# Topological Data Analysis for Cosmology & the String Landscape

#### Gary Shiu University of Wisconsin-Madison

#### String Landscape



#### Swampland

#### String Landscape

10<sup>500</sup> IIB flux vacua [Ashok-Denef-Douglas]

#### Swampland

#### String Landscape

10<sup>755</sup> F-theory vacua [Halverson, Long, Sung]

#### Swampland

#### String Landscape

10<sup>272,000</sup> F-theory vacua

[Taylor-Wang]

#### Swampland

#### String Landscape

10<sup>272,000</sup> F-theory vacua

[Taylor-Wang]

Are there hidden structures?

## **Distribution of String Vacua**



## **Distribution of Large Scale Structure**

Similar **clustering** and **void** features also appear in LSS:



## **Topological Data Analysis**

- When the space of data is huge, we cannot simply "visualize" the structure of data. We need a systematic diagnostic tool.
- Topological data analysis (TDA) is a systematic tool in applied topology to diagnose the "shape" of data.
- To turn a discrete set of data points (point cloud) into a topological space, we need a notion of *persistence*.



Vary simplicial complexes formed by the point cloud with continuous parameters (filtration parameters)

# **Topological Data Analysis**

- TDA is widely used in other fields, e.g., imaging, neuroscience, and drug design. It is well suited for machine learning.
- From the persistent homology of the point cloud, we can test e.g., the effectiveness of drugs. Similarly, we can test:



• A **selector algorithm** is often used due to the huge volume of data. We can test these algorithms on cosmological datasets and the string vacua.

# **TDA and Energy Landscapes**

 Morse theory connects topology of sublevel sets with critical points of a potential:



- TDA can be used to find critical points of energy landscapes.
- Moreover, topology change of sublevel set of energy is often associated with a phase transition.
- Persistent homology can recover known associated topology changes. [Donato, Gori, Pettini, Petri, De Nigris, Franzosi, Vaccarino]

## TDA and the String Landscape

• TDA as a tool for the landscape:



Proposed universal behavior for V > 0 and at large distances in field space [Ooguri, Palti, GS, Vafa]:

 $V(\phi) \sim e^{-a\phi}$ 

- While skepticisms have been raised regarding the existence of *explicit, controlled dS vacua* [Obied, Ooguri, Spodyneiko, Vafa];[Ooguri, Palti, GS, Vafa] (see talks of Andriot and van der Schaar), there seem to be many AdS and Minkowski vacua.
- TDA can be used to find local energy minima in protein conformation space [Haspel, Luo, González].
- Morse theory tells us that the appearance of clusters in sublevel sets is related to local minima of the energy landscape.

## Plan of this talk

- Introduce the basic concepts of topological data analysis: persistent homology, barcodes, and persistence diagrams.
- Applying TDA to constrain **primordial non-Gaussianities**.
- Back to String Data. Computing the persistent homology of string vacua to analyze their structure.
- This talk is based on several projects done in collaboration with



- "Persistent Homology and Non-Gaussianity",
   A. Cole, GS, JCAP **1803**, no. 03, 025 (2018) [arXiv: 1712.08159 [astro-ph.CO]].
- "Topological Data Analysis for the String Landscape",
   A. Cole, GS, arXiv: 1812.06960 [hep-th].

## **Simplicial Complexes**

- In  $\mathbb{R}^3$ , simplices are vertices, edges, triangles, and tetrahedra
- Simplicial complexes are collections of simplices that are:
  - Closed under intersection of simplices
  - Closed under taking faces of simplices
- Combinatorial representations easy calculations for computers



Source: Wikipedia, "Simplicial Complex"

# Simplicial Homology

- Given a simplicial complex, define a boundary operator  $\partial_p$ that maps p-simplices to (p-1)-simplices
  - We want to count independent p-cycles (i.e. p-loops) that are not boundaries of higher-dimensional objects
- Group theoretic:  $Z_p = \ker \partial_p$  ,  $B_p = \operatorname{im}\,\partial_{p+1}$ ,  $\checkmark$

$$H_p \equiv Z_p/B_p$$

- Betti numbers:  $eta_p \equiv \mathrm{rank} H_p$ 
  - 0-th Betti number is number of connected components
  - p-th Betti number is number of independent p-loops
- In practice, homology calculation is a matrix reduction

$$\beta_{0} = 1$$

$$\beta_{1} = 1$$

$$\gamma_{s.}$$

$$\beta_{0} = 1$$

$$\beta_{0} = 1$$

$$\beta_{1} = 0$$

### Persistence

- How to choose simplicial representation of our data?
- Persistent homology: vary simplicial representation  $\Sigma_{\nu}$  of data with some filtration parameter  $\nu$  such that

$$\nu_1 \leq \nu_2 \implies \Sigma_{\nu_1} \subseteq \Sigma_{\nu_2}$$

- Track each distinct feature's lifetime (birth and death)
- Intuition: "real" topological features *persist*, short-lived features are noise
- Procedure is stable against perturbations to data [Cohen-Steiner 2005]



