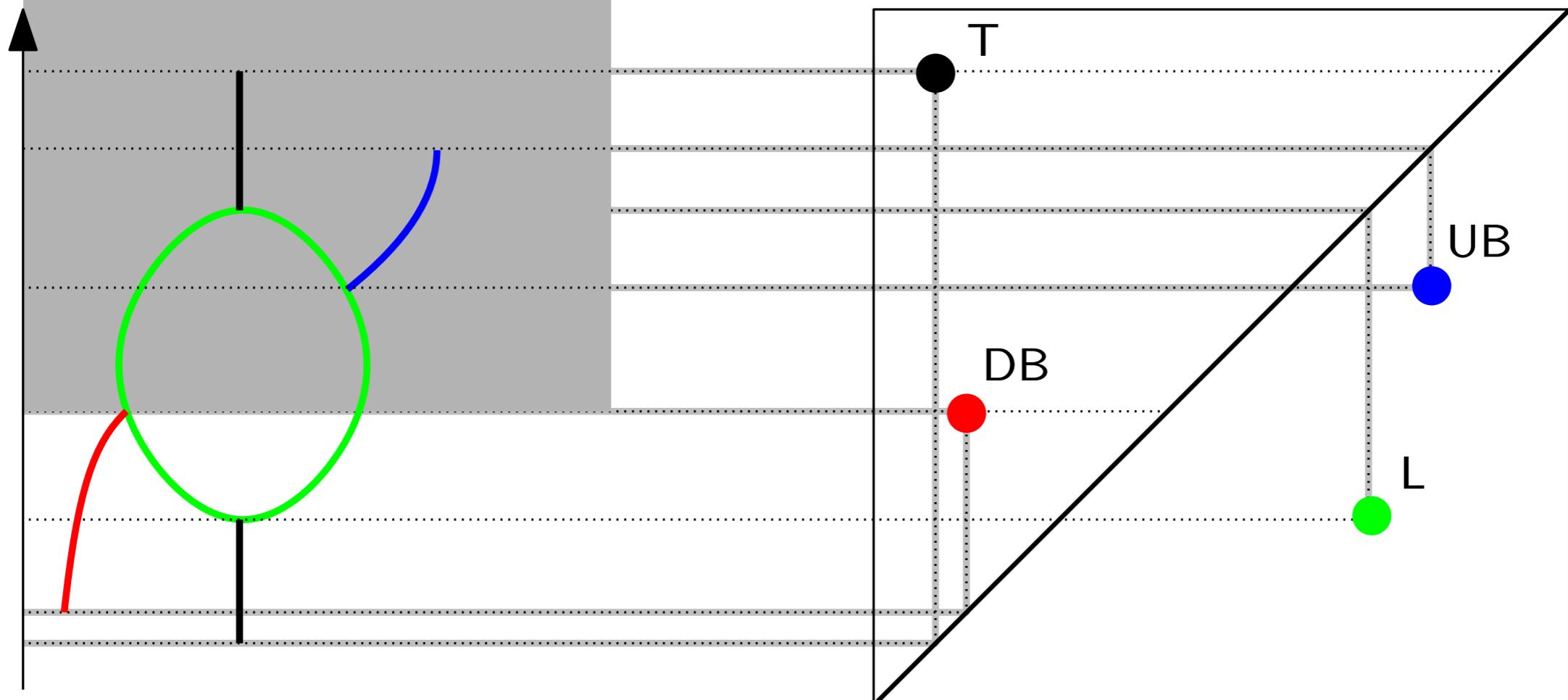
- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Extended Persistence Diagram

Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

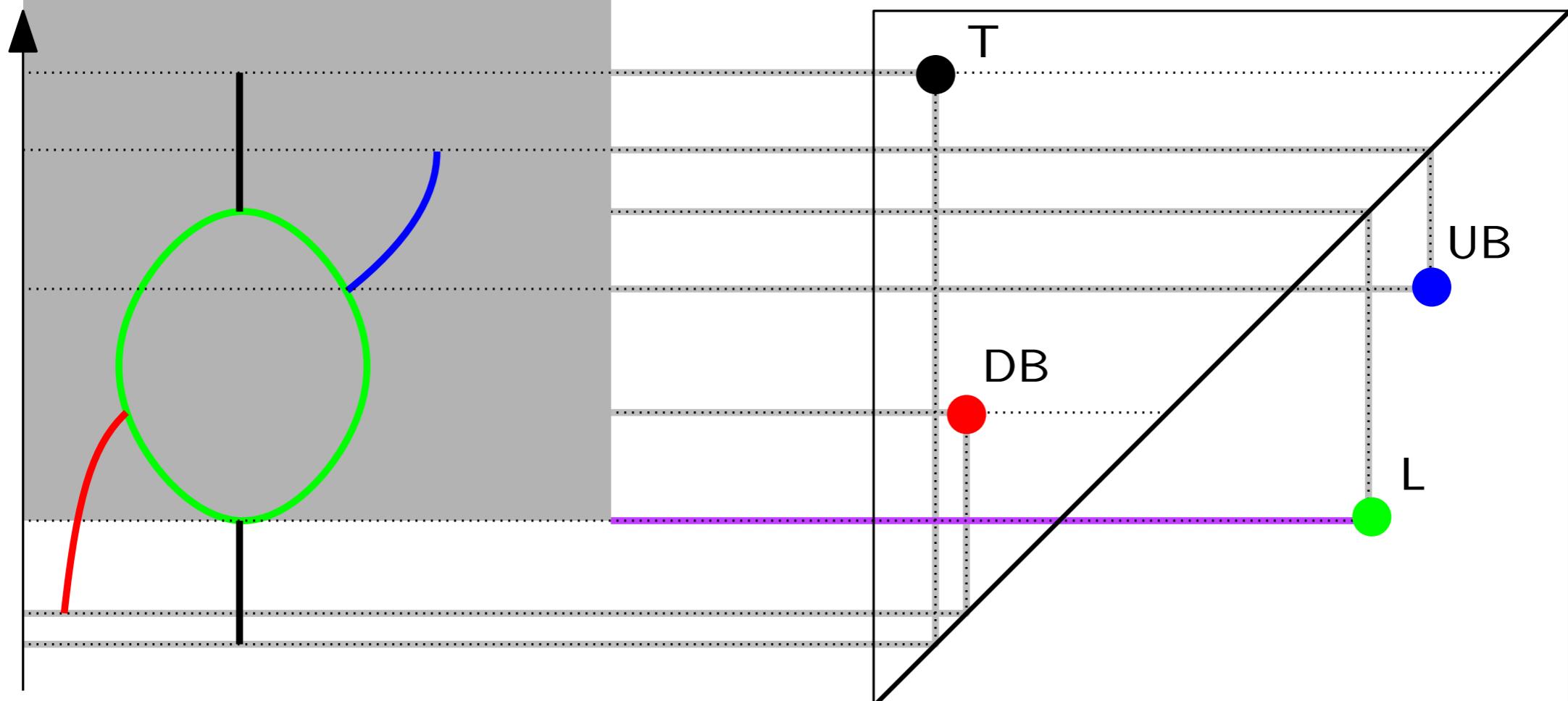
- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Extended Persistence Diagram

Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

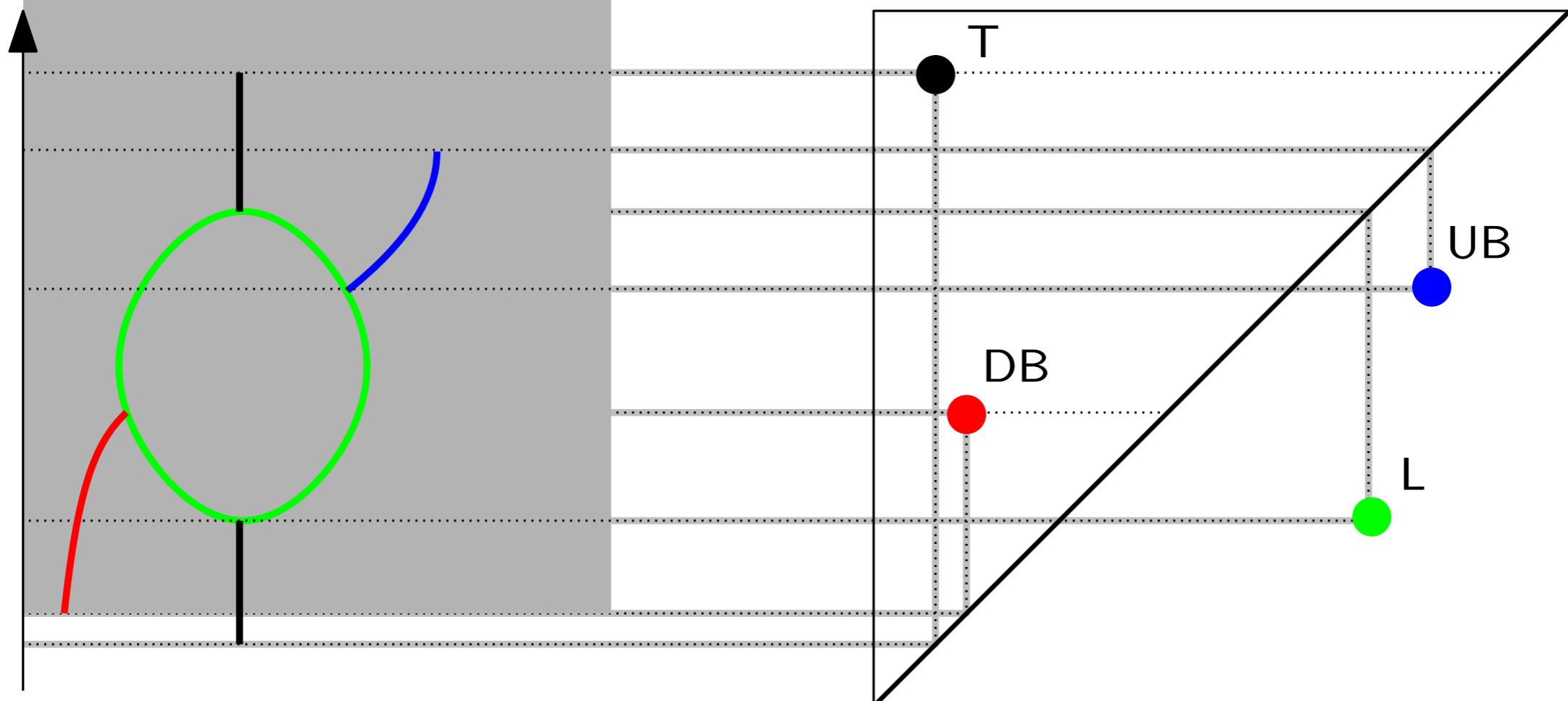
- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Extended Persistence Diagram

Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

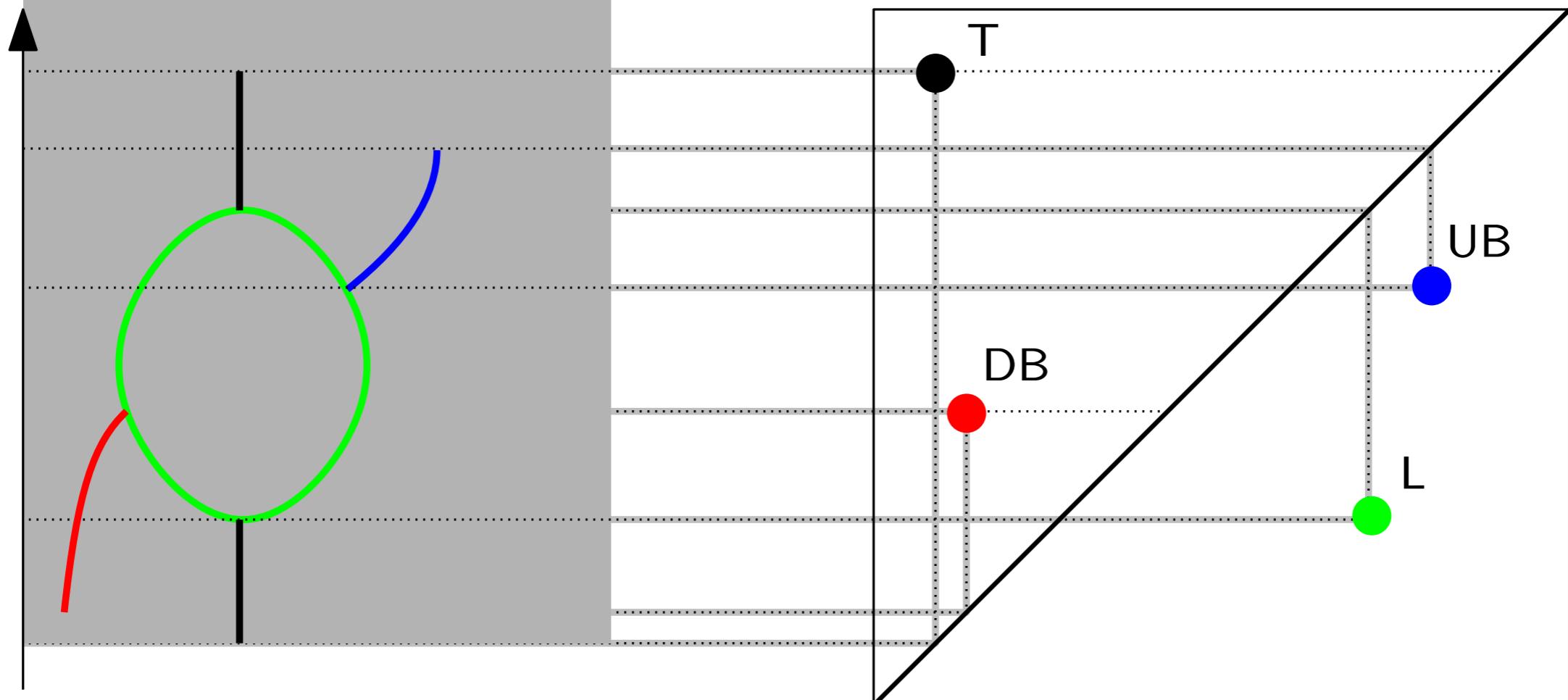
- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Extended Persistence Diagram

Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

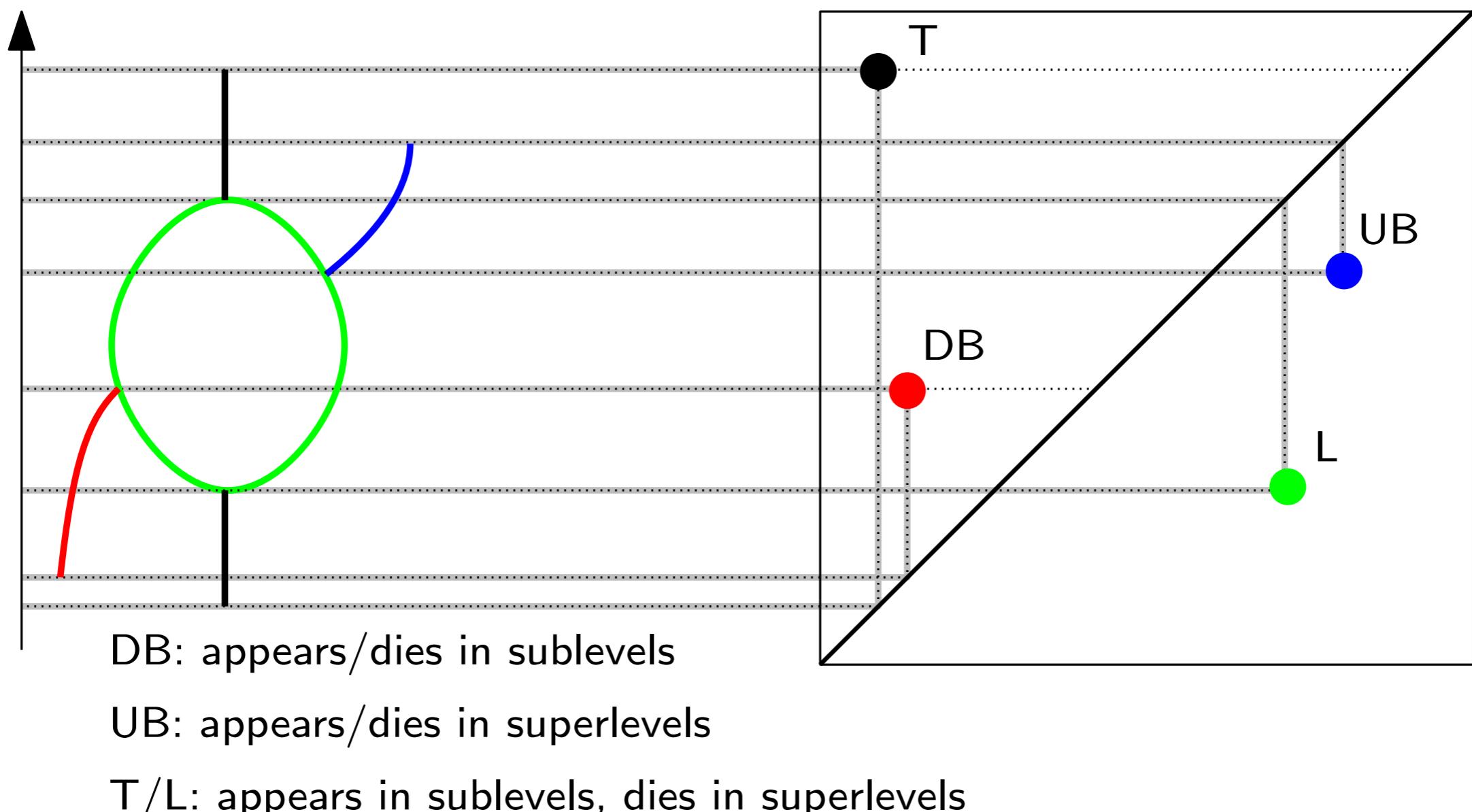
- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Extended Persistence Diagram

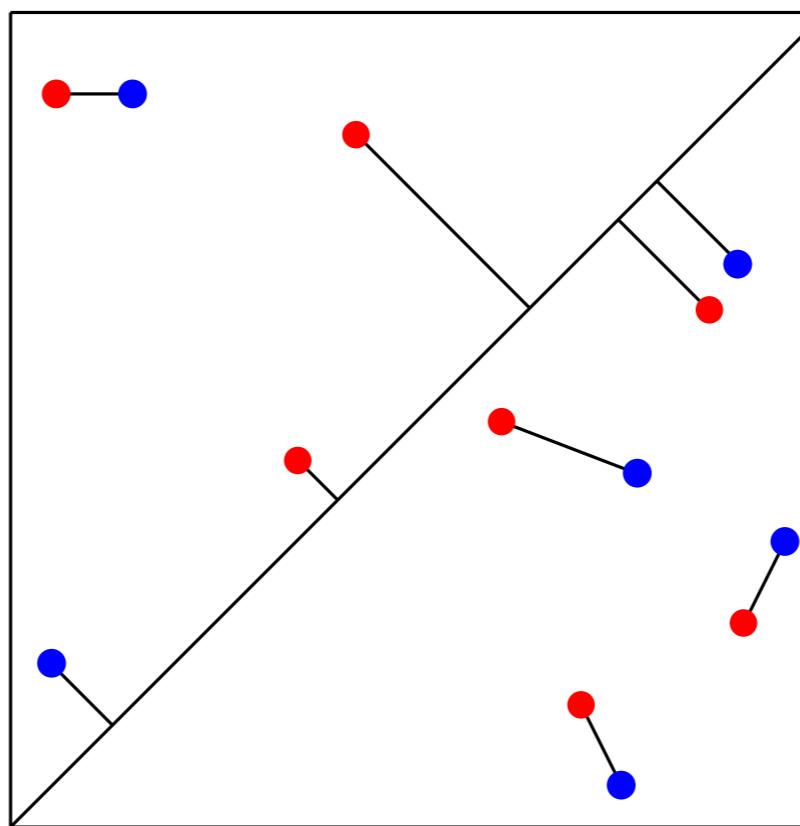
Construction uses **extended persistence**: [Cohen-Steiner, Edelsbrunner, Harer 2008]

- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



# Metric for Extended Persistence Diagrams

---



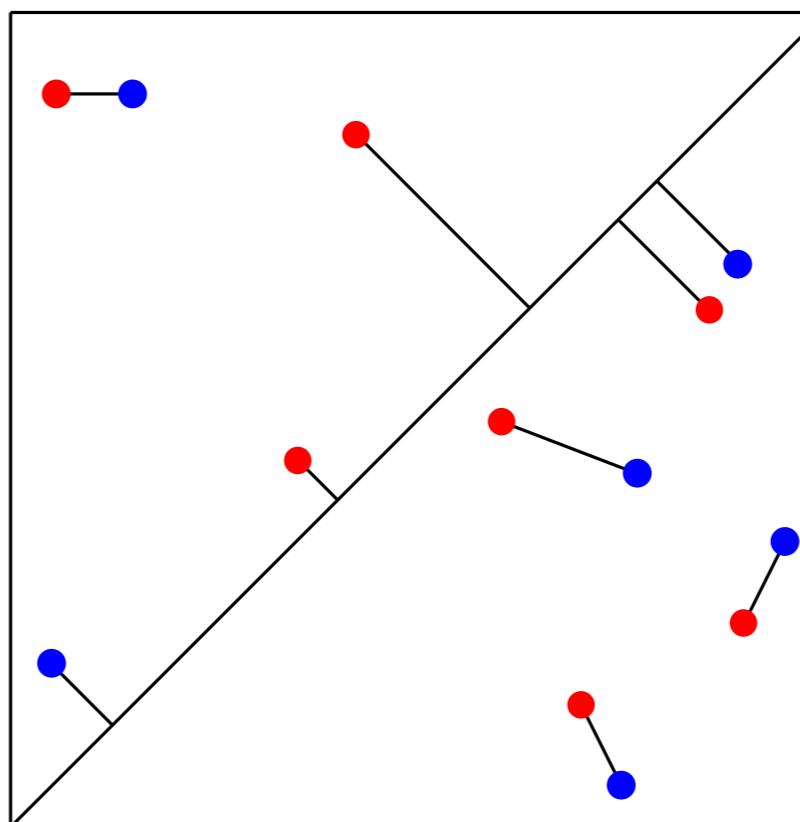
# Metric for Extended Persistence Diagrams

Partial matching  $M : \text{Dg} \leftrightarrow \text{Dg}'$

Matched pair  $(x, y) \in M$ :  $c(x, y) = \|x - y\|_\infty$

Unmatched point  $z \in X \sqcup Y$ :  $c(z) = \|z - \bar{z}\|_\infty$

$$c(M) = \max\{ \max_{(x, y)} c(x, y), \max_z c(z) \}$$



# Metric for Extended Persistence Diagrams

Partial matching  $M : \text{Dg} \leftrightarrow \text{Dg}'$

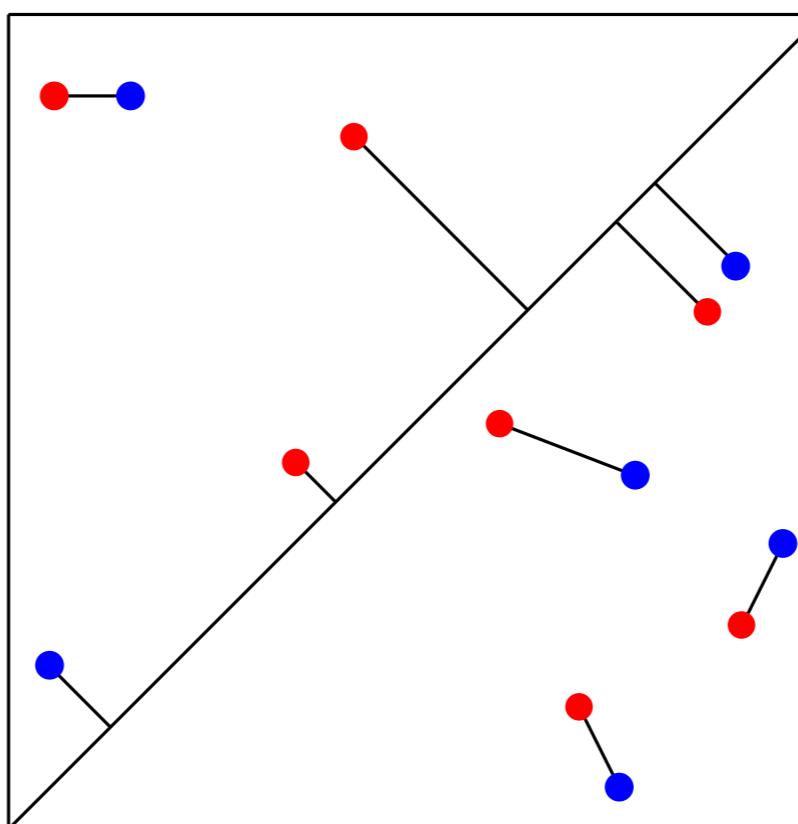
Matched pair  $(x, y) \in M$ :  $c(x, y) = \|x - y\|_\infty$

Unmatched point  $z \in X \sqcup Y$ :  $c(z) = \|z - \bar{z}\|_\infty$

$$c(M) = \max\{ \max_{(x, y)} c(x, y), \max_z c(z) \}$$

**Def:** bottleneck distance:

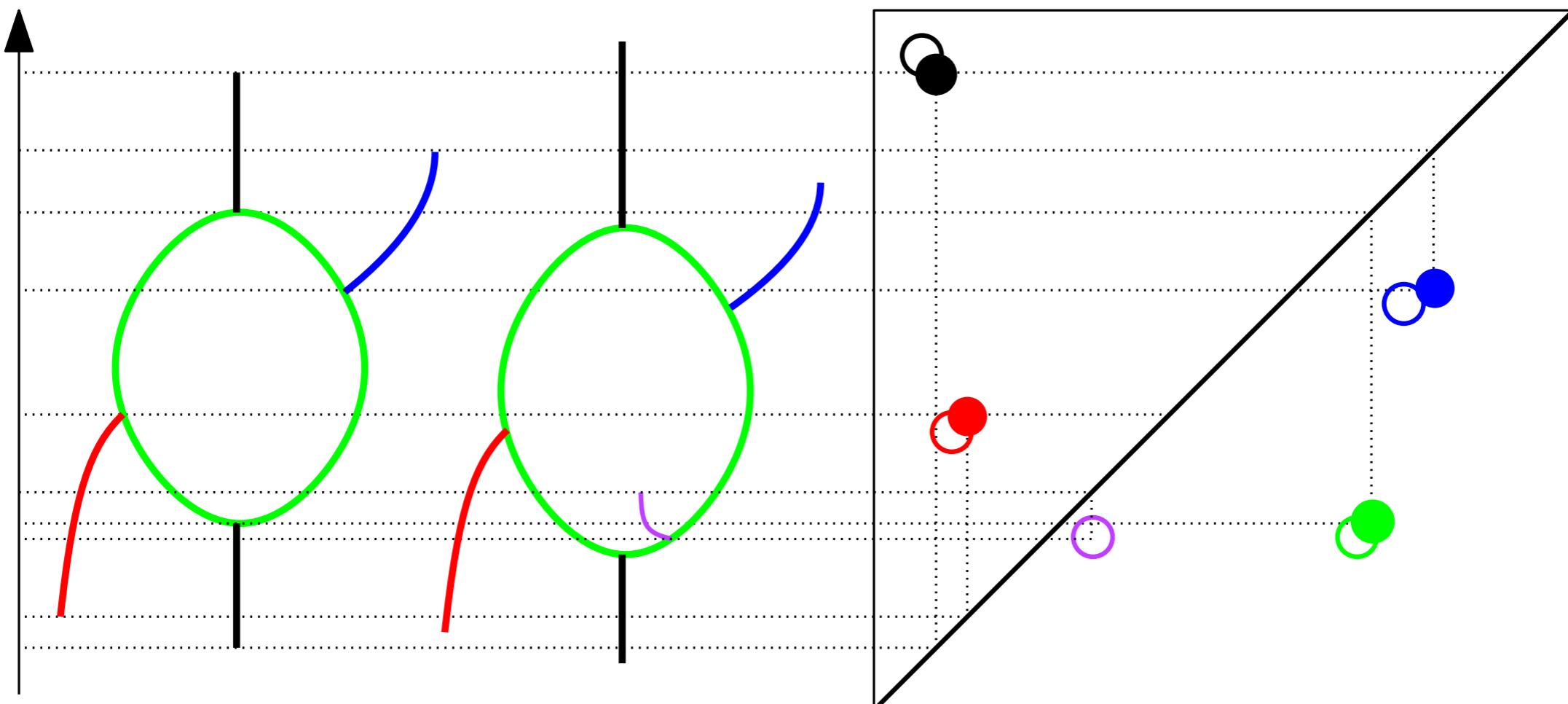
$$d_B(\text{Dg}, \text{Dg}') = \inf_{M: \text{Dg} \leftrightarrow \text{Dg}'} c(M)$$



# Metric Properties

**Thm (stability):** [Bauer, Ge, Wang 2013]

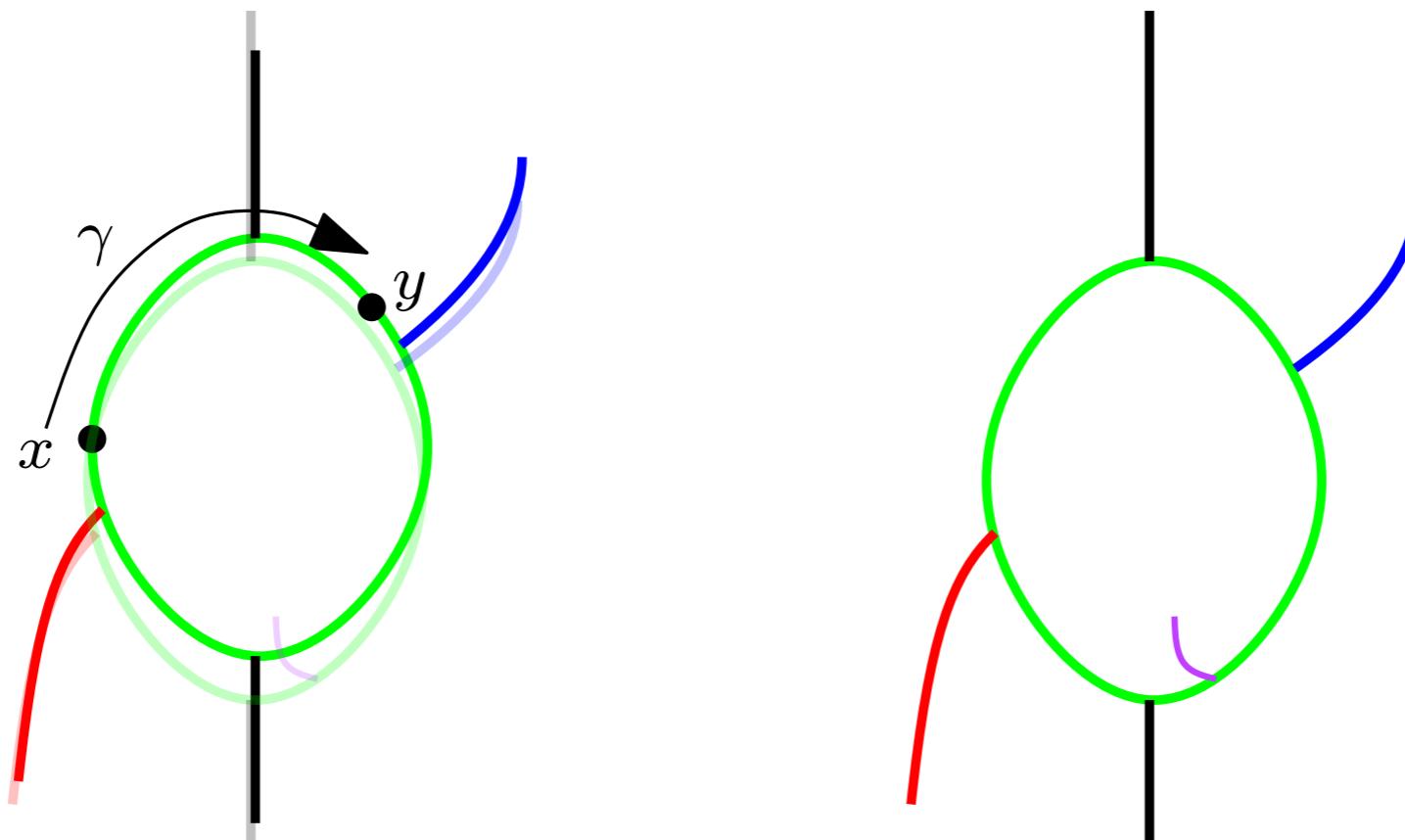
$$d_B(Dg(R_f), Dg(R_g)) \leq 6 d_{GH}(R_f, R_g)$$



# Metric Properties

**Thm (stability):** [Bauer, Ge, Wang 2013]

$$d_B(Dg(R_f), Dg(R_g)) \leq 6 d_{GH}(R_f, R_g)$$

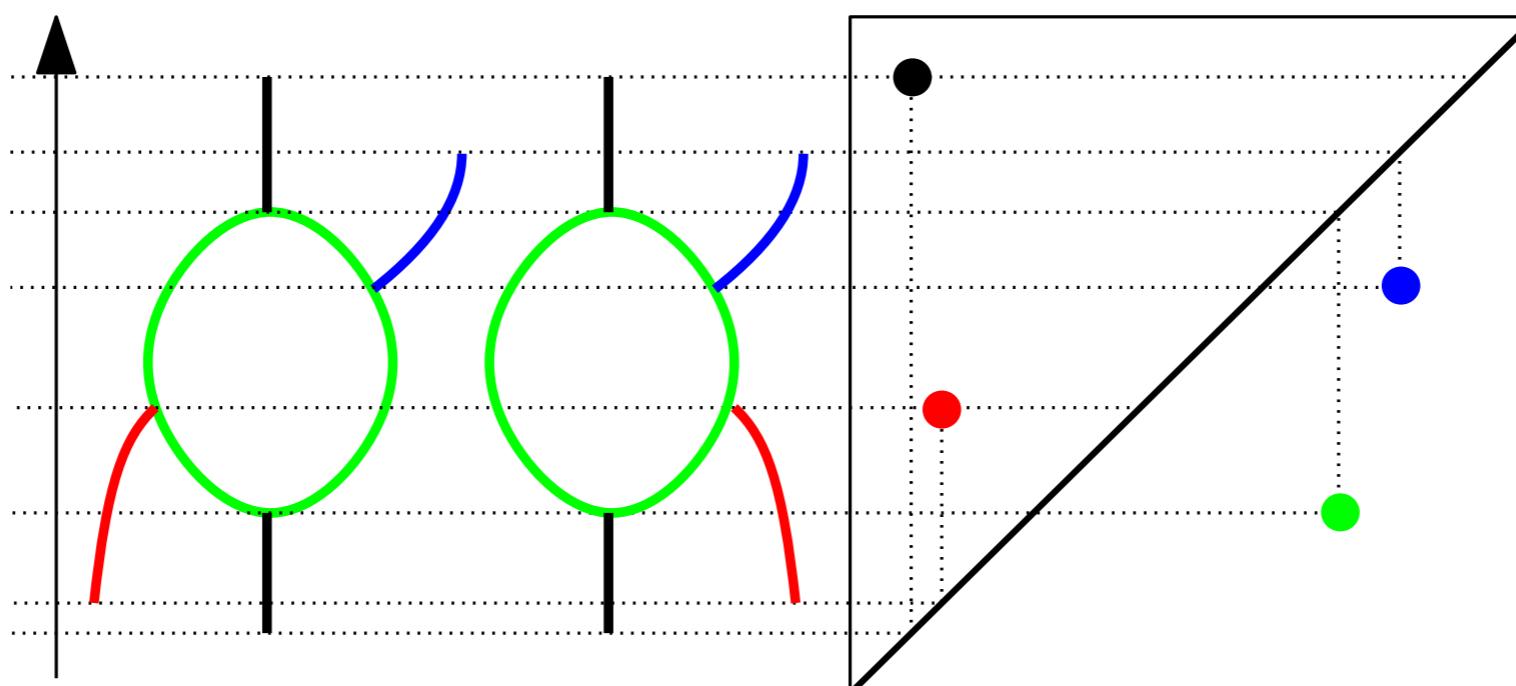


# Metric Properties

**Thm (stability):** [Bauer, Ge, Wang 2013]

$$d_B(Dg(R_f), Dg(R_g)) \leq 6 d_{GH}(R_f, R_g)$$

**Note:**  $d_B(Dg(\cdot), Dg(\cdot))$  is only a pseudometric on Reeb graphs



# Metric Properties

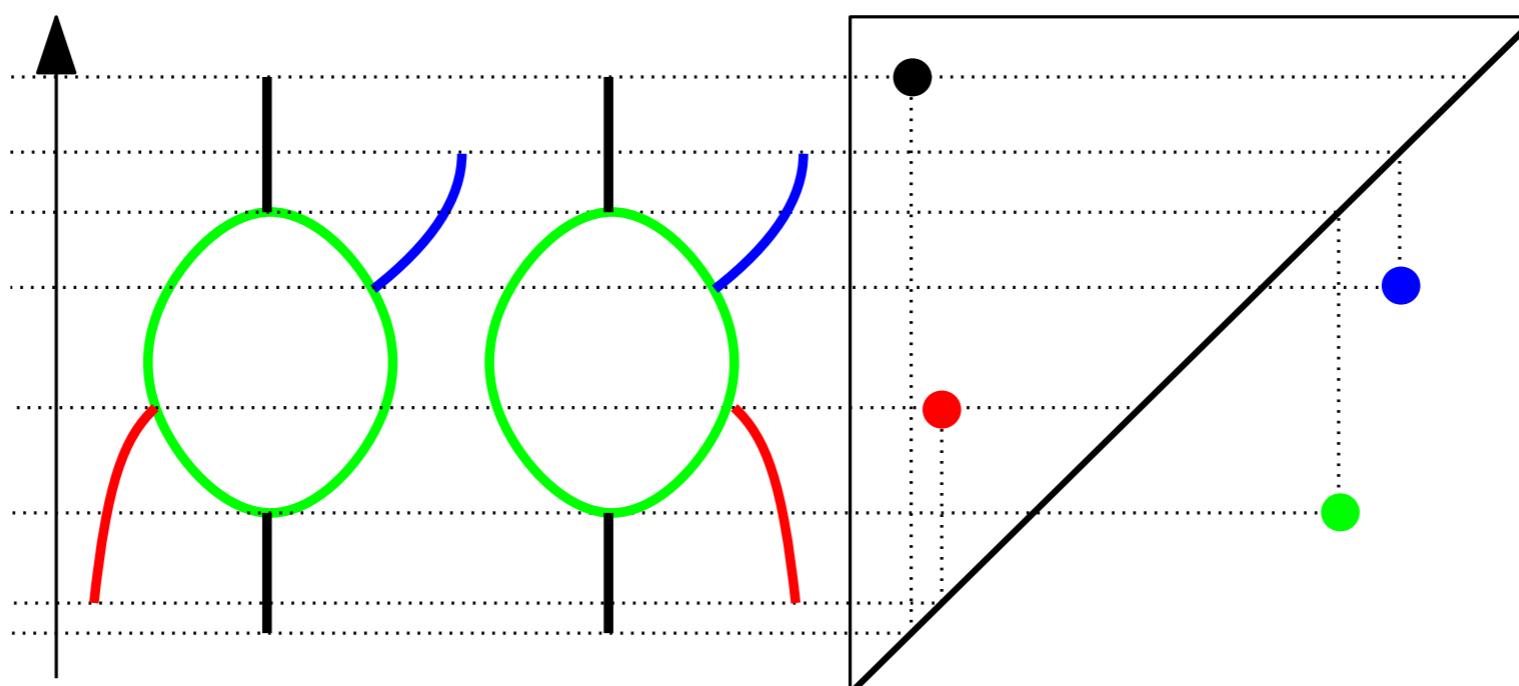
**Thm (stability):** [Bauer, Ge, Wang 2013]

