

Exploring Material Similarity using Graph-Based Crystal Structure

Analysis and Machine Learning

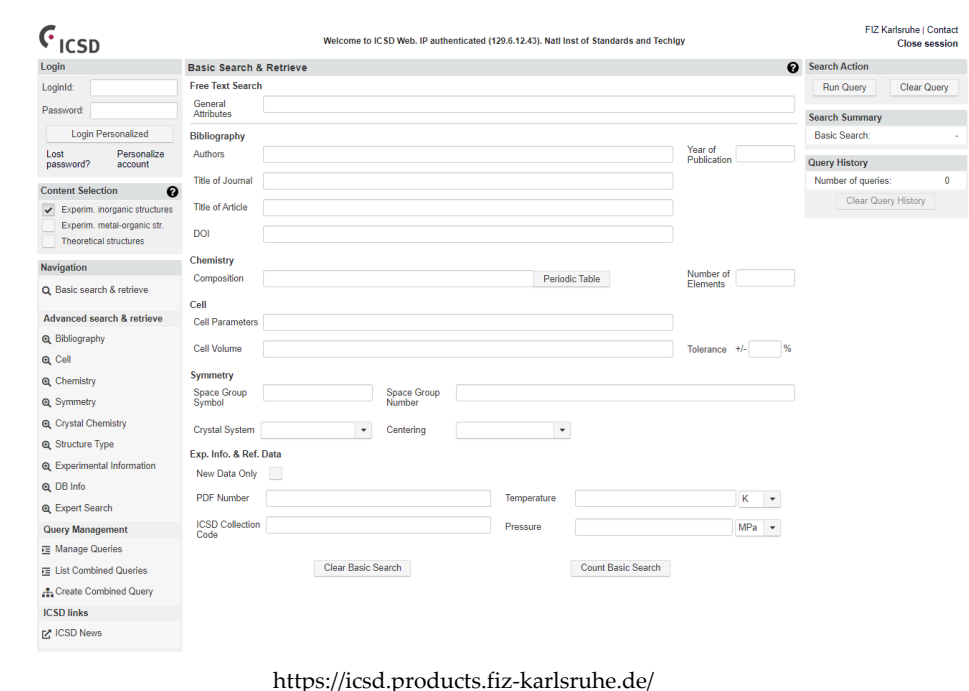
Karen Cao



Background

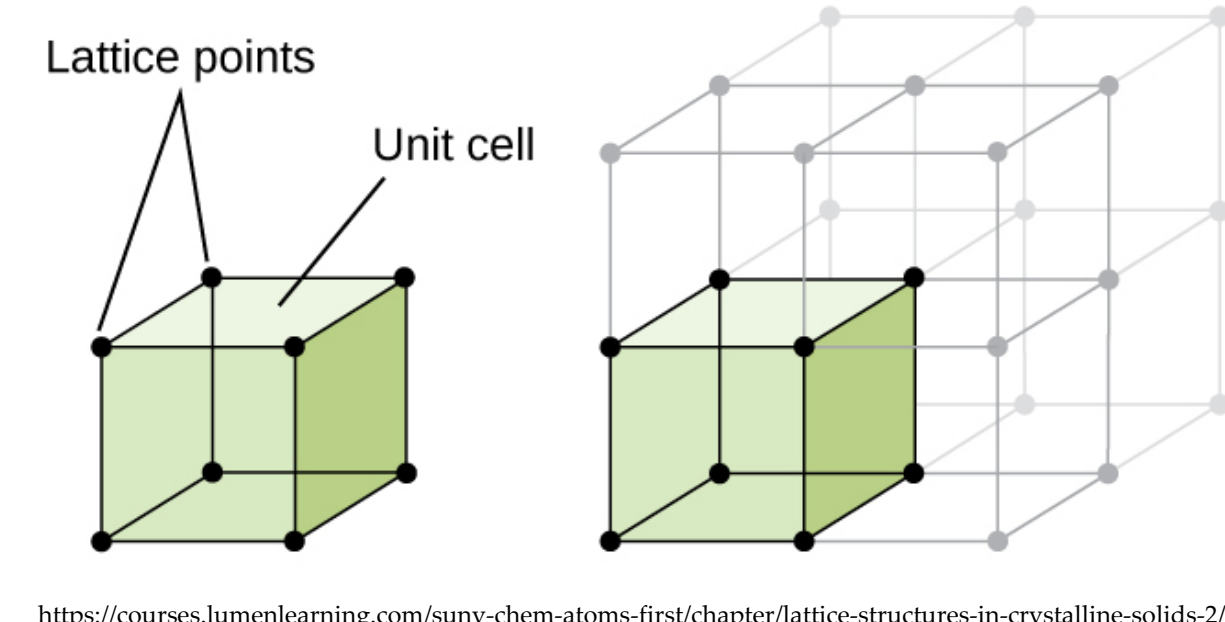
Inorganic Crystal Structure Database (ICSD)

