

1 eplusr: A framework for integrating building energy simulation and
2 data-driven analytics

3 Hongyuan Jia^a, Adrian Chong^{b,*}

4 ^a*SinBerBEST Program, Berkeley Education Alliance for Research in Singapore, Singapore, 138602, Singapore*

5 ^b*Department of Building, School of Design and Environment, National University of Singapore, 4 Architecture Drive,
6 Singapore, 117566, Singapore*

7 **Abstract**

Building energy simulation (BES) has been widely adopted for the investigation of building environmental and energy performance for different design and retrofit alternatives. Data-driven analytics is vital for interpreting and analyzing BES results to reveal trends and provide useful insights. However, seamless integration between BES and data-driven analytics current does not exist. This paper presents eplusr, an R package for conducting data-driven analytics with EnergyPlus. The R package is cross-platform and distributed using CRAN (The Comprehensive R Archive Network). With a data-centric design philosophy, the proposed framework focuses on better and more seamless integration between BES and data-driven analytics. It provides structured inputs/outputs format that can be easily piped into data analytics workflows. The framework also provides an infrastructure to bring portable and reusable computational environment for building energy modeling to facilitate reproducibility research.

8 **Highlights**

- 9 1. Developed an R package that integrates EnergyPlus with data-driven analytics
- 10 2. Structured inputs/outputs format that can be easily piped into data analytics workflows
- 11 3. Facilitates reproducible simulations through Docker
- 12 4. Enables flexible and extensible parametric simulations

13 **1. Introduction**

14 Building energy simulation (BES) is increasingly being used throughout the building's life-cycle for the
15 analysis and prediction of building energy consumption, measurement and verification, carbon evaluation,
16 and cost analysis of energy conservation measures (ECMs) [1,2]. It has played a growing role in the design and
17 operation of low energy, high-performance buildings, and development of policies that drive the achievement
18 of reducing energy use and greenhouse-gas emissions in the buildings sector [3].

19 BES offers an alternative approach that encourages customized, integrated design solutions, and the
20 development of BES tools has been pronounced over the decades [3–5]. The core tools are the whole-
21 building energy simulation programs that provide users with key building performance indicators such as
22 energy use and demand, temperature, humidity, and costs [6]. However, BES, with an iterative nature
23 inside, can produce a large amount of data. The volumes of the data have overwhelmed traditional data
24 analysis methods such as spreadsheets and ad-hoc queries with a large number of factors to be considered
25 [7]. Solutions in most existing software and applications have limited post-processing capacities on BES
26 results. They are not flexible enough to enable a clear understanding and control of how the data is being

*Corresponding Author

Email addresses: hongyuan.jia@bears-berkeley.sg (Hongyuan Jia), adrian.chong@nus.edu.sg (Adrian Chong)

transformed [8,9]. According to a survey of 448 building energy management professionals in the U.S., there is a need to improve the efficacy and integration between data-driven analytics and BES, and efforts should be made to develop integrated tools that are capable of leveraging both methods [10].

As BES becomes more integral to many aspects of architecture design and decision-making processes, computational reproducibility has become increasingly important to researchers, designers and practitioners. Lack of credibility in BES results due to a lack of reproducibility is widely considered a problem by the energy modeling community [11]. Issues in simulation reproducibility are mainly caused by the absence of (1) an integrated workflow between BES and data-driven analytics and (2) a portable and reusable computational environment encapsulating essential software and applications to perform it.

To address these issues, this paper introduces a new framework for integrating BES and data-driven analytics. The framework is different from existing ones because of its data-centric design philosophy. The objectives are (1) to provide better and more seamless integration between BES engine EnergyPlus and R-programming data-driven analytics environment and (2) to build infrastructures for portable and reusable BES computational environment to facilitate reproducibility research in building energy domain. Section 2 describes the state of the art of related research and development. Section 3 introduces the concepts behind the framework, along with its implementation. Section 4 demonstrates the applications of the framework using a medium office building model with four examples, covering various topics, including data exploration, parametric simulation, optimization and calibration.

2. State of the art

2.1. BES software and applications

Over the decades, a wide variety of whole-building energy simulation programs have been developed [6]. In general, they can be classified into three categories [12], including:

1. Applications with integrated simulation engine (e.g. EnergyPlus [13], TRNSYS [14], DeST [15], ESP-r [16], IESVE [17], IDA ICE [18])
2. Software that based on a certain engine (e.g. eQUEST [19], DesignBuilder [20], OpenStudio [21], jEplus [22], Modelkit [23], GenOpt [24])
3. Plugins for other software enabling certain performance analysis (e.g. Ladybug & Honeybee [25,26], eppy [27], MLE+ [28,29], EpXL [30])

Some tools may fall into several categories, such as IESVE and IDA ICE which can also be treated as an interface software toolkit to its engine, and software OpenStudio and DeST which also provides plugins for other software to perform geometry creation and manipulation.

Choosing the appropriate combination of design options using BES is a complex task that requires the management of a large amount of information on the properties of design options and the simulation of their performance [31]. Parametric energy simulation is often needed to take into account the uncertainties and variability of different design variables [32]. However, parametric analysis involves tedious file management tasks, repeated entry of model parameters, the application of design transformations and the execution of large-scale analyses [33], which can be time-consuming and error-prone. Parametric simulation task automation has been proven to be a useful way to reduce human intervention and improve the efficiency of large parametric analysis [34]. Table 1 gives a summary of the characteristics and capabilities of various BES software and plugins for this purpose.

In Table 1, some tools consist of graphical user interfaces (GUIs), while others use general-purpose scripting languages accompanied by a suite of programming features and libraries [34]. Among the tools, EnergyPlus is the most used simulation engine. This may be due to its advantage of free, open-source and cross-platform characteristics.

OpenStudio [21] is a free, open-source software toolkit designed for energy modeling and can be used to efficiently create or modify models, manage individual or multiple simulations, and visualize results. OpenStudio has its own format (.OSM) and schema for EnergyPlus model representation which will eventually be translated into EnergyPlus models. Parametric Analysis Tools (PAT) is a GUI application that is part

138 have constraints on flexibility as the users have to specify exactly what and how features of the design can be
139 manipulated and often are not be able to provide a good workflow for repeating that task across a broader
140 range of situations on different systems. In this case, manual steps have to be performed using other tools,
141 such as a spreadsheet or command-line tools, which introduces additional transcription burden and results in
142 a non-reproducible process [33]. Sometimes, custom solutions have to be created from scratch to automate
143 part portions of the workflows, which may lead to new inefficiencies and potential errors. Currently, no
144 widely adopted solution is able to integrate all processes into one single platform.

145 Moreover, BES often involves the use of multiple applications, software and platforms. To perform crucial
146 scientific processes such as replicating the results, extending the approach or testing the conclusions in other
147 contexts, the indispensable step is to install the software used by the original researchers, which sometimes
148 can become immensely time-consuming if not impossible. It is easy to underestimate the significant barriers
149 raised by a lack of familiar, intuitive, and widely adopted tools for addressing the challenges of computational
150 reproducibility [45].

151 3. Methodology

152 To achieve seamless integration between BES and data-driven analytics, we propose a framework con-
153 sisting of the following (Fig. 1):

- 154 1. I/O processors for structuring BES inputs and outputs for seamless integration with data analytics
155 workflow.
- 156 2. A parametric prototype for conducting flexible and extensible parametric simulations.
- 157 3. A computational environment that is based on Docker containerization [47] to facilitate reproducibility
158 research in the energy simulation domain.

159 The first two components have been packaged into a free, open-source R package *eplusr*¹ which is dis-
160 tributed using CRAN (The Comprehensive R Archive Network). The third component has been encapsulated
161 using Docker containerization and is distributed using Docker Hub².

162 3.1. I/O processors

163 The I/O processors are implemented through three modules shown in Fig. 1, including:

- 164 1. Relational Database module to represent EnergyPlus models and weathers in relational databases,
- 165 2. Object-Oriented Programming (OOP) Model API module for tidy data model modification APIs
- 166 3. Tidy Data Interface module for querying and structuring BES outputs in tidy format.

167 3.1.1. Relational databases

168 The Relational Database module is developed to read, parse and represent EnergyPlus models and weath-
169 ers in relational databases. EnergyPlus Input Data File (IDF) is based on the data schema that is defined in
170 the Input Data Dictionary (IDD). In the proposed framework, data of an IDF and the corresponding IDD
171 are stored as Relational Databases (RD). RD was first proposed by Codd [48] and has become the dominant
172 database model for a number of Relational Database Management Systems (RDMS). It organizes data in a
173 set of rectangular tables with rows and columns. Each table has a primary key which is an unique identifier
174 constructed from one or more columns. A table is linked to another by including the other table's primary
175 key (also called a foreign key).

176 Fig. 2 shows the structure of RD for an EnergyPlus IDF and IDD. The RD data structure follows
177 the idea of database normalization where each variable is expressed in only one place, avoiding any data
178 redundancy. The hierarchy structure of the IDF data schema is retained through various tables. Data
179 integrity is maintained via relations among table variables. Each RD has a **reference** table to store the
180 referencing relations among various field values. To modify an IDF is equal to change the corresponding
181 fields in its RD tables. The RD structure provides the capability to quickly perform data wrangling and
182 fast table joining among entities and variables.

