

GraSPI

Graph-based Structure-Property Identifier

Last updated on November 2020

Motivation: at molecular level a large pool of descriptors/features exists

Dragon molecular descriptor list: 51 pages with 4885 descriptors

List of molecular descriptors calculated by Dragon

ID	Name	Description	Block
1	MW	molecular weight	Constitutional indices
2	AMW	average molecular weight	Constitutional indices
3	Sv	sum of atomic van der Waals volumes (scaled on Carbon atom)	Constitutional indices
4	Se	sum of atomic Sanderson electronegativities (scaled on Carbon atom)	Constitutional indices
44	nCIC	number of rings (cyclomatic number)	Ring descriptors
45	nCIR	number of circuits	Ring descriptors
46	TRS	total ring size	Ring descriptors
47	Rperim	ring perimeter	Ring descriptors
94	DBI	Dragon branching index	Topological indices
95	SNar	Narumi simple topological index (log function)	Topological indices
96	HNar	Narumi harmonic topological index	Topological indices
97	GNar	Narumi geometric topological index	Topological indices
98	Xt	total structure connectivity index	Topological indices
99	Dz	Pogliani index	Topological indices
100	Ram	ramification index	Topological indices
101	BLI	Kier benzene-likeness index	Topological indices
173	MPC04	molecular path count of order 4	Walk and path counts
174	MPC05	molecular path count of order 5	Walk and path counts
175	MPC06	molecular path count of order 6	Walk and path counts
176	MPC07	molecular path count of order 7	Walk and path counts
197	X0	connectivity index of order 0	Connectivity indices
198	X1	connectivity index of order 1 (Randic connectivity index)	Connectivity indices
199	X2	connectivity index of order 2	Connectivity indices
200	X3	connectivity index of order 3	Connectivity indices

Do we need so many?

- Data automation lowers energy barriers

Can we find the smallest set of features

- Feature engineering and data analytics enables the reduction

Challenge: lack of large pool of descriptors @ microstructure level

GraSPI aims to fill the gap

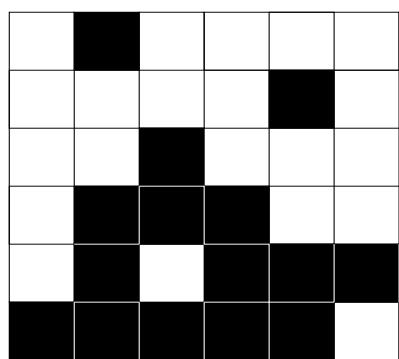
- GraSPI annotates microstructure with the set of descriptors/features.
- Given a segmented micrograph (black and white image), the set of descriptors is calculated.
- The set of descriptors captures size, shape and topological characteristics of the morphology.
- The set of descriptors has been motivated by the application for organic solar cells but can be generalized for other applications
- To enable the generalization, we plan to group the descriptors into mathematical descriptors and physics-encoded descriptors.

What problem GraSPI solves?

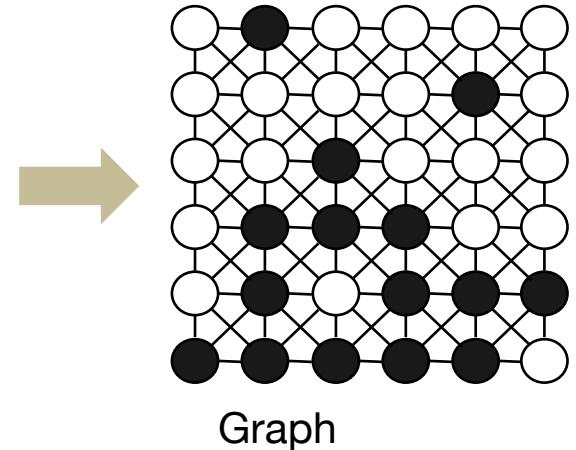
- As of now, given the segmented image, compute set of descriptors (scalars, histograms) and physics-weighted descriptors (scalars and histograms)
- In the future, expand the set of descriptors to include more descriptors, expand the representation to include the skeletons (topological features), expand the physics-based weighting functions, add representation-specific distance measures (beyond the Euclidean distance) and couple it with dimensionality reduction techniques.