- World's largest database for inorganic crystal structures
- 200,000+ compounds
- Provides Crystallographic Information File (CIF) for each crystal structure with unit cell, atomic coordinates, space group, and etc. information



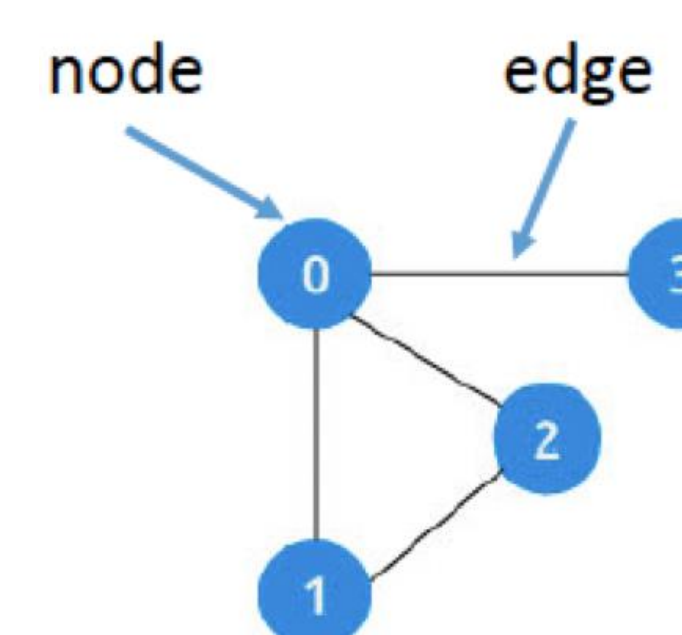
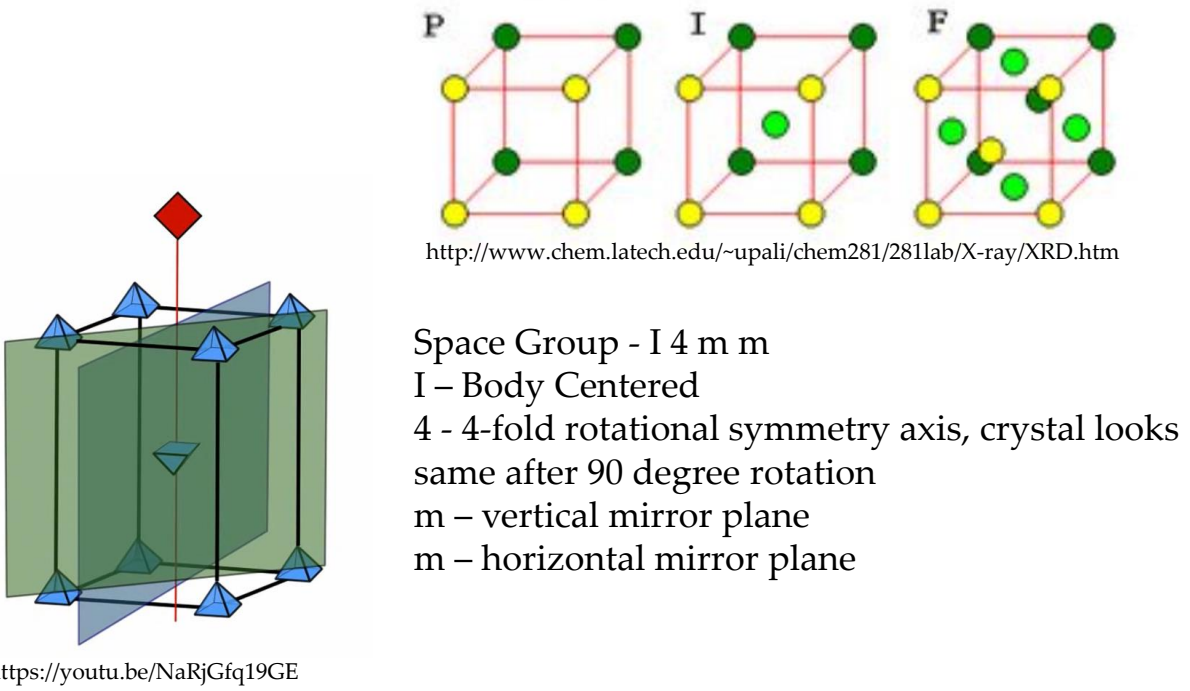
Crystal Structures