$$d_B(Dg(R_f), Dg(R_g)) \leq 6 d_{GH}(R_f, R_g)$$

**Note:**  $d_B(Dg(\cdot), Dg(\cdot))$  is only a pseudometric on Reeb graphs

**Thm:** [C., Oudot 2016]

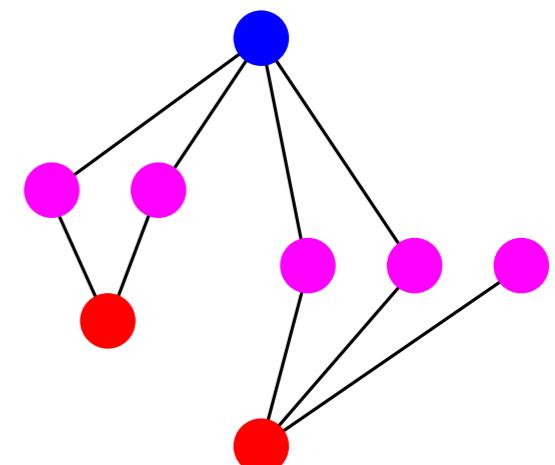
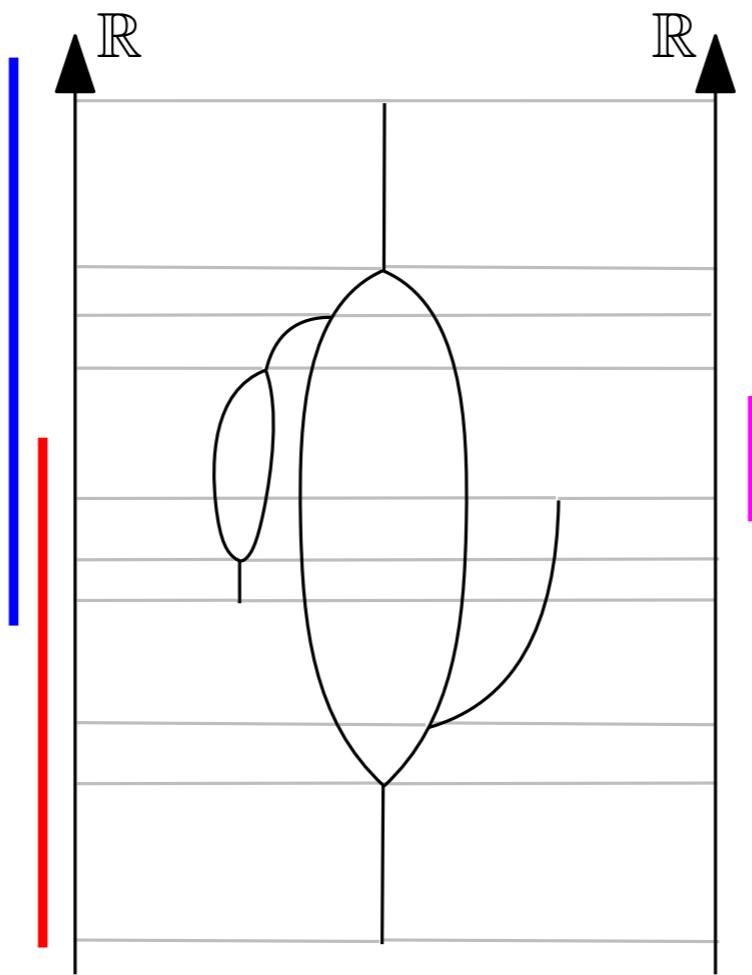
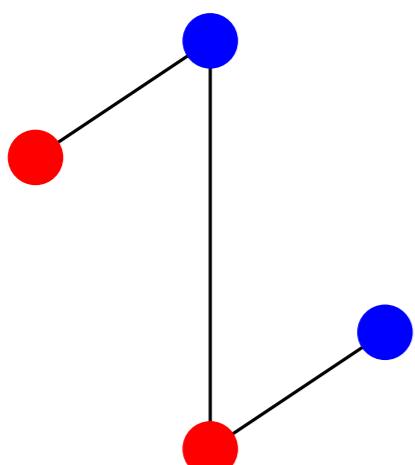
$d_B(Dg(\cdot), Dg(\cdot))$  is *locally* a metric equivalent to  $d_{GH}$



# Structure of Mapper

---

**Reminder:** Mapper  $\equiv$  pixelized Reeb graph



# Structure of Mapper

---

**Def:** Given  $X, f, \mathcal{I}$ :

$$Dg(M_f) = (DB(R_f) \setminus Q_{\mathcal{I}}^{DB}) \cup (UB(R_f) \setminus Q_{\mathcal{I}}^{UB}) \cup (L(R_f) \setminus Q_{\mathcal{I}}^L)$$

# Structure of Mapper

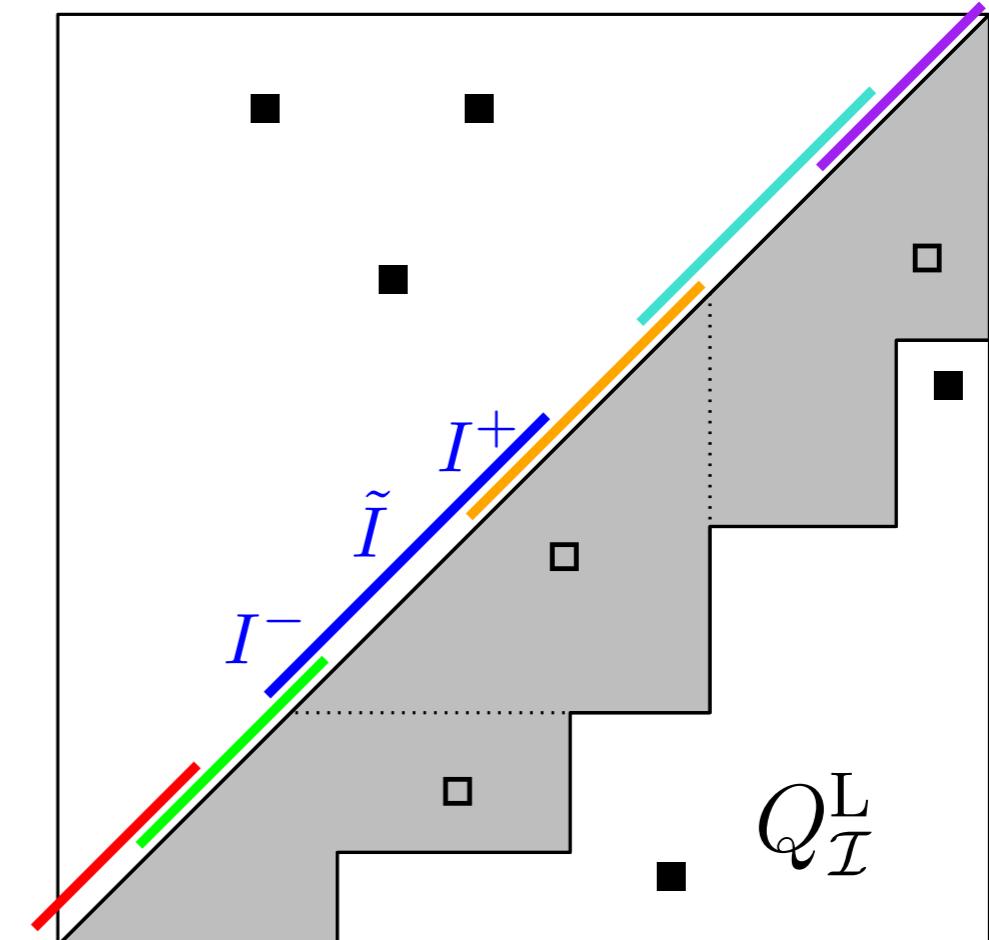
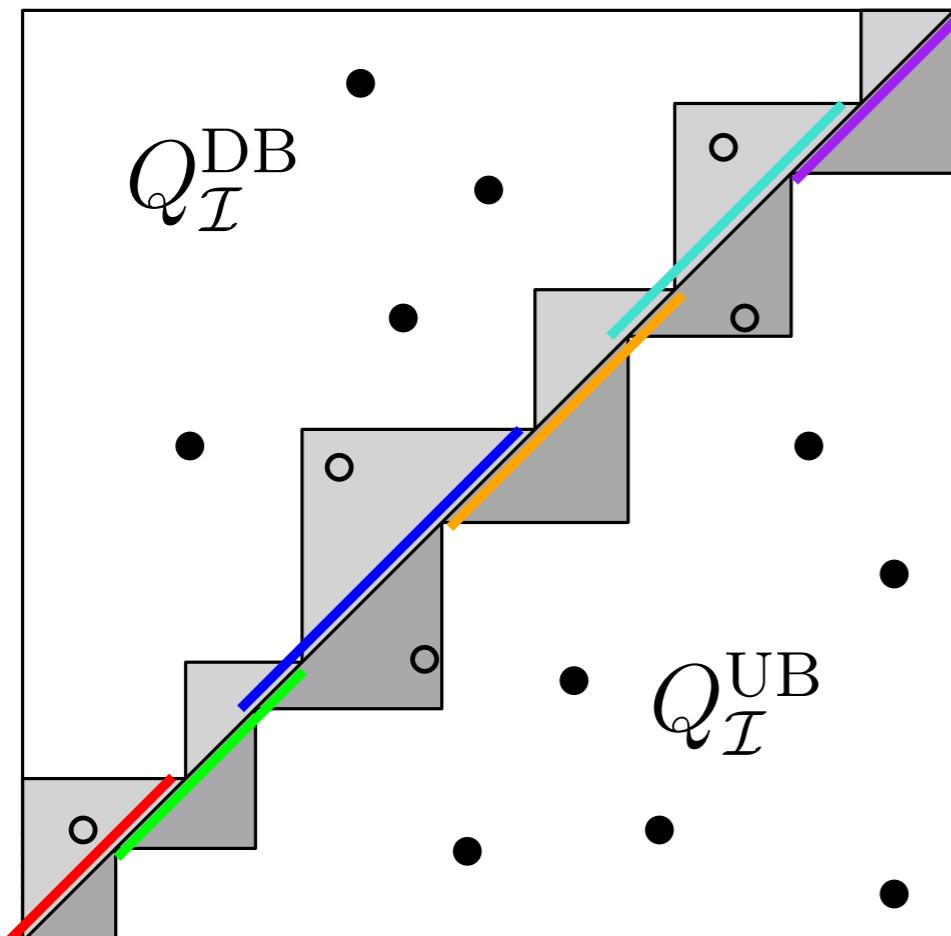
**Def:** Given  $X, f, \mathcal{I}$ :

$$\text{Dg}(\text{M}_f) = (\text{DB}(\text{R}_f) \setminus Q_{\mathcal{I}}^{\text{DB}}) \cup (\text{UB}(\text{R}_f) \setminus Q_{\mathcal{I}}^{\text{UB}}) \cup (\text{L}(\text{R}_f) \setminus Q_{\mathcal{I}}^{\text{L}})$$

$$Q_{\mathcal{I}}^{\text{DB}} = \bigcup_{I \in \mathcal{I}} Q_{\tilde{I} \cup I^+}^+$$

$$Q_{\mathcal{I}}^{\text{UB}} = \bigcup_{I \in \mathcal{I}} Q_{I^- \cup \tilde{I}}^-$$

$$Q_{\mathcal{I}}^{\text{L}} = \bigcup_{\substack{I, J \in \mathcal{I} \\ I \cap J \neq \emptyset}} Q_{I \cup J}^-$$



# Structure of Mapper

**Thm:** [C., Oudot 2016]

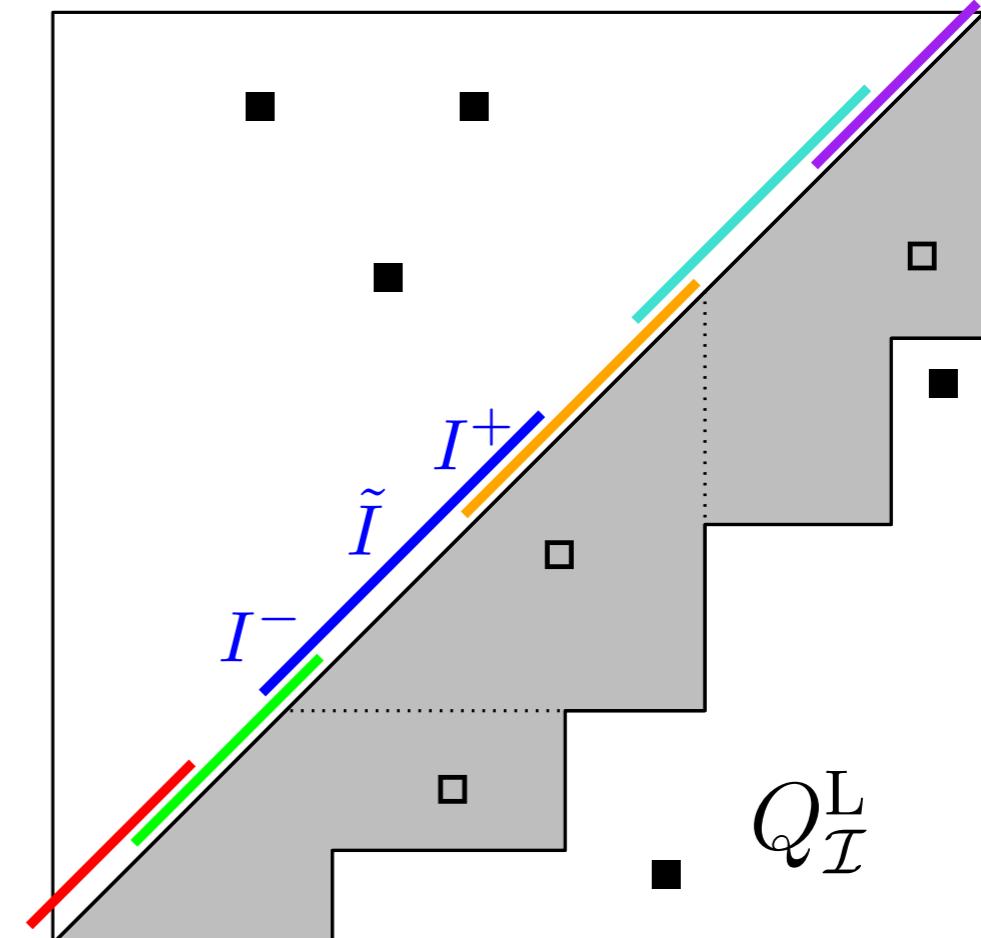
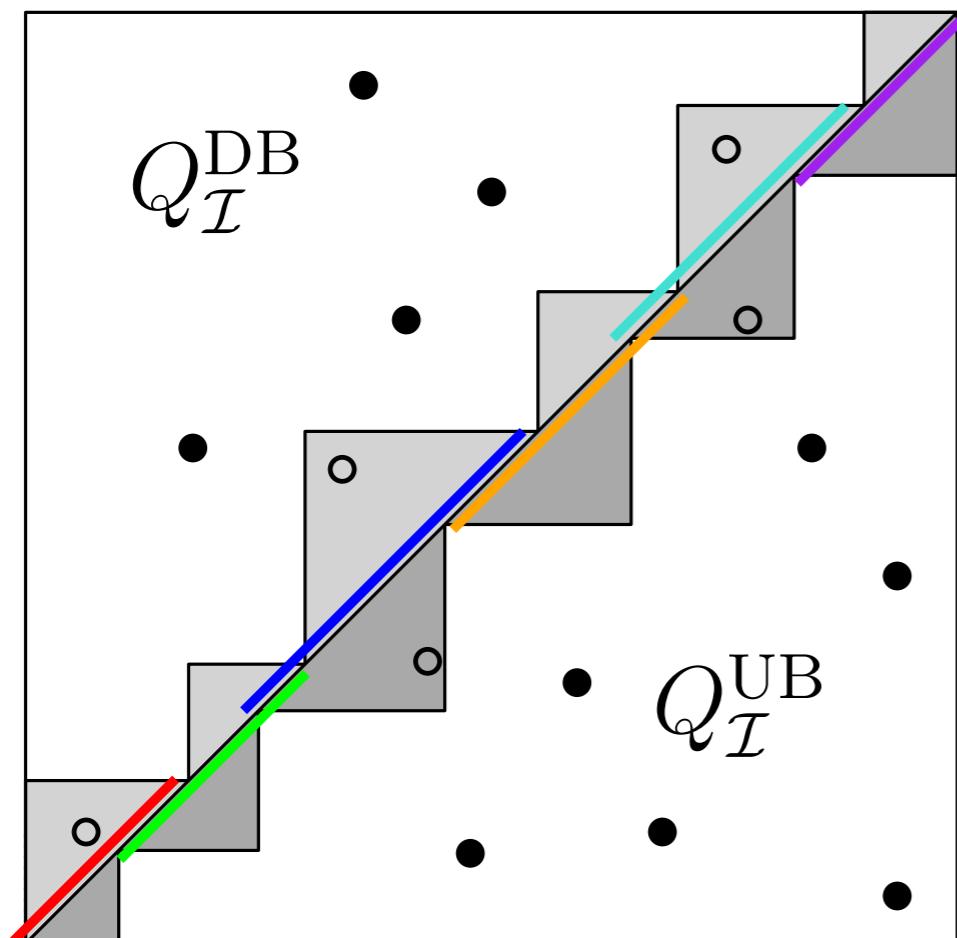
$Dg(M_f)$  provides a **bag-of-features** descriptor for  $M_f(X, \mathcal{I})$ :

$DB \longleftrightarrow$  downward branches

$UB \longleftrightarrow$  upward branches

$T \longleftrightarrow$  trunks (cc)