•





























# Visualizing Persistent Homology

#### Barcodes:

- Each horizontal line represents an independent cycle contributing to a particular Betti number (i.e. a connected component, loop, void...)
- Lines start at birth and end at death
- To calculate Betti number, make vertical slice and count intersections

#### Persistence diagrams:

- Scatter plot, each point representing an independent cycle
- Calculate Betti number by counting "living" cycles



#### Persistence diagrams contain more information than Betti number curves!



# Applying TDA to Cosmology

# Inflation

#### [Starobinsky];[Guth];[Linde];[Albrecht, Steinhardt];...

- Period of accelerated expansion in early universe
  - Solves flatness, horizon, and monopole problems
  - Predicts nearly scale-invariant,
     Gaussian curvature fluctuations
    - Source anisotropies in CMB, inhomogeneities in LSS
- A myriad of models. Taxonomy done mostly through their observables (n<sub>s</sub>, r)





## Anisotropies

 The lowest order correlation we can extract from the anisotropies is the power spectrum

$$\left\langle 0 \left| \hat{\mathcal{R}}_{\mathbf{k_1}} \hat{\mathcal{R}}_{\mathbf{k_2}} \right| 0 \right\rangle = (2\pi)^3 P_{\mathcal{R}}(k_1) \delta(\mathbf{k_1} + \mathbf{k_2}) \qquad \Delta_{\mathcal{R}}^2 = \left( \frac{k^3}{2\pi^2} \right) P_{\mathcal{R}}^2 \propto k^{n_s - 1}$$

- For a Gaussian theory, the power spectrum dictates all higher-pt correlations.
- However, the inflationary fluctuations are not perfectly Gaussian.
- The leading **non-Gaussianity** is the **bispectrum**:

$$\langle 0 | \hat{\mathcal{R}}_{\mathbf{k_1}} \hat{\mathcal{R}}_{\mathbf{k_2}} \hat{\mathcal{R}}_{\mathbf{k_3}} | 0 \rangle = (2\pi)^3 \, \delta^3 (\mathbf{k_1} + \mathbf{k_2} + \mathbf{k_3}) F(\mathbf{k_1}, \mathbf{k_2}, \mathbf{k_3})$$

- Scaling and symmetries imply that F(k<sub>1</sub>, k<sub>2</sub>, k<sub>3</sub>) is fixed by an overall size ~ f<sub>NL</sub> and its ''shape" F(1, k<sub>2</sub>/k<sub>1</sub>, k<sub>3</sub>/k<sub>1</sub>).
- More **powerful discriminator** of inflationary models.

## **Non-Gaussianities**

- The bispectrum for single field slow-roll inflation was computed in [Maldacena, '02];[Acquaviva et al, '02]; its size is f<sub>NL</sub> ~ O(ε,η):
- The bispectrum for general single field inflation was found to be parametrized by 5 parameters [Chen, Huang, Kachru, GS, '06]:



 There is also an "orthogonal shape" but it "looks" qualitatively like the equilateral shape (*challenge for machine learning?*).

## Non-Gaussianities

 More complicated models which involve non-standard initial conditions, features in potential (e.g. axion monodromy), or multiple fields or quasi-single field can give rise to more shapes:



- Like scattering amplitudes in particle physics, non-Gaussianties can reveal interactions governing inflation: *cosmological collider.*
- In collider physics: use *different strategies* for different particles.

# Measuring Non-Gaussianity

 Harmonic space: fits with <u>templates</u> of bispectrum, trispectrum, etc. One can define a "cosine" between distributions:

$$\cos(F_1, F_2) = \frac{F_1 \cdot F_2}{(F_1 \cdot F_1)^{1/2} (F_2 \cdot F_2)^{1/2}}$$

Some shapes are harder to find, e.g.,

Resonant shape (axion monodromy)



- Geometrical/topological: Minkowski functionals (for CMB: area fraction, length of boundaries, and genus of excursion sets)
- Current bound on non-Gaussianity (Planck '15):

$$f_{NL}^{local} = 2.5 \pm 5.7$$
  $f_{NL}^{equil} = -16 \pm 70$ 

(Hotter points are deeper red)





Many distinct components, no loops





 $\nu = 0$ 

Many loops, fewer distinct components

(Sublevel set in black)





One connected component, many loops have been filled in

(Sublevel set in black)



# Sensitivity to Non-Gaussianity

- We first carried out TDA for local NG and with low-resolution maps (/ max~ 1024) as a warmup, more in our pipeline.
- Bined the PDs for different f<sub>NL</sub>, & computed the likelihood function:



- More sensitive statistic than Minkowski functionals or Betti number curves, PDs strengthens topological analysis significantly.
- N.B. Lower resolution maps used here compared to Planck's.
- Other subtle shapes (e.g., resonant non-Gaussianity).

# Applying TDA to String Vacua

## **TDA for String Vacua**



## Toy Example: IIB Flux Vacua on Rigid CY

• **Superpotential**  $W = A\tau + B$  where the flux quanta:

$$A = -h_1 - ih_2, \quad B = f_1 + if_2, \quad h_1, h_2, f_1, f_2 \in \mathbb{Z}$$

subject to **tadpole cancellation**:  $N_{\text{flux}} = f_1 h_2 - h_1 f_2 \leq L_{\text{max}}$ 

• Vacua are mapped to the **fundamental domain** using SL(2,Z).