- 3D arrangement of atoms, molecules, or ions in a crystalline solid
- Spatial arrangement determined by interactions between the atoms such as bonds
- Unit Cell - 3D structure in the form of a parallelepiped that is repeated
- Lattice Point - reference point that is repeated within the crystal structure with identical surroundings



Space Groups

- Composed of Bravais type and symmetry operations that describe the arrangement of atoms in unit cell of a crystal structure
- Hermann-Mauguin (HM) symbols
- Bravais types
 - P - Primitive
 - I - Body Centered
 - F - Face Centered
- Symmetry Operations
 - Rotations
 - Reflections
 - Inversions
 - Glide Plane/Screw Axis
- Each space group is associated with a unique number (230 space groups)

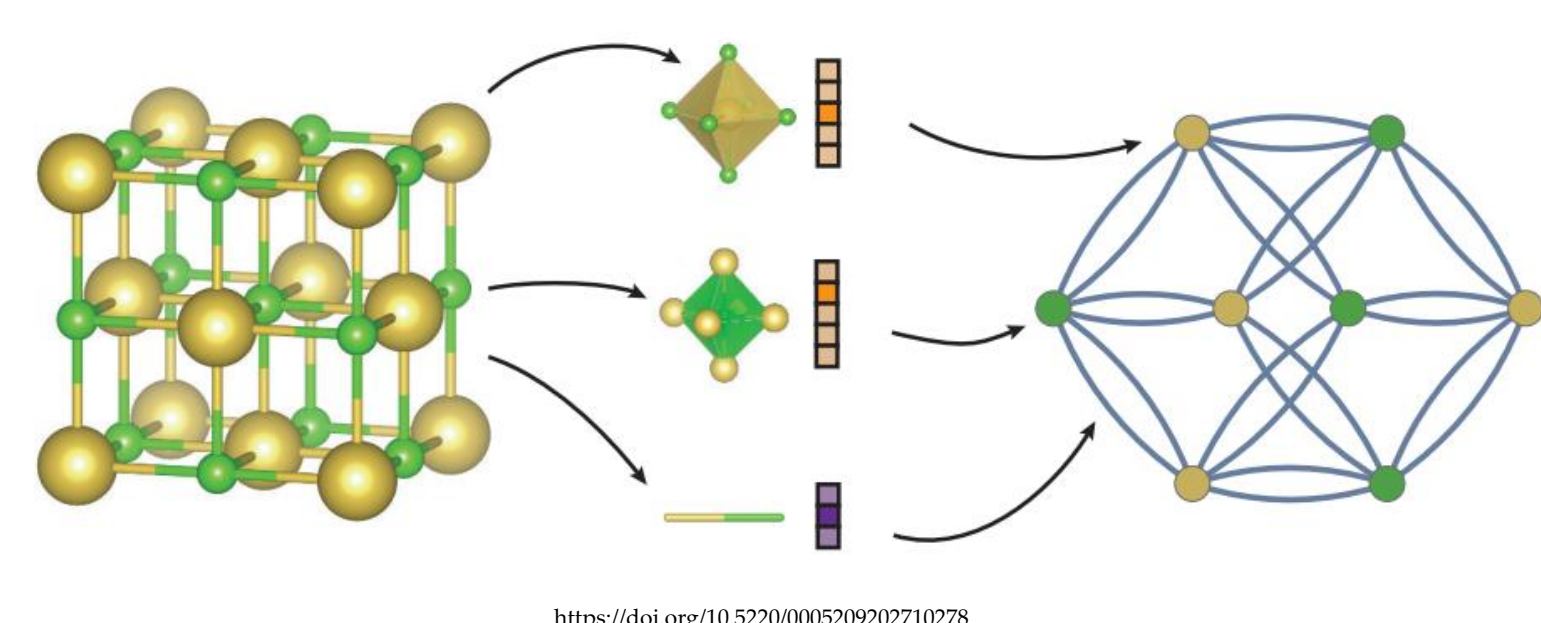


Graphs

- Mathematical structures used to model relationships
- Composed of nodes and edges
- Nodes represent entities or objects
- Edges represent interactions or connections between nodes

Objective

- Investigate the types of crystal structures already discovered and how they are distributed
- Represent crystal structures as graphs in order to create comparisons and place them in clusters or communities

