

# Graph Based Analysis for Gene Segment Interactions In a Scrambled Genome

Masahico Saito  
University of South Florida

BIRS, May 7 – 12, 2017

Coauthors: Mustafa Hajij, Nataša Jonoska,  
Denys Kukushkin

Collaborator: Laura Landweber, Columbia U.

(Supported in part by NIH R01GM109459-01)

## Plan:

- 1 Representing interactions of gene segments by graphs
- 2 Examples of graphs thus obtained
- 3 Applying TDA to a data set consisting of graphs through graph properties
- 4 Outputs

# Outline

Graph Based  
Analysis for  
Gene Segment  
Interactions In  
a Scrambled  
Genome

Masahico  
Saito

Gene segment  
interactions  
and graphs

Examples

Applying TDA  
to graphs

Outputs

**1** Gene segment interactions and graphs

2 Examples

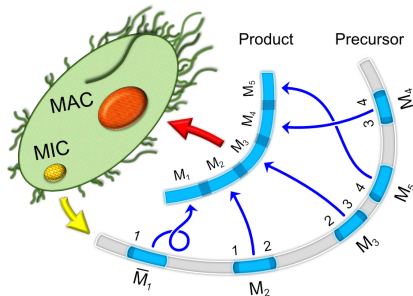
3 Applying TDA to graphs

4 Outputs

# Ciliates

*Oxytricha trifallax*, a species of ciliate, undergoes massive genome rearrangements.

Landweber Lab sequenced the whole genome of *Oxytricha trifallax*.



Micronucleus (MIC) : used for mating/conjugation

Macronucleus (MAC) : functional

Gene assembly : from MIC to MAC

MDS (MAC Destined Sequence) : segments headed to MAC

IES (Internally Eliminated Sequence) : junk DNA, eliminated

# Rearrangement of gene segments in *Oxytricha trifallax*

Graph Based  
Analysis for  
Gene Segment  
Interactions In  
a Scrambled  
Genome

Masahiko  
Saito

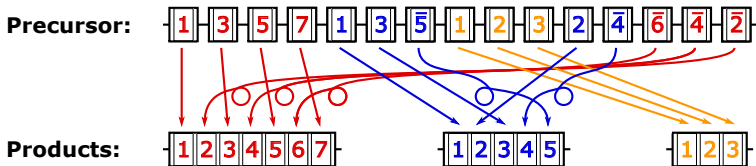
Gene segment  
interactions  
and graphs

Examples

Applying TDA  
to graphs

Outputs

MIC to MAC:



MDSs from different MAC genes interleave in the MIC contig.

# Interaction types of gene segments

Graph Based  
Analysis for  
Gene Segment  
Interactions In  
a Scrambled  
Genome

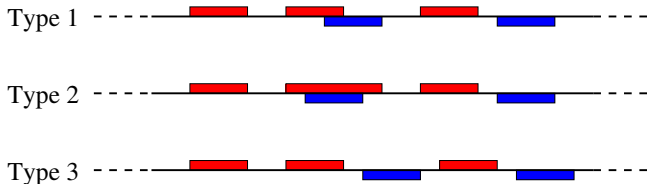
Masahiko  
Saito

Gene segment  
interactions  
and graphs

Examples

Applying TDA  
to graphs

Outputs



Three types of interactions of gene segments:

Type 1: Overlapping

Type 2 : Containment

Type 3 : Interleaving

Seven possible combination types of interaction: Overlapping and interleaving, etc.

# Representing gene interactions by graphs

Graph Based  
Analysis for  
Gene Segment  
Interactions In  
a Scrambled  
Genome

Masahico  
Saito

Gene segment  
interactions  
and graphs

Examples

Applying TDA  
to graphs

Outputs

$H$  : a MIC “contig”  $\iff G = G(H) = (V, E)$  a graph

$g$  : a MAC gene  $\iff$  vertex  $v(g)$

A (colored and directed) edge  $e$  from  $v(g_1)$  to  $v(g_2)$

$\iff g_1$  interacts with  $g_2$

Eg.  $g_1$  is contained in  $g_2 \iff$  an edge from  $v(g_1)$  to  $v(g_2)$

Colors of edges represent a combination type of interactions  
(7 colors)

(About 70 MIC chromosomes, about 283 isomorphism classes  
of graphs, about 16,000 MAC chromosomes.)

# Outline

Graph Based  
Analysis for  
Gene Segment  
Interactions In  
a Scrambled  
Genome

Masahico  
Saito

Gene segment  
interactions  
and graphs

Examples

Applying TDA  
to graphs

Outputs

1 Gene segment interactions and graphs

2 Examples

3 Applying TDA to graphs

4 Outputs



# Examples

Graph Based  
Analysis for  
Gene Segment  
Interactions In  
a Scrambled  
Genome