$L \longleftrightarrow$  loops



# Structure of Mapper

**Thm:** [C., Oudot 2016]

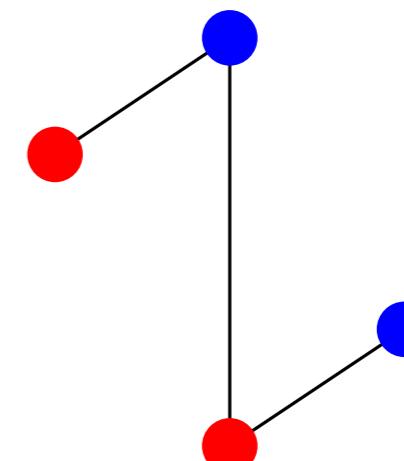
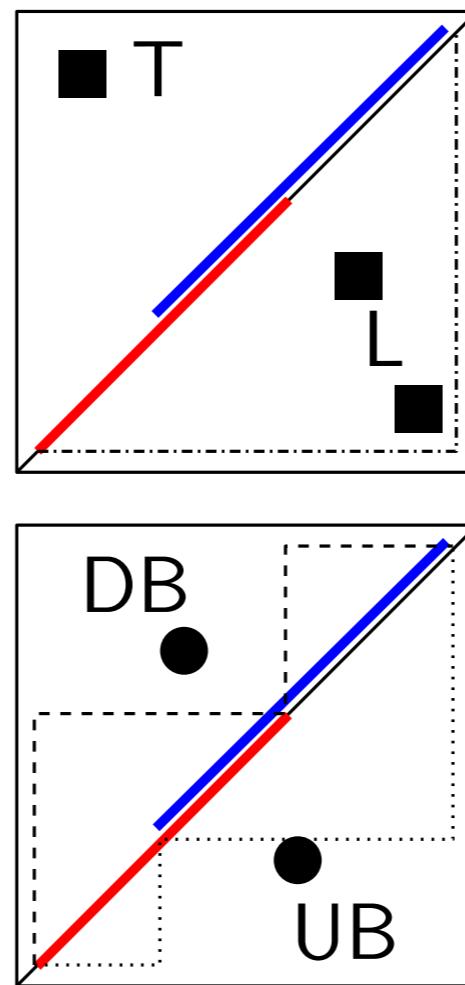
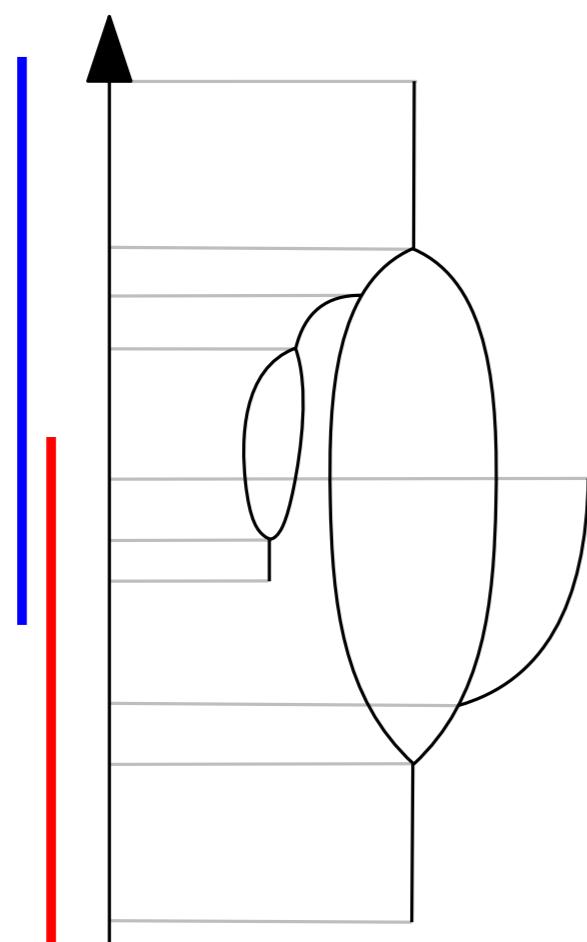
$Dg(M_f)$  provides a **bag-of-features** descriptor for  $M_f(X, \mathcal{I})$ :

$DB \longleftrightarrow$  downward branches

$UB \longleftrightarrow$  upward branches

$T \longleftrightarrow$  trunks (cc)

$L \longleftrightarrow$  loops



# Structure of Mapper

**Thm:** [C., Oudot 2016]

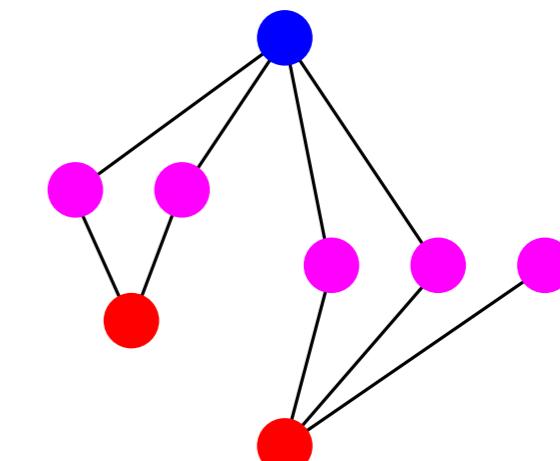
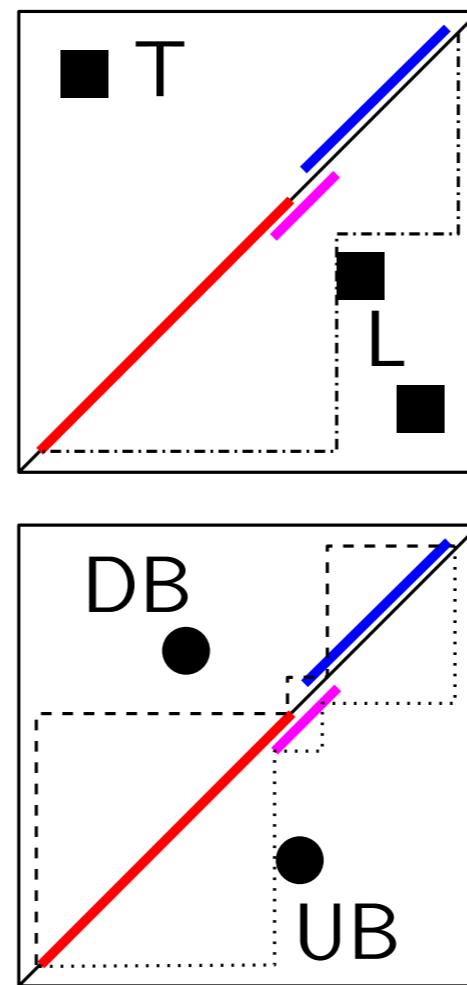
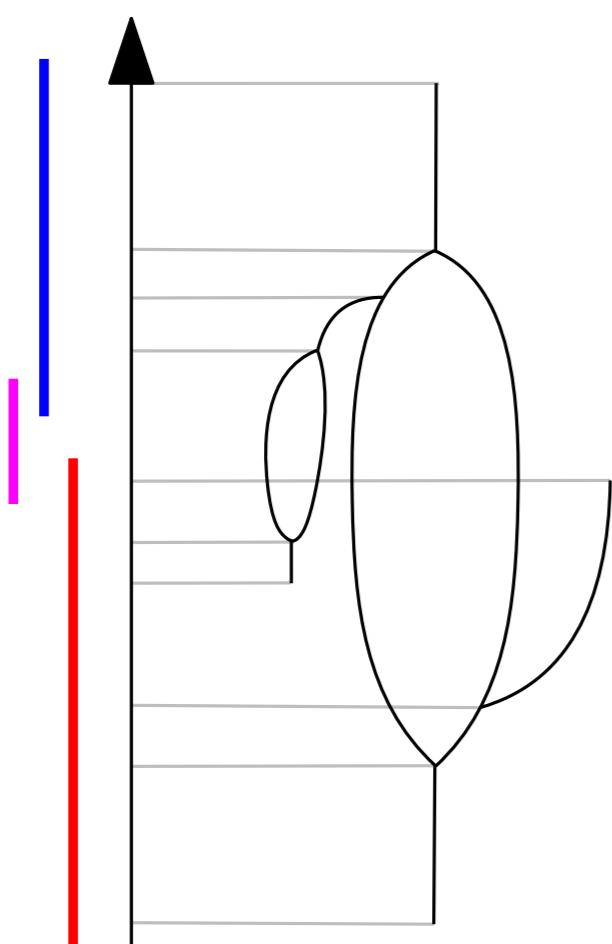
$Dg(M_f)$  provides a **bag-of-features** descriptor for  $M_f(X, \mathcal{I})$ :

$DB \longleftrightarrow$  downward branches

$UB \longleftrightarrow$  upward branches

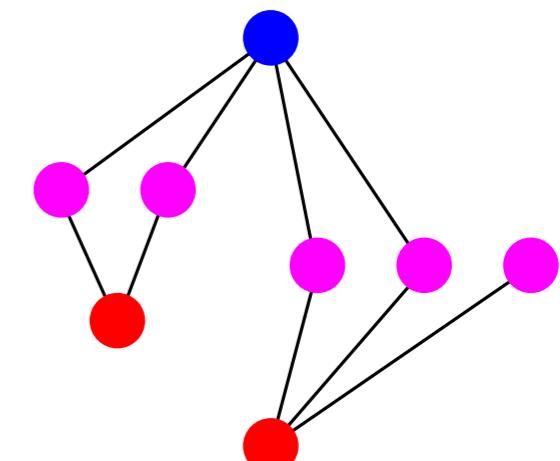
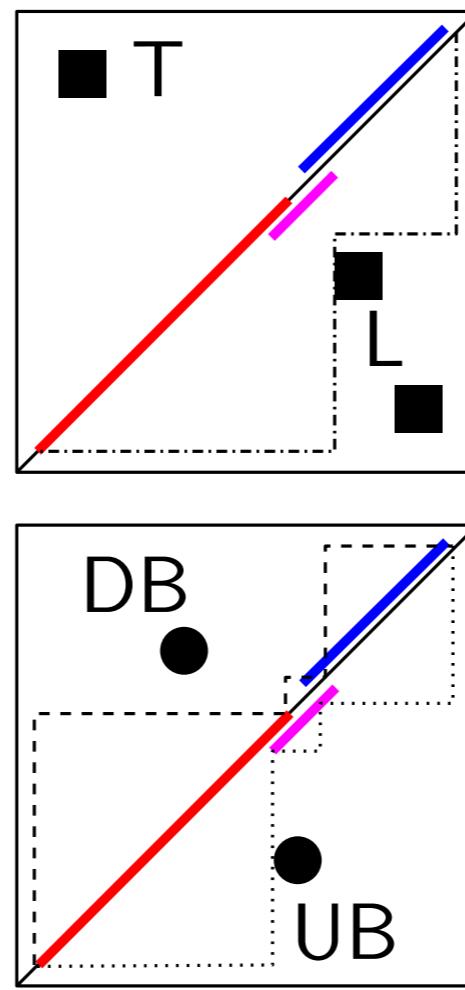
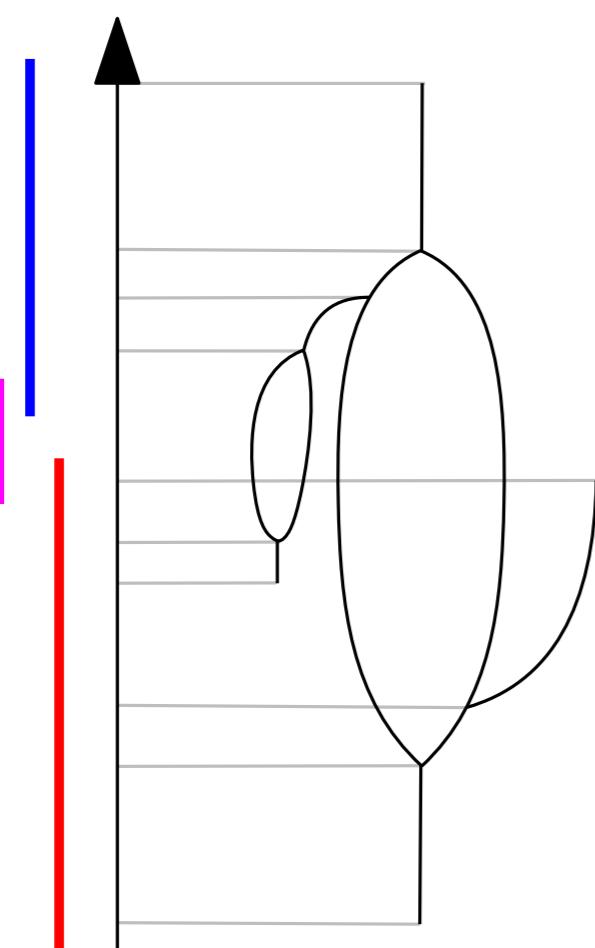
$T \longleftrightarrow$  trunks (cc)

$L \longleftrightarrow$  loops



# Structure of Mapper

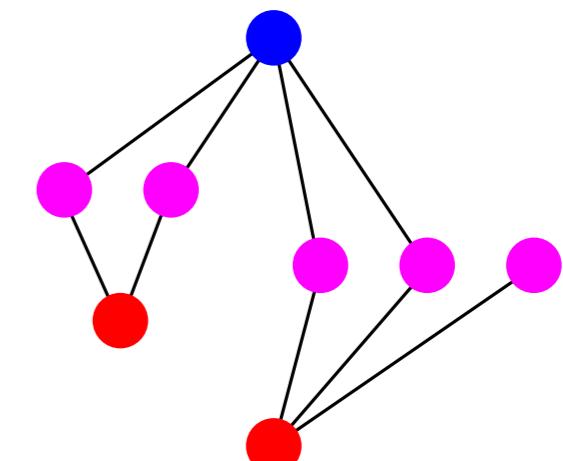
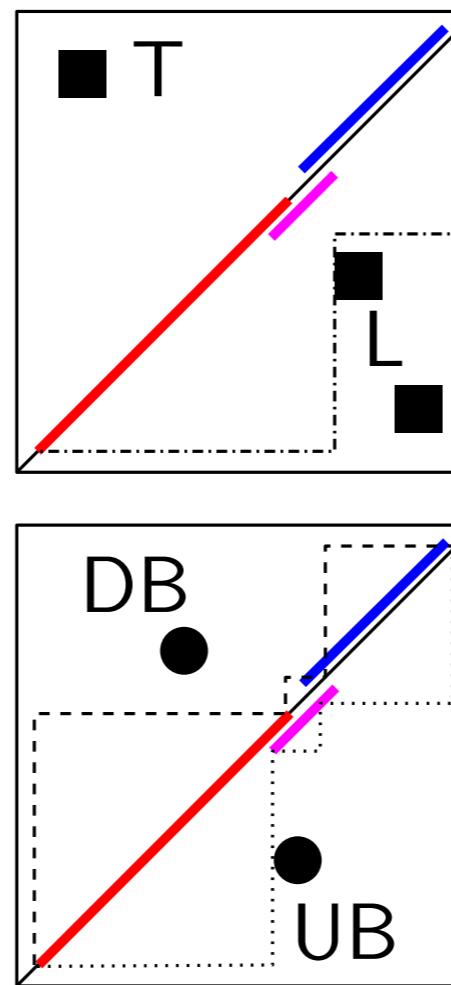
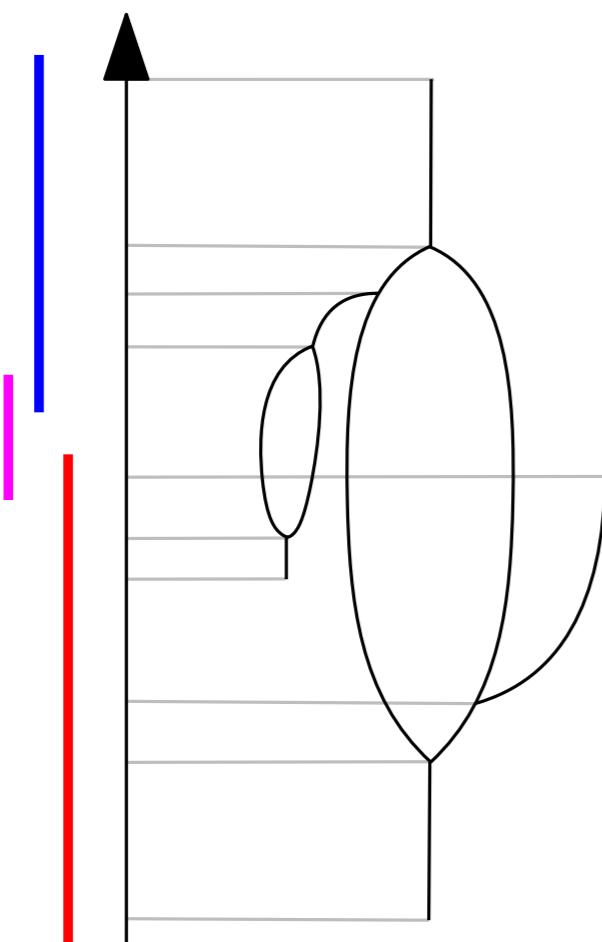
**Cor:**  $Dg(M_f) = Dg(R_f)$  whenever the resolution  $r$  of  $\mathcal{I}$  is smaller than the smallest distance from  $Dg(R_f) \setminus \Delta$  to the diagonal  $\Delta$



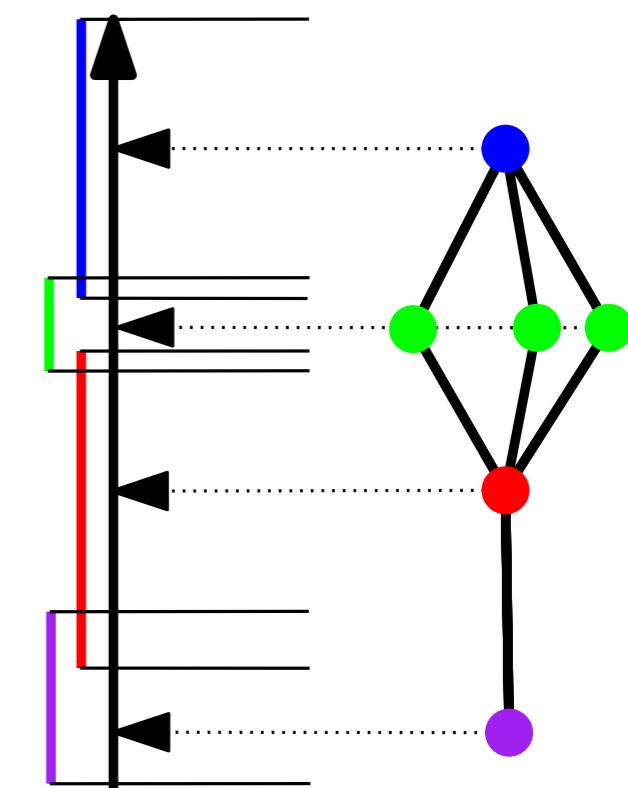
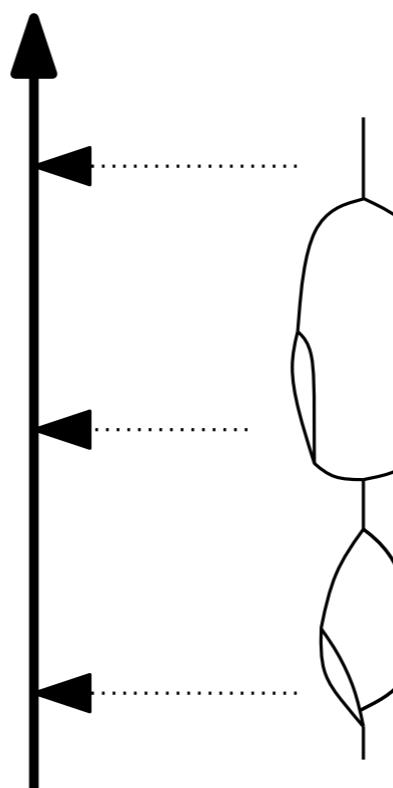
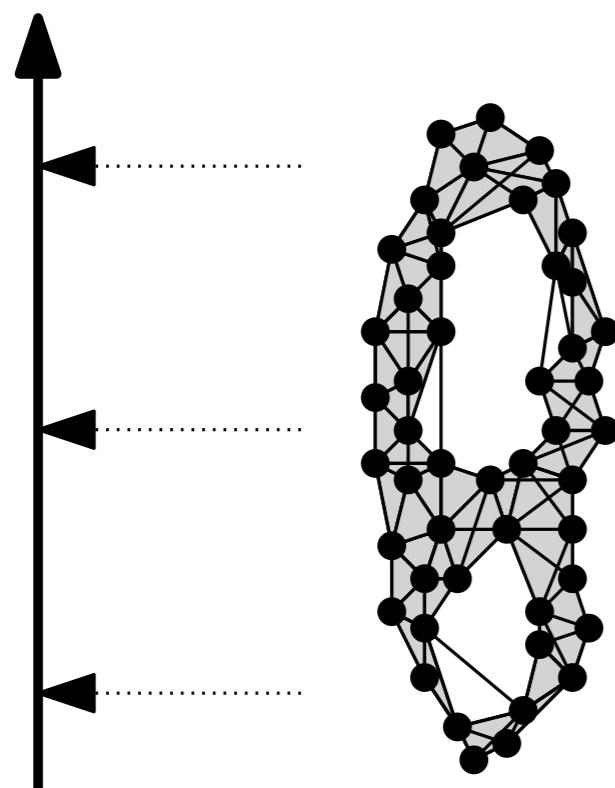
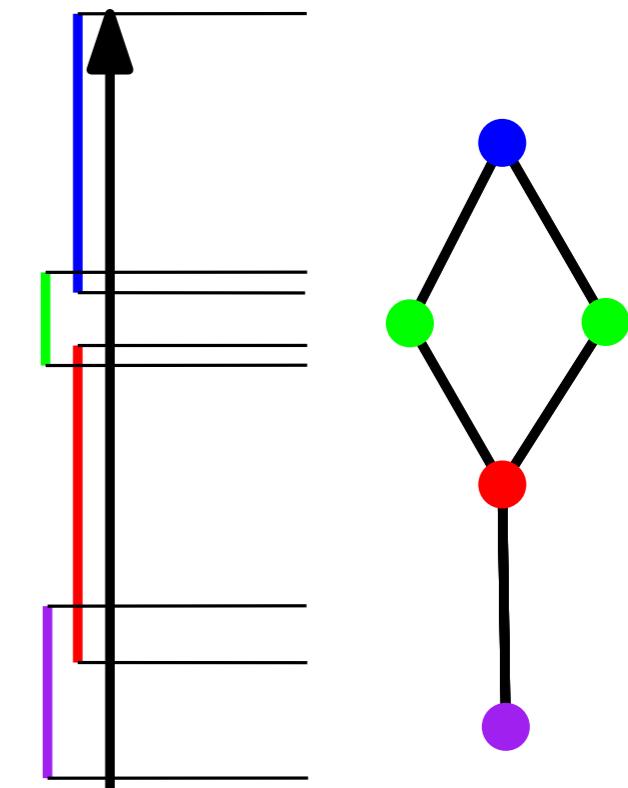
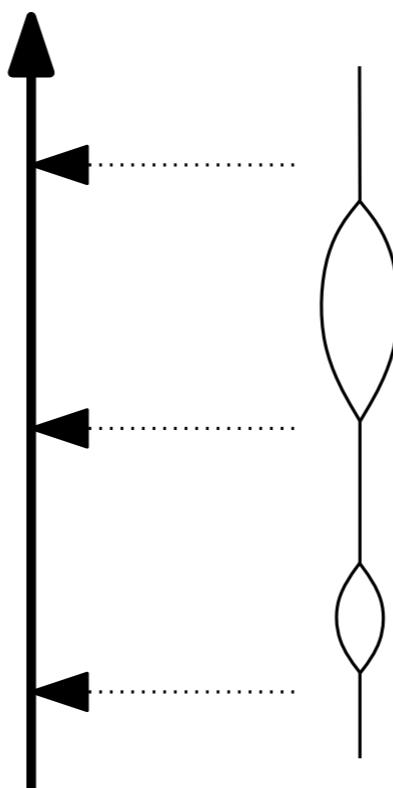
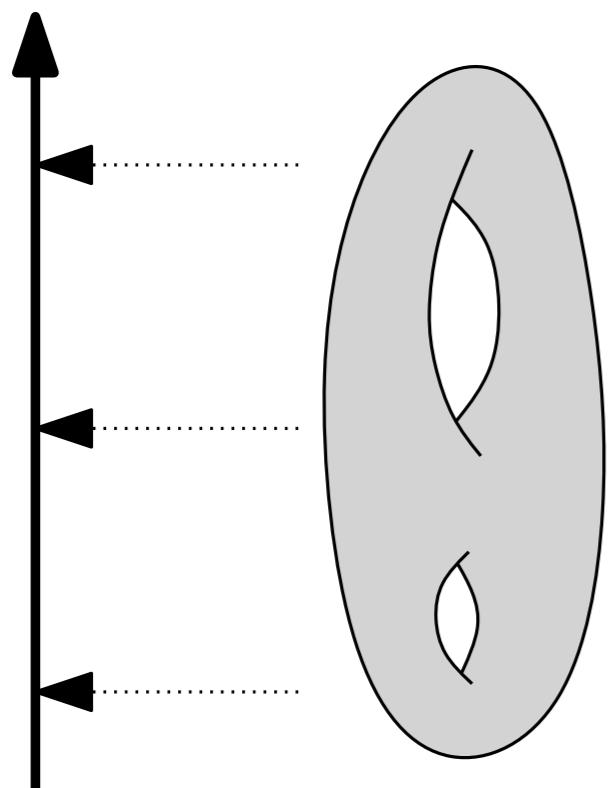
# Structure of Mapper

