# New Basis Set Exchange: An Open, Up-to-Date Resource for the Molecular Sciences Community
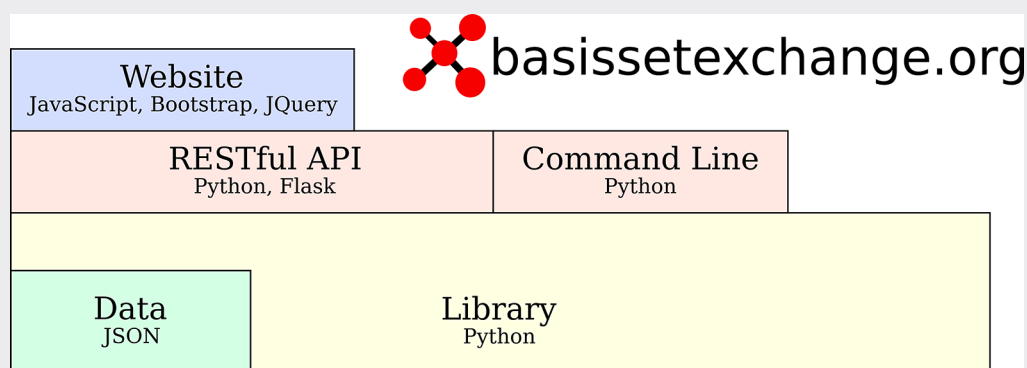
Benjamin P. Pritchard,*,[†],[⊥] Doaa Altarawy,[†],[‡],[⊥] Brett Didier,[§] Tara D. Gibson,[§] and Theresa L. Windus[†],[‖]

[†]Molecular Sciences Software Institute, Virginia Polytechnic Institute and State University, Blacksburg, Virginia 24060, United States
[‡]Department of Computer and Systems Engineering, Alexandria University, Alexandria 21544, Egypt
[§]Pacific Northwest National Laboratory, Richland, Washington 99352, United States
[‖]Iowa State University, Ames, Iowa 50011, United States

**ABSTRACT:** The Basis Set Exchange (BSE) has been a prominent fixture in the quantum chemistry community. First publicly available in 2007, it is recognized by both users and basis set creators as the de facto source for information related to basis sets. This popular resource has been rewritten, utilizing modern software design and best practices. The basis set data has been separated into a stand-alone library with an accessible API, and the Web site has been updated to use the current generation of web development libraries. The general layout and workflow of the Web site is preserved, while helpful features requested by the user community have been added. Overall, this design should increase adaptability and lend itself well into the future as a dependable resource for the computational chemistry community. This article will discuss the decision to rewrite the BSE, the new architecture and design, and the new features that have been added.

## 1. INTRODUCTION AND HISTORY

In computational quantum chemistry, the basis set is a fundamental input to many types of calculations, and the choice of basis set can have dramatic effects on both the accuracy and performance of the simulation. Due to this fundamental nature, basis sets have been developed starting at the beginning of the computational chemistry field, and extending to the present day, with new basis sets taking into account new theories and methods now accessible due to increases in computing power.

Practically, computations are most comparable when they use the same basis set. However, early distribution channels often led to differing incorporation of basis sets in different computational chemistry codes, harming reproducibility and introducing errors in calculations that can go unnoticed. Such errors can be due to many causes, including transcription errors, neglecting publication of the full precision of numerical values, updated values not being incorporated in a timely manner, or development of an "improved" version by different developers that retains the same name as the original basis set.

Therefore, there existed a need to store basis sets in a central location and make them easily and freely accessible by any researcher. The Basis Set Exchange (BSE) is widely considered by the field to be such a resource.

The BSE has its origins in the early 1990s when David Feller was developing databases and reproducibility tools for Gaussian-based computational chemistry simulations. His Computational Chemistry Input Assistant[1] (CCIA) contained databases of basis sets in plain text (the Gaussian Basis Set Database, organized by family) and computational results for relatively small molecules at many different levels of theory (the Methods Performance Database). The computational results included a fairly extensive list of properties such as molecular structure, vibrational modes, moments, transition moments, and absolute and relative energies. Other properties such as dissociation energies, inversion/rotation barriers, electron affinities, ionization potentials, and proton affinities were also made available. The CCIA was a Fortran-77 based

program that ran on a user desktop system and included a graphical interface.