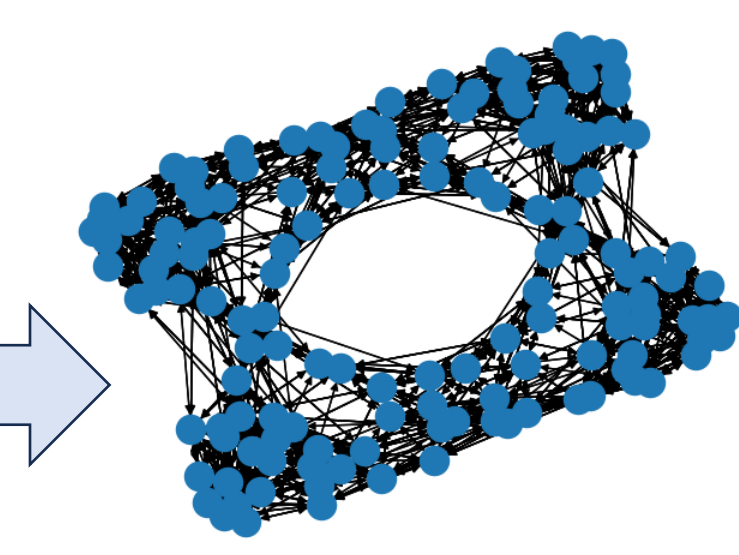


Methods

Transforming crystal structures into graphs

- Takes in a CIF file as input and returns a graph constructed with the JARVIS and DGL python libraries
- Main functions
 - Nearest_neighbor_edges - k-NN (k-nearest neighbor) edge list
 - Radius_graph - edge list based on a specified cutoff radius after several conversions and calculations
- Created by NIST researchers Dr. Kamal Choudhary and Brian DeCost

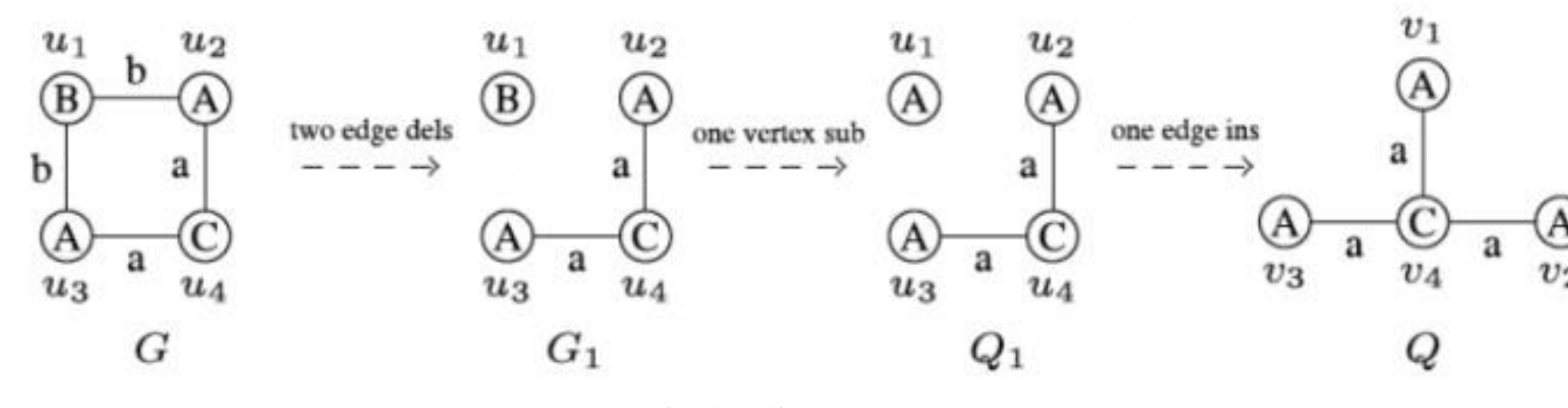
```
#Date: 2020-02-18 17:37:37 +0000 (Thu, 18 Feb 2020) $
#Revision: 170209 $
#URL: http://www.crystallography.net/cif2graph/000000000.cif $
#
# This file is available in the Crystallography Open Database (COD),
# http://www.crystallography.net/. The original data for this entry
# were provided by IUCr Journals, http://journals.iucr.org/.
# The file may be used within the scientific community so long as
# proper attribution is given to the Journal article from which the
# data were obtained.
#
# See: 000000
#
# CIF file name:
# CIF Name: 2020-02-18
# Author: s.
# Name: s.
# Title: s.
# Journal: s.
# Journal Page: s.
# Journal Page First: s.
# Journal Page Last: s.
# Journal Paper: s.
# Journal Title: s.
# Journal Year: s.
# Chemical Formula: s.
# Chemical Formula Weight: s.
# Space Group: s.
# Space Group Number: s.
# Symmetry: s.
# Symmetry Setting: s.
# Symmetry Space Group Name: s.
# Cell Lengths: s.
# Cell Angles: s.
# Cell Volume: s.
# Cell Formula Units: s.
# Cell Formula Units Z: s.
# Top of CIF File
```