We propose a two-level representation

Level 1: voxel-based representation

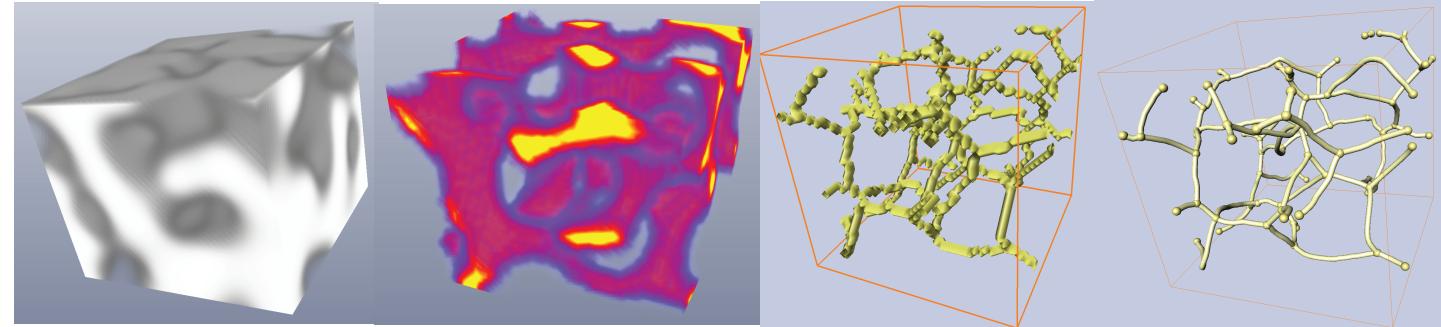


Microstructure



Graph

Level 2: topology preserving representation



Microstructure

Distance map

Medial axis

Graph

Steps involved:

- For each voxel in the microstructure image create vertex in graph
- Connect vertices in the graph following neighborhood of the image (first and second order neighbors)

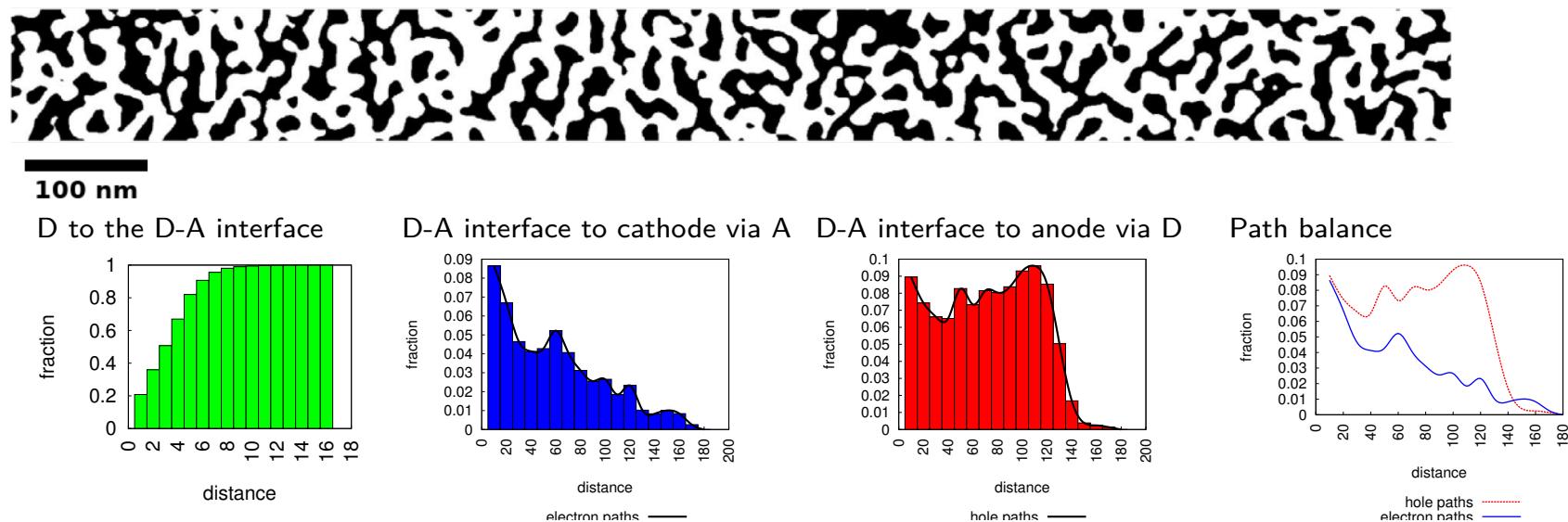
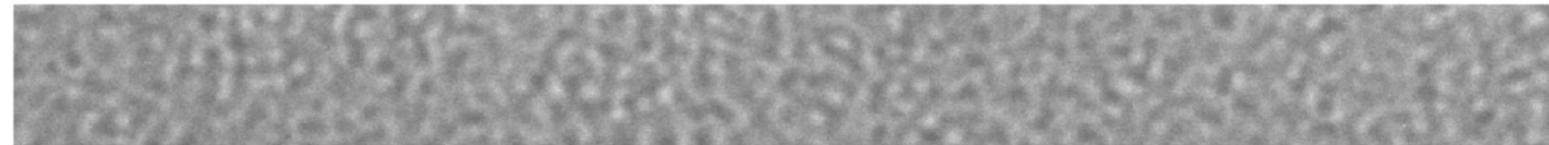
Steps involved:

- calculate distance field to the interface
- Remove voxel based on the distance field to get medial axis
- Convert the medial axis into graph representation

GraSPI: Graph-based Structure descriptor Interrogator

From micrograph to set of statistics and descriptors

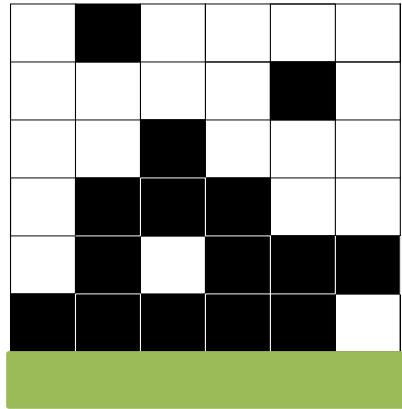
Micrograph



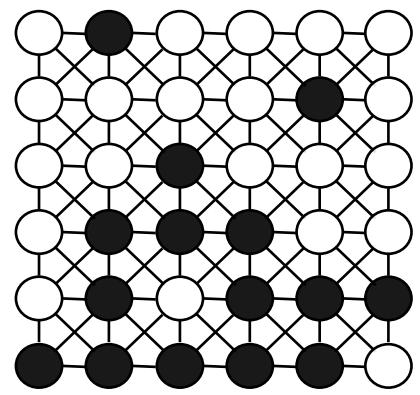
Set of descriptors
/features

- 1 Fraction of light absorbing material: **0.583**.
- 2 Fraction of photoactive material whose distance to the interface is within exciton diffusion length ($d < 10\text{nm}$): **0.994**.
- 3 Fraction of donor domains connected to anode: **0.97**.
- 4 Fraction of acceptor domains connected to cathode: **0.54**.
- 5 Distance from interface to anode via donor: $\mu = 67.82 \text{ nm} \sigma = 38.45 \text{ nm}$.
- 6 Distance from interface to cathode via acceptor: $\mu = 53.61 \text{ nm} \sigma = 39.38 \text{ nm}$.

Graph-based representation is flexible to define and compute wide range of descriptors

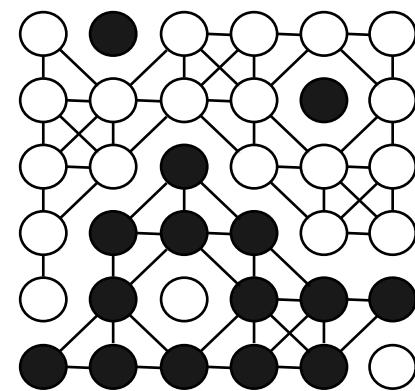
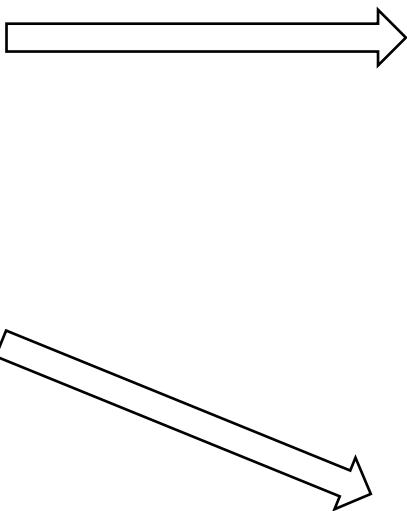