The Gaussian Basis Set Database was also used as the back end to the Gaussian Basis Set Order Form that had a more text-based interface. In the mid 2000s, this Order Form was developed into the popular Basis Set Exchange Web site.[2] This web portal implementation used web services that had the intuitive periodic-table-based interface on top of a rich content management system and a messaging/notification system. An XML Schema was devised to store and retrieve the data and metadata in a Scientific Annotation Middleware layer[3] that also provided capabilities for searching, versioning, locking, access control, and managing provenance. This version of the BSE allowed browsing and downloading of basis sets, user contributions to the BSE, as well as an administrator curation tool. Web services were considered and a Web Services Description Language[4] (WSDL) file was provided for some of the basic services but ultimately the services were not exposed for external access.

The first version of this web-based Basis Set Exchange Web site was created in 2007 at the Pacific Northwest National Laboratory (PNNL) and supported later through the Environmental Molecular Sciences Laboratory (EMSL) at PNNL. The BSE was very popular and was well used throughout the molecular sciences community. It has also become the authoritative source for some basis set data, and many authors directly uploaded their basis sets as a way to publish and distribute their work. It was used by more than 10 000 researchers per month from more than 65 countries. As of August 2019, over 1700 publications cite the suggested reference for the BSE (ref 2), with the original perspective by David Feller (ref 1) garnering over 1500 citations (both according to SciFinder[5]).

The BSE infrastructure was based on a more general collaboratory framework as part of the Collaboratory for Multiscale Chemical Science project[6] and provided substantial functionality that did not get significantly used (such as chat rooms). As the project matured, maintenance of the complex software stack proved to be difficult and the site had significant downtime on a few occasions where it was unavailable to users. In addition, the BSE did not have functionality that the community asked for including the ability to download the whole library (for inclusion in quantum chemistry codes) and an easy, modern callable interface so the data could be connected to other applications; such features would have been very difficult to implement into the existing software stack.

In late 2017, it was decided to retire the original BSE and create a new version, using modern software stacks and best practices. This was to be a joint project between the Molecular Sciences Software Institute[7] (MolSSI) and members of the original team from EMSL and PNNL.

## 2. ARCHITECTURE

The software for the new BSE has been created from scratch, and has been made more modular and lightweight, following generally accepted modern programming design principles. This new architecture is designed to facilitate a wide range of potential use cases, particularly as knowledge of programming and technology increases in the computational chemistry community.

The new BSE is split into three main components. The first is a Python-based library and data repository, and the second is
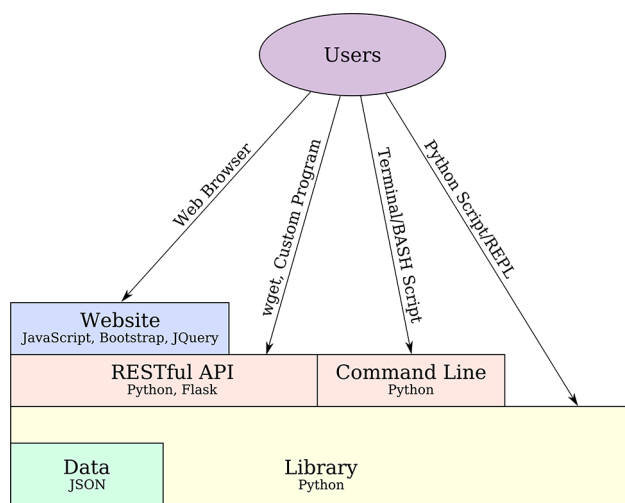


**Figure 1.** Overall architecture of the new BSE and how users would typically access each component. Also given are the programming language and major technologies used. The raw basis data is generally inaccessible to the end user.

a RESTful API (built upon the first component). A publicly accessible Web site is then built upon the library and API (see Figure 1) These components each serve different use cases— while the Web site is convenient for many researchers to browse and download basis sets for one-off calculations, the library can be used by "power users" in a more automated fashion and can even be embedded into computational chemistry codes.

**2.1. BSE Python Library.** The basis set data and manipulation library is the core of the project. Written in Python, it contains all of the raw data for the basis sets, as well as functions to read, manipulate, and convert basis set data. Similarly, reference data (and reference manipulation functions) are also stored in this library.