**Cor:**  $Dg(M_f) = Dg(R_f)$  whenever the resolution  $r$  of  $\mathcal{I}$  is smaller than the smallest distance from  $Dg(R_f) \setminus \Delta$  to the diagonal  $\Delta$

**Thm:** [C., Oudot 2017]  $d_{GH}(M_f(X, \mathcal{I}), R_f(X)) \leq 3r$

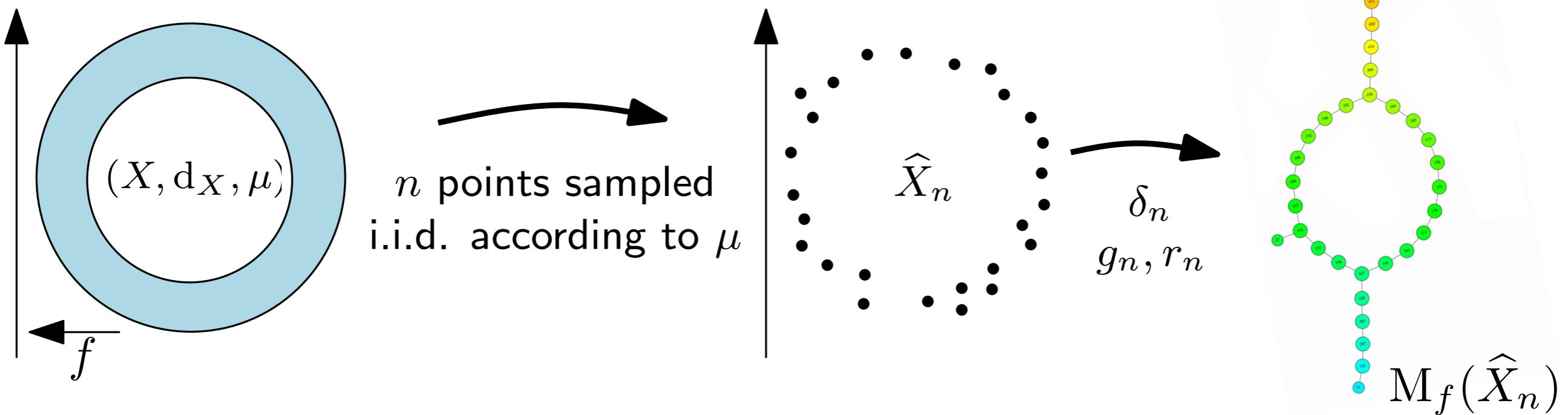


# Statistics for Mapper

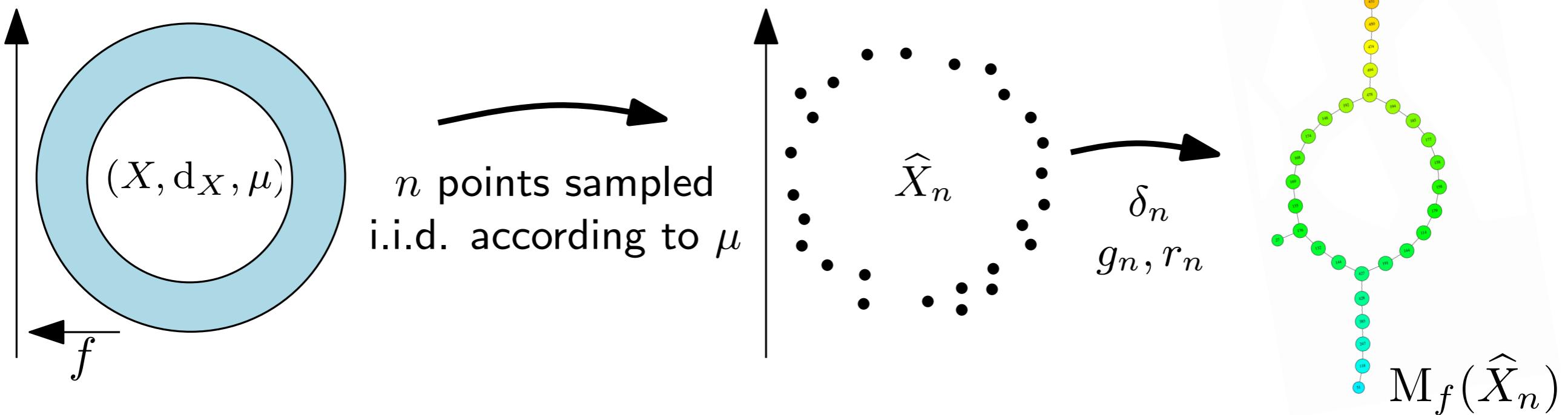


# Setup

---

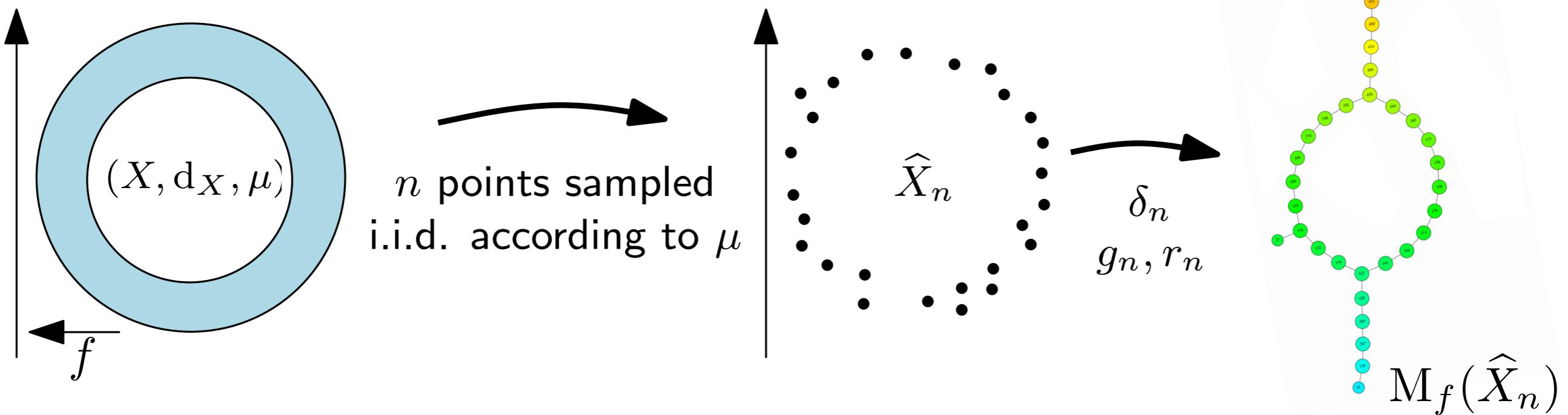


# Setup



**Prop:** [C., Michel, Oudot 2017]  $\hat{M}_n = M_f(\hat{X}_n)$  is measurable

# Setup



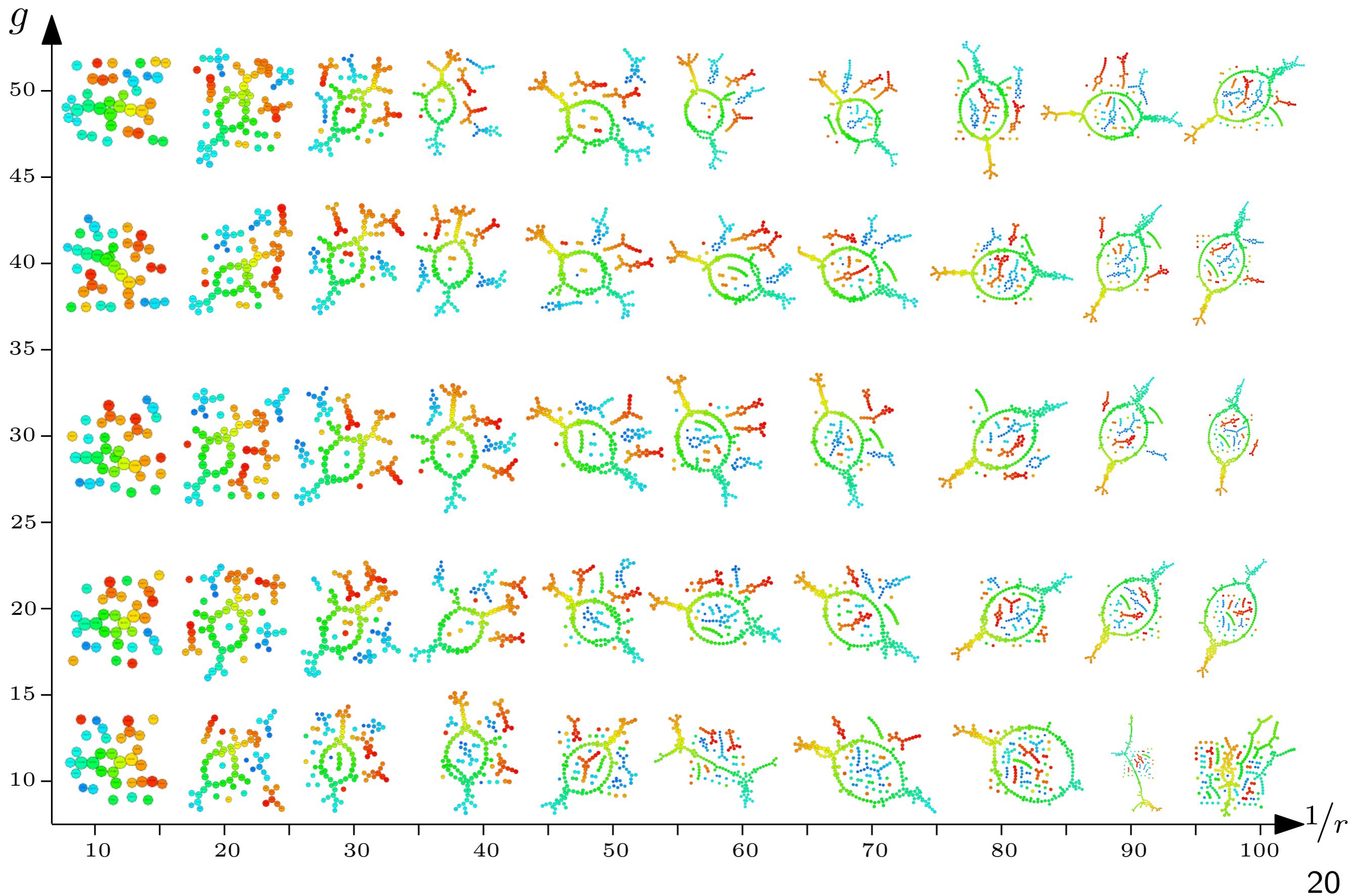
**Prop:** [C., Michel, Oudot 2017]  $\hat{M}_n = M_f(\hat{X}_n)$  is measurable

**Goal:** Find heuristics to compute "good"  $\delta_n, g_n, r_n$

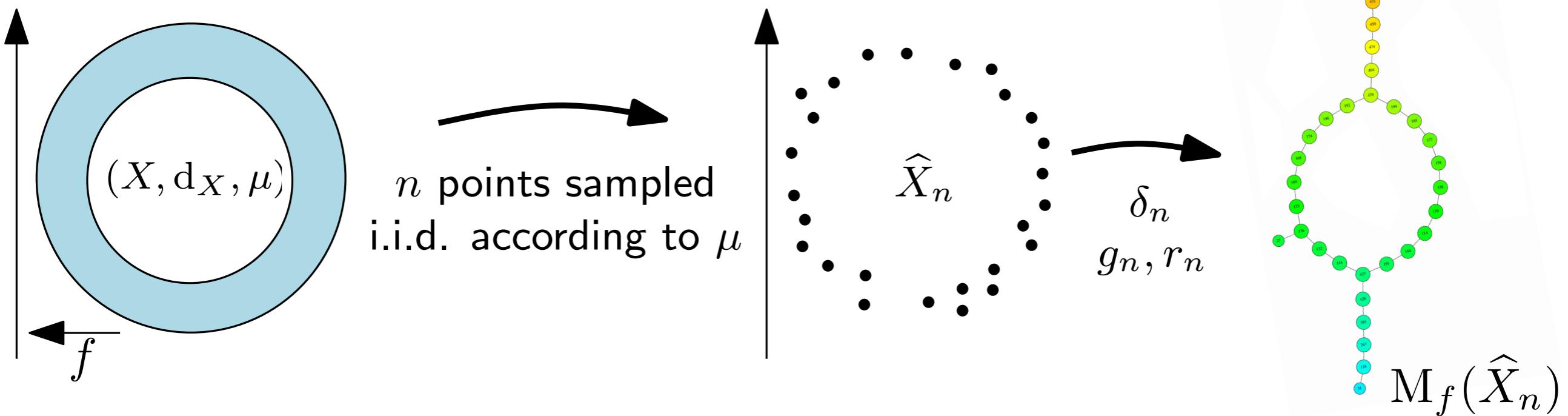
Assess quality through **confidence regions** and **convergence rates**

# Setup

---



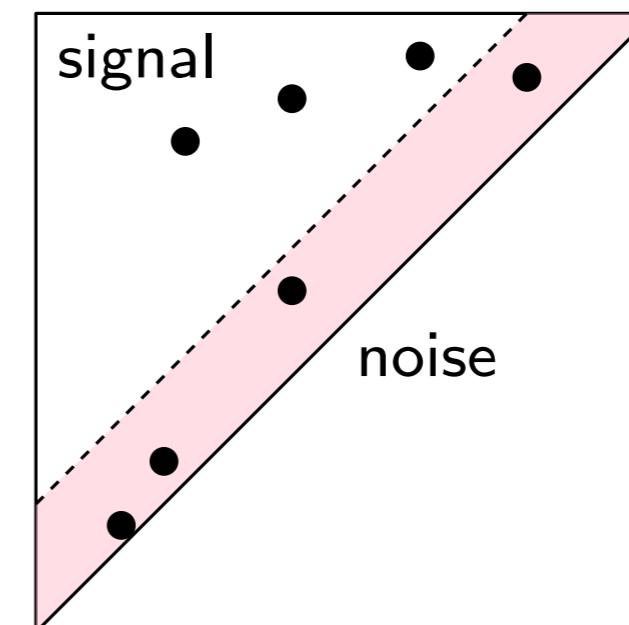
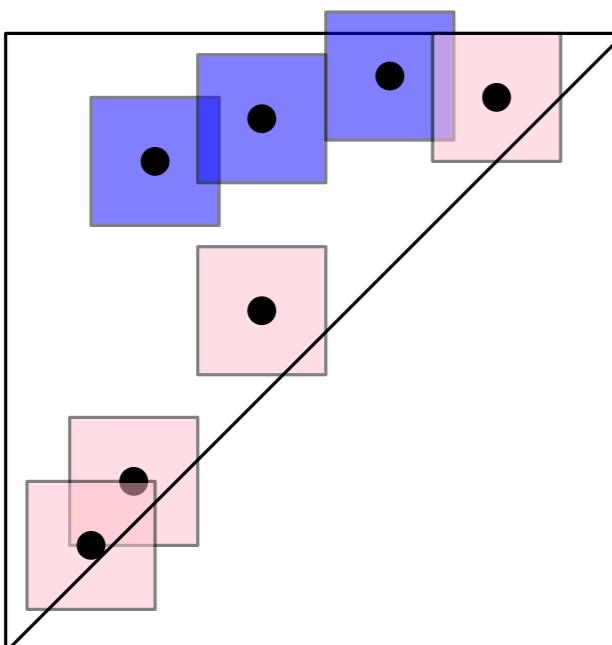
# Setup



