

# volcalc: Calculate predicted volatility of chemical compounds

Kristina Riemer

Eric R. Scott

Laura Meredith

2023-03-21

## Signatories

### Project team

- Kristina Riemer, author/maintainer, Director of Communications and Cybertechnologies Data Science team at University of Arizona
- Eric Scott, contributor, Scientific Programmer and Educator for Communications and Cybertechnologies Data Science team at University of Arizona

### Contributors

- Assistant Professor Laura Meredith developed the original idea for the `volcalc` package along with Kristina Riemer and is supportive of continued development by our team.
- PhD student S. Marshall Ledford has been the main user of early versions of `volcalc` and will continue to provide feedback on the package API and documentation.

### Consulted

Tamás Stirling, maintainer of the `webchem` package (part of `rOpenSci`), was consulted and confirmed that `volcalc` is not replicating the efforts of any similar R packages that we are aware of.

## The Problem

Volatile organic compounds (chemicals that readily evaporate under ambient conditions) are important in a number of fields and contexts including involvement in plant defense against herbivores, as routes of microbial communication, and as important indoor pollutants, to name a few. Yet measures of volatility are time consuming to calculate experimentally and not available for the vast majority of chemical compounds in chemical information databases. However, methods exist for predicting measures of volatility from chemical structure (Pankow and Asher 2008). The `volcalc` package aims to automate the following steps for a given compound: 1) downloading data on chemical structure, 2) parsing those data to discover chemical functional groups, 3) applying the SIMPOL algorithm to predict volatility from functional groups and molecular weight. The current draft version of `volcalc` does all this, but is limited to working with compounds present in the Kyoto Encyclopedia of Genes and Genomes (KEGG) database. This proposal focuses on expanding the scope of `volcalc` and preparing it for a larger and more diverse user base so anyone interested in volatile organic compounds can integrate it into their workflow.

## The proposal

### Overview

- 1-2 sentence overview of what work is proposed.

This is more background, move to next section?:

Version \_\_\_\_ of `volcalc` was created in \_\_\_\_ as part of a [data science incubator](#) project in collaboration between Dr. Kristina Riemer and Dr. Laura Meredith at University of Arizona. `volcalc` is the first project, to our knowledge, to implement the SIMPOL method for predicting chemical vapor pressures and enthalpies of vaporization (Pankow and Asher 2008) in an R package. This current version of the `volcalc` package has been successfully used to calculate volatility estimates for *all* compounds in the Kyoto Encyclopedia of Genes and Genomes (KEGG) database. However, in its current form, it is limited to only working with chemical compounds in KEGG. This project will generalize this functionality of `volcalc` to work with any chemical—not just those in the KEGG database. `volcalc` currently downloads chemical information as molfiles from the KEGG API as a starting point. Molfile is an open format and various tools exist to translate other standard representations of chemical structure such as SMILES and InChIKey to molfiles (e.g. using the [OpenBabel](#) command line utility or in R with [ChemmineOB](#)). Refactoring the code in `volcalc` to work with *any* molfile and preparing the package for wider distribution will make this powerful tool accessible to researchers across a variety of domains.

## Detail

The current version of `volcalc` focuses specifically on estimating volatility of compounds in the Kyoto Encyclopedia of Genes and Genomes (KEGG) database. The main function in `volcalc`, `calc_vol()` downloads chemical structure data using the KEGG API as a `.mol` file. It then reads that `.mol` file in and parses it to find functional groups. It then applies an algorithm (published in Pankow and Asher (2008)) to predict vapor pressure, and outputs to the user a relative measure of volatility. Other functionality in the package includes ...

`volcalc` is developed on GitHub and distributed under an MIT license. Project repository: <https://github.com/Meredith-Lab/volcalc>

## Project Goals

Our goals for the proposed project fall into two main categories: 1) to make `volcalc` useful for applications beyond estimating relative volatility for compounds in the KEGG database, and 2) polishing the package in preparation for an initial submission to CRAN.

For the first goal, we will focus on decoupling the data access and volatility estimation functionality of existing code in the `volcalc` package. This code is already written and the main work here is in mindfully re-factoring. The minimum viable product here is a function that can calculate volatility when provided a path to a `.mol` file and a vignette demonstrating how to couple this with chemical data sources such as the `webchem` package.

The second goal will involve improving test coverage (although current coverage is high at 96.88% as reported by `covr::package_coverage()`), establishing continuous integration with GitHub actions, improving documentation, and satisfying R CMD check. The minimum viable product here is a package that has gone through the steps suggested by `usethis::use_release_issue()` and is ready to submit to CRAN.

A reach goal is to allow users to input other chemical structure representations besides molfiles. The `ChemmineOB` package can translate a variety of chemical structure representations to molfiles and is already an indirect dependency of `volcalc` through its dependency on `ChemmineR`. Adding an argument to our `volcalc` function to specify the input format, and passing it to `ChemmineOB`'s translation function would be a way to expand the usability of `volcalc` even further.

## Project plan

### Dates

- Project start date: June 1, 2023
- Project end date: January 31, 2024

## Start-up phase

I just realized there's nothing about getting package CRAN-ready on here

### Milestone 1: July 1, 2023

- Implement CI with GitHub actions
- Check code coverage with `codecov` package
- Use GitHub Issues or Discussions to brainstorm eventual API (i.e. function names, argument names, how many exported functions, etc.)

Estimated work: 10 hours

## Technical delivery

### Milestone 2: September 1, 2023

- Re-factor `calc_vol()` code to split KEGG download and SIMPOL calculation functionality
- Deprecate arguments and functions appropriately as necessary
- Update documentation to reflect new function usage

Estimated work: 40 hours

### Milestone 3: November 1, 2023

- Create a vignette demonstrating both KEGG usage and more general usage (i.e. providing a path to a .mol file) for volatility estimation
- Improve package documentation by adding citations, details, and additional examples where appropriate
- Create a `pkgdown` website for `volcalc`
- Create a CITATION.cff file, make a GitHub release, and archive code on Zenodo

Estimated work: 40 hours

### Milestone 4: January 31, 2024

- Add functionality to supply other chemical representations besides molfiles as input
- Add to vignette(s) examples of integrating `volcalc` with `webchem` to estimate volatility for an arbitrary set of compounds (i.e. not from KEGG), and `volcalc` with Biocyc

Estimated work: 40 hours

## Other aspects

Dissemination plan:

- After the initial re-factor (milestone 2 above), we plan to share the package with **webchem** contributors via our rOpenSci Slack channel for feedback & suggestions. We will encourage them to share the project with their networks of collaborators as well.
- At the project conclusion we will:
  - prepare blog post for <https://datascience.cct.arizona.edu/>
  - prepare a twitter announcement to share from ([cct\\_datascience?](#))
  - prepare a short demonstration video to be published to our [YouTube channel](#) that can be shared on the README of our repository
  - We will also work with our collaborator Laura Meredith to identify potential domain-specific venues to promote the use of **volcalc**

Estimated work: 20 hours

## Requirements

**People**

**Processes**

**Tools & Tech**

**Funding**

**Summary**

**Success**

**Definition of done**

**Measuring success**

**Future work**

**Key risks**

Pankow, J F, and W E Asher. 2008. "SIMPOL.1: A Simple Group Contribution Method for Predicting Vapor Pressures and Enthalpies of Vaporization of Multifunctional Organic

Compounds.” *Atmos. Chem. Phys.* <https://doi.org/10.5194/acp-8-2773-2008>.