# An Efficient Method for Indexing Temporal Gene Expression Datasets

Shawn Oliai, Casimir Tokarski, Guenter Tusch, PhD

Health Informatics and Bioinformatics Graduate Program, School of Computing and Information Systems

Grand Valley State University, Grand Rapids, MI, USA

## Our Project

**Purpose:** Identifying Temporal High Throughput Gene Expression Data Sets for Comparative Transcriptome Analysis.

**BACKGROUND:** Comparative transcriptome analysis of high throughput gene expression temporal experiments helps with understanding of the complexities of living organisms with great value for diagnosis, treatment, and prevention of human diseases.

**PROCEDURES:** This study explores the feasibility of searching for temporal patterns based on knowledge-based temporal abstractions. Those imply conversion of expression values into an interval-based qualitative representation expressing amount of change over time in terms of trends, as "increasing", "decreasing", "constant" etc. We searched 20,000 articles between 2008-2018, and mined for keywords such as *'high throughput', 'gene expression', 'transcriptome analysis', 'temporal', 'time series', 'longitudinal, 'chronological'* and other synonym terms.

**METHODS AND MATERIALS**: To identify those temporal patterns, these studies need to be indexed appropriately since much of the publicly available high throughput expression data is of non-temporal nature. The index needs to be based on the MIAME standard as used for example for abstracting NCBI's Gene Expression Omnibus (GEO) datasets. A simple keyword search of abstracts from NCBI GEO will only result in a large number of false positives.

**RESULTS:** We essentially see this research as a text mining process to find the correct set of pertinent articles for the topic at hand. Using random samples of keyword search results, we repeatedly refined the search query to obtain better indexing of the datasets from temporal studies. We utilized a series of appropriate words for the algorithm to find in the search and thus were able to significantly improve on false positive and false negative search results.

## Conclusion

**CONCLUSIONS:** After finding 9,694 articles with our selected terms, the were 212 articles including all of the terms and 'temporal' words. Of those 41 had one or more related GSE files for our final consideration.
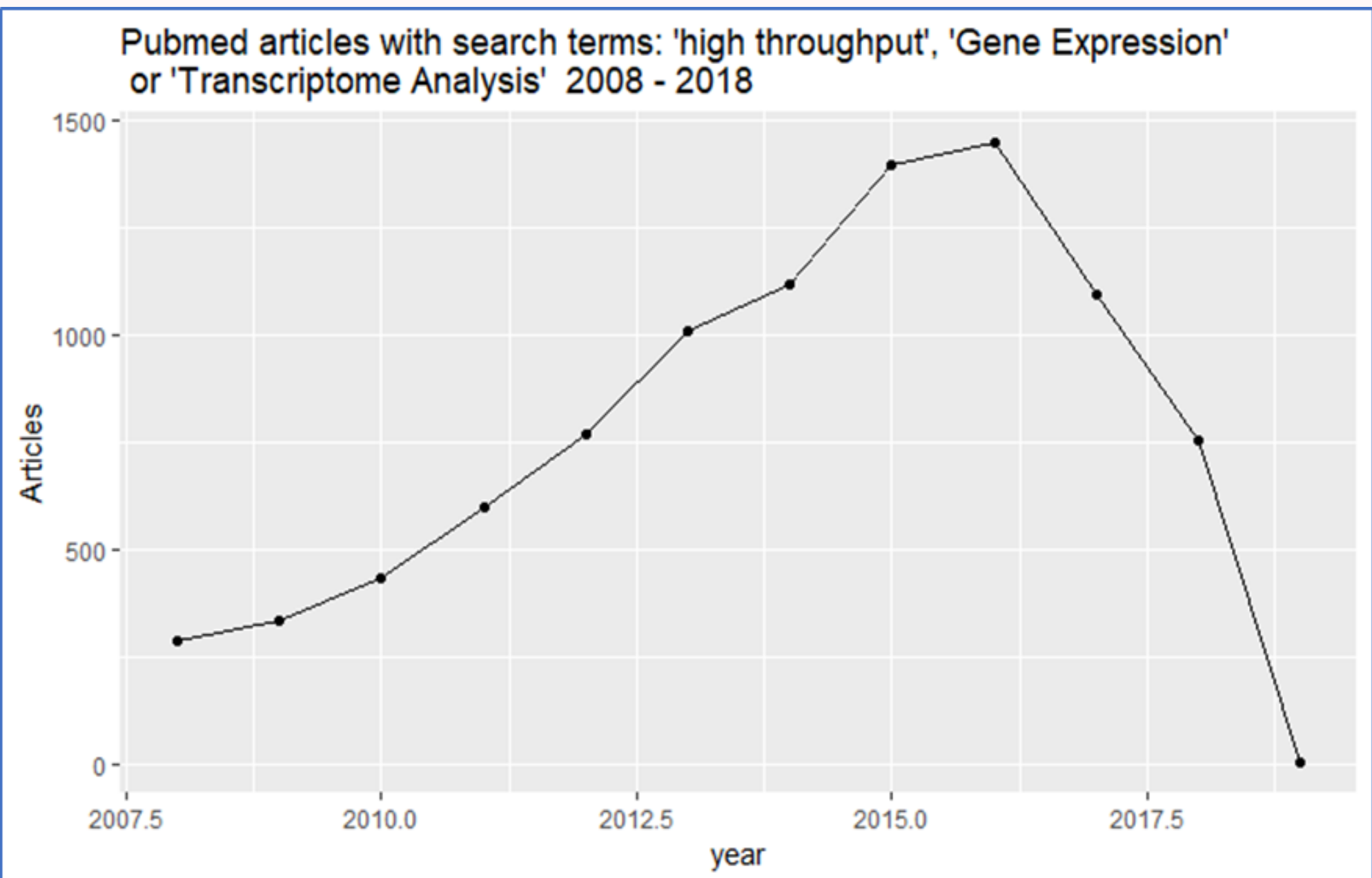
Our promising model can assist researchers in bioinformatics, genetics, medicine and scientific investigations, to find more appropriate selection of scientific data online, while saving various resources (time, capital, talent). The refined indexing algorithm can be used in future studies to compare patterns of gene expression in bioinformatics.