The BSE Python library is hosted on GitHub (https://github.com/MolSSI-BSE/basis_set_exchange). Users are able to install the library through typical Python distribution channels and use the library (via the API or the command line interface) to obtain basis set data and information. For example, a user may use the library to create many basis set files for a range of quantum chemistry software packages. A particularly powerful use of the library is for a software developer to integrate the library into their own quantum chemistry package, negating the need to manage their own library of basis set data. This would also lead to a consistent definition of basis sets across chemistry codes.

All basis set data (exponent, coefficients, references, etc.) is stored in JavaScript Object Notation format[8] (JSON) that is easily readable from Python and many other popular programming languages. Numerical data is stored as strings to allow for exact matching with published data, which may have more or less than double precision. Data is stored semantically such that future migration to other formats should be relatively simple, should that need arise.

The data for a basis set is divided into several files; the purpose of this is to facilitate construction of basis sets which are formed via combination of individual components, while reducing data duplication. For example, the 6-31G* basis set is formed via a file containing the 6-31G basis and a separate file containing just the polarization functions.
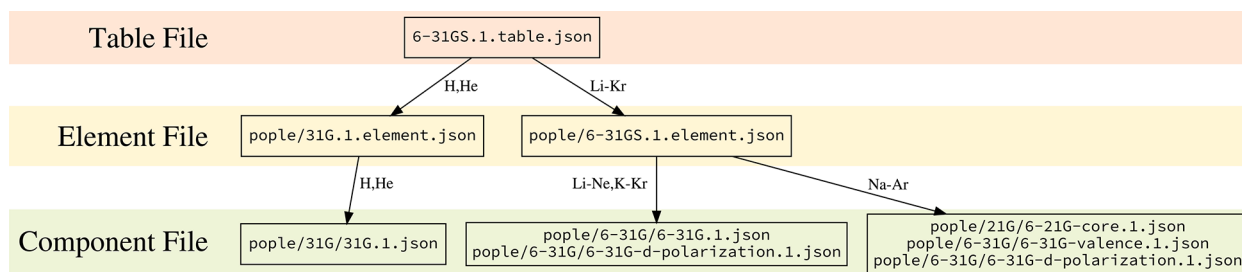
**Figure 2.** Composition of the 6-31G* basis set from component and elemental basis files. Data from the component files is combined within the element files, which is then referenced from within table files.
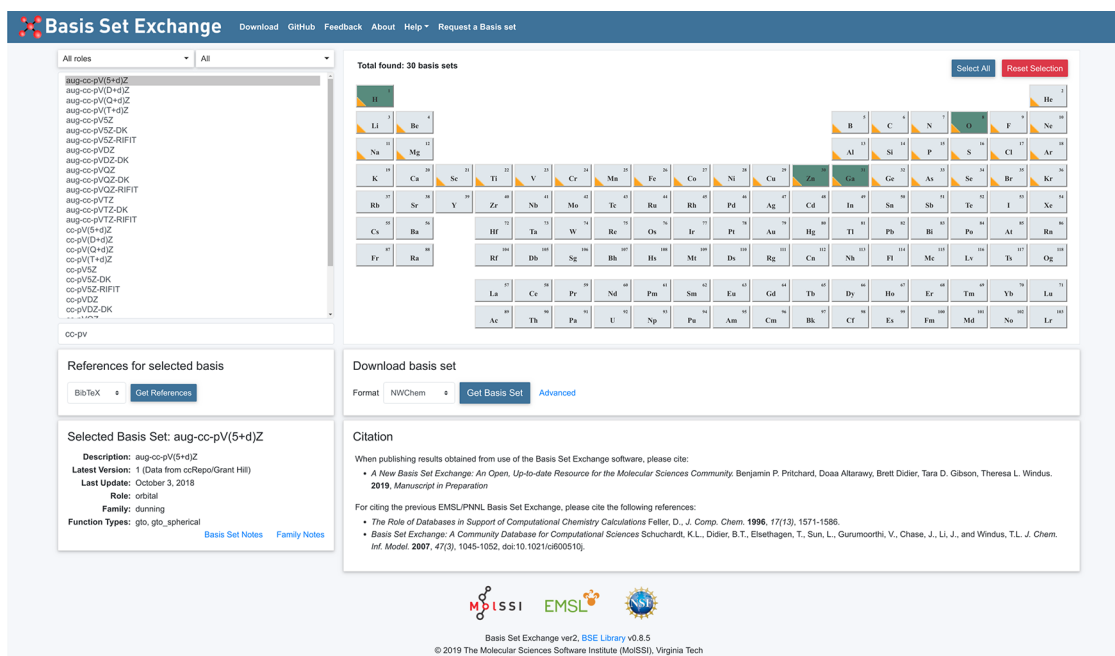