Microstructure



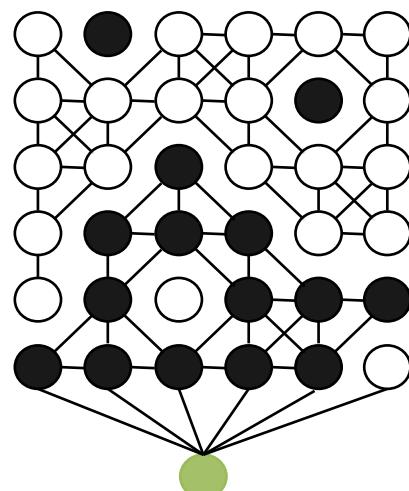
Graph= $\{V, E\}$

Filter graph



Query filtered graph
Using standard
graph-based
algorithms

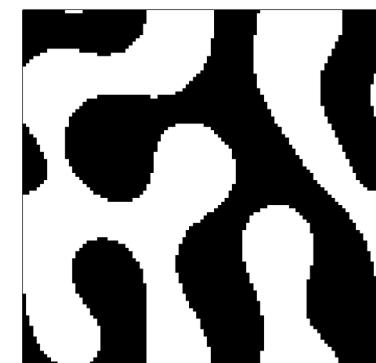
- connected component algorithm
- Shortest path



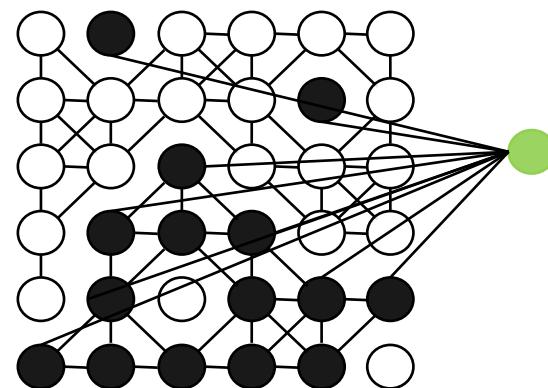
Steps involved:

- For each voxel in the microstructure image create vertex in graph
- Connect vertices in the graph following neighborhood of the image (first and second order neighbors)
- Filter graph and query the graph

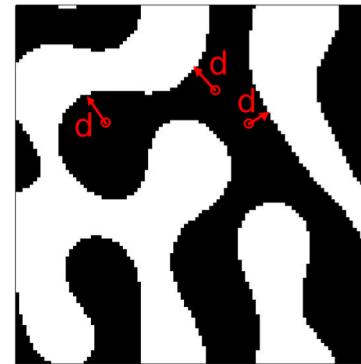
Graph-based representation is flexible to define and compute math- and physics-based descriptors



Represent microstructure as a graph and add meta vertices

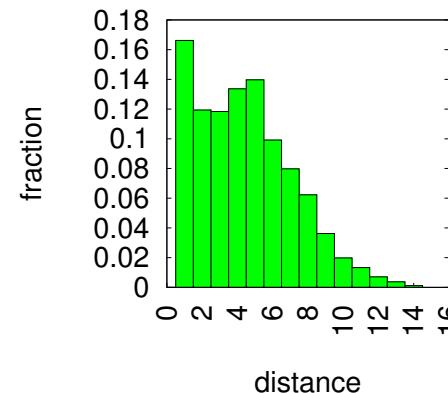


Compute all distances from black pixels to the interface (green meta vertex)



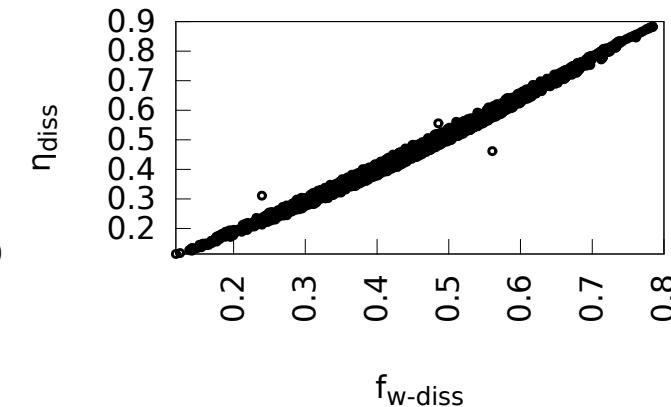
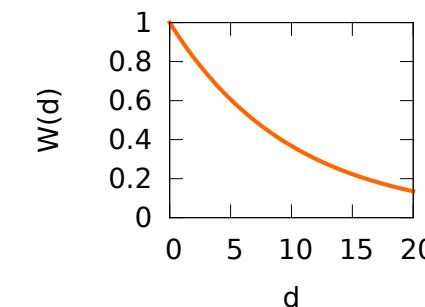
Filter graph
Make graph query:
Find shortest paths from green vertex to all black vertices