¹Github Repository: <https://github.com/hongyuanjia/eplusr>

²Docker Hub Link: <https://hub.docker.com/r/hongyuanjia/eplusr>

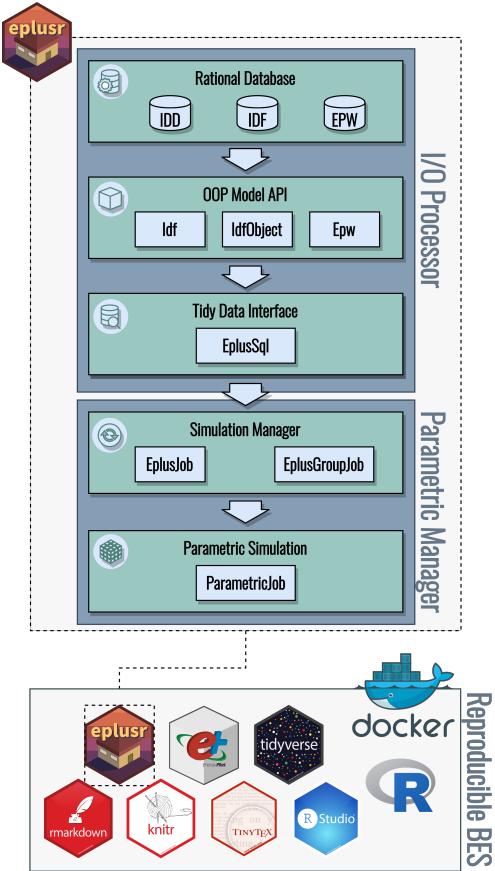


Figure 1: An architecture overview of the proposed framework which includes three main components: (1) I/O processors, (2) Parametric prototype and (3) Computational environment for reproducible BES using Docker containerization

183 3.1.2. *Object-oriented programming model API*

184 The Object-oriented programming (OOP) Model API module enables users to perform queries and
 185 modifications on EnergyPlus models programmatically. OOP [49] is a programming paradigm that focuses
 186 on the objects to manipulate rather than the logic required to manipulate them. It provides a clear modular
 187 structure for programs and is good for defining abstract data types. OOP hides implementation details and
 188 makes it possible to develop a clearly defined interface for each abstraction.

189 Fig. 3 gives an overview of the OOP Model API module. It introduces three groups of classes, including
 190 (1) `Idd` class and `IddObject` class for a whole and part of an IDD, (2) `Idf` class and `IdfObject` class for
 191 a whole and part of an IDF, and (3) `Epw` class for an EPW (EnergyPlus Weather). Each class provides a
 192 number of methods to manipulate the encapsulated data. An extensive rule-based data model validator has
 193 been developed to check the integrity of data before any modifications.

194 `Idf` class exposes flexible interfaces to modify field values in different scope levels, including single-object
 195 level, grouped-object level, and whole-class level, enabling to alter a number of objects at the same time.
 196 Both `Idf` and `IdfObject` class provide a `to_table()` method to extract certain or all parts of a model into
 197 one `data.table` object, which is an extension of R's table representation but extremely optimized for fast
 198 computation [50]. The `load()` and `update()` methods in `Idf` class can take any model data in table format
 199 as input, create and modify large number of objects accordingly. Section 4 demonstrates some of the APIs
 200 in this module.

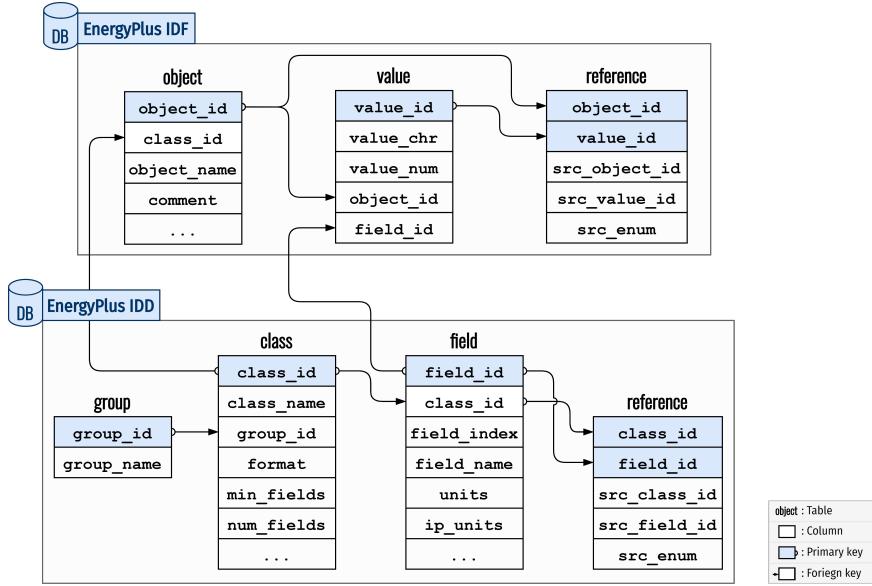


Figure 2: Structure of relational databases for an EnergyPlus IDF and IDD

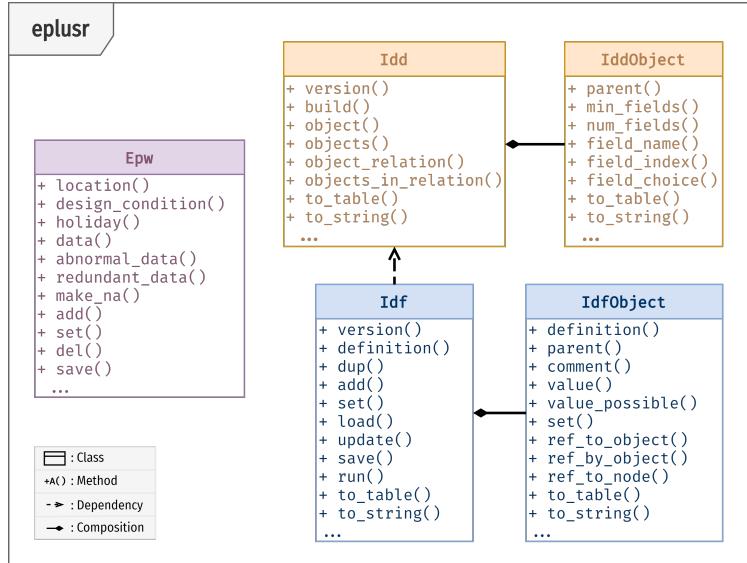


Figure 3: Overview of OOP Model API

3.1.3. Tidy data interface

The Tidy Data Interface module is designed to extract and represent EnergyPlus simulation results from the SQLite output into tidy tables. The concept of tidy data format was first proposed by Wickham [51], as a standard way of mapping the meaning of a dataset to its structure. It means that each variable forms a column, each observation forms a row, and each type of observational unit forms a table (see Table (b) in Fig. 4). This structure makes it intuitive for an analyst or a computer to extract needed variables. It is particularly suited for vectorized programming languages like R. The layout ensures that values of different variables from the same observation are always paired [51,52] and is well fitted for data analyses using the *tidyverse* R package ecosystem [53].

Table (a) in Fig. 4 shows an example of the standard format from EnergyPlus CSV table output, while

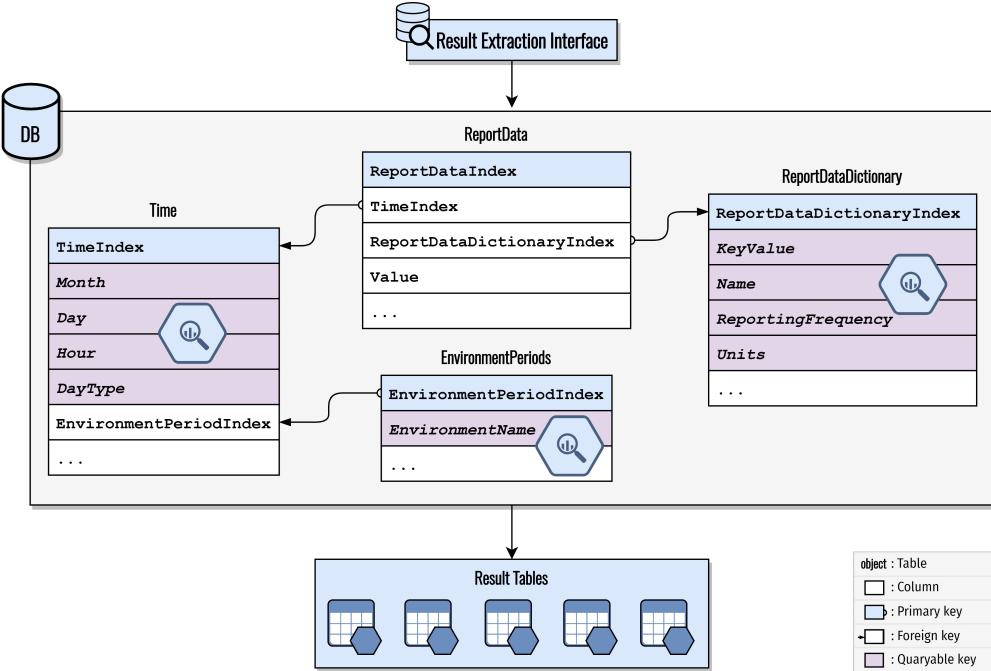


Figure 5: Overview of the tidy data interface for variable and meter outputs

239 3.2. *Parametric prototype*

240 The parametric prototype in the framework provides a set of abstractions to ease the process of parametric
 241 model generation, design alternative evaluation, and large parametric simulation management. An overview
 242 of the parametric prototype implementation is shown in Fig. 6.

243 A parametric simulation is initialized using a seed model and a weather file. Design alternatives are
 244 specified by applying a *measure* function to the seed model. The concept of *measure* in the prototype is
 245 inspired by a similar concept in OpenStudio [21] but tailored for flexibility and extensibility. A measure is
 246 simply an R function that takes an *Idf* object and any other parameters (e.g. t_1 to t_5 in Fig. 6) as input,
 247 and returns a set of modified *Idf* objects as output, making it possible to leverage other modules in the
 248 framework and apply statistical methods and libraries existing in R to generate design options. After a
 249 measure is defined, the method *apply_measure()* takes it and other parameter values specified to create a
 250 set of models. The *run()* method will run all parametric simulations in parallel and place each simulation
 251 outputs in a separate folder. All simulation metadata will keep updating during the whole time and can be
 252 retrieved using the *status()* method for further investigations.

253 The *ParametricJob* class leverages the tidy data interface to retrieve parametric simulation results in
 254 a tidy format. Despite that, a number of methods are also provided to read various output files, includ-
 255 ing simulation errors (*eplusout.err*), report data dictionary (*eplusout.rdd*), and meter data dictionary
 256 (*eplusout.mdd*). For all resulting tidy tables, an extra column containing the simulation job identifiers is
 257 prepended in each table. It can be used as an index or key for further data transformations, analyses and
 258 visualization to compare results of different simulated design options.