-0.5

 $\tau$ -plane

0.5

### **Persistence** Pairing



- In general, not possible to visualize a *higher dim.* data space.
- For example, flux vacua of IIB orientifold on CY hypersurface:

$$\sum_{i=1}^{4} x_i^8 + 4x_0^2 - 8\psi x_0 x_1 x_2 x_3 x_4 = 0, \quad x_i \in \mathbf{WP}^4_{1,1,1,1,4}$$

has  $h^{1,1} = 1$ ,  $h^{2,1} = 149$  and discrete symmetry  $\Gamma = Z_8^2 \times Z_2$ . The only  $\Gamma$ -invariant moduli: complex structure modulus  $\psi$  & axio-dilaton  $\tau$ .





- To identify cluster, apply density cutoff (excises cluster, results in identifiable void)
- Does this cluster/void exist in the full four-dimensional space? (Might not if clustering correlates with structure in axiodilaton.) Are there significant higher dimensional features?
  - These questions can be answered with persistent homology

- To identify cluster, apply density cutoff (excises cluster, results in identifiable void)
- We found a long-lived 1cycle in the full four-dim.
   space and only observe short-lived higher dimension features (sampling noise)



- To identify cluster, apply density cutoff (excises cluster, results in identifiable void)
- We found a long-lived 1cycle in the full four-dim.
   space and only observe short-lived higher dimension features (sampling noise)



long-lived 1-cycle With Density Filter Vdeath To identify cluster, apply density cutoff (excises 0.15 cluster, results in identifiable blue:0-cycles void) orange:1-cycles 0.10 green:2-cycles We found a long-lived 1red:3-cycles cycle in the full four-dim. space and only observe short-lived higher dimension features (sampling noise)

0.02

0.04

0.06

D.08

0.10

Work in progress [Cole,GS]

0.12

Vbirth

## Flux Vacua on Symmetric T<sup>6</sup>

- Factorizable  $T^6 = (T^2)^3$  with equal complex structure  $z_1 = z_2 = z_3 = z_2$
- Two complex moduli: complex structure modulus z and axio-dilaton  $\tau$ .
- Number-theoretical methods were used to find distributions of vacua with W=0 and with discrete symmetries [DeWolfe, Giryavets, Kachru, Taylor]

#### Generic vacua on z-plane





 How do "cuts" like restricting to W=0 vacua (e.g., discrete R-symmetry, motivated by [Nelson, Seiberg]) change the topology of distribution?

## Flux Vacua on Symmetric T<sup>6</sup>

Reasonable expectation:	generic moduli dist topology	"topology" of cut	resulting topology
	trivial	trivial	trivial
	nontrivial	trivial	reduced complexity
	trivial	nontrivial	nontrivial
	nontrivial	nontrivial	complicated

• Comparing persistent homology:



 W=0 cut adds complexity! Long-lived higher dimensional topological features differs from that for generic vacua.

# Sampling in TDA

- We can't realistically include all  $10^{500}$  vacua as vertices
- Can sample the topology via the witness complex:
  - From the entire point cloud Z, choose a *landmark set L* as the complex's vertices. Often chosen randomly or via sequential maxmin algorithm
  - Let  $m_k(z)$  be the distance from some  $z \in Z$  to the (k+1)-nearest landmark point. Then, given filtration parameter  $\mathcal{V}$ , the simplex  $[l_0 l_1 \dots l_k]$  is included in the witness complex if  $\max \{d(l_0, z), d(l_1, z), \dots, d(l_k, z)\} \le \nu + m_k(z)$









Conclusions

## Conclusions

- Applications of TDA to cosmological datasets and string vacua.
- Persistence diagrams strengthen constraints on local non-Gaussianities, and potentially other shapes & other observables.
- Techniques we developed can be applied to analyze the structure of string vacua. We performed initial study of simple flux vacua.
- Next step is to examine the topology of string vacua point clouds with desired features, supplementing earlier work on statistics:
  - Enhanced symmetries [DeWolfe, Giryavets, Kachru, Taylor], ...
  - Particle physics features [Marchesano, GS, Wang];[Dienes];[Gmeiner, Blumenhagen, Honecker, Lust, Weigand], [Douglas, Taylor], ...
- Topology of Energy Landscape of String theory?
- String Landscape vs the Swampland?

## Conclusions



- Applications of TDA to cosmological datasets and string vacua.
- Persistence diagrams strengthen constraints on local non-Gaussianities, and potentially other shapes & other observables.
- Techniques we developed can be applied to analyze the structure of string vacua. We performed initial study of simple flux vacua.
- Next step is to examine the topology of string vacua point clouds with desired features, supplementing earlier work on statistics:
  - Enhanced symmetries [DeWolfe, Giryavets, Kachru, Taylor], ...
  - Particle physics features [Marchesano, GS, Wang];[Dienes];[Gmeiner, Blumenhagen, Honecker, Lust, Weigand], [Douglas, Taylor], ...
- Topology of Energy Landscape of String theory?
- String Landscape vs the Swampland?