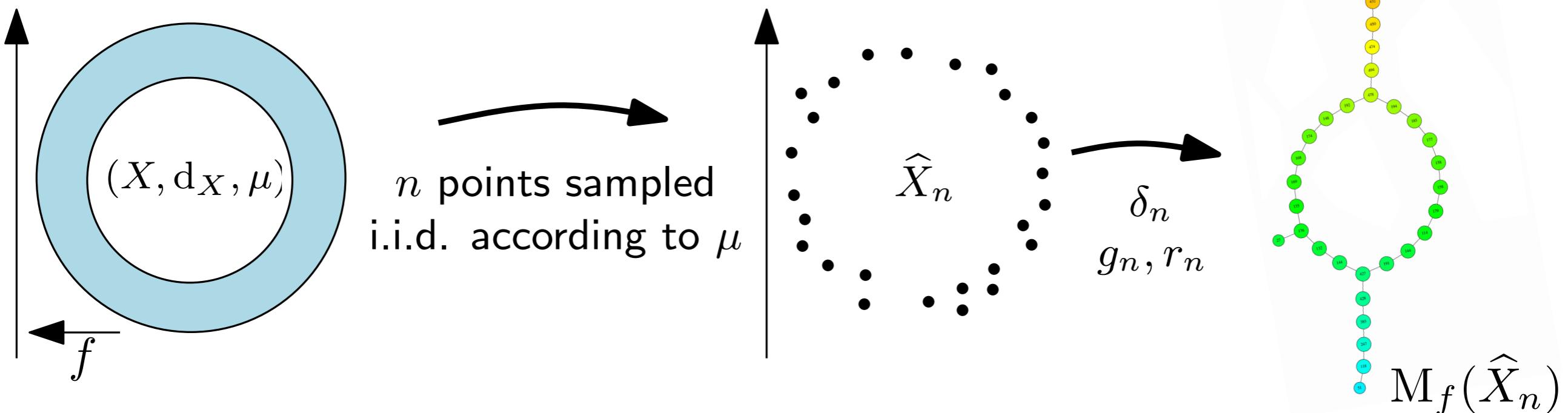
**Confidence regions:** given  $\alpha \in (0, 1)$ , find  $c_n(\alpha) \geq 0$  s.t.:

$$\limsup_{n \rightarrow \infty} \mathbb{P} \left( d_B \left( \hat{M}_n, R_f(X) \right) > c_n(\alpha) \right) \leq \alpha$$

→  $d_B$ -ball of radius  $c_n(\alpha)$  around  $Dg(\hat{M}_n)$



# Setup



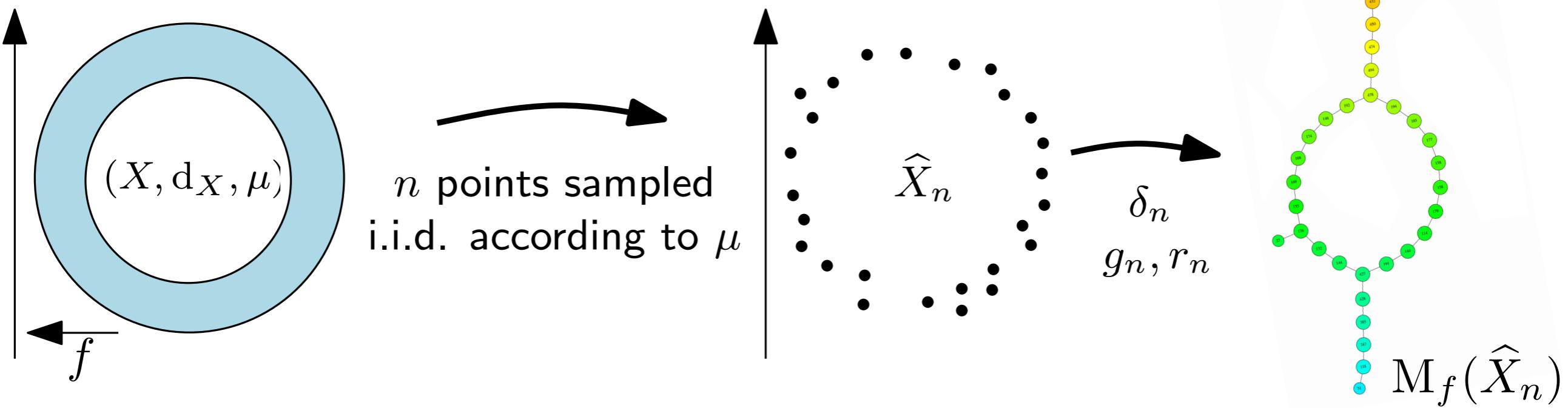
**Confidence regions:** given  $\alpha \in (0, 1)$ , find  $c_n(\alpha) \geq 0$  s.t.:

$$\limsup_{n \rightarrow \infty} \mathbb{P} \left( d_B \left( \hat{\mathbf{M}}_n, \mathbf{R}_f(X) \right) > c_n(\alpha) \right) \leq \alpha$$

→  $d_B$ -ball of radius  $c_n(\alpha)$  around  $\text{Dg}(\hat{\mathbf{M}}_n)$

**Convergence Rate:** estimate  $\mathbb{E} \left[ d_B(\hat{\mathbf{M}}_n, \mathbf{R}_f(X)) \right]$  w.r.t.  $n$

# Rate of Convergence



**Regularity of the filter function:** (exact) modulus of continuity of  $f$

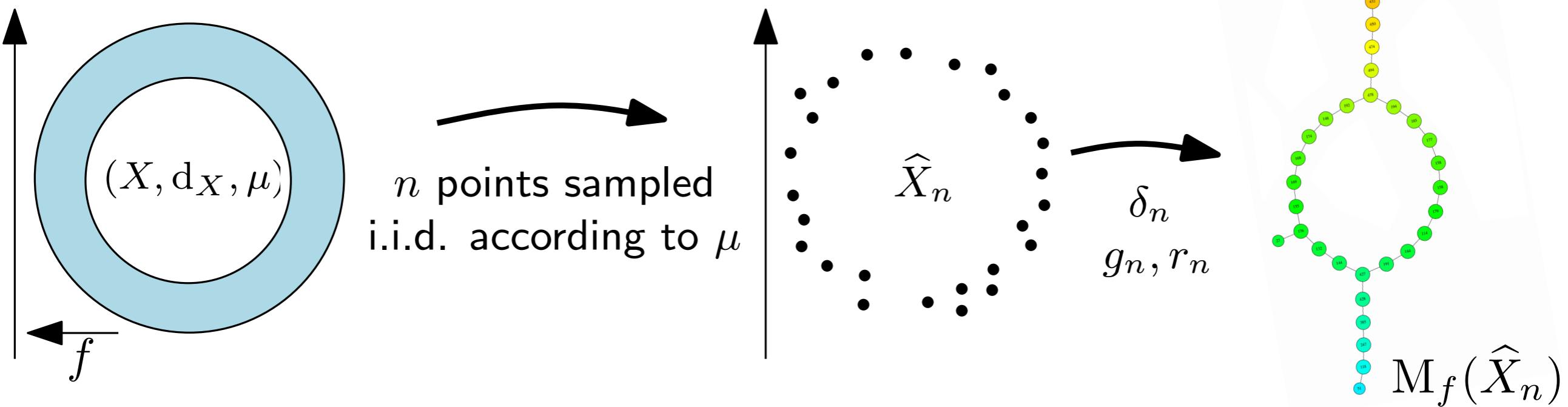
$$\omega(\delta) = \sup_{\|x-x'\| \leq \delta} |f(x) - f(x')|$$

**Approximation inequality:** [C., Michel, Oudot 2017]

Let  $\hat{X}_n \subset X$ . Under some *regularity assumptions* on  $X, f, \delta, r, g$ , one has:

$$d_B \left( R_f(X), M_f(\hat{X}_n) \right) \leq r + 2\omega(\delta)$$

# Rate of Convergence



$$4d_H(\hat{X}_n, X) \leq \delta \leq C(X)$$

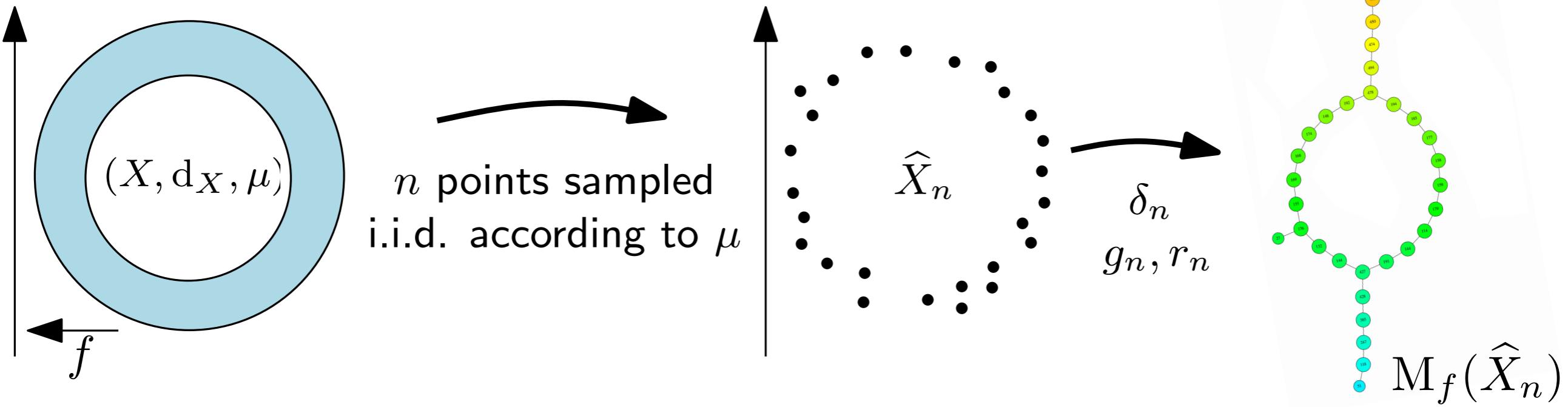
$$\max\{|f(X_i) - f(X_j)| : \|X_i - X_j\| \leq \delta\} < gr$$

**Approximation inequality:** [C., Michel, Oudot 2017]

Let  $\hat{X}_n \subset X$ . Under some **regularity assumptions** on  $X, f, \delta, r, g$ , one has:

$$d_B(R_f(X), M_f(\hat{X}_n)) \leq r + 2\omega(\delta)$$

# Rate of Convergence



$$4d_H(\hat{X}_n, X) \leq \delta \leq C(X)$$

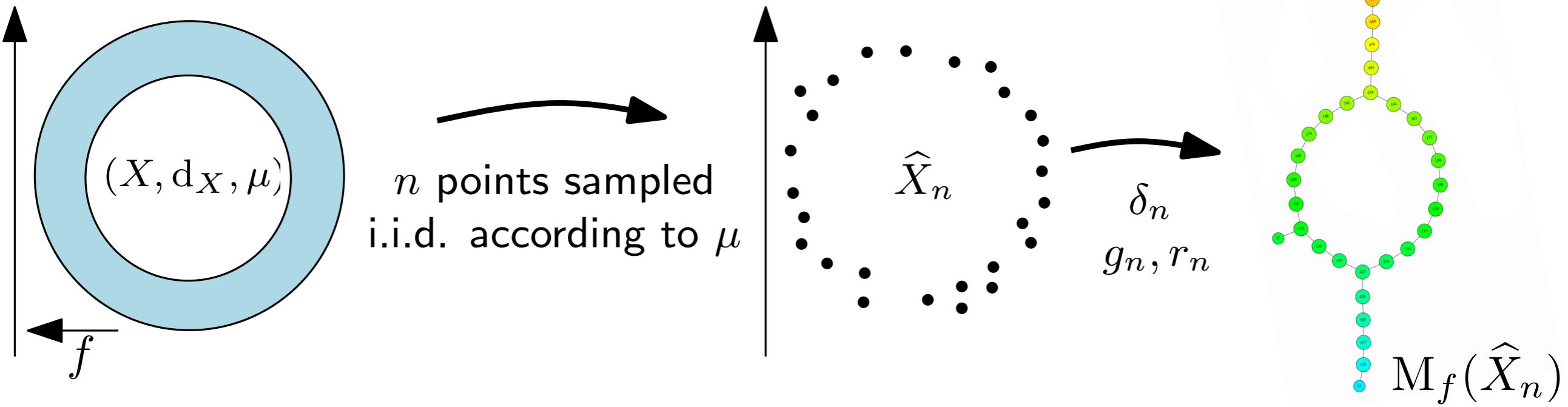
$$\max\{\max\{|f(X_i) - f(X_j)|, |\hat{f}(X_i) - \hat{f}(X_j)|\} : \|X_i - X_j\| \leq \delta\} \leq rg$$

**Approximation inequality:** [C., Michel, Oudot 2017]

Let  $\hat{X}_n \subset X$ . Under some **regularity assumptions** on  $X, f, \delta, r, g$ , one has:

$$d_B(R_f(X), M_{\hat{f}}(\hat{X}_n)) \leq 2r + 2\omega(\delta) + \max\{|f(X_i) - \hat{f}(X_i)|\}$$

# Rate of Convergence



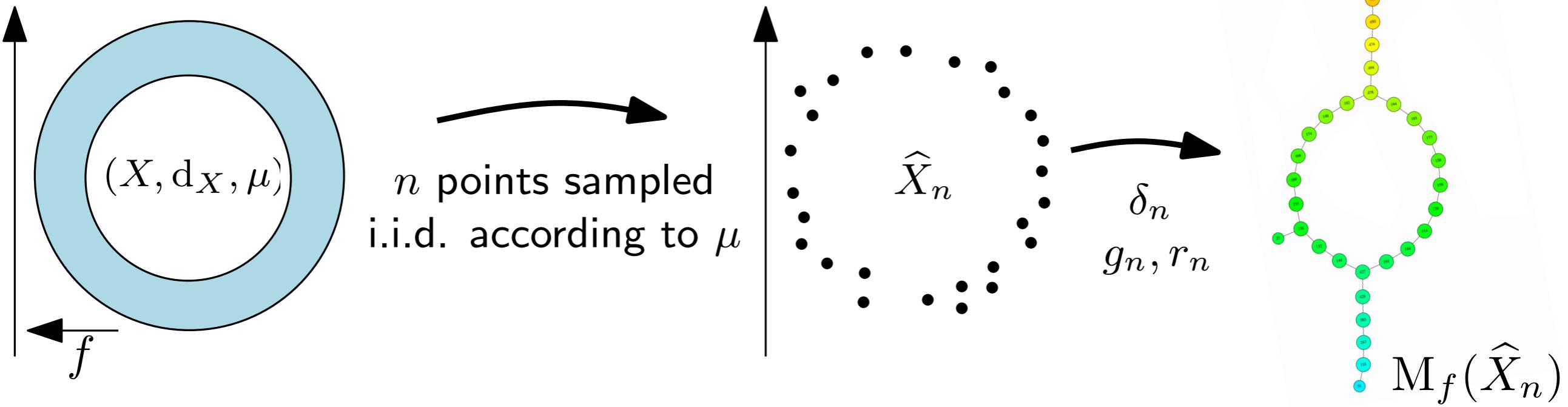
$$V_n = \max\{f(X_i) - f(X_j) : \|X_i - X_j\| \leq \delta_n\}$$

**Thm:** [C., Michel, Oudot 2017]

If  $\mu$  is  $(a, b)$ -standard, then for  $\delta_n = 4 \left( \frac{2 \log n}{an} \right)^{1/b}$ ,  $g_n \in \left( \frac{1}{3}, \frac{1}{2} \right)$ ,  $r_n = \frac{V_n}{g_n}$ , one has:

$$\sup_{\mu \in \mathcal{P}} \mathbb{E} \left[ d_B \left( M_f(\hat{X}_n), R_f(X) \right) \right] \lesssim \omega \left( \frac{\log n}{n} \right)^{1/b}$$

# Rate of Convergence



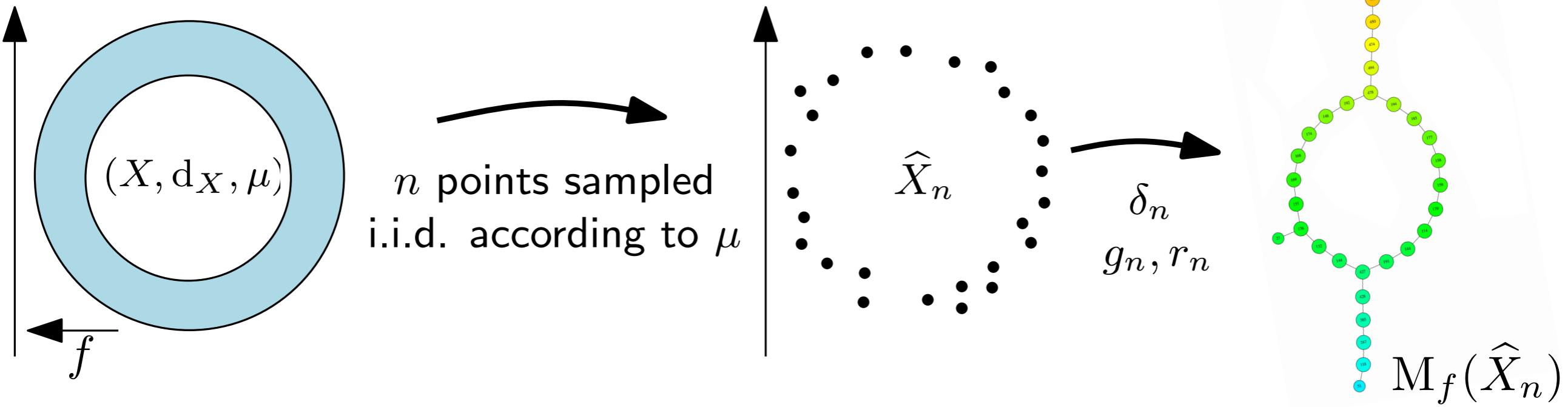
$$V_n = \max\{f(X_i) - f(X_j) : \|X_i - X_j\| \leq \delta_n\}$$

**Thm:** [C., Michel, Oudot 2017]

If  $\mu$  is  $(a, b)$ -standard, then for  $\delta_n = 4 \left( \frac{2 \log n}{an} \right)^{1/b}$ ,  $g_n \in \left( \frac{1}{3}, \frac{1}{2} \right)$ ,  $r_n = \frac{V_n}{g_n}$ , one has:

$$\sup_{\mu \in \mathcal{P}} \mathbb{E} \left[ d_B \left( M_f(\hat{X}_n), R_f(X) \right) \right] \lesssim \omega \left( \frac{\log n}{n} \right)^{1/b}$$

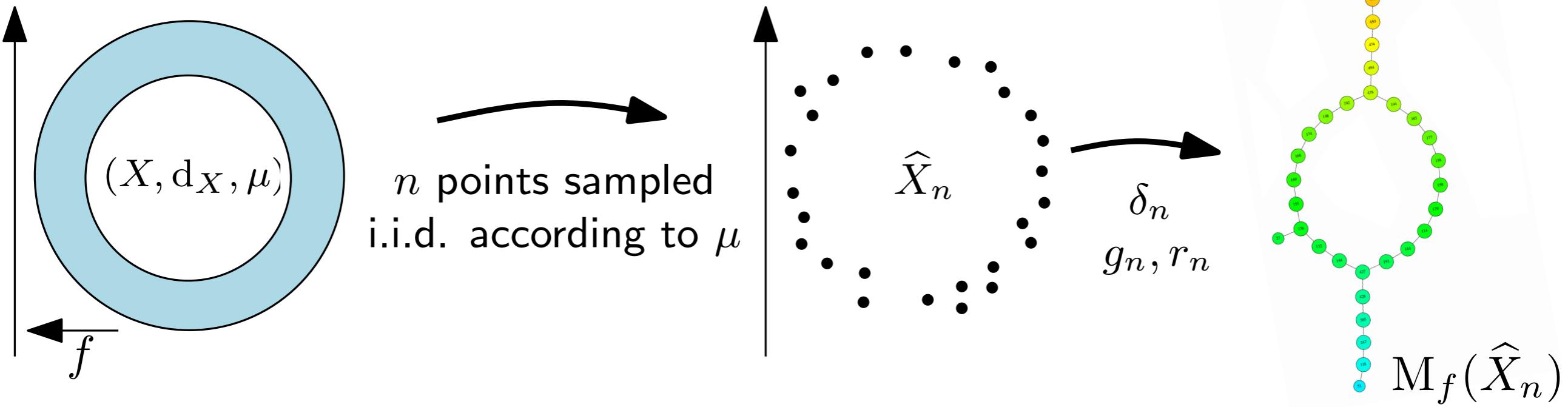
# Rate of Convergence



Subsampling to tune  $\delta_n$ : let  $\beta > 0$  and take  $s(n) = n \log(n)^{-(1+\beta)}$