259 The proposed parametric prototype is designed to be simple yet flexible and extensible. One good
 260 example of the extensibility of this framework is the *epluspar*⁴ R package, which provides new classes
 261 for conducting specific parametric analyses on EnergyPlus models, including sensitivity analysis using the
 262 Morris method [55] and Bayesian calibration using the method proposed by Chong [1]. All the new classes
 263 introduced are based on the *ParametricJob* class. The main difference mainly lies in the specific statistical

⁴Github Repository: <https://github.com/hongyuanjia/epluspar>

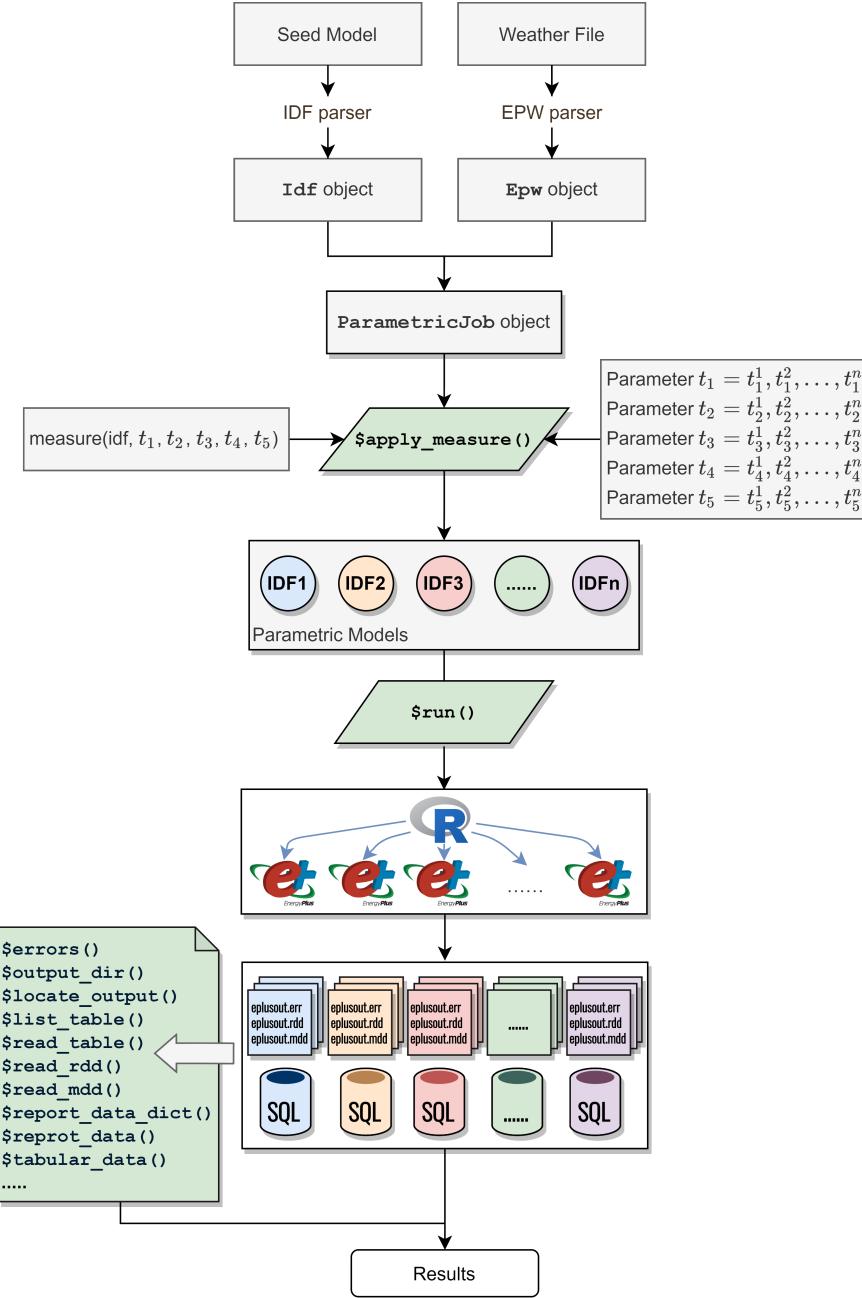


Figure 6: Workflow of a parametric simulation

264 method used for sampling parameter values when calling `apply_measure()` method. Few examples of this
 265 application have been provided in Section 4.

266 3.3. Computational environment for reproducible BES

267 The Docker containerization for BES aims to provide infrastructures to bring portable and reusable
 268 computational environments to facilitate reproducible BES applications. Peng [44] summarized two major
 269 components to successful reproducible research: (1) data, i.e. the availability of raw data from the experiment,
 270 and (2) code, i.e. the availability of the statistical code and documentation to reproduce the results.

271 In the context of BES, these two component will be (1) the building energy models and (2) the code to
272 perform simulations and following data-driven analytics. However, the complex and rapidly changing nature
273 of computer environments makes it immensely challenging to reproduce the same workflow and results even
274 with the original data and code are available. To address this issue, a reproducible BES computational
275 environment has been developed based on the Docker containerization technology, which captures the full
276 software stack, including all software dependencies in a portable and reusable image.

277 Docker [47] is a popular open-source tool for containerization and has shown its potential to improve
278 computational reproducibility [45,56]. The Rocker Project was launched in 2014 as a collaboration to provide
279 high-quality Docker images containing the R environment and has seen both considerable uptakes in the R
280 community and substantial development and evolution [57]. The proposed reproducible BES computational
281 environment is built upon the `rocker/verse` images. It contains four groups of toolchains needed for
282 common BES and data-driven analytics workflows using the eplusr framework:

- 283 1. Statistical computing environment, including the latest R environment and RStudio Server, a web-
284 based integrated development environment for R programming
- 285 2. BES engine, including EnergyPlus of specified version and the eplusr R package
- 286 3. Data analytics toolkits, including a collection of tidyverse [53] R packages for data import, tidying,
287 manipulation, visualization and programming
- 288 4. Literate programming environment, including R Markdown related packages for dynamic document
289 generation

290 The first three have been described in previous sections. Literate programming is a programming
291 paradigm introduced by Knuth [58] in which the explanation of a computer program is given, together
292 with snippets of source code. Recently, there have been significant efforts to develop literate programming
293 infrastructure to reproducibly perform and communicate data analyses, including R Markdown [59], Jupyter
294 notebook [60], just to name a few.

295 The R Markdown format is powered by the knitr R package [61] and Pandoc [62]. Knitr executes
296 the computer code written in various programming languages embedded, and converts R Markdown to
297 Markdown. Pandoc processes the resulting Markdown and render it to various output formats, including
298 PDF, HTML, Word, etc. The R Markdown format has been a widely adopted authoring framework for
299 data science. It can be used to both save and execute code and generate high-quality reports that can
300 be shared with an audience. Together with `rmarkdown` [63] and `tinytex` [64] packages, the proposed BES
301 computational environment can be easily adapted to any R-centric workflows and enables researchers in the
302 BES field to build and archive reproducible analytics.

303 The source files of Docker configuration were written in several text files so-called Dockerfiles and are
304 publicly available and hosted via GitHub⁵. Further evolutions can be taken to make the computational
305 environment tailored to different audiences and use purposes. The docker approach is suited for moving
306 between local and cloud platforms when a web-based integrated development environment is available, such
307 as RStudio Server [45], providing the scalability potential for large cloud-based BES computation.

308 4. Applications

309 To show how the eplusr framework can be used, examples are presented in four topics: (1) data exploration,
310 (2) parametric simulation, (3) multi-objective optimization (MOO) using Genetic Algorithm (GA),
311 and (4) Bayesian calibration. For all examples, the U.S. Department of Energy (DOE) medium office
312 reference building model in compliance with Standard ASHRAE 90.1 – 2004 [65] is used. Fig. 7 shows a
313 3D view of the building geometry. It is a 3-story, 15-zone medium office building with a total floor area
314 of 4982 m². A central packaged air conditioning unit with a gas furnace is equipped on each story. The
315 air distribution systems are Variable Air Volume (VAV) terminal boxes with electric reheating coils. The
316 typical meteorological year 3 (TMY3) weather data of Chicago was used for all the simulations.

⁵Github Repository: <https://github.com/hongyuanjia/eplusr-docker>

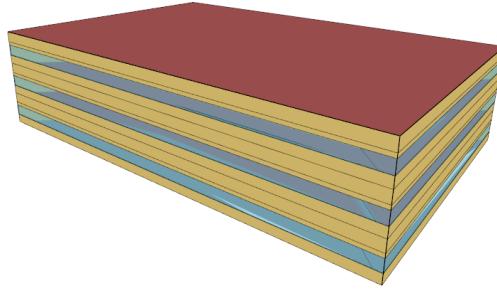


Figure 7: 3D view of DOE medium office reference building

317 *4.1. Data exploration*

318 Data exploration is an essential aspect of BES. It is often used to reduce large volumes of simulation
 319 data to a manageable size so that efforts can be focused on analyzing the most relevant data. This example
 320 demonstrates the data exploration process of obtaining (1) energy use intensity (EUI) and (2) heating and
 321 cooling demand profile using annual simulation results. The energy use intensity (EUI) is one key indicator
 322 for building energy performance [66], and its breakdown can provide potential directions of where ECMS
 323 should be applied to reduce energy usage. When evaluating the feasibility of free-cooling applications in
 324 buildings, the heating and cooling demand profile plays an essential role in the determination of the potential.
 325 This example showcases the basic features of the proposed framework with the main focus on how the tidy
 326 data interface can provide a seamless workflow to extract BES output, feed it into data analysis pipelines
 327 and turn the results into understanding and knowledge.

328 Fig. 8 shows an overview of a typical data exploration workflow using BES output extracted by the tidy
 329 data interface described in Section 3.1.3. Listing 1 shows the R code to achieve it.