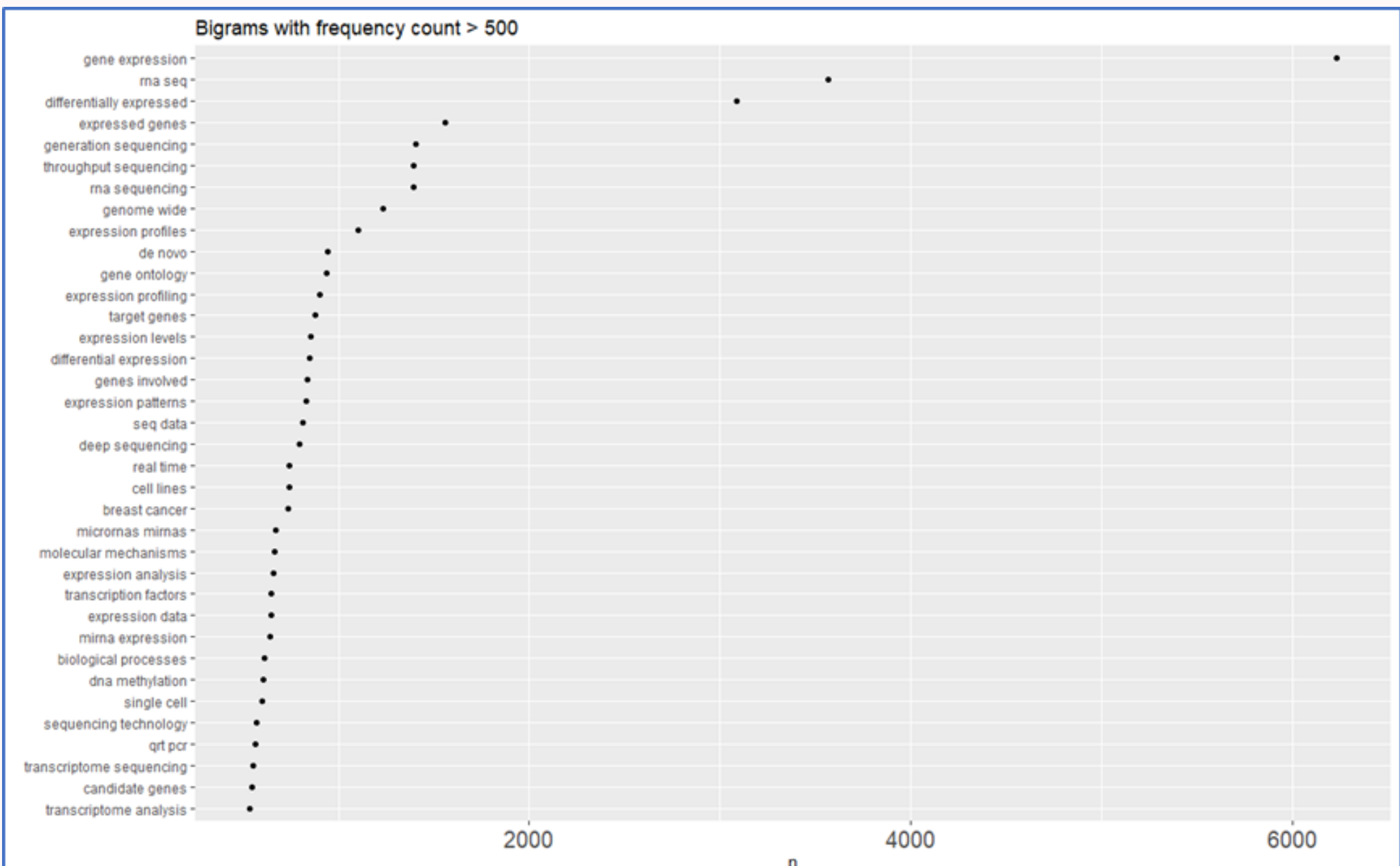**FUTURE EXPLORATION:** There are several appropriate considerations on this project:

- Examine the methodology in other biological science fields and databases beyond PubMed, and compare the results

- Compare/contrast results with the PubMed's internal text mining tools

- Develop a Topic Network for further exploration, along with deeper insights into LDA modelling

### Acknowledgement:

PubMed Articles with Selected Key Terms



Bigram for high throughput



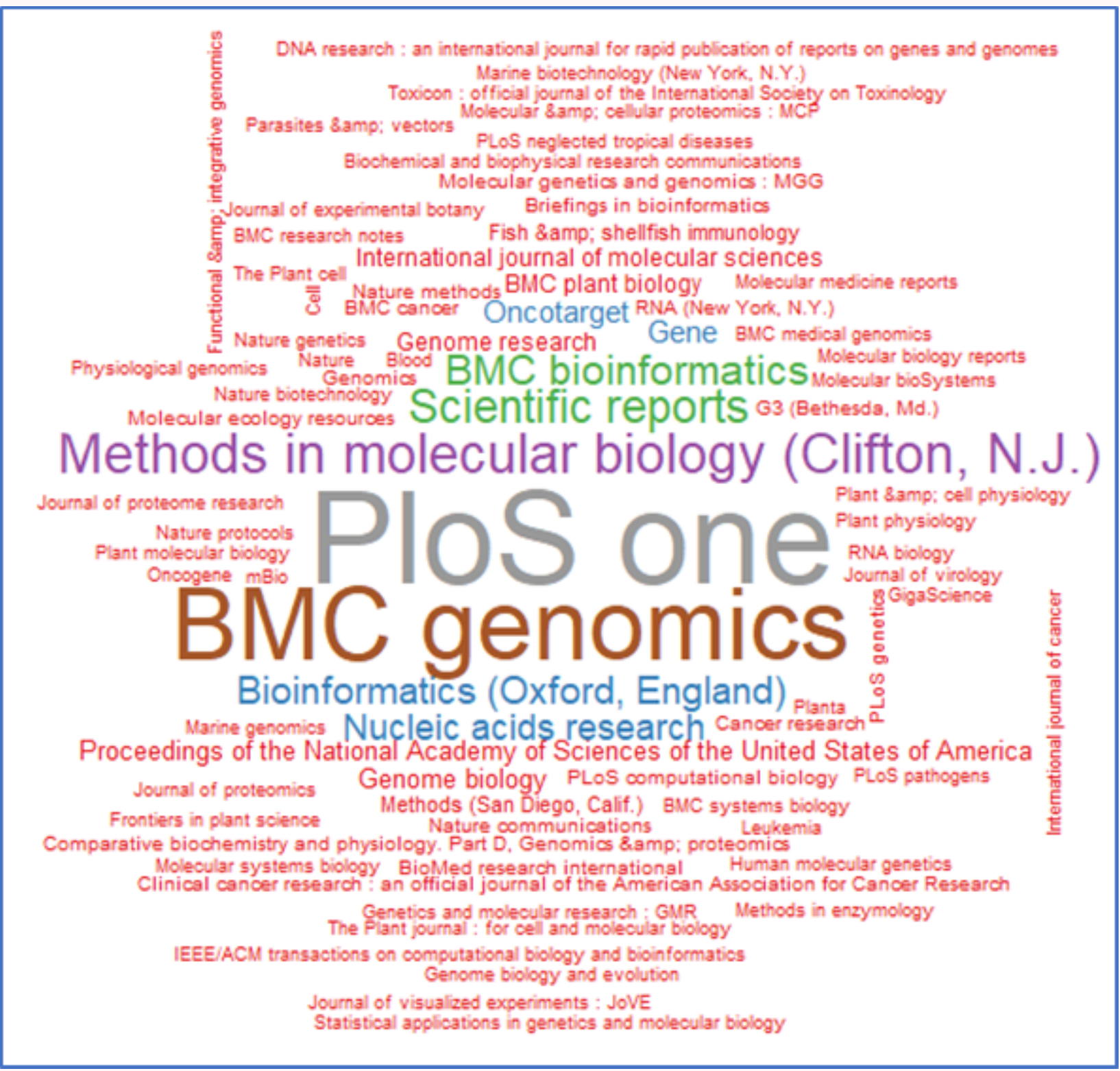Articles based on Any of the Keyword Term Sets



Final Results of Appropriate Articles with GSEs



Word Cloud for top 1500, repeated at least 750 times



Bigram of words with n>500 high throughput



Journals with at least 10 articles on the topic words

| | 2008-2018 | March 2018 |
|---|---|---|
| | Entire Set | Sample Set |
| **Run Time** | 384 s | 42 s |
| **Raw Data** | 9694 | 1077 |
| **Clean (!na)** | 9277 | 919 |
| | | |
| **'gene expression'** | 3520 | 324 |
| **'high throughput'** | 729 | 55 |
| **'temporal'** | 212 | 20 |
| **'time series'** | 36 | 1 |
| **'longitud\*'** | 21 | 0 |
| **'chronol\*'** | 6 | 0 |
| **Article + GSE** | 40 | 3 |