Graph visualization with NetworkX Python Library

Measuring Similarity between Graphs

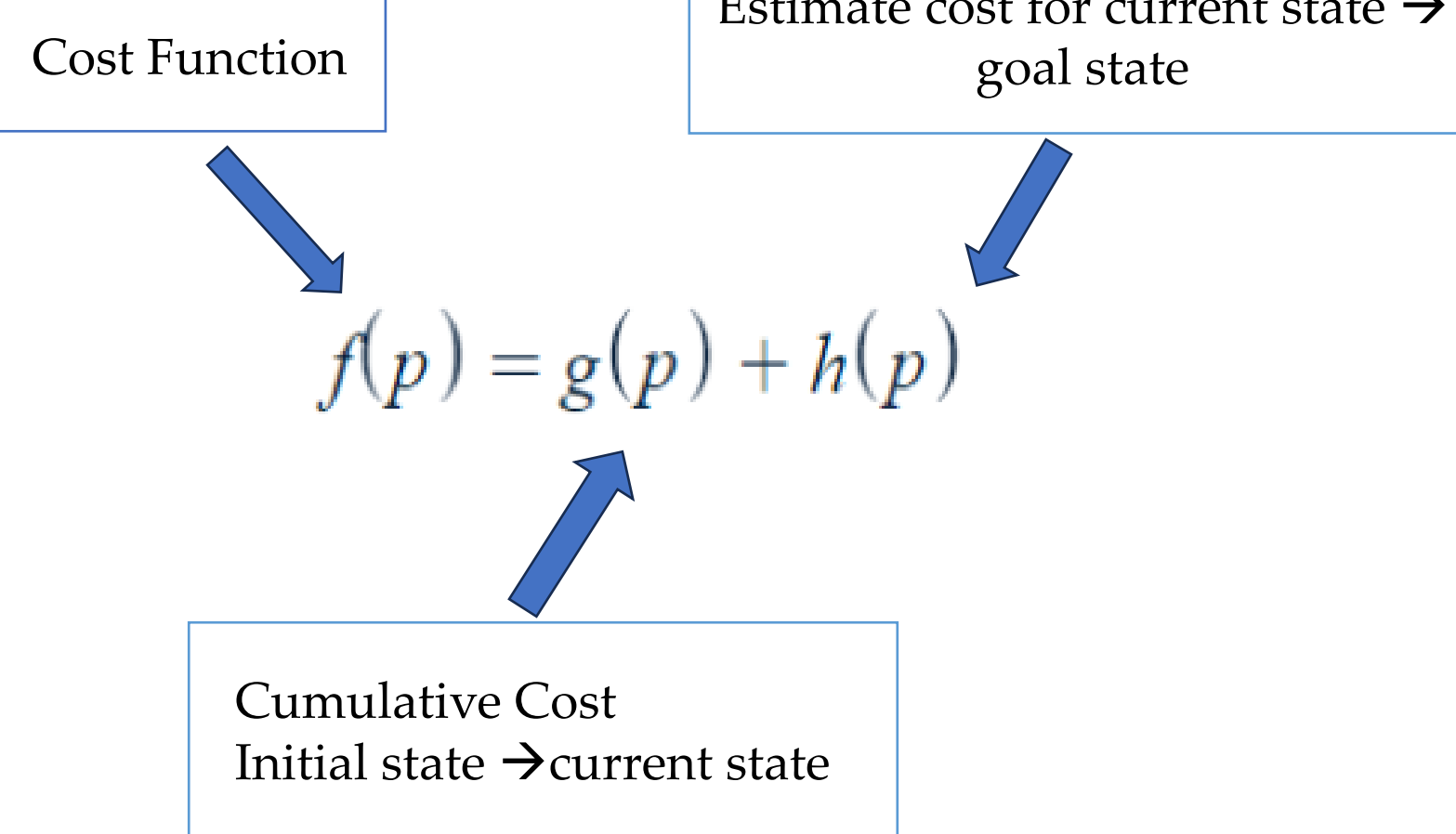
- Graph Edit Distance (GED) - minimum cost to transform one graph into another through a sequence of edit operations
- Edit operations - node and edge insertion, deletion, substitution
- The idea of using graph edit distance was inspired by NIST researcher Dr. Debra Audus