$\delta_n = d_H(\hat{X}_n^{s(n)}, \hat{X}_n)$  where  $\hat{X}_n^{s(n)} \subset \hat{X}_n$  of size  $s(n)$

# Rate of Convergence



Subsampling to tune  $\delta_n$ : let  $\beta > 0$  and take  $s(n) = n \log(n)^{-(1+\beta)}$

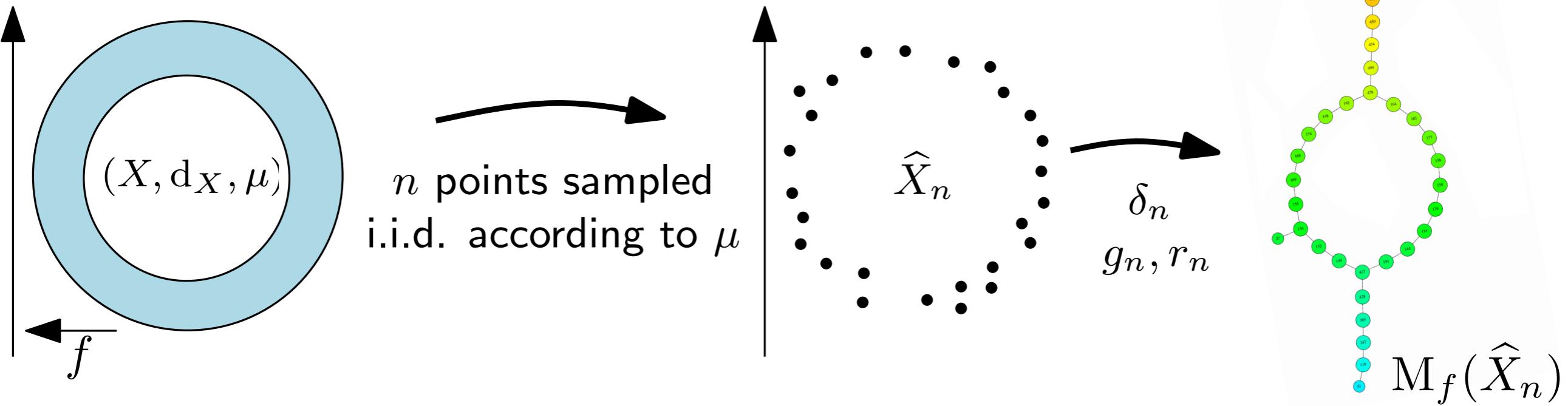
$\delta_n = d_H(\hat{X}_n^{s(n)}, \hat{X}_n)$  where  $\hat{X}_n^{s(n)} \subset \hat{X}_n$  of size  $s(n)$

**Thm:** [C., Michel, Oudot 2017]

If  $\mu$  is  $(a, b)$ -standard, then for  $\delta_n, g_n \in (\frac{1}{3}, \frac{1}{2})$ ,  $r_n = \frac{V_n}{g_n}$ , one has

$$\sup_{\mu \in \mathcal{P}} \mathbb{E} \left[ d_B \left( M_f(\hat{X}_n), R_f(X) \right) \right] \lesssim \omega \left( \frac{\log(n)^{2+\beta}}{n} \right)^{1/b}$$

# Rate of Convergence



Subsampling to tune  $\delta_n$ : let  $\beta > 0$  and take  $s(n) = n \log(n)^{-(1+\beta)}$

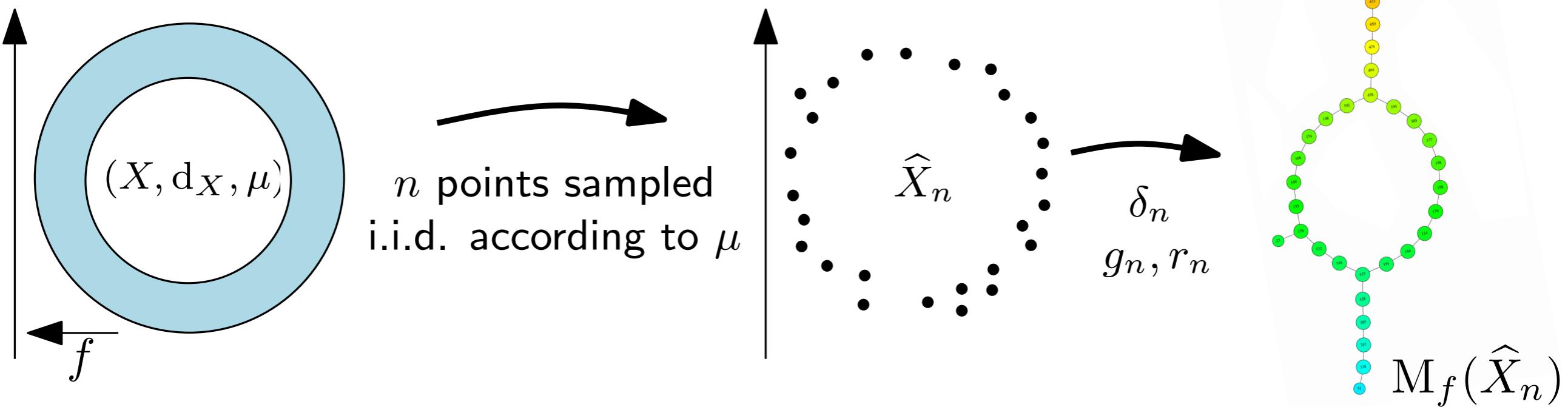
$\delta_n = d_H(\hat{X}_n^{s(n)}, \hat{X}_n)$  where  $\hat{X}_n^{s(n)} \subset \hat{X}_n$  of size  $s(n)$

**Thm:** [C., Michel, Oudot 2017]

If  $\mu$  is  $(a, b)$ -standard, then for  $\delta_n, g_n \in (\frac{1}{3}, \frac{1}{2})$ ,  $r_n = \frac{\max\{V_n, \widehat{V}_n\}}{g_n}$ , one has

$$\sup_{\mu \in \mathcal{P}} \mathbb{E} \left[ d_B \left( M_{\hat{f}}(\hat{X}_n), R_f(X) \right) \right] \lesssim \omega \left( \frac{\log(n)^{2+\beta}}{n} \right)^{1/b} + \mathbb{E} \left[ \max |f(X) - \hat{f}(X)| \right]$$

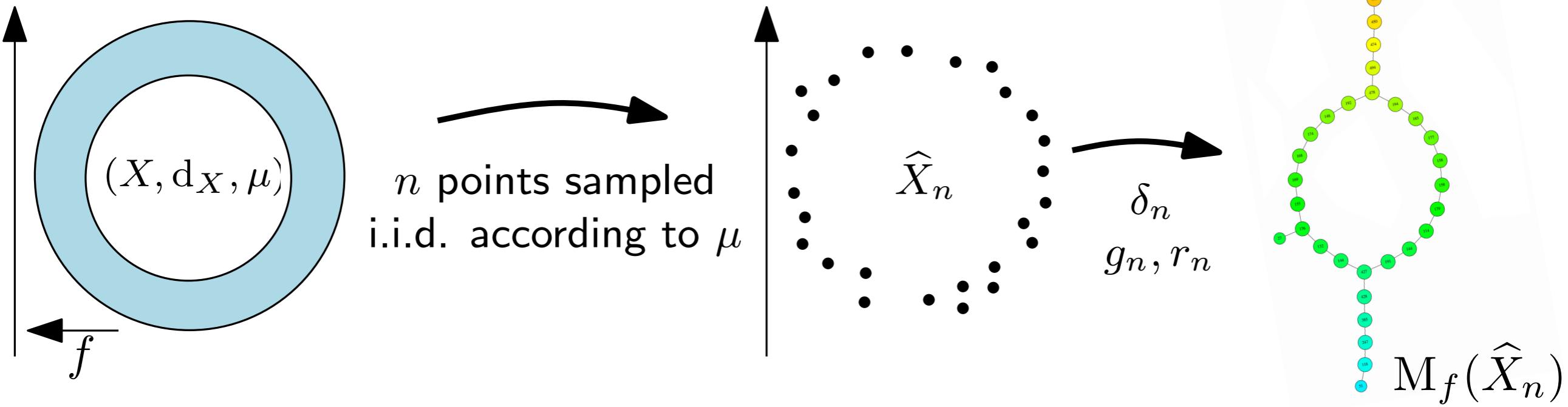
# Rate of Convergence



**Minimax Optimality:** [C., Michel, Oudot 2017] for any estimator  $\hat{R}_n$ ,

$$\omega \left( \frac{1}{n} \right)^{1/b} \lesssim \sup_{\mu \in \mathcal{P}} \mathbb{E} \left[ d_B \left( \hat{R}_n, R_f(X) \right) \right]$$

# Rate of Convergence



**Ex : PCA filter**

$\Pi_1$ : 1st principal direction of covariance operator

$\hat{\Pi}_1$ : 1st principal direction of empirical covariance operator

Using [Biau et. al. 2012]:

$$\mathbb{E} \left[ d_B \left( R_{\Pi_1}(X), M_{\hat{\Pi}_1}(\hat{X}_n) \right) \right] \lesssim \left( \frac{\log(n)^{2+\beta}}{n} \right)^{1/b} \vee \frac{1}{\sqrt{n}}$$

# Confidence Regions

---

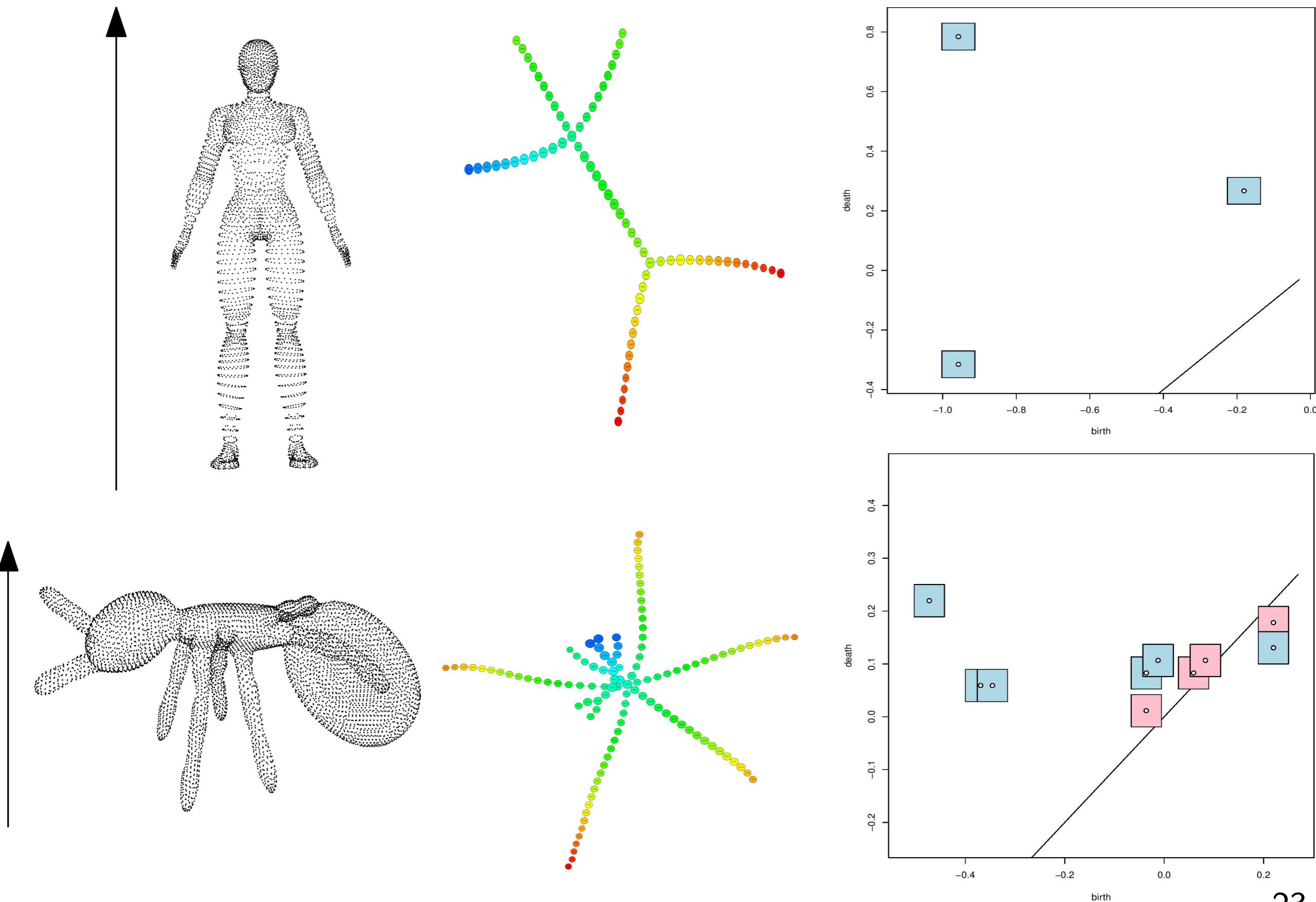
**Either from proof of previous result with:**

$$\mathbb{E} \left[ d_B \left( M_f(\hat{X}_n), R_f(X) \right) \right] = \int_{\alpha} \mathbb{P} \left( d_B \left( M_f(\hat{X}_n), R_f(X) \right) \geq \alpha \right) d\alpha$$

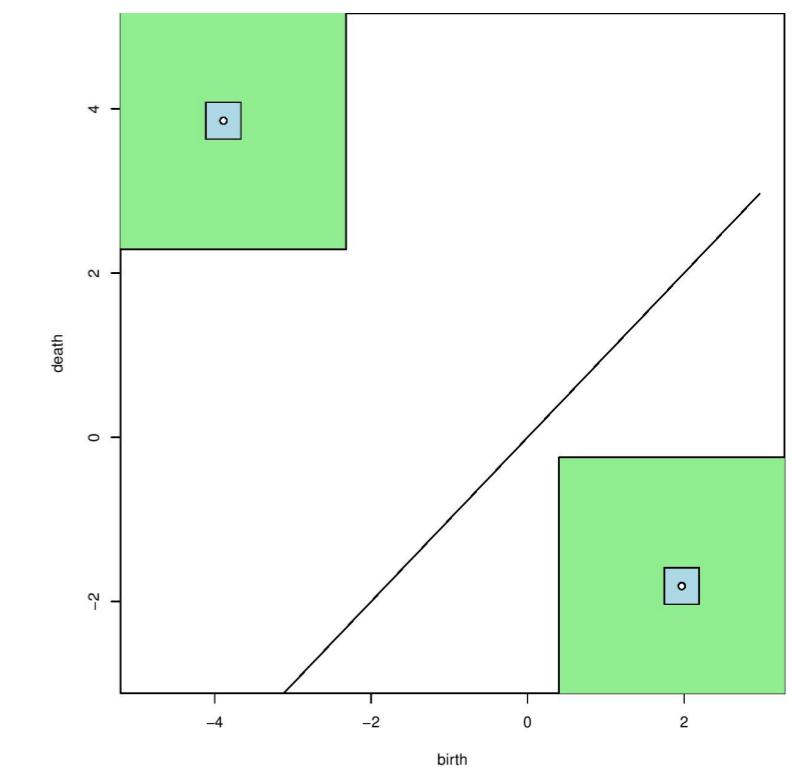
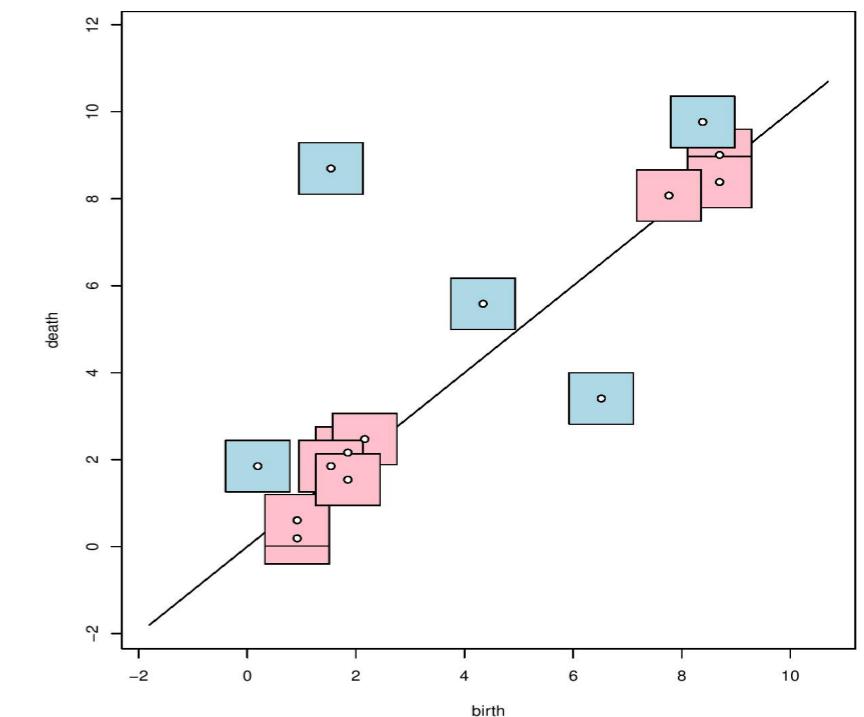
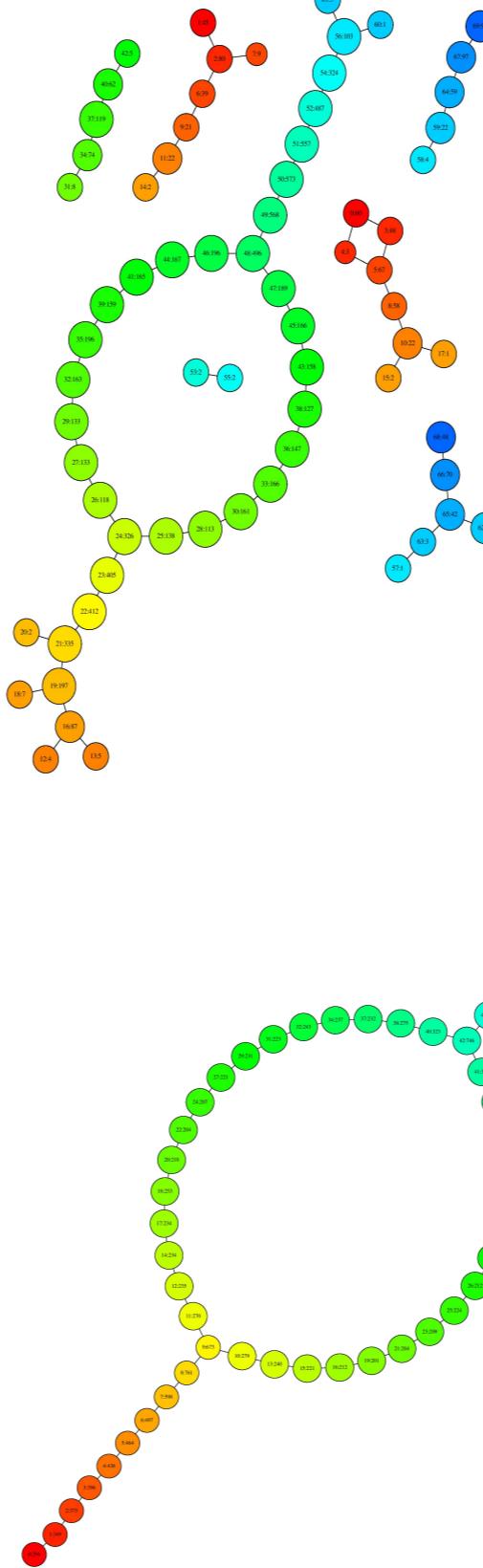
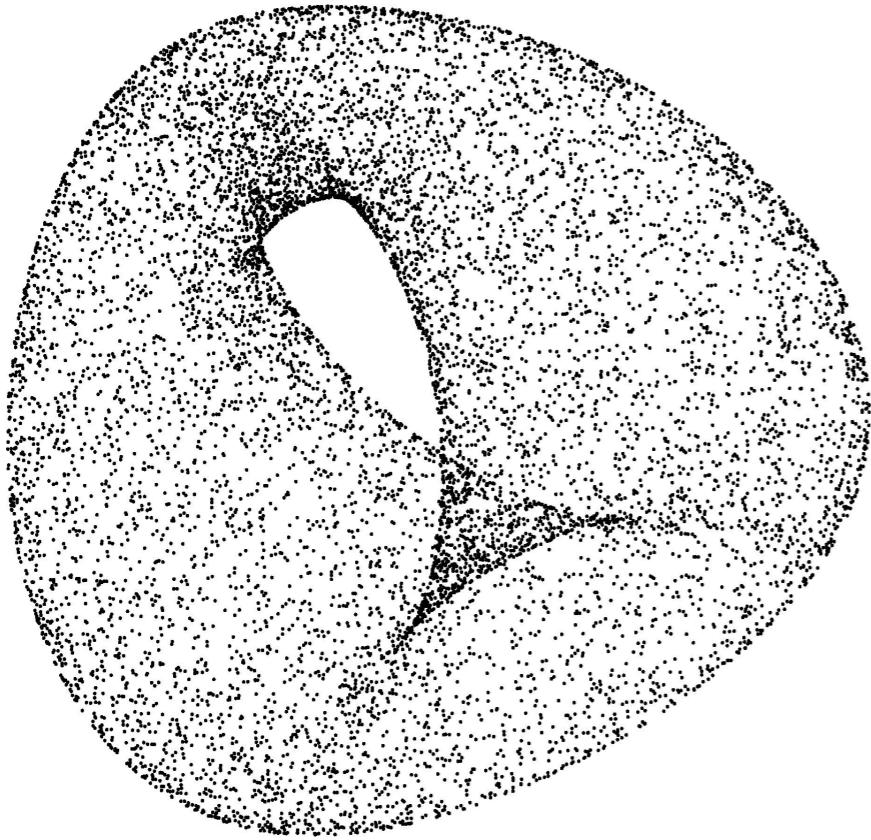
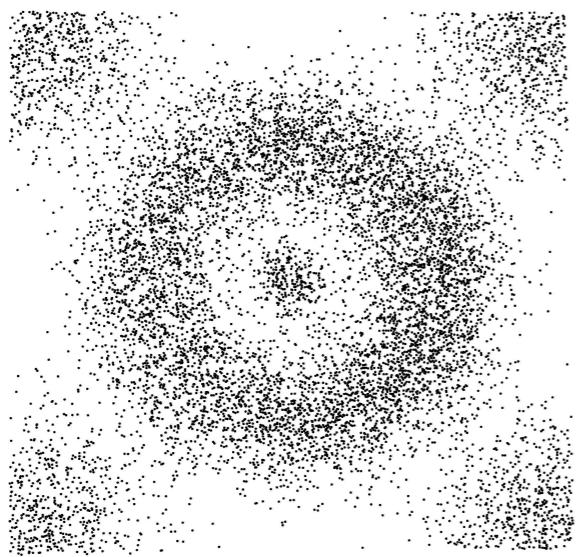
**Or bootstrap (only empirical):**