Compute histograms



Simple post processing
Make queries to extract more specific descriptors

Add the weighting function that capture some aspect of the physics

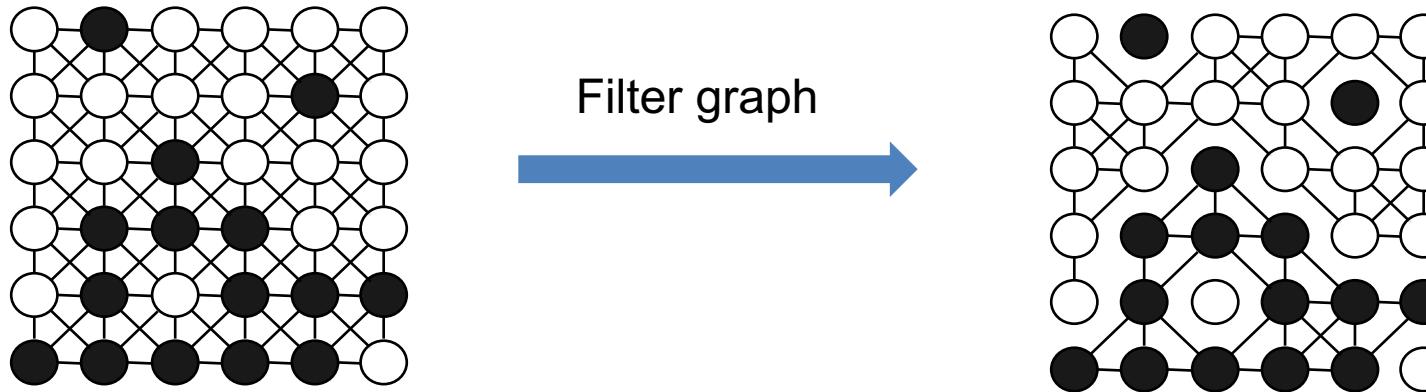


This physics-based descriptor is meaningful and serves as a surrogate for PDE-based model.
It captures non-local effects.

How it works

- Represent morphology as graph
- Add meta vertices to facilitate the computations of the descriptor (top/bottom surface, interface)
- Learn how to filter the graph to facilitate graph queries for various features
- Execute basic algorithms: e.g., compute connected components, compute the shortest paths

Descriptor=graph filtering (meta vertices)+graph algorithms



- ① Construct undirected graph $G = (V, E)$, where V is a set of vertices and E is a set of edges.
- ② Define function $c : V \rightarrow \{0, \dots, 3\}$, which assigns a color to each vertex in V , where:
 - 0 represents electron-donor material (black),
 - 1 represents electron-accepting material (white),
 - 2 represents cathode (blue),
 - 3 represents anode (red).
 - 4 represents interface (green).
- ③ Construct $G' = (V, E')$, where $E' = \{e = (u, v) \in E \mid c(u) = c(v)\}$, i.e. E' is a set of edges connecting vertices of the same color.
- ④ Identify connected components in G' .

```
1 void DetermineConnectedComponents( gt::graph_t* G, const std::vector<COLOR>& color,
2 std::vector<int>& components){
3     edge_same_color_predicate p(*G, color);
4     boost::filtered_graph<gt::graph_t, edge_same_color_predicate> FG(*G, p);
5     boost::connected_components(FG, &components[0]);
6 }

1 class edge_same_color_predicate {
2 public:
3     edge_same_color_predicate() : G_(0), color_(0) { }
4     edge_same_color_predicate(const gt::graph_t& G, const std::vector<COLOR>& color)
5         : G_(&G), color_(&color) { }
6     bool operator()(const gt::edge_t& e) const {
7         return (*color_)[boost::source(e, *G_)] == (*color_)[boost::target(e, *G_)];
8     }
9 private:
10    const gt::graph_t* G_;
11    const std::vector<COLOR>* color_;
12};
```

How to use it?

