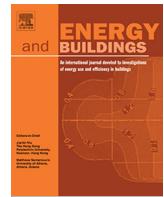




ELSEVIER

Contents lists available at ScienceDirect



Data science for building energy efficiency: A comprehensive text-mining driven review of scientific literature



Mahmoud M. Abdelrahman, Sicheng Zhan, Clayton Miller, Adrian Chong *

Department of Building, School of Design and Environment, National University of Singapore, 4 Architecture Drive, Singapore 117566, Singapore

ARTICLE INFO

Article history:

Received 16 December 2020

Revised 19 February 2021

Accepted 2 March 2021

Available online 13 March 2021

Keywords:

Reference mining
Natural language processing
Data science
Built environment
Building energy efficiency
Word embeddings

ABSTRACT

The ever-changing data science landscape is fueling innovation in the built environment context by providing new and more effective means of converting large raw data sets into value for professionals in the design, construction and operations of buildings. The literature developed due to this convergence has rapidly increased in recent years, making it difficult for traditional review approaches to cover all related papers. Therefore, this paper applies a natural language processing (NLP) method to provide an exhaustive and quantitative review. Approximately 30,000 scientific publications were retrieved from the Elsevier API to extract the relationship between data sources, data science techniques, and building energy efficiency applications across the life cycle of buildings. The text-mining and NLP analysis reveals that data sciences techniques are applied more for operation phase applications such as fault detection and diagnosis (FDD), while being under-explored in design and commissioning phases. In addition, it is pointed out that more data science techniques that are to be investigated for various applications. For example, generative adversarial networks (GANs) has potential in facilitating parametric design; transfer learning is a promising path to promoting the application of optimal building operation;

© 2021 Elsevier B.V. All rights reserved.

1. Introduction

With advances in information technology, buildings today are collecting ever-larger amount of real-time data from various heterogeneous sources [14,83]. The vast amount of data also have led to increased data awareness and data science applications [99]. These innovations have led to an explosion of research in this field, resulting in thousands of publications in this area (e.g., Fig. 3). It is now the case that researchers are in a position in which there are significantly more research publications available than what can be processed and digested by human [40]. Numerous literature reviews are also being produced to aggregate research literature; however, this is also not a trivial process due to the volume of research in this area.

1.1. Using data science to quantify the impact of data science on buildings

In order to address this challenge, the concept of using text-mining methods to analyze scientific literature has gained traction. The academic knowledge is exponentially expanding; thousands of

research articles are authored by domain experts every day [77]. Thanks to the recent advancements in Natural Language Processing (NLP), it has been viable to extract knowledge from a large corpus of such structured text. Typically, information from the literature is extracted via traditional narrative literature review/survey of a finite number of articles (hundreds) [94], besides other methods such as questionnaire surveys and expert interviews. However, there are some challenges in applying these methods on a large scale [8]. Specifically, conducting a manual literature review on a large number of papers requires huge effort. It is even more challenging if the literature review is cross-disciplinary such as extracting relations between various data.

1.2. Similar studies

Several conventional literature reviews have been completed in recent years to capture the innovation occurring due to the convergence of data science and building energy performance research during different lifecycle phases [49,109,60,21,48,132,69]. Wang and Srinivasan explored the use of single versus ensemble-based models for building energy prediction [153]. Roth et. al explored the use of various data-driven techniques in the context of benchmarking building [134]. Colm et. al. [51] investigated Machine learning methods for maximizing measurement and verification

* Corresponding author.

E-mail address: adrian.chong@nus.edu.sg (A. Chong).

(M&V) accuracy with an application on a real building. This application concluded sufficient accuracy despite some limitations such as poor data quality and insufficient metering. In the operation phase, data science methods were found promising to tackle the challenges in building system control [96]. Fault detection and diagnosis (FDD) is another important application of improving building energy performance, where data science methods are commonly used [168]. Furthermore, energy audit and commissioning of buildings using data analytics has been investigated by Rohloff et. al. to minimize the performance testing hours and maximize the value of the test results [132]. Beyond single buildings, data-driven methods are also useful for demand response and smart grid applications [50,112,47]. On the urban scale, energy efficiency applications also have grasped the interest of researchers. For example, many researchers investigated district heating and cooling systems [88,128] and Urban Building Energy Modelling (UBEM) [11,65,127]. Most of these reviews have indicated the potentials of using big data (such as sensing data from IoT and urban building energy modelling data) and machine learning.

These reviews cover the specific application of data science to various facets of the building energy paradigm however, they are constrained by the ability of human-driven analysis to make qualitative relationships between a relatively small number of papers. Each review is only able to analyze between 100-120 publications. An emerging field of analysis of scientific literature is seeking to extract insights from quantities of publications in the tens of thousands instead of only the hundreds. These studies have been completed in fields, such as the humanities [131], bio-medicine [143], and frameworks have been built for more general text mining purposes [147]. In the building domain, some studies adopted bibliometric reviews of the global trends in different building-related issues such as BIM [137,91], Green Buildings [167], Life-cycle assessment [52], Building maintenance [130] among other aspects.