- draw  $\hat{X}_n^* = X_1^*, \dots, X_n^*$  iid from  $\mu_{\hat{X}_n}$  (empirical measure on  $\hat{X}_n$ )
- compute  $d^* = d_B \left( M_f(\hat{X}_n^*), M_f(\hat{X}_n) \right)$
- repeat  $N$  times to get  $d_1^*, \dots, d_N^*$
- let  $q_\alpha$  be the  $(1 - \alpha)$  quantile of  $\frac{1}{N} \sum_{i=1}^N \mathbb{I}(\sqrt{n} d_i^* \geq t)$
- take  $c_n(\alpha) = \frac{q_\alpha}{\sqrt{n}}$

# Experiments

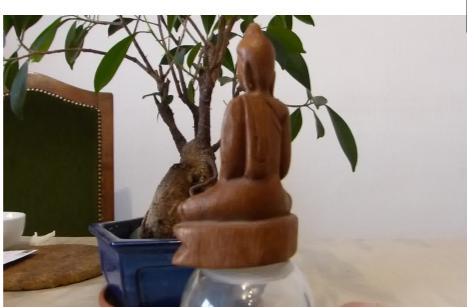
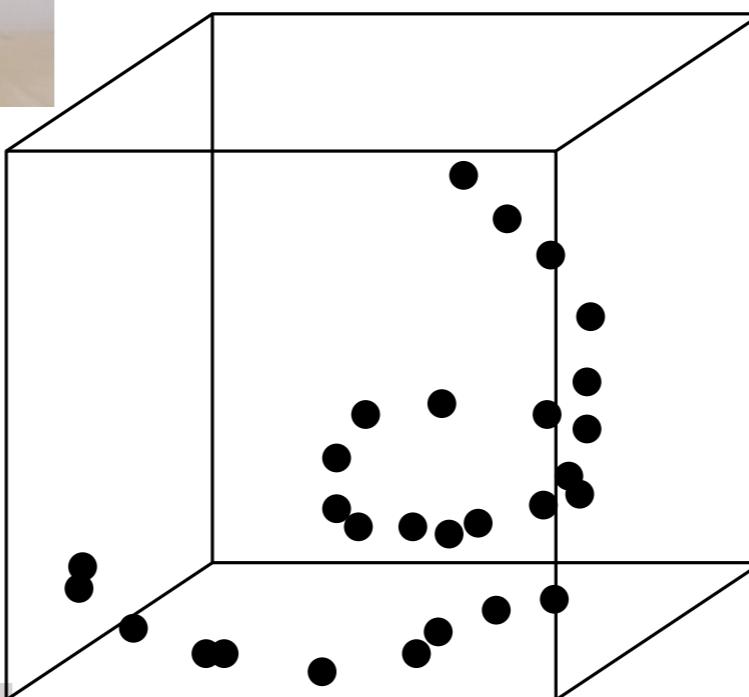


# Experiments



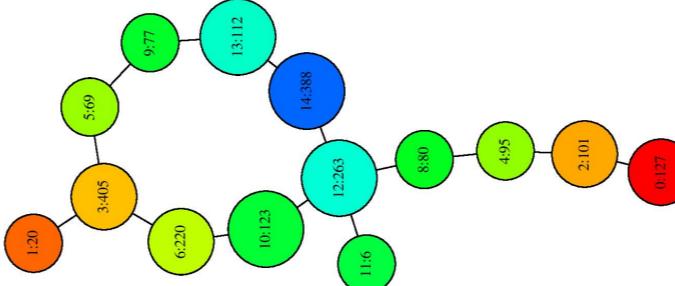
# Experiments

---

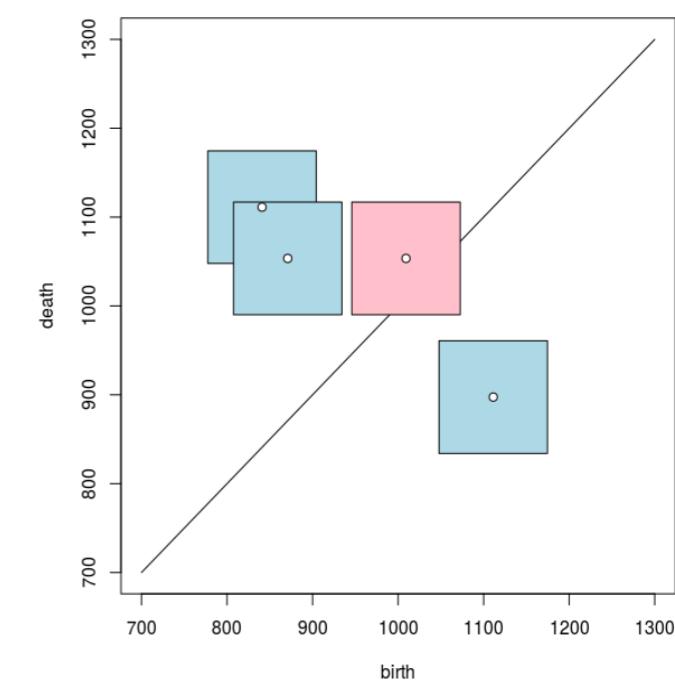
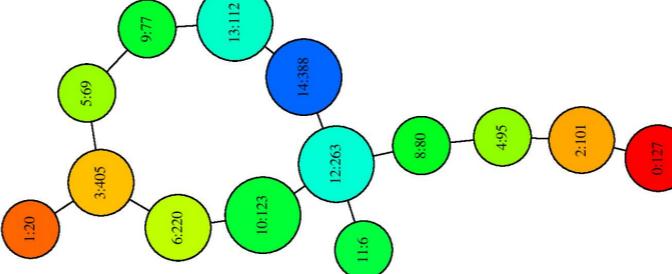


# Experiments

---



# Experiments



# Learning

---

Learning is also possible *directly on Persistence Diagrams*

Kernel methods can compensate 'bad' properties of  $d_B$  and  $d_W^p$

# Learning

---

Learning is also possible *directly on Persistence Diagrams*

Kernel methods can compensate 'bad' properties of  $d_B$  and  $d_W^p$

**Finite-dimensional embedding:** [C., Oudot, Ovsjanikov 2015]

Concatenate pairwise distances of  $Dg$  into single vector  $\Phi(Dg)$

**Thm:**  $\|\Phi(Dg) - \Phi(Dg')\|_\infty \leq 2 d_B(Dg, Dg')$

Learning is also possible *directly on Persistence Diagrams*

Kernel methods can compensate 'bad' properties of  $d_B$  and  $d_W^p$

**Finite-dimensional embedding:** [C., Oudot, Ovsjanikov 2015]

Concatenate pairwise distances of  $Dg$  into single vector  $\Phi(Dg)$

**Thm:**  $\|\Phi(Dg) - \Phi(Dg')\|_\infty \leq 2 d_B(Dg, Dg')$

**Gaussian Kernel:** [C., Cuturi, Oudot 2017]

Modify  $d_W^1$  into  $SW$  and use  $\exp\left(-\frac{SW(\cdot, \cdot)}{2\sigma^2}\right)$

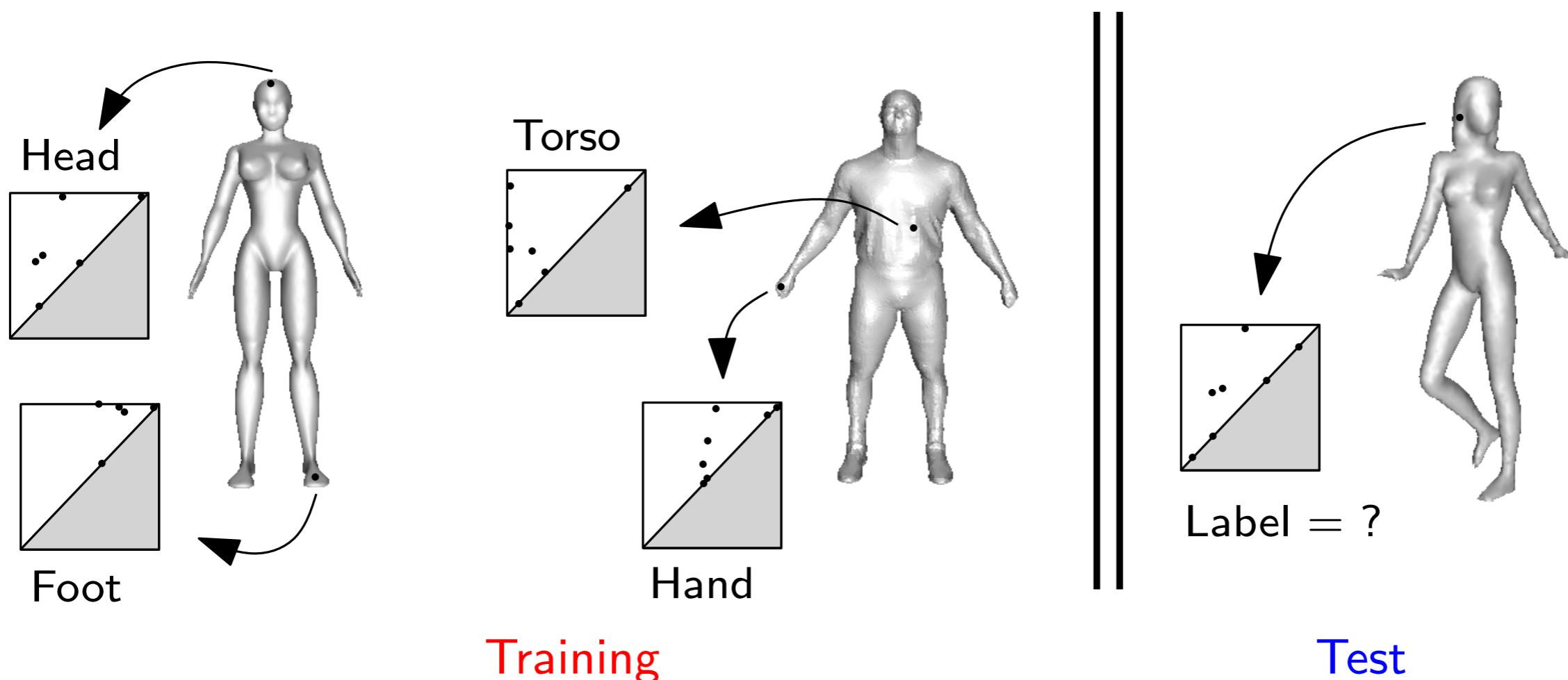
**Thm:**  $C(N) d_W^1(Dg, Dg') \leq SW(Dg, Dg') \leq 2\sqrt{2} d_W^1(Dg, Dg')$

# Application to supervised shape segmentation

**Goal:** segment 3d shapes based on examples

**Approach:**

- train a (multiclass) classifier on PDs extracted from the training shapes
- apply classifier to PDs extracted from query shape



# Application to supervised shape segmentation

**Goal:** segment 3d shapes based on examples

**Approach:**

- train a (multiclass) classifier on PDs extracted from the training shapes
- apply classifier to PDs extracted from query shape

**Accuracies (%):**

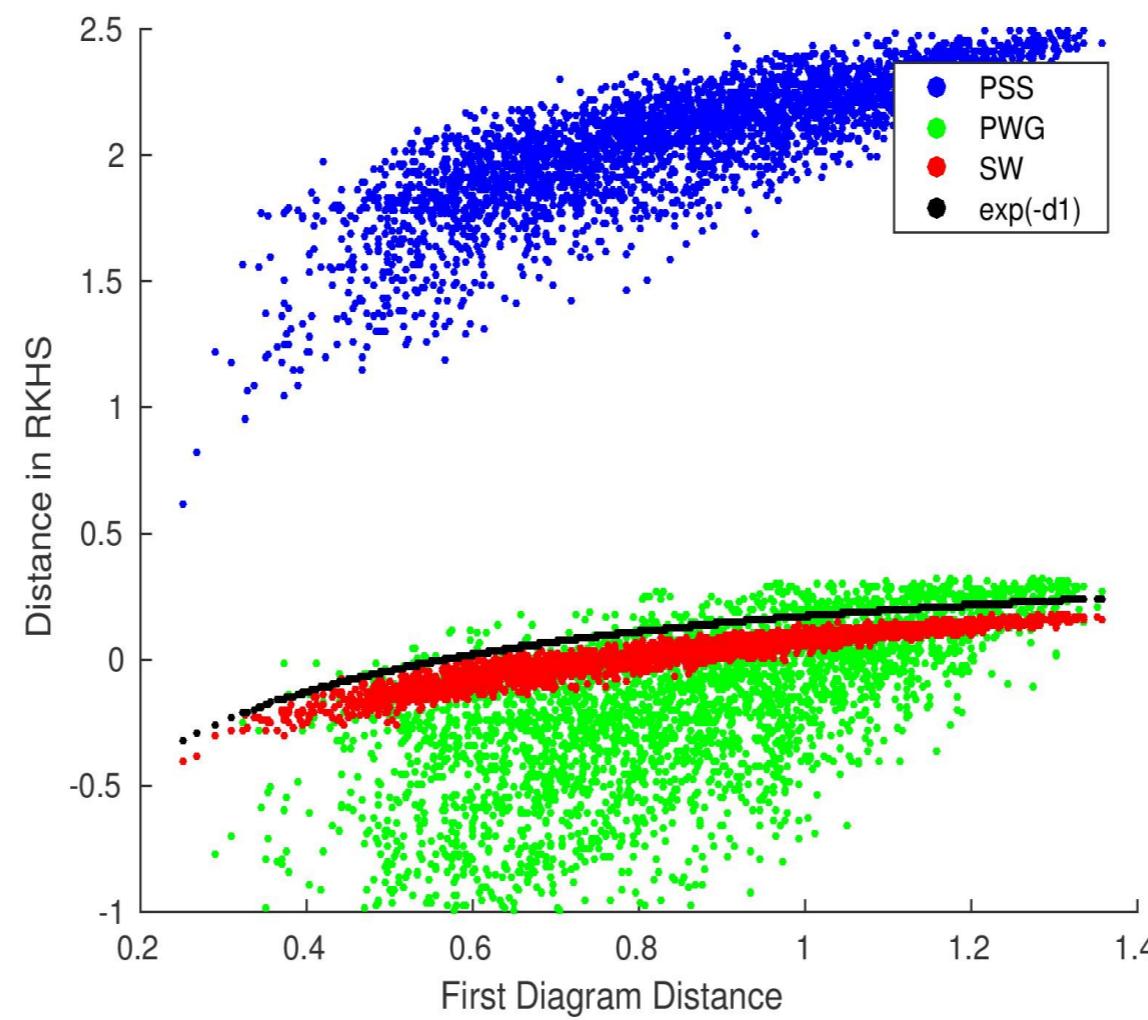
|          | TDA  | geometry | TDA + geometry |
|----------|------|----------|----------------|
| Human    | 74.0 | 78.7     | <b>88.7</b>    |
| Airplane | 72.6 | 81.3     | <b>90.7</b>    |
| Ant      | 92.3 | 90.3     | <b>98.5</b>    |
| FourLeg  | 73.0 | 74.4     | <b>84.2</b>    |
| Octopus  | 85.2 | 94.5     | <b>96.6</b>    |
| Bird     | 72.0 | 75.2     | <b>86.5</b>    |
| Fish     | 79.6 | 79.1     | <b>92.3</b>    |

# Application to supervised shape segmentation

**Goal:** segment 3d shapes based on examples

**Approach:**

- train a (multiclass) classifier on PDs extracted from the training shapes
- apply classifier to PDs extracted from query shape



# Conclusion

---

Structure and Stability of Mapper

Parameter Selection for Mapper

Kernel Methods for Persistence Diagrams

## Extensions:

Multivariate function  $f : X \rightarrow \mathbb{R}^n$

Space of Reeb graphs (curvature, barycenters, interpolation...)

Other types of statistics for Persistence Diagrams/Mappers

# Many thanks to:

---



Maks Ovsjanikov



Bertrand Michel



Marco Cuturi



Steve Oudot

ÉCOLE DOCTORALE

Sciences et technologies  
de l'information  
et de la communication (STIC)

ગુઢી **GUDHI** DAAD