Masahico  
Saito

Gene segment  
interactions  
and graphs

Examples

Applying TDA  
to graphs

Outputs

## Sample graphs in the largest cluster

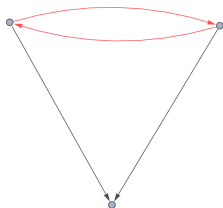


Figure :  
ctg7180000087289

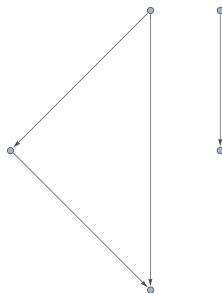


Figure :  
ctg7180000069936

# Examples

Graph Based  
Analysis for  
Gene Segment  
Interactions In  
a Scrambled  
Genome

Masahico  
Saito

Gene segment  
interactions  
and graphs

Examples

Applying TDA  
to graphs

Outputs

## Sample graphs in the largest cluster (cont.)

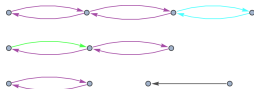


Figure :  
ctg7180000087650

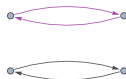
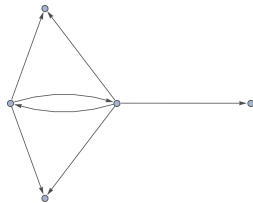


Figure :  
ctg7180000069209

# Examples

Graph Based  
Analysis for  
Gene Segment  
Interactions In  
a Scrambled  
Genome

Masahico  
Saito

Gene segment  
interactions  
and graphs

Examples

Applying TDA  
to graphs

Outputs

## Two in a small component

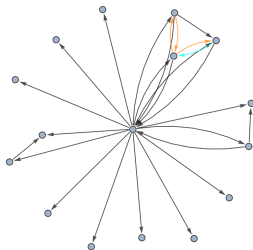


Figure :  
ctg7180000088096



Figure :  
ctg7180000067742

# Examples

Graph Based  
Analysis for  
Gene Segment  
Interactions In  
a Scrambled  
Genome

Masahico  
Saito

Gene segment  
interactions  
and graphs

Examples

Applying TDA  
to graphs

Outputs

## Singleton components

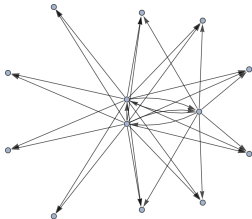


Figure :  
ctg7180000067761  
joins the large  
component at  $d = 11$

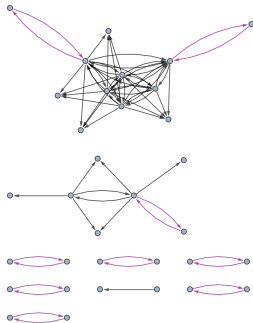


Figure :  
ctg7180000067223  
joins the large  
comonent last at  
 $d = 24$

# Outline

Graph Based  
Analysis for  
Gene Segment  
Interactions In  
a Scrambled  
Genome

Masahico  
Saito

Gene segment  
interactions  
and graphs

Examples

Applying TDA  
to graphs

Outputs

1 Gene segment interactions and graphs

2 Examples

**3 Applying TDA to graphs**

4 Outputs

# From graphs to point-cloud to filtration

Graph Based  
Analysis for  
Gene Segment  
Interactions In  
a Scrambled  
Genome

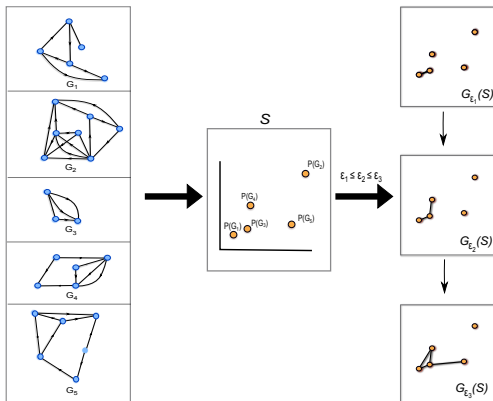
Masahico  
Saito

Gene segment  
interactions  
and graphs

Examples

Applying TDA  
to graphs

Outputs



Left : the set of graphs that represents the contigs.

Middle : Represent graphs as points in a  $c$  Euclidean space.

Right : Construct a filtration on the point-cloud.

# From graphs to point-cloud

Graph Based  
Analysis for  
Gene Segment  
Interactions In  
a Scrambled  
Genome

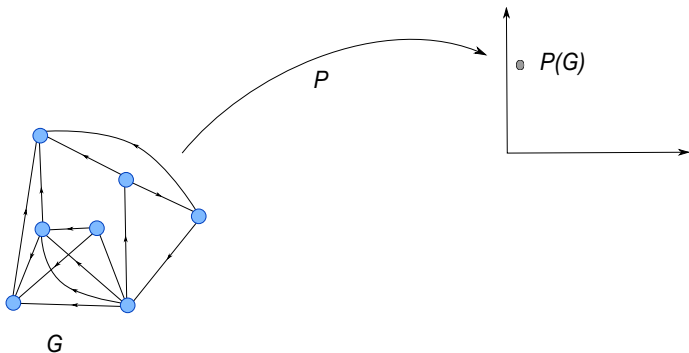
Masahico  
Saito

Gene segment  
interactions  
and graphs

Examples

Applying TDA  
to graphs

Outputs



Associate to every graph  $G$  a feature vector  $P(G)$ , a point in a Euclidean space, that represents the graph  $G$ .