Different tools and approaches of text-mining have been used in literature to conduct literature reviews. [23,41,138,158] used VOS-viewer [149] to create bibliometric networks and density map between articles in different fields. Other researchers used CiteNet-Explorer [150] to track the citation relations across articles in scientific research [38,142] among others. Other tools such as CiteSpace, BibExcel, SciMAT, Sci² Tool have been extensively reviewed by [110]. However, all these tools come with a graphical user interface (GUI) which limits the user ability to extend it beyond its embedded algorithms. Additionally, these tools only use the articles' metadata (title, abstract, authors, keywords, references, date ..etc) not the article body full text. Therefore, many researchers used open-sourced libraries such as the Natural Language Toolkit NLTK [95], Glove [120], Python/scikit-learn [119], word2Vec [102,100,103] to develop a model that performs a specific task.

This paper aims to address the challenges and deficiencies of typical literature reviews and capture the full extent of the relationships between data science and building energy performance. Given these circumstances, the current study adopts text mining survey and natural language processing to extract different segments of building data usability and their relevant users. This effort is the first text-mining and NLP review of its kind in the building energy performance research domain.

The paper is organized as follows. Section 2 provides an outline of the data extraction from the publisher's API, the text mining process, and the quantification of relationships between the different concepts being compared. Section 3 illustrates the overview graphics extracted from the mining process that show the diversity of data science techniques applied to buildings. Section 4 provides a high-level analysis of the trends and gaps found in the literature with respect to data science for building performance analysis.

Finally, Section 5 concludes the analysis and provides insight on reproducibility and further analysis using the data set.

2. Methodology

The current study follows three types of research designs, namely, text-mining survey, natural language processing (NLP) semantic analysis, and relation graph extraction. Each one of these three designs is distributed across a five-phase method of data collection, preprocessing, and processing. These five phases, summarized in Fig. 1, are: 1) Identifying the querying keywords of each category, 2) Extracting the relevant articles with their corresponding metadata using ELSEVIER api, 3) Pre-processing the data, 4) Applying the NLP algorithms, 5) Extracting the relationships and creating the relation graph network.

2.1. Keyword identification

Four distinct categories of keywords are identified that were used for querying the articles for this analysis. Specifically, the categories are data, data science, energy efficiency, and phase. These keywords are meant to constitute a relational network to extract the use of different data-points, techniques, algorithms, and applications during the building life cycle phase as illustrated in Fig. 2. The analysis of the relationships between these concepts forms the foundation to understand what techniques and data sources are popular in the building energy performance domain and which ones are underutilized.

2.1.1. Definitions

To set the context, the following are more detailed definitions of each of these concept categories:

Def.1 Data: (data) refers to different types of data used in buildings, including design specifications' data such as thermal comfort and indoor environmental quality; metered data such as temperature, humidity, energy consumption, and chilled water flow rates; and spatial data such as building geometry, spaces and zones.

Def.2 Data Science: (data_science) refers to models and algorithms used by different users during different building life-cycle phases. For example, the use of energy simulation, data mining and visualization, machine learning models would be included in this category.

Def.3 Energy Efficiency: (energy_efficiency) refers to the various categories of potential application of data science in the building energy analysis domain. These techniques range from conventional approaches such as automated fault detection and diagnostics (AFDD) to more contemporary innovations such as urban-scale district energy modelling.

Def.4 Phase: (phase) refers to the building life-cycle phase/stage. We defined 5 phases found in the literature: design phase, commissioning, operation and maintenance, and retrofit.

Each of these categories consists of manually defined initial keywords. We obtained these keywords by conducting a preliminary survey over the existing literature.

2.1.2. Keywords acquisition

A preliminary literature survey was conducted to obtain the keywords of each of the categories. For example, keywords that are related to the data category include: meter readings, energy consumption, load profile, thermal mass, electricity pricing, schedule, thermal comfort, etc.