- C++: command line program
(input morphology using one of two formats: array or graph format)
- Python: cythonized grapi
(input morphology as numpy array)

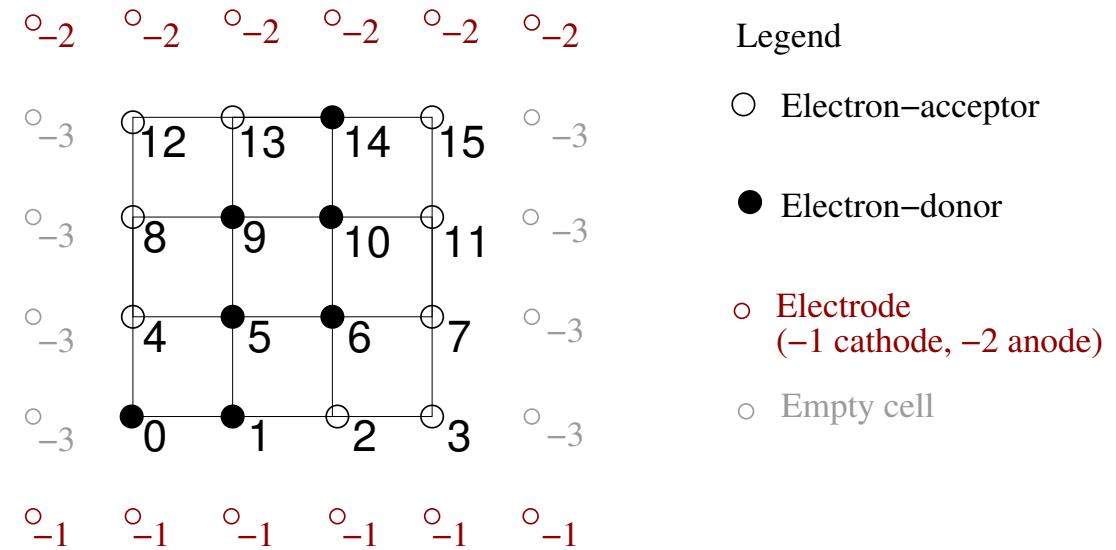
Inputs, useful scripts, example, outputs

Input: two accepted formats

```
4 4 1  
1 1 0 0  
0 1 1 0  
0 1 1 0  
0 0 1 0
```

Array:
nx ny nz
[array of labels: size nx*ny*nz]

We use a row-wise order to store the matrix of pixels/vertices.



Graph:
nVertices
[indexOfVertex labelOfVertex (triplet: indexOnNeighbor weight/distance orderOfNeigbor)] nVertices Lines

```
16  
0 1  4 1 f 5 1.42 s 1 1 f -1 1 f  
1 1  5 1 f 6 1.42 s 2 1 f -1 1 f 0 1 f 4 1.42 s  
2 0  6 1 f 7 1.42 s 3 1 f -1 1 f 1 1 f 5 1.42 s  
...  
14 1 -2 1 f 15 1 f 11 1.42 s 10 1 f 9 1.42 s 13 1 f  
15 0 -2 1 f 11 1 f 10 1.42 s 14 1 f  
16 0 0 s 1 0 s 2 0 s 3 0 s  
17 15 0 s 14 0 s 13 0 s 12 0 s
```