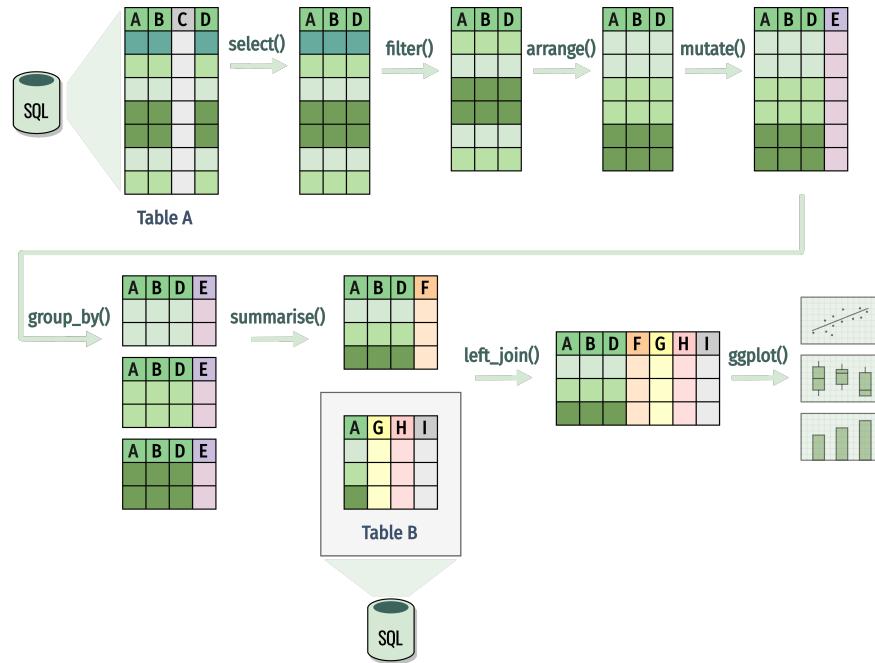


Figure 8: Workflow overview of data exploration using BES output extracted by the tidy data interface

330 Lines 38 – 53 in Listing 1 shows how to use methods `tabular_data()`, `read_table()` and `report_data()`
 331 provided by the tidy data interface to extract building area and building energy consumption, zone meta-

332 data, and cooling and heating demands, with all formatted in a tidy representation. Note that instead of
 333 presenting the simulated date and time as strings, the `report_data()` adds a time-series column `datetime`
 334 in `POSIXct` based on a derived year value using the algorithm described in Section 3.1. Moreover, the
 335 tidy data interface also provides a number of additional columns shown in Fig. 5, which makes it quite
 336 convenient and straightforward to directly perform further data transformations. Lines 59 – 124 in Listing
 337 1 demonstrate the benefits of the tidy format in selecting columns using `select()`, subsetting rows using
 338 `filter()`, sorting rows using `arrange()`, adding new variables using `mutate()`, summarizing data using
 339 a combination of `group_by()` and `summarize()`, joining tables using `left_join()`, and data visualization
 340 using `ggplot()`.

341 Based on the building energy consumption data (line 42 in Listing 1) and the building area (line 39 in
 342 Listing 1), the electricity consumption breakdown from various end-use categories was calculated, and a pie
 343 chart was created (shown in Fig. 9) using only 15 lines of codes (line 60 – 82 in Listing 1). From Fig. 9,
 344 we can see that most of the energy has been consumed by interior electric equipment, followed by indoor
 345 lighting. It indicates that ECMS which help to reduce the plug loads and lighting power density (LPD)
 346 may have a promising potential in improving the overall energy performance. We will perform further
 347 investigations on this in Section 4.2.

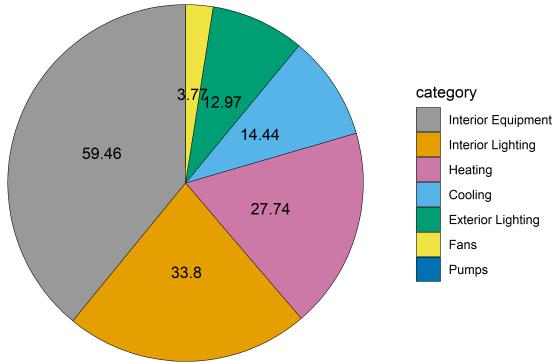


Figure 9: Annual electricity consumption breakdown in kWh/(m²a)

348 Fig. 10 shows the profile of monthly heating and cooling demands in a unit of MJ/m². It is calculated
 349 based on the zone metadata (line 45 in Listing 1) and hourly air system heating and cooling energy outputs
 350 (lines 49 – 53 in Listing 1). With the tidy format and additional metadata columns, these two data fit well in
 351 the data pipeline, making it straightforward and intuitive to perform data transformation and visualization.
 352 Fig. 10 is a result of only around 30 lines of code (line 41 – 53 and 84 – 124 in Listing 1). In Fig. 10, we
 353 can see that during the transition seasons, including March, April, October, and November, the heating and
 354 cooling demands are relatively small compared to summer and winter seasons, indicating the potential of
 355 free-cooling applications.

356 4.2. Parametric simulation

357 This example demonstrates the process of performing parametric simulation analyses using the proposed
 358 parametric prototype. The main focuses are on showcasing the capabilities of (1) creating parametric models
 359 by applying measures and (2) easing the comparative analysis by reusing code snippets developed in data
 360 exploration process.

361 As shown in Fig. 9, the plug loads and interior lighting systems consumed more than 60% of total
 362 electricity. In this example, we will investigate the energy-saving potentials of ECMS on reducing the plug
 363 loads and LPD. Fig. 11 gives an overview of the workflow and Listing 2 shows the actual R code.

364 Measures are functions that describe how an energy model should be modified based on input parameter
 365 values. Lines 8 – 18 in Listing 2 shows a simple measure that modifies the LPD. The core code is line 14
 366 that assigns all related fields in a whole class to input values, taking advantage of the flexibly-scoped OOP
 367 model API. Lines 20 – 40 in Listing 2 shows a measure that modifies the off-work schedule values of plug

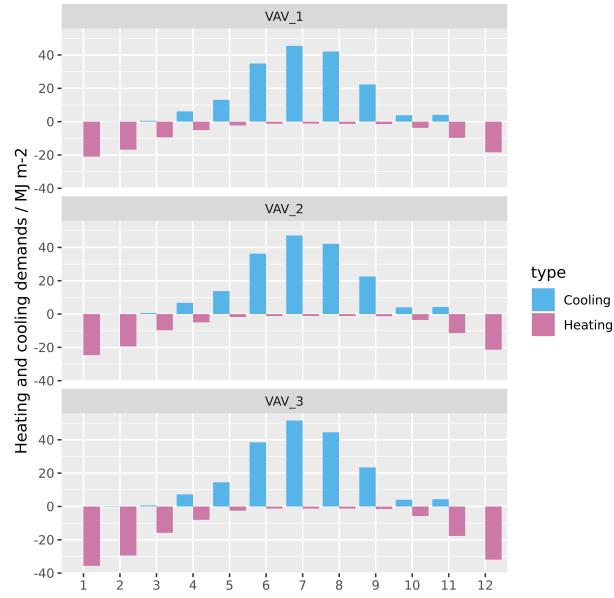


Figure 10: Monthly heating and cooling demand profile

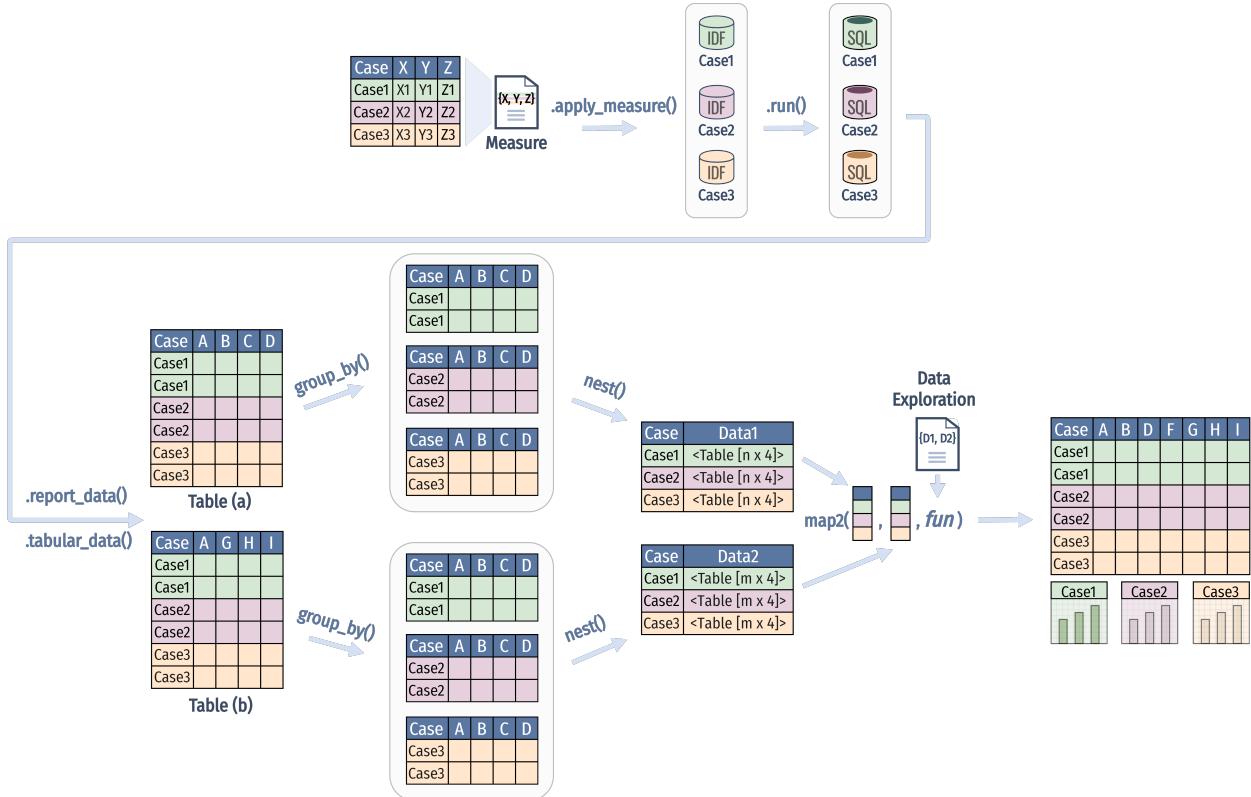


Figure 11: Workflow overview of parametric simulation analysis using the proposed parametric prototype

404 using the Genetic Algorithm (GA). The `GAOptimJob` class leverages the proposed framework in terms of
 405 data structure and parametric simulation management and is capable of defining any number of arbitrary
 406 customized objective functions. It implements flexible general-purpose GA interfaces to solve BES-based
 407 single- or multi-objective optimization problems. This example demonstrates how to use the `GAOptimJob`
 408 class to solve a MOO problem, i.e. reducing carbon emissions and discomfort hours of the medium office
 409 reference building at the same time, by varying (1) indoor heating and cooling setpoint temperatures, (2)
 410 window-to-wall ratio (WWR) and (3) exterior wall insulation thickness. Fig. 13 gives an overview of a
 411 typical GA-based MOO process using the `GAOptimJob` class. Listing 3 shows the actual R code. The
 412 process shown in Fig. 13 can be divided into four main parts:

- 413 1. Specify optimization parameters
 414 2. Create optimization objective functions
 415 3. Set GA operators
 416 4. Gather results and perform further analyses

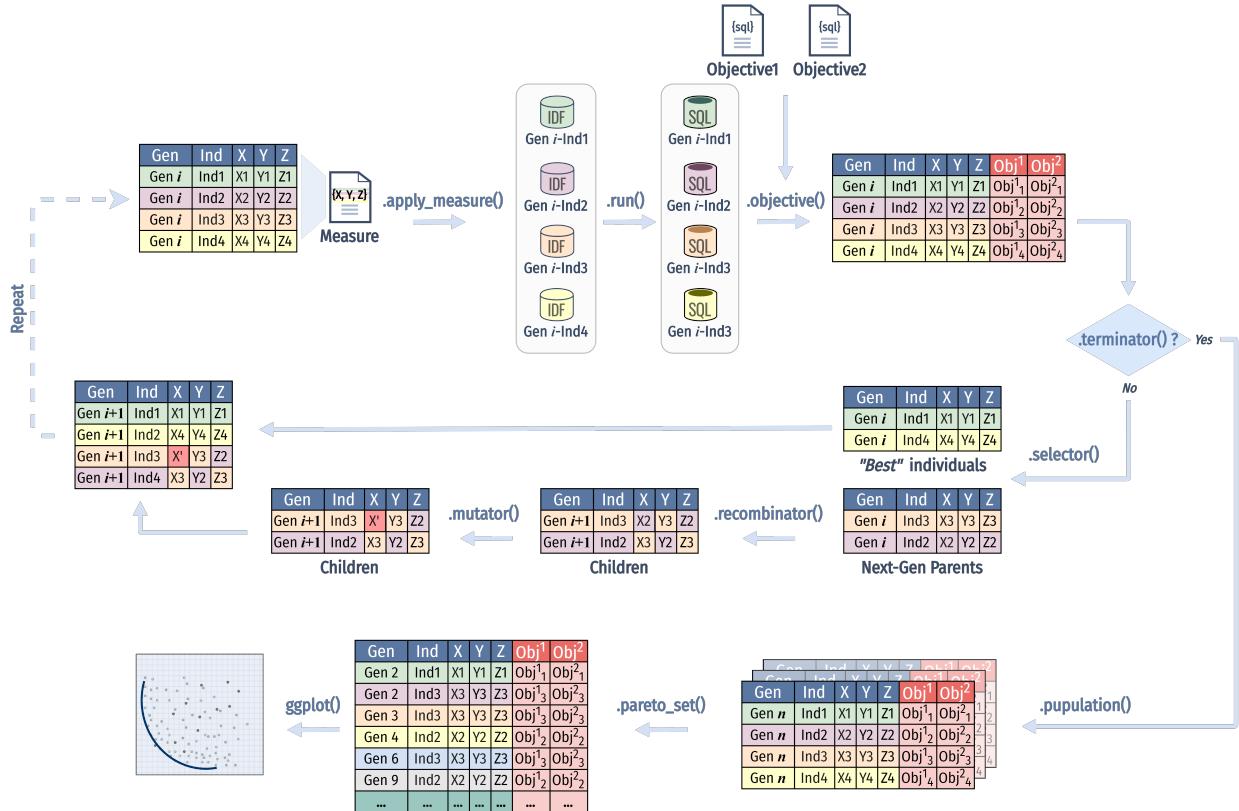


Figure 13: Workflow overview of performing multi-objective optimization on an EnergyPlus model using the `epluspar` package that is based on the proposed parametric prototype

417 Built on top of the parametric prototype, `GAOptimJob` class provides a similar interface for parametric
 418 model generation in the `apply_measure()` method. It can take the same measure functions to describe
 419 how optimization parameters should be modified. Lines 11 – 62 in Listing 3 define three measure functions
 420 to accept various design options in terms of (1) indoor heating and cooling set point, (2) window-to-wall
 421 ratio (WWR) and (3) exterior wall insulation thickness. The actual optimization parameter values in each
 422 generation are automatically calculated based on the GA operators and provided to `apply_measure()`
 423 method for parametric model generation.

424 GAOptimJob provides the flexibility to define objective functions of an optimization problem using any
 425 results from the simulation outputs. The `objective()` method takes objective definitions, evaluates them
 426 after simulations, and extracts the fitness together with optimization parameter values into a tidy table for
 427 post-processing using GA operators. In this example, Lines 88 – 95 in and lines 97–104 in Listing 3 define
 428 functions to extract the annual total carbon emissions and discomfort hours counted based on the Standard
 429 ASHRAE 55 – 2004 from the standard reports using tidy data interface. The `objective()` method in Line
 430 107 in Listing 3 takes these two objective functions and tells the algorithm the minimization optimization
 431 direction.

432 GAOptimJob class has three key genetic operators (methods): (1) `selector()` (to select individuals to
 433 breed a new generation), (2) `mutator()` (to alter parts of one solution randomly), and (3) `recombinator()`
 434 (also called crossover, to swap parts of the solution with another), providing detailed procedures and steps
 435 on how to generate children from parent solutions. Lines 114 – 118 in Listing 3 directly specify those three
 436 operators with the default values that are tweaked to directly perform MOO using the Non-Dominated
 437 Sorting Genetic Algorithm (NSGA-II). The `terminator()` method is used to specify conditions to terminate
 438 the computation. In this example, we set it to stop when one hundred generations have been evaluated.

439 With all objectives, parameters, and operators specified, the optimization will start with the `run()`
 440 method. In this example, we have twenty individuals per generation, resulting in a total of two thousand
 441 annual energy simulations. Once one of the conditions specified in `terminator()` is met, all populations
 442 and Pareto set can be extracted into two tidy tables for further analyses, using the `population()` and
 443 `pareto_set()` method (Lines 130 and 133 in Listing 3). Fig. 14 shows the Pareto front of discomfort hours
 444 and total carbon emissions generated using lines 135 – 143 in Listing 3. The final Pareto font contained 20
 445 unique solutions.

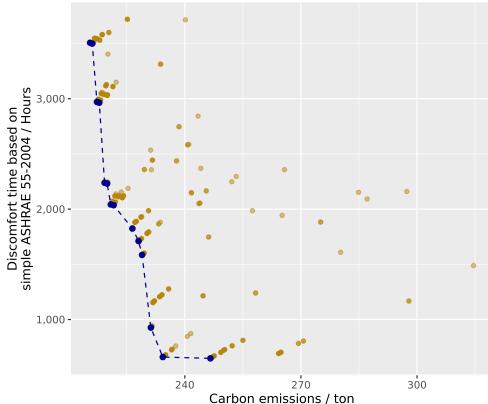


Figure 14: Pareto front of discomfort hours and carbon emissions

446 Fig. 15 shows the parallel coordinates charts of the Pareto set. The carbon emissions have seen a
 447 significant reduction from the original value of 290ton. However, there were 10 out of 20 solutions in the
 448 Pareto set that performed worse in terms of providing a satisfactory indoor thermal environment. One
 449 possible solution to avoid this is to add a constraint when evaluating the fitness of the `discomfort_hours`
 450 objective, making sure all solutions that have larger discomfort hours should be abandoned.

451 4.4. Bayesian calibration

452 Model calibration is an essential process to achieve greater confidence in BES results. In recent years
 453 there has been an increasing application of Bayesian approaches for BES calibration [1,69,70]. Bayesian
 454 calibration is carried out following the statistical formulation proposed by Kennedy and O'Hagan [71]. This
 455 example demonstrates the model calibration workflow using the epluspar R package. The epluspar R package
 456 implements the Bayesian calibration algorithm proposed by Chong [1] and guidelines proposed by Chong

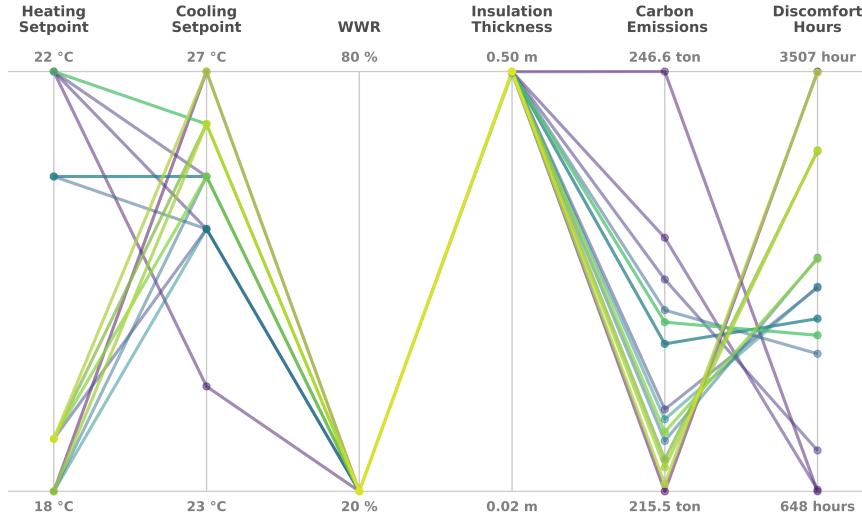


Figure 15: Parallel coordinates chart of the Pareto set

and Menberg [71], and encapsulates them into the `BayesCalibJob` class. The `BayesCalibJob` class inherits from the parametric prototype and thus can leverage all the parametric simulation management capabilities.

Fig. 16 gives an overview of a typical workflow of Bayesian calibration using the `BayesCalibJob` class in the epluspar package. Specifically, Listing 4 shows the workflow of calibrating one VAV fan total efficiency in the medium office reference model using observed fan air flow rate and electrical power.