- Edit Operations:
- Two edge deletions
 - One vertex substitution
 - One edge insertion
- Each operation costs 1 → Total cost of 4

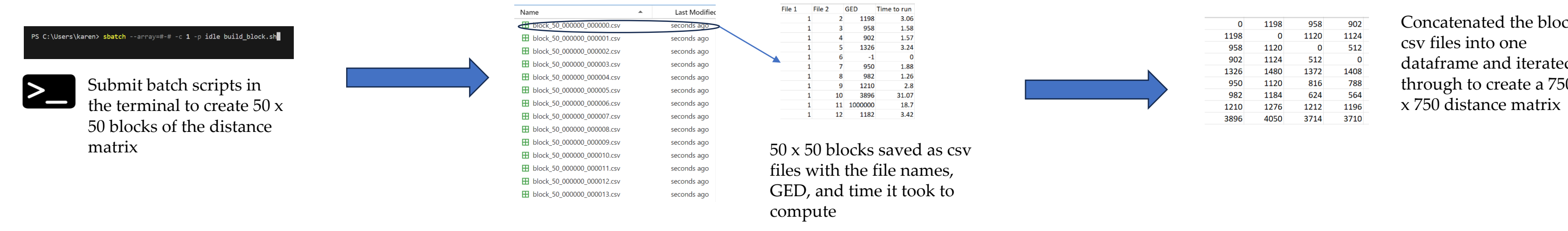
Calculating Graph Edit Distance (GED)

- Based on the A* Search Algorithm which finds shortest paths very efficiently
- Priority queue to store nodes and edges
- Elements are retrieved based on their "priority value" or order in which they are added (First in first out - FIFO)
- Apply each possible operation to create a new graph state and calculate the cost function f(p)
- Remove node with smallest f(p) from priority queue
- Used the NetworkX Python Library to calculate GED
- Returns the best GED calculation within a maximum number of seconds to execute
- Calculating GED is computationally expensive and time consuming so limiting run time was necessary
- Only able to calculate GEDs between 750 crystal structures



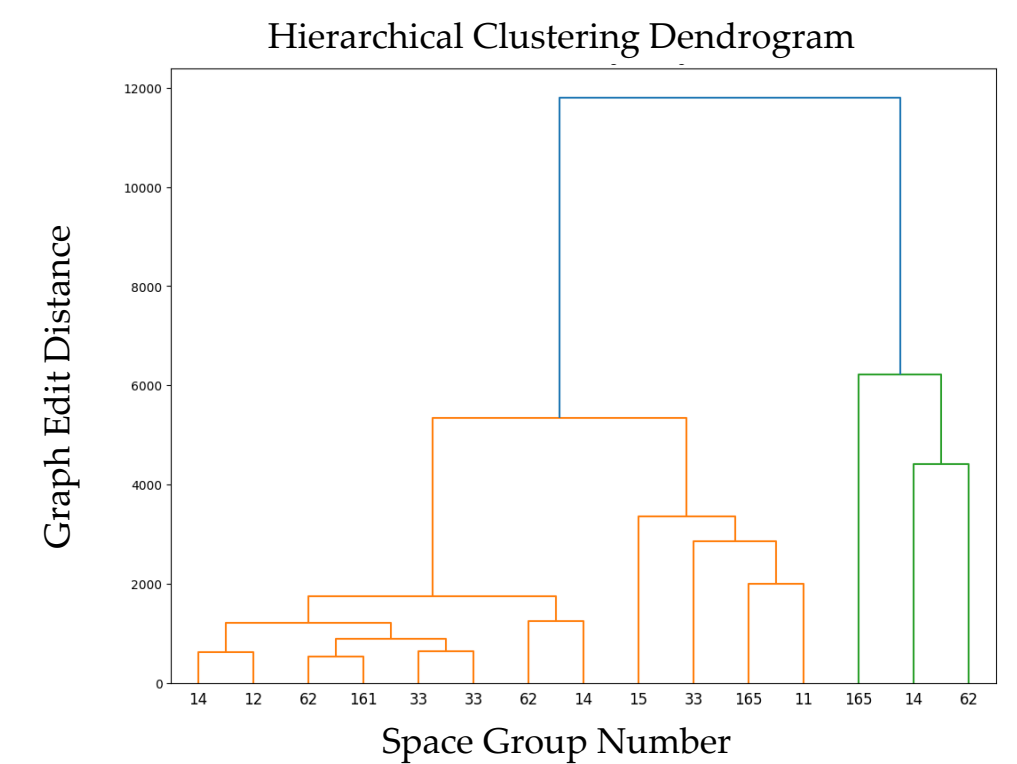
Creating a Distance Matrix

Ran the processes parallel on multiple nodes of a clustered server to increase efficiency of GED calculations



