




# Graph Layout by Path-Guided Stochastic Gradient Descent

S. Heumos<sup>1\*</sup> , A. Guarracino<sup>2\*</sup> , and E. Garrison<sup>3,4</sup> 

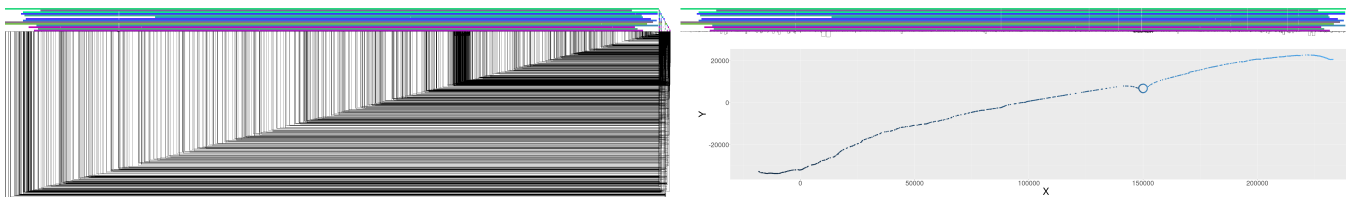
<sup>1</sup>Quantitative Biology Center (QBiC) Tübingen, University of Tübingen, Tübingen, Germany

<sup>2</sup>University of Rome Tor Vergata, Via della Ricerca Scientifica 1, Rome, Italy

<sup>3</sup>Genomics Institute, University of California Santa Cruz, Santa Cruz, CA, USA

<sup>4</sup>Biomolecular Engineering and Bioinformatics, University of California Santa Cruz, Santa Cruz, CA, USA

\*Contributed equally



**Figure 1:** Left: 1D visualization of a pangenome graph built from raw sets of alignments using the *edyeet* aligner and the variation graph inducer *seqwish*. The colored bars represent the binned, linearized renderings of the embedded paths versus the pangenome sequence. The black lines under the paths represent the topology of the graph. Right top: the same 1D visualization, but of the sorted graph by applying our algorithm in one dimension. Right bottom: 2D visualization of the 1D-sorted graph using our algorithm in two dimensions. Each dot represents a node. The nodes' *x*-coordinates are on the *x*-axis and the *y*-coordinates are on the *y*-axis, respectively.

## Abstract

Pangenome graphs built from raw sets of alignments may have complex structures generated by common patterns of genome variation. These structures can introduce difficulty in downstream analyses, visualization, mapping, and interpretation. Graph sorting aims to find the best node order for a 1D and 2D layout to simplify these complex regions. Pangenome graphs embed pangenomic sequences as paths in the graph, but to our knowledge, no algorithm takes into account this biological information in the sorting. Moreover, existing 2D layout methods struggle to deal with large graphs. For these reasons, we present a new layout algorithm that orders the nodes of a pangenome graph using a path-guided stochastic gradient descent (SGD) approach. The methodology is inspired by Zheng *et al.*, and is applicable in 1D and in 2D. In our implementation, the algorithm moves a single pair of nodes at a time, optimizing the disparity between the layout distance of a node pair and the actual nucleotide distance of a path traversing these nodes: 1. The first node of a pair is a uniform path position pick from all nodes. 2. The second node of a pair is sampled from the same path following a Zipfian distribution. 3. The path nucleotide distance of the nodes in the pair guides the actual layout distance update of these nodes. The magnitude of the update depends on the current learning rate of the SGD.

Figure 1 (left) shows in 1D an unsorted pangenome graph of a human variable number repeat region. The presence of many and long links which connect distant regions of the pangenome highlights that the node order is not optimal in one dimension. Figure 1 (right top) displays the same graph, but sorted by applying our 1D path-guided SGD algorithm. The graph presents shorter links, allowing a clearer visualization. Figure 1 (right bottom) shows a 2D visualization of the previously 1D sorted graph. The coloring of the nodes follows perfectly the node position gradient, reflecting the goodness of the 1D sorting. Moreover, the nodes are well distributed on the 2D plane, without overlapping. The knot highlights the presence of a copy number variation (CNV), also visible as a dense region of links in the 1D plot. This means both visualizations match each other. Our multi-threaded implementation presents a working prototype that is based on succinct graph data structures. In progress is the exploration of the path-guided SGD parameter space, in order to get the best layout as quickly as possible. During this process, we also evaluate path-guided metrics in order to measure a graph's stress level. In the future, we want to find out performance boundaries applying the algorithms up to gigabase-scale pangenome graphs. We also plan to compare our proposed 2D graph layouting concept with existing pangenome graph visualization tools. The 1D path-guided SGD implementation is a key step in general pangenome analyses such as the pangenome graph linearization and simplification pipeline *smoothxg*.

**Acknowledgements:** We thank Vincenza Colonna for organizing the Crusco Hackathon and the Forentum Ritrovato museum for hosting it. S.H. acknowledges funding from the Central Innovation Programme (ZIM) for SMEs of the Federal Ministry for Economic Affairs and Energy of Germany.

## CCS Concepts

• Human-centered computing → Graph drawings; • Applied computing → Bioinformatics; Population genetics;