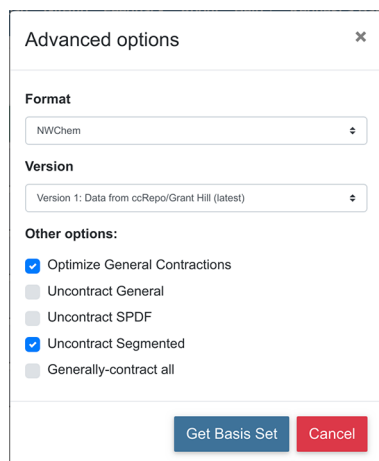**Figure 3.** Screenshot of the new BSE Web interface.



**Figure 4.** Advanced options for downloading basis sets.

Such a scheme requires four separate types of files. At the bottom are *component* files, which contain actual basis set data (exponents, coefficients, effective core potentials, references, etc.). This basis data is stored per-element, and a component file can contain multiple elements. Also contained in these component files is metadata containing provenance information for the data in that file.

Next, this component data is combined via *element* files. An element file for 6-31G* would contain entries pointing to the component file containing 6-31G and the component file containing the polarization. The 6-31G* basis for hydrogen and helium are handled separately and would not contain the polarization functions (see Figure 2).

At the next level, this elemental data is combined for all elements with a *table* basis file. The table basis contains links to an elemental basis set for each element, forming what is generally considered a complete "basis set". This file also contains metadata about the particular version of the basis set described by that file (see section 2.1.2).

The split between elemental and table basis sets can be somewhat awkward. However, this layout is convenient to describe overall basis sets that contain data from several elemental sets. A common example would be 6-31G*, which would contain the hydrogen and helium data from the 31G elemental basis, with the rest of the elements from the 6-31G* elemental basis. The 6-31G** basis set would be similar, except that H and He would be incorporated via the 31G** elemental basis.

The three kinds of files previously described are versioned by incorporating a version number into the filename. This allows the library to handle updates and improvements to basis sets, while keeping old data accessible for future reproducibility
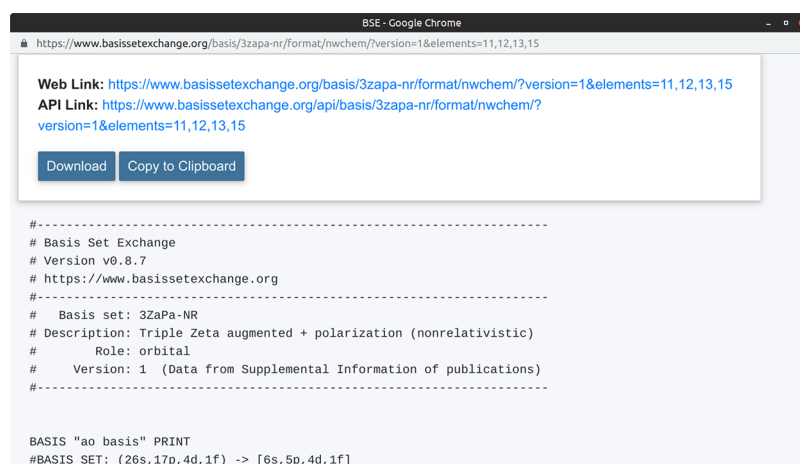
4816

**Figure 5.** Download basis set page for a specific basis set in the NWChem format.

**Table 1. Overview of the Basis Set Exchange REST APIs**[a]

| API URL | Description |
|---|---|
| /api/metadata | get all the metadata describing the basis sets and their elements. |
| /api/formats | get a list of the available download formats of the basis sets |
| /api/reference_formats | get a list of the available download formats of the references |
| /api/basis/<basis>/format/<fmt> | download a basis set in a specific format. Elements can be specified using GET parameters (see REST APIs docs) |
| /api/references/<basis>/format/<fmt> | download references for a specific basis set |
| /api/notes/<basis> | get notes of the given basis set |
| /api/family_notes/<family> | get notes of the given basis set family |

[a]The URL parameters are as follows: <basis> is the name of a basis set, <fmt> is the desired format of the basis set or reference, <family> is the name of a basis set family.

and repeatability purposes. This backward compatibility is something that was widely seen as lacking in the previous BSE.