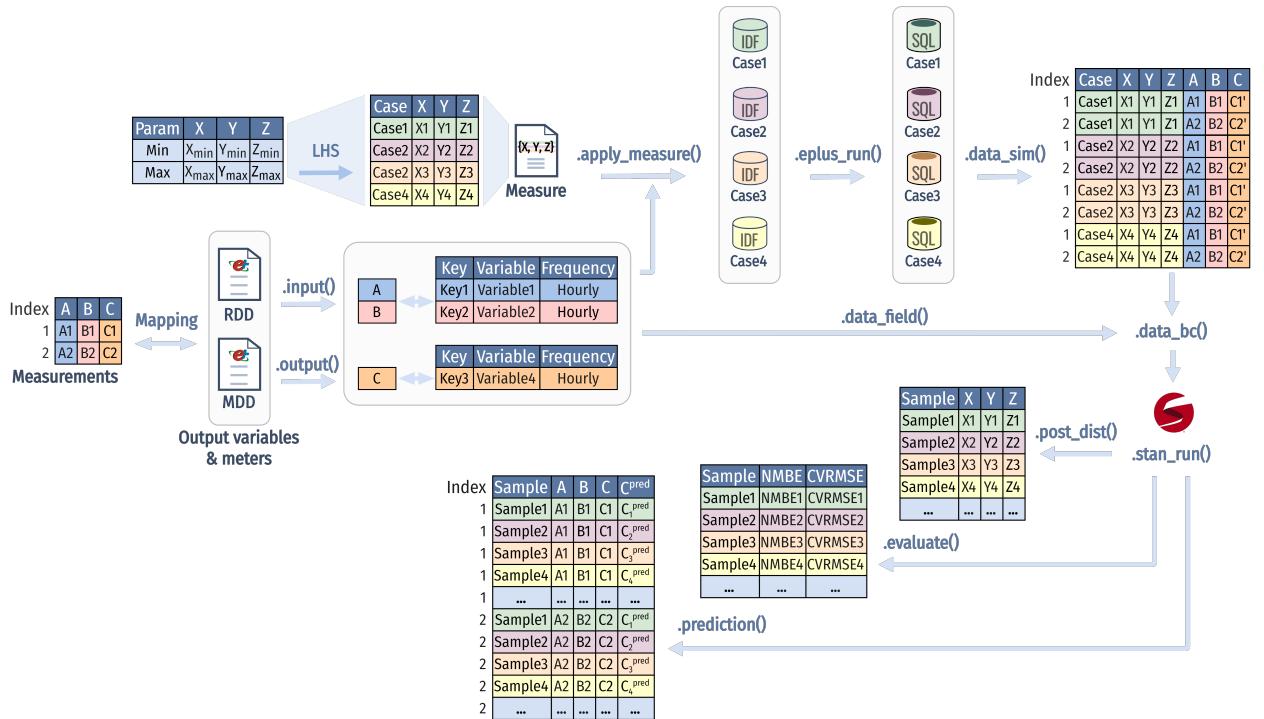


Figure 16: Workflow overview of performing Bayesian calibration on an EnergyPlus model using the epluspar package that is based on the proposed parametric prototype

The initial step for a Bayesian calibration is to collect data for observable input and output. In this example, since the reference model represents a virtual building with no measured data, we created some synthetic data for the examined period of July 1st to July 3rd using simulations (lines 1 – 27 in Listing 4). Also, the TMY3 weather data was used, instead of the Actual Meteorological Year (AMY) weather data. In real practice, the actual measurable variables may not be directly representable in EnergyPlus. In this case, an essential mapping process has to be performed to transform measured variables into EnergyPlus output variables (listed in RDD) and meters (listed in MDD), and connect the transformed measured values with the model using schedule files or other techniques. The next step is to specify the observable input and output variables for the calibration using the `input()` and `output()` methods in `BayesCalibJob` class (lines 30 – 40 in Listing 4).

Following the Bayesian calibration guidelines described in [71], the `epluspar` R package also introduces a `SensitivityJob` class based on the parametric prototype to perform calibration parameter screening with the Morris method. In `BayesCalibJob` class, the calibration parameters, together with the number of EnergyPlus simulations, can be described using either the `param()` method (lines 42 – 46 in Listing 4) or the `apply_measure()` method. Each calibration parameter should be given a lower and upper bound value. Once the calibration parameters are given, the Latin Hypercube Sampling (LHS) algorithm will be used to generate parameter sample values based on the lower and upper bound and the simulation number (lines 48 – 49 in Listing 4). The benefit of LHS is that it will try to cover as much as possible in the multi-dimensional space of the calibration parameters.

After completing all simulations using the `eplus_run()` method (lines 51 – 52 in Listing 4), the `data_sim()` method gathers all simulated data of calibration input and output variables and aggregates the data into the same time frequency as the actual measured data (lines 54 – 55 in Listing 4). Method `data_bc()` will combine measured data and simulated data, and transform them into a list with the proper format for the Bayesian calibration algorithm (lines 60 – 61 in Listing 4) written in probabilistic programming language Stan. With all data required specified, the `stan_run()` method is used to call the Stan program (lines 63 – 64 in Listing 4). Once completed, the posterior distributions of calibration parameters and predicted output variable values can be retrieved using method `post_dist()` and `prediction()`, respectively (lines 66 – 70 in Listing 4). The uncertainty statistical indicators, including Normalized Mean Biased Error (NMBE) and Coefficient of Variation of the Root Mean Squared Error (CVRMSE) can be retrieved using the `evaluate()` method (lines 72 – 73 in Listing 4). Each method returns a tidy table that is well fit for the tidyverse data pipeline for data transformation and visualization.

Fig. 17 gives a density plot showing the posterior distributions of calibrated fan total efficiency, created using lines 75 – 80 in Listing 4. The mean value of the posterior distribution was 0.60, which is quite close to the actual value of 0.59 specified in the original model. Fig. 18 gives a box plot showing the distribution of CVRMSE and NMBE per Markov Chain Monte Carlo (MCMC) sampling, created using lines 82 – 86 in Listing 4. The mean NMBE value was quite close to zero, and the average CVRMSE is around 3.0%. Both values met the thresholds of $\text{CVRMSE} \leq 15\%$ and $\text{NMBE} \leq 5\%$ set by ASHRAE [72]. The satisfactory results are expected since we use synthetic data. However, the overall workflow shown in this example can be applied to Bayesian calibration on building energy models with real measured data.

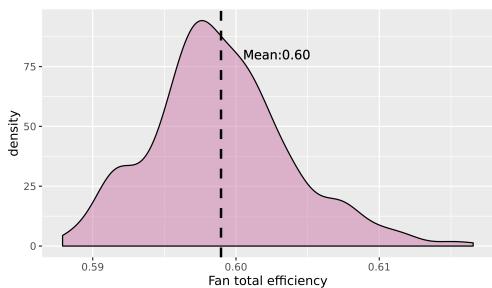


Figure 17: Posterior distribution of fan total efficiency

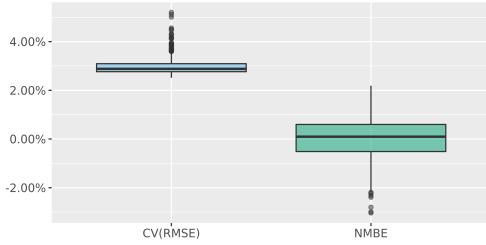


Figure 18: Distribution of CV(RMSE) and NMME per MCMC sample

501 5. Limitations

502 The main limitation of the proposed framework lies in its R-oriented workflows, which currently may not
 503 be widely adopted in industry fields. Prospective users of the framework who do not know R must spend
 504 time learning how to use it. This drawback may be compensated by the growing user community.

505 Since the eplusr framework is mainly focused on modifying existing BES models, instead of creating
 506 new ones from scratch, currently it has limited capacities to perform sophisticated geometry transformation,
 507 including surface matching and rotation. Operations such as creating or replacing one whole HVAC system
 508 may also be time-consuming processes.

509 6. Conclusion

510 Building energy simulation (BES) has been widely adopted for the investigation of environmental and
 511 energy performance for different design and retrofit alternatives. The absence of seamless integration of BES
 512 and data-centric analysis raises problems in both the productivity and also the credibility of BES studies.
 513 This paper proposed a holistic framework called ‘eplusr’ to bridge the gap between the building energy
 514 simulation and data science domains.

515 Eplusr differs from existing frameworks by its data-centric design philosophy. It provides a tidy data
 516 interface for BES that matches the semantics of the simulation results with the data representation. The
 517 tidy data interface provides the possibility to query BES outputs with various types of specifications, which
 518 makes it easy and straightforward to extract simulation results from any time period and for any variables
 519 in a consistent manner. The tidy-formatted results can be easily fed to various data-centric analytics using
 520 existing tools in R.

521 The parametric prototype developed in this framework provides a set of abstractions to ease the process
 522 of parametric model generation, design alternative evaluation, and large parametric simulation management.
 523 It is capable of defining various analyses using any algorithms available in R. The flexibility and extensibility
 524 of the parametric simulation prototype in this framework are demonstrated by its easy adoption to perform
 525 multi-objective optimization and Bayesian calibration.

526 The need for reproducibility in BEM is growing significantly, together with the ongoing trend of the in-
 527 creasing complexity of BEM projects. The eplusr framework provides a solution by developing a portable and
 528 reusable BES computational environment based on Docker containerization, encapsulating the toolchains
 529 for statistical computing, building energy modeling, data analytics, and literate programming. The open-
 530 source nature of the proposed framework will advocate the BES domain to embrace the tools essential for
 531 maintaining a reproducible workflow.

532 Supplementary materials

533 The supplementary files include code and datasets used in this article. More is available at <https://github.com/ideas->
 534 lab-nus/eplusr-paper.