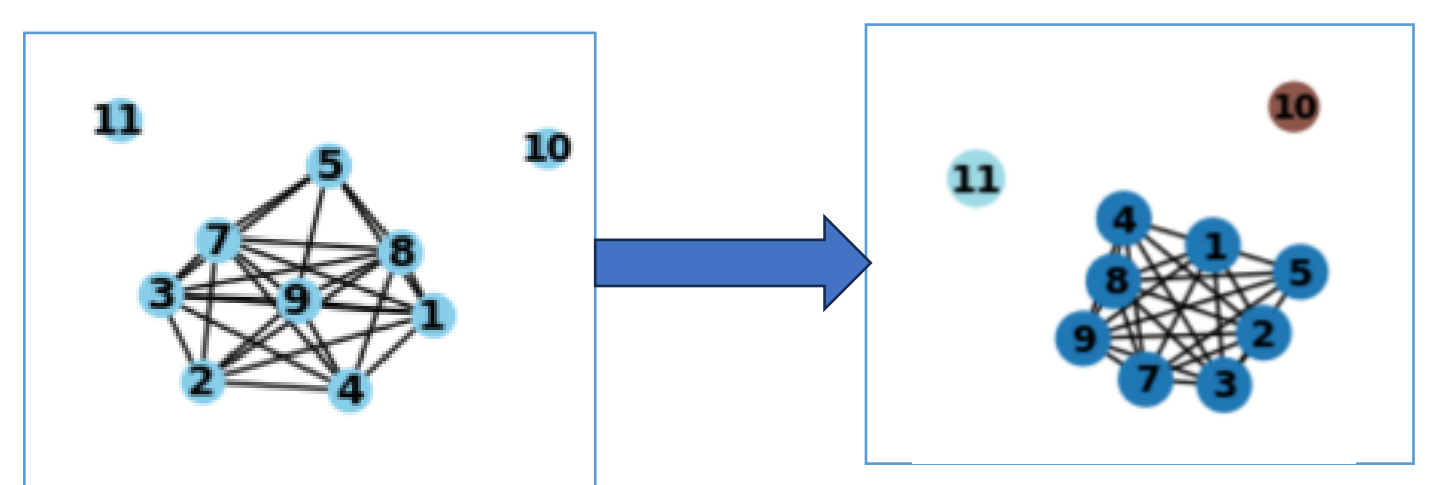
Hierarchical Clustering

- Groups similar crystal structures into groups/clusters
- How it works:
 - Treats each graph as a separate cluster
 - Repeatedly merges two clusters that are closest
 - Iterates until all clusters merged together