Metadata (names, descriptions, corresponding auxiliary basis sets, etc.) are stored in the fourth kind of file. It is expected that this metadata would be shared among all versions, and therefore does not contain any version information. In this metadata, basis sets are tagged with a *role*. These roles include "orbital", "jkfit", "jfit", and "rifit", among others. The corresponding auxiliary basis set to be used with a given orbital basis is also stored and can be queried through the library.

Along with the metadata, plain text files containing notes for an individual basis set may also be present. Notes for an entire basis set family ("dunning", for example) are also available. These notes may contain a discussion of the uses, provenance, and data from reference calculations for a basis set or family. Also in these notes, discussion of variations and differences of the basis set in different codes may also be present.

*2.1.1. Manipulation.* The library described in the previous section also contains functionality for manipulating and converting basis sets. These manipulations are made accessible to the user both as part of the main API that is used when obtaining a basis set, as well as standalone functions accessible from Python code, including the REPL environment and Jupyter notebooks.

One main manipulation function is to compose a basis set (for a given set of elements) and then convert this basis set into a format usable by a particular quantum chemistry software package. The top-level API of the library provides functions for this use case, as well as functions for obtaining reference data, notes, and for filtering the list of basis sets. In addition, the option for "Optimized General Contractions"[9]

from the previous BSE is still available to the user. New functionality, such as uncontracting or recontracting basis sets, has also been added.

*2.1.2. Versioning and Testing.* As mentioned previously, data within the library is versioned, such that updating data is possible while retaining previous versions for backward compatibility and reproducibility. Data that is equivalent to the original BSE is pinned to version 0 (zero). Changes to the basis set increment this version number. Basis sets contain a brief description of each version, noting where the data for that version came from or what has changed. A more extensive discussion of what has changed and why it has changed may be present in the notes or family notes.

The policy of what constitutes a "change" in a basis set is somewhat imprecise. We have decided that only modifications which may alter the results of previous calculations should be considered a new version. Therefore, changing exponents or coefficients would certainly force a new version; updating metadata, references, descriptions, or notes should not require a new version. Similarly, adding basis functions for elements previously missing from the basis set should not increment the version, as the results from repeating previous calculations would remain identical. The one exception to this is that version 0 (representing the data from the original BSE) will contain only the elements present in the original BSE.

To help keep the data consistent, and to ensure reliability in the manipulation pipelines, extensive testing of the library and data is implemented. All basis set data (formatted for the various codes) has been downloaded from the original BSE, and output from the new BSE is regularly compared against this previous output using continuous integration. When new data is added from an authoritative source (for example, from
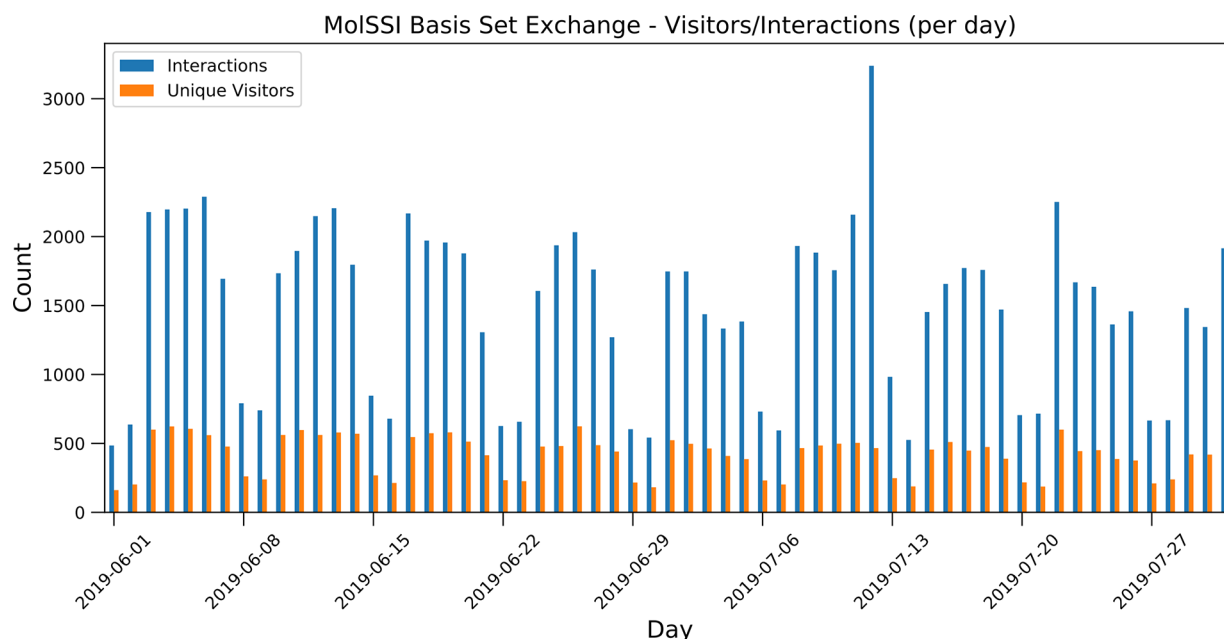
**Figure 6.** Unique users and interactions with the BSE Web site for June and July, 2019.