- [59] Y. Xie, J.J. Allaire, G. Grolemund, R Markdown: The Definitive Guide, Chapman and Hall/CRC, Boca Raton, Florida, 2018. <https://bookdown.org/yihui/rmarkdown/> (accessed May 28, 2020).
- [60] T. Kluyver, B. Ragan-Kelley, F. Pérez, M. Bussonnier, J. Frederic, J. Hamrick, J. Grout, S. Corlay, P. Ivanov, S. Abdalla, C. Willing, Jupyter Notebooks: A publishing format for reproducible computational workflows, in: Positioning and Power in Academic Publishing: Players, Agents and Agendas, 2016: pp. 87–90. <https://doi.org/10.3233/978-1-61499-649-1-87>.
- [61] Y. Xie, Dynamic Documents with R and knitr, Second Edition, 2 edition, Routledge, Boca Raton, 2015.
- [62] A. Krewinkel, R. Winkler, Formatting Open Science: Agilely creating multiple document formats for academic manuscripts with Pandoc Scholar, PeerJ Computer Science. 3 (2017) e112. <https://doi.org/10.7717/peerj-cs.112>.
- [63] J.J. Allaire, Y. Xie, J. McPherson, J. Luraschi, K. Ushey, A. Atkins, H. Wickham, J. Cheng, W. Chang, R. Iannone, A. Dunning, A. Yasumoto, B. Schloerke, C. Dervieux, F. Aust, J. Allen, J. Seo, M. Barrett, R. Hyndman, R. Lesur, R. Storey, R. Arslan, S. Oller, RStudio Inc, Rmarkdown: Dynamic Documents for R, (2020). <https://CRAN.R-project.org/package=rmarkdown> (accessed May 28, 2020).
- [64] Y. Xie, TinyTeX: A lightweight, cross-platform, and easy-to-maintain LaTeX distribution based on TeX Live, TUGboat. (2019) 30–32. <http://tug.org/TUGboat/Contents/contents40-1.html>.
- [65] K. Field, M. Deru, D. Studer, Using DOE Commercial Reference Buildings for Simulation Studies: Preprint, National Renewable Energy Lab. (NREL), Golden, CO (United States), 2010. <https://www.osti.gov/biblio/988604> (accessed May 29, 2020).
- [66] S.-H. Yoon, C.-S. Park, Objective Building Energy Performance Benchmarking Using Data Envelopment Analysis and Monte Carlo Sampling, Sustainability. 9 (2017) 780. <https://doi.org/10.3390/su9050780>.
- [67] Z. Yang, B. Becerik-Gerber, A model calibration framework for simultaneous multi-level building energy simulation, Applied Energy. 149 (2015) 415–431. <https://doi.org/10/f7fnf7>.
- [68] J. Bossek, Ecr 2.0: A modular framework for evolutionary computation in R, in: Proceedings of the Genetic and Evolutionary Computation Conference Companion, ACM, Berlin Germany, 2017: pp. 1187–1193. <https://doi.org/10.1145/3067695.3082470>.
- [69] W. Tian, S. Yang, Z. Li, S. Wei, W. Pan, Y. Liu, Identifying informative energy data in Bayesian calibration of building energy models, Energy and Buildings. 119 (2016) 363–376. <https://doi.org/10.1016/j.enbuild.2016.03.042>.
- [70] D.H. Yi, D.W. Kim, C.S. Park, Parameter identifiability in Bayesian inference for building energy models, Energy and Buildings. 198 (2019) 318–328. <https://doi.org/10.1016/j.enbuild.2019.06.012>.
- [71] A. Chong, K. Menberg, Guidelines for the Bayesian calibration of building energy models, Energy and Buildings. 174 (2018) 527–547. <https://doi.org/10/gd5s5b>.
- [72] ASHRAE, Guideline 14-2014, Measurement of Energy and Demand Savings, (2014).

721 Appendix A. Code for data exploration

```

1 # install packages
2 install.packages(c("eplusr", "tidyverse"))
3
4 # load package
5 library(eplusr)
6 library(tidyverse) # for data-driven analytics
7
8 # get EnergyPlus v9.1 installation directory
9 dir <- eplus_config(9.1)$dir
10
11 # use example model and weather file distributed with EnergyPlus v9.1
12 path_model <- file.path(dir, "ExampleFiles/RefBldgMediumOfficeNew2004_Chicago.idf")
13 path_weather <- file.path(dir, "WeatherData/USA_IL_Chicago-OHare.Intl.AP.725300_TMY3.epw")
14
15 # read model
16 idf <- read_idf(path_model)
17
18 #####
19 # Model API #
20 #####
21
22 # make sure weather file input is respected
23 idf$SimulationControl$Run_Simulation_for_Weather_File_Run_Periods <- "Yes"
24
25 # make sure energy consumption is presented in kWh
26 idf$OutputControl_Table_Style$Unit_Conversion <- "JtoKWH"
27
28 # save the modified model into a temporary folder
29 idf$save(file.path(tempdir(), "MediumOffice.idf"))
30
31 # run annual simulation
32 job <- idf$run(path_weather)
33
34 #####
35 # Tidy data interface #
36 #####
37
38 # read building area from Standard Reports
39 area <- job$tabular_data(table_name = "Building Area", wide = TRUE)[[1L]]
40
41 # read building energy consumption from Standard Reports
42 end_use <- job$tabular_data(table_name = "End Uses", wide = TRUE)[[1L]]
43
44 # read zone metadata from Standard Input and Output
45 zones <- job$read_table("Zones")
46
47 # read hourly air-conditioning system output with all additional metadata for
48 # the annual simulation from Variable Output
49 aircon_out <- job$report_data(
50   name = sprintf("air system total %s energy", c("heating", "cooling")),
51   environment_name = "annual",
52   all = TRUE
53 )
54

```

```

55 #####
56 # Data-driven analytics #
57 #####
58
59 # calculate Energy Use Intensity (EUI) for electricity
60 eui <- end_use %>%
61   # only select columns of interest
62   select(category = row_name, electricity = `Electricity [kWh]`) %>%
63   # get rid of category with empty energy consumption
64   filter(electricity > 0.0) %>%
65   # order by value
66   arrange(-electricity) %>%
67   # calculate EUI
68   mutate(eui = round(electricity / area$'Area [m2]'[1], digits = 2)) %>%
69   # calculate proportion of each category
70   mutate(proportion = round(eui / eui[1] * 100, digits = 2)) %>%
71   # remove electricity column
72   select(-electricity)
73
74 # plot a pie chart to show EUI breakdown
75 p_eui <- eui %>%
76   filter(category != "Total End Uses") %>%
77   mutate(category = as_factor(category),
78         ypos = rev(cumsum(rev(proportion))- 0.5*rev(proportion))) %>%
79   ggplot(aes("", proportion, fill = category)) +
80   geom_bar(stat = "identity", width = 1, color = "black", size = 0.2) +
81   geom_text(aes(y = ypos, label = {eui[eui < 0.1] <- NA; eui})) +
82   coord_polar("y", start = 0)
83
84 # calculate air-conditioned floor area per storey
85 storey <- zones %>%
86   # exclude plenum zones
87   filter(is_part_of_total_area == 1) %>%
88   # group by centroid height
89   group_by(centroid_height = round(centroid_z, digits = 4)) %>%
90   # calculate total floor area
91   summarise(floor_area = sum(floor_area)) %>%
92   ungroup() %>%
93   # add storey index
94   arrange(centroid_height) %>%
95   mutate(storey = seq_len(n()), air_system = paste("VAV", storey, sep = "_")) %>%
96   select(air_system, floor_area)
97
98 # get monthly heating and cooling demands per served area
99 aircon_out_mon <- aircon_out %>%
100   # only consider weekdays
101   filter(!day_type %in% c("Holiday", "Saturday", "Sunday")) %>%
102   # add an identifier column to indicate cooling and heating condition
103   mutate(type = case_when(
104     str_detect(name, "Heating") ~ "Heating",
105     str_detect(name, "Cooling") ~ "Cooling"
106   )) %>%
107   # add floor area served by each air-conditioning system
108   left_join(storey, c("key_value" = "air_system")) %>%
109   # calculate the monthly averaged heating and cooling demands in MJ/m2
110   group_by(month, type, air_system = key_value) %>%

```

```

111 summarise(system_output = sum(value) / 1e6 / floor_area[1]) %>%
112 ungroup()
113
114 # plot a pie chart to show the heating and cooling demand profile
115 p_aircon_out <- aircon_out_mon %>%
116   mutate(month = as.factor(month)) %>%
117   mutate(system_output = case_when(
118     type == "Heating" ~ -system_output,
119     type == "Cooling" ~ system_output
120   )) %>%
121   ggplot() +
122   geom_col(aes(month, system_output, group = type, fill = type), position = "dodge") +
123   facet_wrap(vars(air_system), ncol = 1) +
124   labs(x = "", y = "Heating and cooling demands / MJ m-2")

```

722 Listing 1: Data exploration on the EUI and the heating and cooling demands

723 Appendix B. Code for parametric simulation

```

1 # create a parametric prototype of given model and weather file
2 param <- param_job(idf, path_weather)
3
4 #####
5 # Create Measures #
6 #####
7
8 # create a measure for modifying LPD
9 set_lpd <- function (idf, lpd = NA) {
10   # keep the original if applicable
11   if (is.na(lpd)) return(idf)
12
13   # set 'Watts per Zone Floor Area' in all 'Lights' objects as input LPD
14   idf$set(Lights := list(watts_per_zone_floor_area = lpd))
15
16   # return the modified model
17   idf
18 }
19
20 # create a measure for reducing plug loads during off-work time
21 set_nightplug <- function (idf, frac = NA) {
22   # keep the original if applicable
23   if (is.na(frac)) return(idf)
24
25   # extract the plug load schedule into a tidy table
26   sch <- idf$to_table("bldg_equip_sch")
27
28   # modify certain schedule value specified using field names
29   sch <- sch %>%
30     mutate(value = case_when(
31       field %in% paste("Field", c(4,14,16,18)) ~ sprintf("%.2f", as.numeric(value) * frac),
32       TRUE ~ value
33     ))
34
35   # update schedule object using the tidy table

```

```

36     idf$update(sch)
37
38     # return the modified model
39     idf
40 }
41
42 # combine two measures into one
43 ecm <- function (idf, lpd, nightplug_frac) {
44   idf %>% set_lpd(lpd) %>% set_nightplug(nightplug_frac)
45 }
46
47 #####
48 # Apply Measures #
49 #####
50
51 # apply measures and create parametric models
52 param$apply_measure(ecm,
53   lpd = c(    NA,    7.0,    5.0,           NA,           NA,      5.0),
54   nightplug_frac = c(    NA,     NA,     NA,       0.6,       0.2,     0.2),
55   # name of each case
56   .names = c("Ori", "T5", "LED", "0.6Frac", "0.2Frac", "LED+0.2Frac")
57 )
58
59 # run parametric simulations in parallel
60 param$run()
61
62 #####
63 # Data-driven analytics #
64 #####
65
66 # read building energy consumption from Standard Reports
67 param_end_use <- param$tabular_data(table_name = "End Uses", wide = TRUE)[[1L]]
68
69 # calculate EUI breakdown
70 param_eui <- param_end_use %>%
71   select(case, category = row_name, electricity = `Electricity [kWh]`) %>%
72   filter(electricity > 0.0) %>%
73   arrange(-electricity) %>%
74   mutate(eui = round(electricity / area$'Area [m2]'[1], digits = 2)) %>%
75   select(case, category, eui) %>%
76   # exclude categories that did not change
77   filter(category != "Pumps", category != "Exterior Lighting")
78
79 # extract the seed model, i.e. "Ori" case as the baseline
80 ori_eui <- param_eui %>% filter(case == "Ori") %>% select(-case)
81
82 # calculate energy savings based on the baseline EUI
83 param_savings <- param_eui %>%
84   right_join(ori_eui, by = "category", suffix = c("", "_ori")) %>%
85   mutate(savings = (eui_ori - eui) / eui_ori * 100) %>%
86   filter(case != "Ori")
87
88 # plot a bar chart to show the energy savings
89 p_param_savings <- param_savings %>%
90   mutate(case = factor(case, names(param$models()))) %>%
91   ggplot(aes(case, savings, fill = category)) +