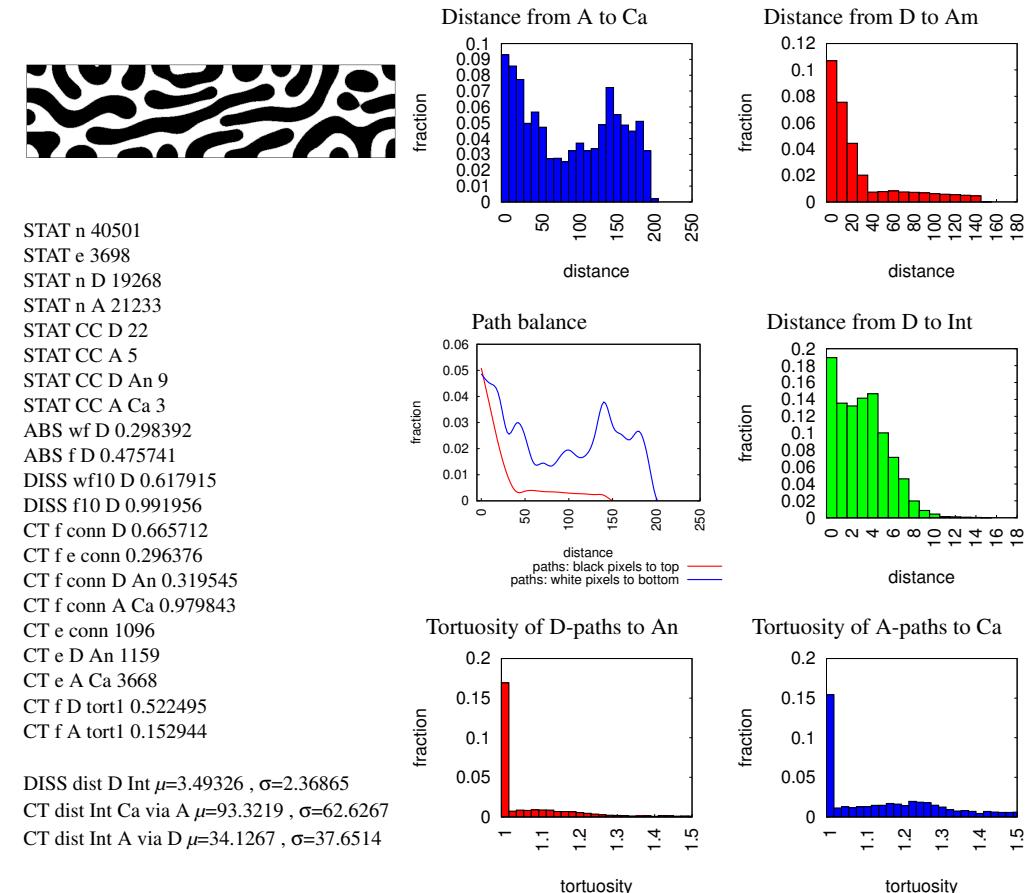
This is a custom format of adjacency list that additionally encodes information about the weights and order of the edge

Example of usage

1. Go to folder folder
`examples/2phaseMorphologies/thinFilm`
Input: set of morphologies in the array format
2. Execute script `RunTests.sh`
Outcome: set of descriptors in log file and set of files with distances
3. Execute script `PrepareReport.sh`
Outcome: pdf with the report.



1 Morphology: data_0.528_3.8_000160



Output

GraSPI outputs the descriptors to the standard output that can be redirected to another stream (e.g., log files). It also determines several sets of shortest distances. The corresponding distances/path lengths are saved following files

- DistancesGreenToRedViaBlack.txt
- DistancesGreenToBlueViaWhite.txt
- DistancesBlackToRed.txt
- DistancesWhiteToBlue.txt
- DistancesBlackToGreen.txt
- TortuosityBlackToRed.txt
- TortuosityWhiteToBlue.txt

The names of the file indicates the conditions used to filter the graph. For example file DistancesBlackToGreen.txt stores all the shortest distances between any donor/black/0 voxel and the interface/green.

Useful scripts

GraSPI has been set up to execute the high throughput analysis.

In the folder examples, two scripts have been included: RunTests.sh and PrepareReport.sh (both in bash – the command language)

RunTests.sh shows how to execute GraSPI for all input files

PrepareReport.sh shows how to generate histograms and aggregate descriptors with histogram into one report. [Dependencies: latex, gnuplot]

Computational complexity

GraSPI computes the set of descriptors at the low cost:

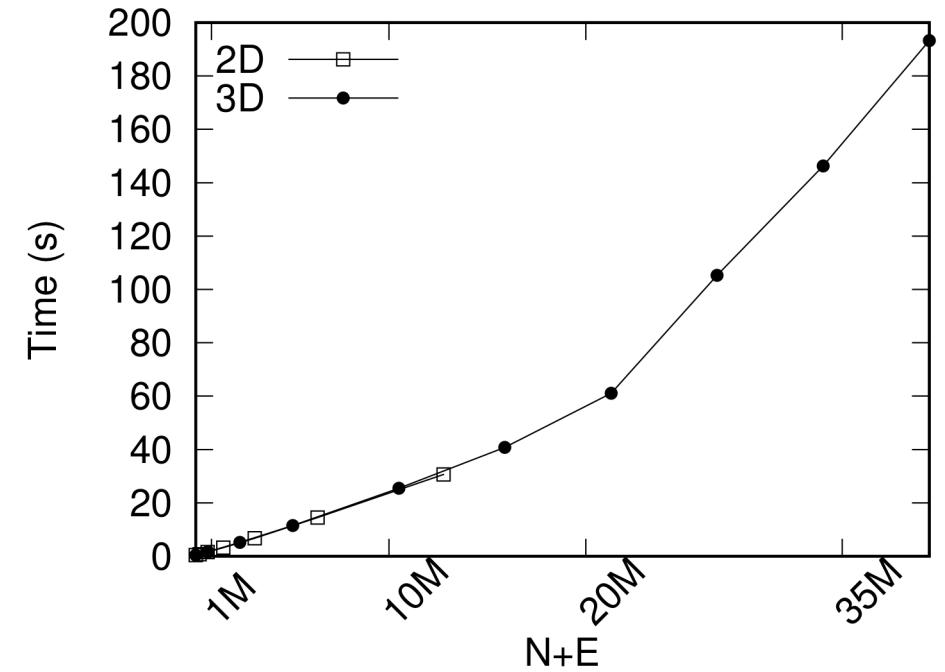
2D examples – order of seconds

3D examples – order of minutes

Note linear complexity of the execution time

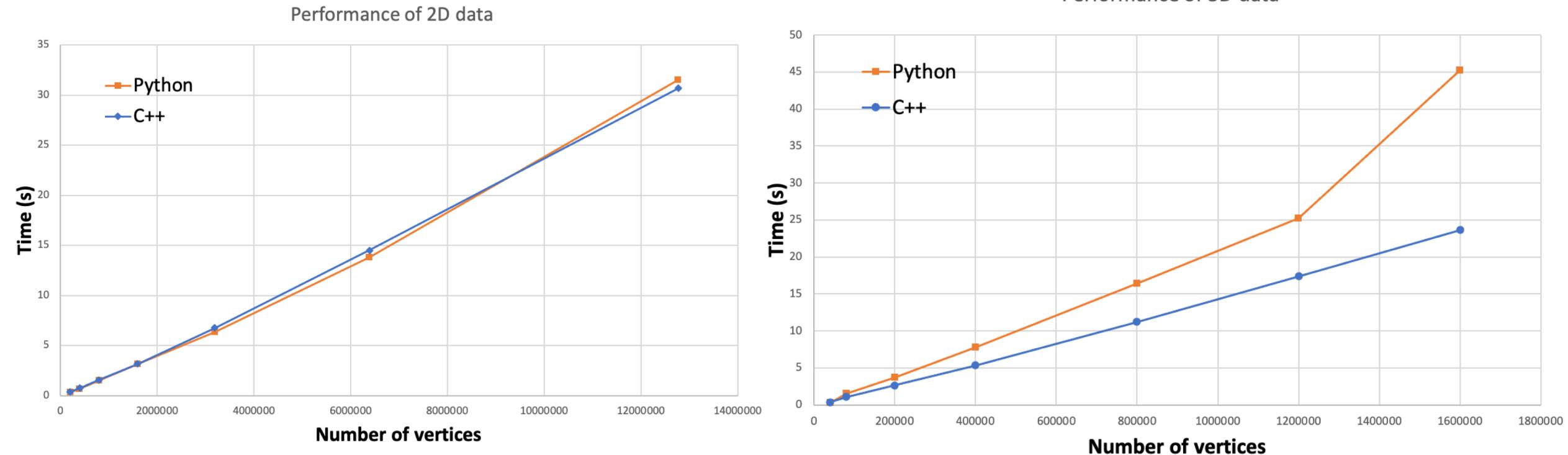
The largest problem solved on the desktop machine (35M) under 3 minutes

Technical details: Coded in C/C++ using boost graph library



N - number of vertices
E - number of edges

Computational complexity: comparison between c++ and python



Note a small overhead due to the cythonization

Resources

See our paper for the definition of descriptors used

<https://arxiv.org/abs/1106.3536v2>

How to cite this package:

Wodo, Olga, Srikanta Tirthapura, Sumit Chaudhary, and Baskar Ganapathysubramanian. "A graph-based formulation for computational characterization of bulk heterojunction morphology." *Organic Electronics* 13, no. 6 (2012): 1105-1113