Supporting Information of a paper or from data files sent directly from the author), this reference data is included in the test set as well. In this fashion, accidental mistakes or modifications can be caught in an automated fashion, without requiring extensive manual review.

Other tests are implemented as well, such as detecting when data is present that is not part of a basis set, and making sure manipulations work properly. Data is also tested against a schema. These tests ensure that the data remains clean and reliable.

*2.1.3. Command-Line Interface.* The library also contains a command-line interface (CLI), allowing for easy use of commonly used functionality (such as obtaining a basis set or references in a particular format). This allows a user to explore the data and manipulations without requiring programming knowledge. This also allows for integration of the library into workflows that do not directly use Python. TAB completion is present, making it easier for users to input commands or awkward basis set names.

**2.2. BSE Web Site and RESTful APIs.** The new Basis Set Exchange is built using modern Web development technologies and applying Model-view-controller (MVC) principles. Application of MVC means that data is separated from its presentation which allows more flexibility and extensibility. The Web site uses the BSE Python library (section 2.1) to extract the metadata in order to populate the web interface and browse the basis sets. This means that the BSE Web site can always be up-to-date with any new changes in the library (such as addition of new basis sets or implementation of new formats).

The Web site is built using Flask,[10] a popular Python web framework, which relies on the Jinja[11] templating system to render web pages. The front end is designed using Jinja templates, JavaScript, and the Bootstrap CSS[12] library for formatting HTML and allowing a uniform look and feel. The Web site is designed to be responsive, meaning that it works on any device from desktops to mobile phones. The Web site back end uses MongoEngine[13] and MongoDB,[14] a NoSQL database

system, for storing information about the usage of the Web site and basis sets.

With the new BSE, we tried to maintain the overall structure of the previous BSE that thousands of researchers worldwide are familiar with. As shown in Figure 3, the main screen includes the periodic table on the right (updated with new elements), and a list of available basis sets on the left. The user can perform filtering on the basis sets using various filters and then download a specific basis set in some chosen format.

*2.2.1. Web Site Features.* The BSE Web site provides a convenient and easy way to access the BSE Python library data using any device from anywhere in the world. The new BSE is designed to load quickly with low response time—data is fetched from the backend server only when needed. Any interaction, other than downloads, is performed on the client side without unnecessary remote server access.

One of the main features of the Web site is the ability to filter the list of available basis sets using different options such as roles and whether a basis set is all-electron or contains ECPs. Additionally, basis sets can be filtered by typing any part of its name in the search box below the basis sets list. The user can select/deselect elements from the periodic table by clicking them to further filter the available basis sets. Two of the new features added to the new BSE is the ability to "Select All" of the elements in the periodic table for which a basis set is defined and to clear the selection; these features were requested by many users of the old Web site. The number of available basis sets after all the combined filtering options are applied is shown on the top of the periodic table.

After filtering and selecting a specific basis set from the list, the user can download the basis set data. This can be done quickly by selecting the download format from the list (such as NWChem,[15] Gaussian,[16] Psi4,[17] etc.) and then clicking the "Get Basis Set" button (similar to the previous BSE). New basis sets formats are added to the Web site automatically when new formats are available in the BSE Python library. Alternatively, the user can choose the advanced download option (Figure 4). This option is new in the BSE and gives more customized options for the user such as selection of a

specific version of the basis set, optimization of general contractions, and various manipulations of contracted/uncontracted basis sets. Finally, the basis set can be viewed, copied to the clipboard, or downloaded as a file (Figure 5). The URL given when viewing the basis set contains all of the information required to obtain the data, and can be shared with collaborators or the community at large.

