

# Visualizing Multiscale Complexity and Features of Biological Sequence Datasets

Leonid Zaslavsky<sup>1\*</sup>, Boris Fedorov<sup>1</sup>, Azat Badretdin<sup>1</sup>, Yiming Bao<sup>1</sup>,  
James R. Brister<sup>1</sup>, Stacy Ciufo<sup>1</sup>, William Klimke<sup>1</sup>, Kathleen O'Neill<sup>1</sup>,  
Alexander Souvorov<sup>1</sup> and Tatiana Tatusova<sup>1</sup>

<sup>1</sup>National Center for Biotechnology Information, National Library of Medicine,  
National Institutes of Health, 45 Center Drive, Bethesda, Maryland, 20894, USA

\*Corresponding author: [leonid.zaslavsky@nih.gov](mailto:leonid.zaslavsky@nih.gov)

## Abstract

The number of genomic and protein sequences in public databases have grown tremendously in recent years. Sophisticated tools are required to provide meaningful overview and request-specific visual representation of data.

Users exploring NCBI sequence databases have different needs and benefit from request-centric data representations with multiple scales of resolution. This allows an overview of large-scale homology relationships and focusing on a relevant scale of resolution, with views of rearrangements and variations between the sequences. Thereby adaptively focused visualization is refined to the sequences most relevant to the request with the rest of the data represented more coarsely. Aggregation techniques also provide a mechanism to deal with the great variation in sequence population density in a database.

Aggregated representation can also be applied to metadata. In some cases (e.g., virus variation), an annotation is based on either season/year or geographic location, and can be easily aggregated for a group of sequences. In other cases, more complex approach is required, such as comparison to another tree with named hierarchy of nodes (e.g., genome tree or taxonomy) and keyword extraction.

Our algorithm for adaptive visual representation of large trees along with metadata aggregation has been implemented within the NCBI Virus Variation Resources and the NCBI 16S rRNA Resource, with ongoing work on a online visualization framework allowing for visualization of more general sets of protein and genomic sequences.