

Comvi - Comparative Visualization of Molecular Surfaces using Similarity-based Clustering

Wilhelm Buchmüller, Shoma Kaiser, Damir Ravilija, Enis ...

Abstract— The goal of this paper is to show the reader the abstract methods and concrete applications that were used to extract and compare features and rank the similarity of the molecular protein maps. Further we present a new method of how the won data can be visualized on high resolution and large displays with a The paper describes the process and the approaches that were taken to solve this task.

Index Terms—Clustering, Similarity, feature extraction, Visualization, high-resolution display, Powerwall, MegaMol, VISUS

1 INTRODUCTION

TODO: add sections to tasks in this section Over the span of 6 months we, the authors of this paper, have researched and implemented a comparative clustering of molecular maps. The task consisted of several parts: This work was based on the MegaMol project. MegaMol is a simulation tool developed by the Universitt Stuttgart and the TU Dresden(?) **TODO: cite megamol**. It can be used to visualize particle data, simulations on atomic scale and other molecular processes. Due to its modular nature, it can be extended with modules to interact with other modules. In this paper we guide you through the **TODO: verbose** MSMCLUSTER plugin, its capabilities and inner workings.

The next task was to retrieve molecular image data through existing MegaMol plugins **TODO: cite molecular maps**. For this a special binary of megamol was compiled and will be released in a separate project **TODO: make github link tocomvi public**. The next task was to extract a expressive feature vector from those images and to find a metric to cluster them by similarity.

The last task which was also developed in a plugin in MagaMol was the visualization on the VISUS POWERWALL. The POWERWALL is a very high definition display that can be used to visualize large data(sets). Due to its size its possible to display much more information than on a regular screen.

The POWERWALL also supports a tracker device that can transmit 6 degrees of freedom, so for the interaction step we had more freedom to work with than with traditional human interaction devices **TODO: HID abbreviation correct (?)**, for this last step we also researched the possible interactions with the tracking devices on the POWERWALL.

Over the cours of the next few pages you will learn how we approached these challenges and how we (attempted) solved them, what worked and what didn't.

2 RELATED WORK

TODO: cite kolesar for clustering Clustering proteins by similarity or at least comparing individual proteins has been subject of existing work.

The paper [1] already had similar approaches to our results. Kolesar et al. used a 10 dimensional feature vector based on invariant image moments defined by hu **TODO: cite hu image moments**

Another approach for 3D protein data were 3D zernicke moments explored by **TODO: 3d surfer**. The approach is basically the same as in Kolesar et al. but the TODO used Zernicke moments instead of traditional image moments and extended them to three dimensions.

2.1 Initial Challenges and encountered problems

It is clear that the task required from us that we learn how to compare the images, measure the distances between the images, and cluster these images. The given task required that we use a similarity based clustering algorithm.

Initially we were given the choice we could either chose to find similarities and cluster the proteins in the .pdb format or given as bitmap image generated by the MolecularMaps plugin in MegaMol **TODO: cite megamol protein image**.

Since handling image files which give information about the protein in two dimensions was easier than dealing with the pdb file format which results in three dimensional visualizations we decided to start with a two dimensional approach.

2.2 Approaches to the Clustering-Problem

Right of the start we had several ideas of how we could approach this problem. With the recent trend in machine learning we had a couple of ideas of how we could determine a similarity metric between two images or classify an image into a more usable vector of data.

We ended up using a higher dimensional feature vector described in **TODO: give section label** to determine the similarity between two protein maps because we didnt manage to train a custom model in the given timeframe, due to inexperience **TODO: this can be said better** and non existance of labeled data.

But our relatively spartan results with a pretrained Imagenet **TODO: cite imagenet** **TODO: cite darknet publication** model let us to believe that given the knowledge on the subject and humanly labeled data (based on known featuresit **TODO: LEFT OFF HERE** should be definety possible for this specified task to find a machine learning solution using neural networks/autoencoders.

2.3 Approaches to the Visualization-Task

Our approaches to visualizing the given clusters were the following, the reader is reminded that we are not just visualizing the clusters on a "normal machine" but rather the POWERWALL, a projected display with effectively 6-24 times the resolution of a consumer grade display. Details on the POWERWALL can be found **TODO: cite powerwall publicatoin here if available**.

If we are given so much pixel real estate we are given the freedom to draw smaller pixels since we still will be able to see them on the POWERWALL.

During the duration of the project the idea of a 3D visualization was discussed among the team, but we settled for a 2D visualization. This had a couple of reasons.

Wilhelm Buchmüller Shoma Kaiser
buch.willi@googlemail.com example@example.com
Enis ... Damir Rwilja
example@example.com example@example.com

Manuscript received 31 March 2011; accepted 1 August 2011; posted online 23 October 2011; mailed on 14 October 2011.
For information on obtaining reprints of this article, please send email to: tvcg@computer.org.

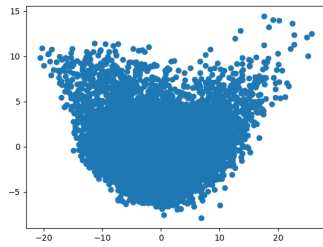


Figure 1: Early test with the Oxford flower dataset **TODO: cite oxford flower dataset**

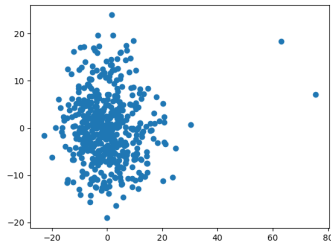


Figure 2: Early test with the BMW car dataset **TODO: cite stanford car dataset**

While the interactions with the data in 3D would have been more fun since we had more degrees of freedom to work with. But we could not find a way to present the data in a way such that with just a glance the user could intuitively interpret the data that would be displayed on the screen.

... So we decided to settle for a 2D visualization. Our approach is rather boring but it works. On startup we display nothing, if the user chooses his supplied image data and the algorithms used to cluster the pictures he gets the option to start the visualization.

The visualization consists of displaying the image of a cluster representative with a simple rectangle. This way the user knows exactly what to expect to be in the cluster.

If the user wishes to have a closer look at the images in the given cluster he can click onto the representative and will get a view of all the images in that cluster.

We decided after testing with toy test and real datasets that one level of subclustering is enough, after 2 levels of subclustering the clusters get **TODO: find better expression** noisy and ambiguous.

In both the main and subcluster view the representatives are visualized with "force directed layout", to avoid

Table 1: lorem ipsum tabulated

dataset	full performance (fps)	half performance (ms)
balls	1,243	0.1
buckets	23	23
bolts	23,312,134.3	22.1

2.8 Finding the

2.9 Discussion of results

2.9.1 Dolor

ACKNOWLEDGMENTS

We would like to thank our supervisors Michael Krone and Florian Fries as well as our project examiner Prof. Ertl . We are grateful for the experience and knowledge that working on this project has given us. This work was partially funded by cake and cookies.

REFERENCES

- [1] I. Kolesár, J. Byška, J. Parulek, H. Hauser, and B. Kozlíková. Unfolding and interactive exploration of protein tunnels and their dynamics. In *Proceedings of the Eurographics Workshop on Visual Computing for Biology and Medicine*, pages 1–10. Eurographics Association, 2016.

2.3.1 Approaches to the Interaction with the Powerwall **TODO: cite correct pub and use correct name**

2.4 Finding a feature vector to cluster the images

2.5 Finding the best performing similarity measure

2.6 Finding the best performing clustering algorithm

2.7 Testing the feature vector with other datasets