One of the advantages of the new BSE Web site is the ability to download the newly curated reference list of a selected basis set. Those references have been revised and updated in the BSE Python library and are available for download through the Web site as plain text, JSON, or BibTeX. More formats are expected to be added in the future. Lastly, the Web site now provides the ability to download the complete collection of basis sets in a specific basis format (such as NWChem or Gaussian).

Users and basis set developers can request new basis sets and new formats through the Web site by filling out a web form, which is then submitted to the BSE support team.

*2.2.2. REST APIs.* In addition to the Web site and the Python interfaces, we provide read-only access to the BSE data through RESTful APIs which can be used from any programming language. REST (REpresentational State Transfer) APIs are a uniform way to access data that is usually used over HTTP. Our REST APIs are stateless, meaning that each request is independent from another and contains all the necessary information to retrieve the data. The REST APIs have a uniform interface to access basis sets, providing a standard way for the user to communicate with the server to query and download the desired data.

Table 1 shows an overview of the REST APIs of the BSE Web server. They can be called using any programming language that supports HTTP requests. A full explanation of how to use the BSE REST APIs can be found in the documentation which is linked in the help menu of the Web site. Each URL should start with the Web site URL, i.e., https://www.basissetexchange.org. For example, to get the metadata, use the API URL https://www.basissetexchange.org/api/metadata. The HTTP request method for all the APIs in Table 1 is GET since the data is read only.

## 3. USAGE STATISTICS

The new BSE collects anonymous statistics related to how the Web site is used. The previous BSE was retired on June 1, 2019 after which all traffic was redirected to the new Web site. Since then, the BSE receives about 9000 to 10 000 unique visitors a month, performing about 45 000 interactions (viewing the homepage, downloading a basis set or reference, etc). Figure 6 shows the daily unique users and interactions for this time period. The BSE has users from all over the world, with most users coming from the United States and China. Consistently, Microsoft Windows is the operating system most commonly used to view the Web site (approximately 59%), with Apple MacOS and various Linux distributions making up most of the rest. Although the BSE is not targeted at mobile platforms, about 5% of the visits use Apple iOS.

Table 2 shows the most popular basis sets and basis set formats during the same time period. The Ahlrichs (def2) and Dunning families are very popular, while historically important basis sets such as STO-3G and 6-31G also remain in the top ten. The Gaussian format is by far the most popular, accounting for more than half of all downloads.

**Table 2. Most Popular Basis Sets and Basis Set Formats (by Number of Downloads) for a 61-Day Period (June and July 2019)**

| Basis set | Downloads | |
|---|---|---|
| def2-TZVP | 1816 | |
| LANL2DZ | 1304 | |
| def2-SVP | 1059 | |
| aug-cc-pVTZ | 791 | |
| STO-3G | 760 | |
| def2-TZVPP | 728 | |
| cc-pVDZ | 721 | |
| cc-pVTZ | 647 | |
| 6-31G | 608 | |
| aug-cc-pVDZ | 564 | |
| **Basis set format** | **Downloads** | |
| Gaussian | 27543 | 56.55% |
| NWChem | 8176 | 16.79% |
| GAMESS US | 5208 | 10.69% |
| Molpro | 2337 | 4.80% |
| Turbomole | 1386 | 2.85% |
| CFour | 876 | 1.80% |
| Dalton | 645 | 1.32% |
| Psi4 | 635 | 1.30% |
| Molcas | 594 | 1.22% |
| Gamess UK | 409 | 0.84% |

## 4. CONCLUSION

The BSE has served the computational chemistry community very well since it debuted more than a decade ago. However, its legacy codebase was seen as a maintenance burden and could not adapt to a changing environment and userbase, where new features could not be easily implemented. The BSE has been rewritten using more modern programming languages and techniques, emphasizing modularity and usability as primary features.

It is hoped that this new architecture remains adaptable to the future needs of the community and that the project will be seen as the standard, authoritative source for basis set information well into the future.