```

```

92 |     geom_bar(position = "dodge", stat = "identity", width = 0.6, color = "black",
93 |                 show.legend = FALSE) +
94 |     facet_wrap(vars(category), nrow = 2) +
95 |     labs(x = NULL, y = "Energy savings / %") +
96 |     coord_flip()

```

724

Listing 2: Parametric simulation of ECMs on plug loads and LPD

725 Appendix C. Code for multi-objective optimization using Genetic Algorithm

```

1 # load package
2 library(eppluspar)
3
4 # create a GA optimization job
5 ga <- gaoptim_job(idf, path_weather)
6
7 #####
8 # Optimization variables #
9 #####
10
11 # define a measure to change heating setpoint
12 set_heating_setpoint <- function (idf, sp) {
13   sp <- as.character(sp)
14   idf$set(htgsetp_sch = list(field_6 = sp, field_16 = sp, field_21 = sp))
15   idf
16 }
17
18 # define a measure to change cooling setpoint
19 set_cooling_setpoint <- function (idf, sp) {
20   sp <- as.character(sp)
21   idf$set(clgsetp_sch = list(field_6 = sp, field_13 = sp))
22   idf
23 }
24
25 # define a measure to change the window-to-wall ratio
26 set_wwr <- function (idf, wwr) {
27   # extract data of all windows
28   win <- idf$to_table(class = "FenestrationSurface:Detailed", wide = TRUE, string_value = FALSE)
29
30   # extract data of all parent walls
31   wall <- idf$to_table(win[["Building Surface Name"]], wide = TRUE,
32                         string_value = FALSE, group_ext = "index"
33   )
34
35   # calculate new X and Y coordinates for windows
36   cols <- sprintf("Vertex %s-coordinate", c("X", "Y", "Z"))
37   ratio <- c(0.999, 0.999, wwr)
38   cal_coords <- function (coords, ratio) {
39     list(round((coords[[1]] - mean(coords[[1L]])) * ratio + mean(coords[[1L]]), 3))
40   }
41   wall[, .SDcols = cols, by = "id", c(cols) := mapply(
42     cal_coords, coords = .SD, ratio = ratio, SIMPLIFY = FALSE
43   )]
44

```

```

45 # update coordinates of windows
46 coords <- wall[, lapply(.SD, unlist), .SDcols = cols, by = "id"]
47 coords <- lapply(coords[, -"id"], function (x) as.data.frame(t(matrix(x, nrow = 4))))
48 for (axis in c("X", "Y", "Z")) {
49   cols <- sprintf("Vertex %i %s-coordinate", 1:4, axis)
50   win[, c(cols) := coords[[sprintf("Vertex %s-coordinate", axis)]]]
51 }
52
53 idf$update(dt_to_load(win))
54
55 idf
56 }
57
58 # define a measure to change the insulation thickness of the exterior wall
59 set_insulation <- function (idf, thickness) {
60   idf$set(`Steel Frame NonRes Wall Insulation` = list(thickness = thickness))
61   idf
62 }
63
64 # combine all measures into one
65 design_options <- function (idf, htg_sp, clg_sp, wwr, insulation_thickness) {
66   idf <- set_heating_setpoint(idf, htg_sp)
67   idf <- set_cooling_setpoint(idf, clg_sp)
68   idf <- set_wwr(idf, wwr)
69   idf <- set_insulation(idf, insulation_thickness)
70   idf
71 }
72
73 # specify design space of parameters
74 ga$apply_measure(design_options,
75   htg_sp = choice_space(seq(18, 22, 0.5)),
76   clg_sp = choice_space(seq(23, 27, 0.5)),
77   wwr = float_space(0.2, 0.8),
78   insulation_thickness = float_space(0.02, 0.5)
79 )
80
81 # validate to make sure all measures and objective functions work properly
82 ga$validate(ddy_only = FALSE)
83
84 #####
85 # Optimization objectives #
86 #####
87
88 # define an objective function to get carbon emissions
89 carbon_emissions <- function (idf) {
90   as.double(idf$last_job()$tabular_data(
91     report_name = "emissions data summary",
92     row_name = "Annual Sum or average",
93     column_name = "carbon equivalent:facility"
94   )$value)
95 }
96
97 # define an objective function to get discomfort hours
98 discomfort_hours <- function (idf) {
99   as.double(idf$last_job()$tabular_data(
100     table_name = "comfort and setpoint not met summary",

```

```

101     row_name = "time not comfortable based on simple ASHRAE 55-2004",
102     column_name = "facility"
103   )$value)
104 }
105
106 # set optimization objectives
107 ga$objective(carbon_emissions, discomfort_hours, .dir = "min")
108
109 #####
110 # GA operators #
111 #####
112
113 # specify how to mix solutions
114 ga$recombinator()
115 # specify how to change parts of one solution randomly
116 ga$mutator()
117 # specify how to select best solutions
118 ga$selector()
119 # specify the conditions when to terminate the computation
120 ga$terminator(max_gen = 100L)
121
122 # run optimization
123 ga$run(mu = 20)
124
125 #####
126 # Gather results and perform further analyses #
127 #####
128
129 # get all population
130 population <- ga$population()
131
132 # get Pareto set
133 pareto <- ga$pareto_set()
134
135 # plot Pareto front
136 p_pareto <- ggplot() +
137   geom_point(aes(carbon_emissions, discomfort_hours), population, color = "darkgoldenrod", alpha = 0.5) +
138   geom_line(aes(carbon_emissions, discomfort_hours), pareto, color = "darkblue", linetype = 2) +
139   geom_point(aes(carbon_emissions, discomfort_hours), pareto, color = "darkblue", size = 2) +
140   scale_x_continuous("Carbon emissions / ton", labels = scales::number_format(scale = 0.001)) +
141   scale_y_continuous("Discomfort time based on\nsimple ASHRAE 55-2004 / Hours",
142     labels = scales::number_format(big.mark = ","))
143 )

```

726

Listing 3: Multi-objective optimization using Genetic Algorithm

727 Appendix D. Code for Bayesian calibration

```

1 #####
2 # NOTE: for demonstration, we use the seed model to generate some synthetic data
3 # clone the original model
4 tmp <- idf$clone()
5 # remove all existing run periods
6 tmp$RunPeriod <- NULL

```

```

7 # add a new run period from Jul 1st to Jul 3rd
8 tmp$add(RunPeriod = list("test", 7, 1, NULL, 7, 3))
9 # add variables of interest to output
10 tmp$Output_Variable <- NULL
11 tmp$add(Output_Variable = list("VAV_1_Fan", "Fan Electric Power", "Hourly"))
12 # get rid of design day variable output data
13 tmp$SimulationControl$set(run_simulation_for_sizing_periods = "No")
14 # save the model to a temporary file
15 tmp$save(tempfile(fileext = ".idf"))
16 # run simulation
17 job <- tmp$run(path_weather)
18 # extract fan electric power in 6-hourly frequency
19 fan_power <- job$report_data(all = TRUE) %>%
20   epluspar:::report_dt_aggregate("6 hour") %>%
21   eplusr:::report_dt_to_wide()
22 # insert Gaussian noise
23 fan_power <- fan_power %>%
24   select(-`Date/Time`) %>%
25   rename(power = everything()) %>%
26   mutate(power = power + rnorm(length(power), sd = 0.05 * sd(power))) %>%
27   mutate(power = case_when(power < 0.0 ~ 0.0, TRUE ~ power))
28 ##########
29
30 # load library
31 library(epluspar)
32
33 # create a `BayesCalibJob` object:
34 bc <- bayes_job(idf, path_weather)
35
36 # specify parameters that can be measured
37 bc$input("VAV_1 Supply Equipment Outlet Node", "System Node Mass Flow Rate", "Hourly")
38
39 # specify the parameter to predict
40 bc$output("VAV_1_Fan", "Fan Electric Power", "Hourly")
41
42 # specify parameters to calibrate
43 bc$param(
44   VAV_1_Fan = list(fan_total_efficiency = c(0.4, 0.8)),
45   .num_sim = 30, .names = "FanEfficiency"
46 )
47
48 # get sample parameter values generated using Latin Hypercube Sampling (LHS)
49 bc$samples()
50
51 # run simulations from Jul 1st to Jul 3rd
52 bc$eplus_run(run_period = list("example", 7, 1, NULL, 7, 3))
53
54 # gather simulated data in 6-hour time frequency
55 bc$data_sim("6 hour")
56
57 # set field data
58 bc$data_field(fan_power)
59
60 # get input data for Stan
61 bc$data_bc()
62

```

```

63 # run Bayesian calibration using Stan
64 res <- bc$stan_run(iter = 50, chains = 2)
65
66 # extract posterior distributions of calibration parameter
67 dist <- bc$post_dist()
68
69 # extract prediction values
70 pred <- bc$prediction()
71
72 # evaluate the uncertainties including NMSE and CV(RMSE)
73 uncert <- bc$evaluate()
74
75 # draw a density plot for the posterior distributions of calibration parameters
76 p_dist <- dist %>%
77   pivot_longer(-sample) %>%
78   ggplot() +
79   geom_density(aes(value, fill = name), alpha = 0.5) +
80   geom_vline(aes(xintercept = mean(value)), linetype = 2, size = 1)
81
82 # draw a boxplot to show the distributions of uncertainty indices
83 p_uncert <- uncert %>%
84   pivot_longer(-sample) %>%
85   ggplot() +
86   geom_boxplot(aes(name, value, fill = name), alpha = 0.5)

```

728

Listing 4: Bayesian calibration