# From graphs to point-cloud

The vector  $P(G)$  is defined as follows.

**Global Features Vector:**  $P_g(G) = \langle V(G), E(G), CN(G) \rangle$

$V(G)$  : # of vertices,  $E(G)$  # of edges  $P_g(G)$

$CN(G)$  : the size of the largest clique in  $G$ .

**Valence Features Vector:**  $P_v(G)$  : the valency of the vertices ordered decreasingly.

**The Clique Vector:**  $P_c(G)$  : # of cliques containing the vertex, in the same order of vertices of  $P_v(G)$ .

$d = \max(\text{valency})$

Concatenate 0s if  $|P_v(G)| < d$

$P(G) \in \mathbb{R}^{2d+3}$  : concatenation of  $P_g(G), P_v(G), P_c(G)$ .



Graph Based  
Analysis for  
Gene Segment  
Interactions In  
a Scrambled  
Genome

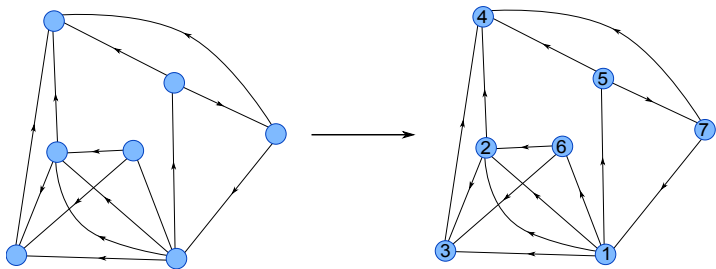
Masahico  
Saito

Gene segment  
interactions  
and graphs

Examples

Applying TDA  
to graphs

Outputs



$$P_V(G) = \langle 6, 5, 4, 4, 3, 3, 3 \rangle$$

Graph Based  
Analysis for  
Gene Segment  
Interactions In  
a Scrambled  
Genome

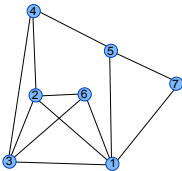
Masahico  
Saito

Gene segment  
interactions  
and graphs

Examples

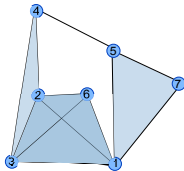
Applying TDA  
to graphs

Outputs

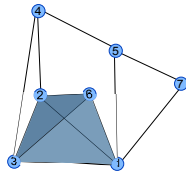


(1,1,1,1,1,1,1)

(5,4,4,3,3,3,2)



(4,3,3,1,1,3,1)



(1,1,1,0,0,1,0)

$$P_c(G) = \langle 11, 9, 9, 5, 5, 8, 4 \rangle$$

# Outline

Graph Based  
Analysis for  
Gene Segment  
Interactions In  
a Scrambled  
Genome

Masahico  
Saito

Gene segment  
interactions  
and graphs

Examples

Applying TDA  
to graphs

Outputs

1 Gene segment interactions and graphs

2 Examples

3 Applying TDA to graphs

**4** Outputs

# Barcodes of connected components

Graph Based  
Analysis for  
Gene Segment  
Interactions In  
a Scrambled  
Genome

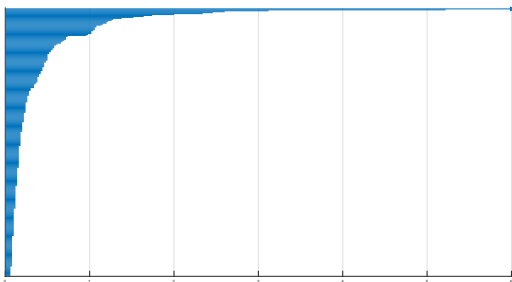
Masahico  
Saito

Gene segment  
interactions  
and graphs

Examples

Applying TDA  
to graphs

Outputs



The barcode diagram describing the birth and death of the connected components.

# Tree diagram of merging components

Graph Based  
Analysis for  
Gene Segment  
Interactions In  
a Scrambled  
Genome

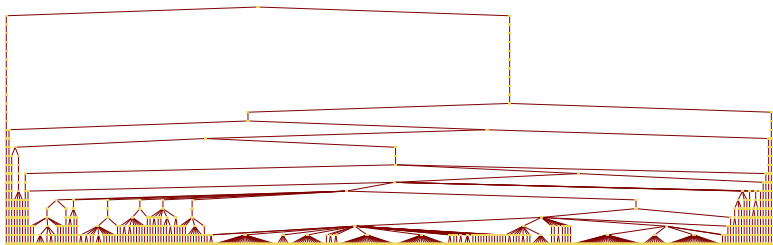
Masahico  
Saito

Gene segment  
interactions  
and graphs

Examples

Applying TDA  
to graphs

Outputs



Tree diagram representing merging components

Graph Based  
Analysis for  
Gene Segment  
Interactions In  
a Scrambled  
Genome

Masahico  
Saito

Gene segment  
interactions  
and graphs

Examples

Applying TDA  
to graphs

**Outputs**

Thank you !