## ■ AUTHOR INFORMATION

**Corresponding Author**
*E-mail: bpp4@vt.edu.
**ORCID** ⊕
Benjamin P. Pritchard: 0000-0003-2136-0606
Doaa Altarawy: 0000-0002-7795-4422
Theresa L. Windus: 0000-0001-6065-3167

**Author Contributions**
B.P.P. is the lead developer and responsible for data curation, migration of data from the old BSE, and the basis set Python library. D.A. is the main developer of the new BSE web application (the web interface, REST APIs, and the statistics database). B.D., T.D.G., and T.L.W. contributed to the requirements and design and managed the collaborative details. All authors read and approved the manuscript.

**Author Contributions**
[⊥]B.P.P. and D.A. contributed equally.

**Notes**
The authors declare no competing financial interest.

## ■ REFERENCES

(1) Feller, D. The role of databases in support of computational chemistry calculations. *J. Comput. Chem.* **1996**, *17*, 1571−1586.

(2) Schuchardt, K. L.; Didier, B. T.; Elsethagen, T.; Sun, L.; Gurumoorthi, V.; Chase, J.; Li, J.; Windus, T. L. Basis set exchange: A community database for computational sciences. *J. Chem. Inf. Model.* **2007**, *47*, 1045−1052.

(3) Myers, J.; Chappell, A.; Elder, M.; Geist, A.; Schwidder, J. ReIntegrating the Research Record. *Comput. Sci. Eng.* **2003**, *5*, 44−50.

(4) Roberto Chinnici, R.; Moreau, J.-J.; Ryman, A.; Weerawarana, S. *Web Services Description Language (WSDL)* Version 2.0 Part 1: Core Language; W3C Recommendation, 2007.

(5) *SciFinder.* https://scifinder.cas.org/, [Online; accessed August 2019].

(6) Myers, J.; Allison, T.; Bittner, S.; Didier, B.; Frenklach, M.; Green, W.; Ho, Y.; Hewson, J.; Koegler, W.; Lansing, C.; et al. A Collaborative Informatics Infrastructure for Multi-Scale Science. *Cluster Comput* **2005**, *8*, 243−253.

(7) Krylov, A.; Windus, T. L.; Barnes, T.; Marin-Rimoldi, E.; Nash, J. A.; Pritchard, B.; Smith, D. G. A.; Altarawy, D.; Saxe, P.; Clementi, C.; et al. Perspective: Computational chemistry software and its advancement as illustrated through three grand challenge cases for molecular science. *J. Chem. Phys.* **2018**, *149*, 180901.

(8) Bray, T. The JavaScript Object Notation (JSON) Data Interchange Format; *RFC 8259,* **2017**.

(9) Hashimoto, T.; Hirao, K.; Tatewaki, H. Comment on Dunning's correlation-consistent basis sets. *Chem. Phys. Lett.* **1995**, *43*, 190−192.

(10) Ronacher, A. *Flask: a lightweight web application framework.* 2010−; https://flask.palletsprojects.com/, [Online; accessed September 2019].

(11) Ronacher, A. Jinja: a full-featured template engine for Python. 2010−; https://jinja.palletsprojects.com/, [Online; accessed September 2019].

(12) *Bootstrap: The most popular HTML, CSS, and JS library in the world.* 2010−; https://getbootstrap.com/, [Online; accessed August 2019].

(13) Marr, H. *MongoEngine: a Document-Object Mapper for working with MongoDB from Python.* 2009−; http://mongoengine.org/, [Online; accessed September 2019].

(14) MongoDB, I. *MongoDB: The database for modern applications.* 2009−; https://www.mongodb.com/, [Online; accessed September 2019].

(15) Valiev, M.; Bylaska, E.; Govind, N.; Kowalski, K.; Straatsma, T.; Dam, H. V.; Wang, D.; Nieplocha, J.; Apra, E.; Windus, T.; et al. NWChem: A comprehensive and scalable open-source solution for large scale molecular simulations. *Comput. Phys. Commun.* **2010**, *181*, 1477−1489.

(16) Frisch, M. J.; et al. *Gaussian 16*, revision C.01; Gaussian Inc.: Wallingford CT, 2016.

(17) Parrish, R. M.; Burns, L. A.; Smith, D. G. A.; Simmonett, A. C.; DePrince, A. E.; Hohenstein, E. G.; Bozkaya, U.; Sokolov, A. Y.; Di Remigio, R.; Richard, R. M.; et al. Psi4 1.1: An Open-Source Electronic Structure Program Emphasizing Automation, Advanced Libraries, and Interoperability. *J. Chem. Theory Comput.* **2017**, *13*, 3185−3197.