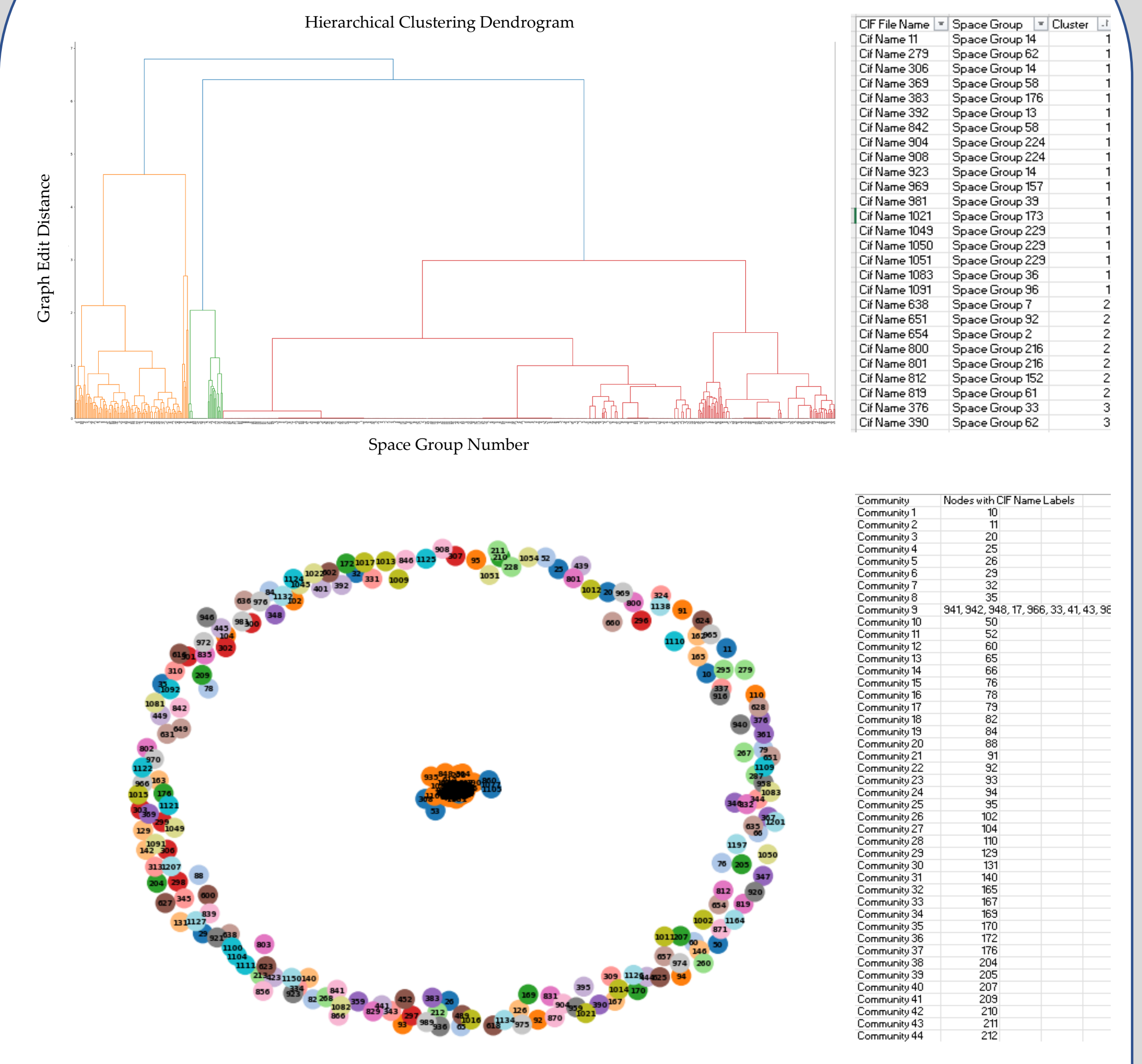


Community Detection

- Locates tightly connected nodes in a graph
- Graph with nodes for each crystal and edges constructed based on cutoff GED
- How it works:
 - Assigns different community to nodes
 - Considers each neighboring community for placing nodes
 - Node placed in neighbor community based on modularity
 - Modularity -



Results



Future Work

- Closely investigate the clusters and communities to understand the types of structures which have not been investigated as much, helping facilitate the generation of new materials
- How has this distribution has changed over time?
- Continue calculating the graph edit distance for a complete 200,000 x 200,000 distance matrix

References

Abu-Aisheh, Z., Raveau, R., Ramel, J., & Martineau, P. (2015). An Exact Graph Edit Distance Algorithm for Solving Pattern Recognition Problems. *International Conference on Pattern Recognition Applications and Methods*. <https://doi.org/10.5220/0005209202710278>

Jayawickrama, T. D. (2021, December 29). Community Detection Algorithms. *Towards Data Science. Medium*. <https://towardsdatascience.com/community-detection-algorithms-9bd8951e7dae>

Patlolla, C. R. (2022, August 15). Understanding the concept of Hierarchical clustering Technique. *Towards Data Science. Medium*. <https://towardsdatascience.com/understanding-the-concept-of-hierarchical-clustering-technique-c6e8243758ec>

Sands, D. E. (2012). *Introduction to Crystallography*. Courier Corporation.

Xie, T., & Grossman, J. C. (2018). Crystal Graph convolutional neural networks for an accurate and interpretable prediction of material properties. *Physical Review Letters*, 120(14). <https://doi.org/10.1103/physrevlett.120.145301>

2.3. Clustering. (n.d.). Scikit-learn. <https://scikit-learn.org/stable/modules/clustering.html#hierarchical-clustering>

Communities - NetworkX 3.1 documentation. (n.d.). <https://networkx.org/documentation/stable/reference/algorithms/community.html>

graph_edit_distance - NetworkX 3.1 documentation. (n.d.). https://networkx.org/documentation/stable/reference/algorithms/generated/networkx.algorithms.similarity.graph_edit_distance.html

Acknowledgements

- Mentor: Dr. William Ratcliff
- Special thanks to Mr. Paul Kienzle
- NCNR Mentors: Dr. Paul Butler, Mr. Jeff Krzywon, Dr. Yun Liu
- CHRNS Sponsors: Julie Borchers, Leland Harriger