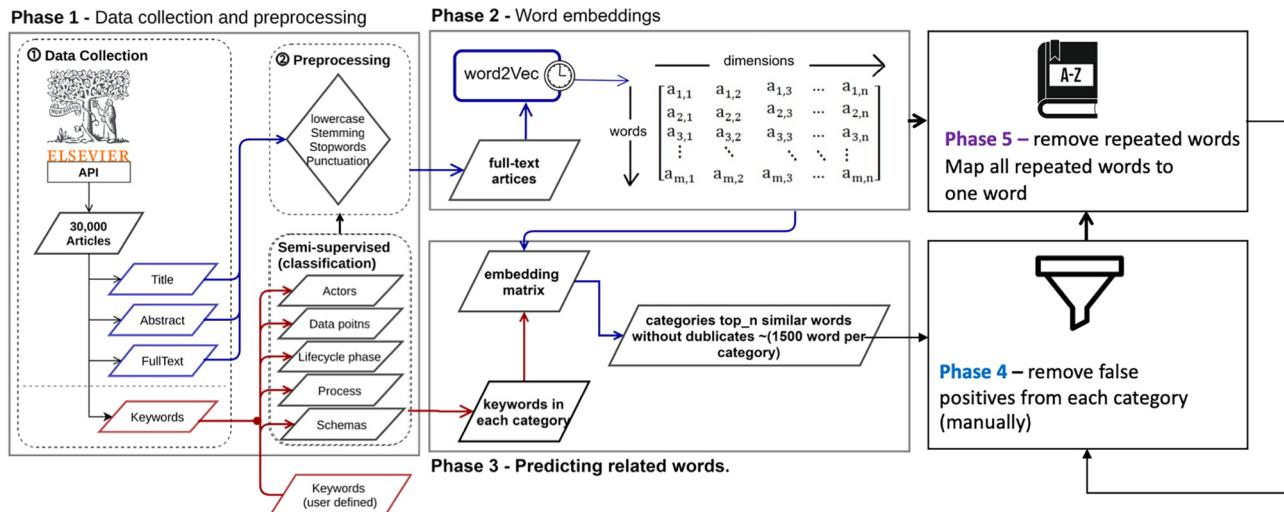


Fig. 1. The flowchart shows the methodology used in this research 1) Identifying the querying keywords of each category, 2) Extracting the relevant articles with their corresponding metadata using ELSEVIER api, 3) Pre-processing the data, 4) Applying the NLP algorithms, 5) Extracting the relationships and creating the relation graph network.

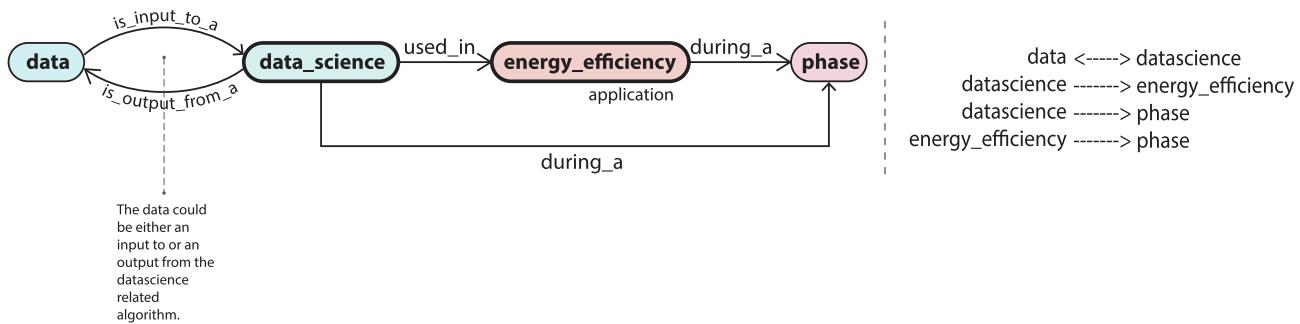


Fig. 2. Overview of the categories of concepts analysed in this text-mining analysis and their relationships with each other.

Each of the keywords has been paired with words to restrict the search query to the built environment. These restrictive words are “building”, built environment, and buildings. For example, using the word “Haystack” which indicates a building schema results in an irrelevant output such as “...Finding a needle in a haystack”.

2.2. Text mining survey

ELSEVIER is one of the largest scientific publishing and aggregation organizations. They first introduced an API for the public for text-mining research in 2014 [151]. By opening their database, researchers can extract full texts and metadata from more than 11 million research items using ELSEVIER API. In this research, the same approach was used to obtain full versions of about 30,000 papers by querying the keywords extracted from the previous step. The articles come alongside their corresponding metadata, such as date of publishing, authors and affiliation, journal (container), title, abstract, keywords, amongst others. In this analysis, we use the publications extracted from this API as a representative sample from the building energy research domain as these journals are the highest cited in energy and buildings.

2.2.1. Article filtering

The initial query process has resulted in 45,000 articles from more than 1000 journals. However, many of these articles are duplicated. Thus, after removing the duplicates, the accumulative

number of articles reached around 30,000 articles. All of these articles come with a rich amount of metadata including publishing date, authors and their affiliations, keywords, number of citations, besides abstract and title. Fig. 3 illustrates the top number of papers per journal and the number of published papers per year. From this result, it is observed that the majority of the articles come from building, energy, and sensor-related journals. At this stage, the extracted articles were ready for preprocessing and preparation.

2.3. Text prepossessing

The preprocessing phase aims at preparing the extracted full text for the data mining process. The data preparation includes removing unwanted words from the articles, making the words consistent, and tokenization of words or group of words. Firstly, there are two types of unwanted words can be identified: 1) titles, subtitles, and annotations such as introduction, literature review, figure, table. These words are repeated in every article and may cause bias in the subsequent processes. 2) stop words; the term “stop-words” refers to words that are frequently repeated yet not meaningful for the context such as the, a, in, of. If these stop-words are included, they will cause bias in the NLP models. Many tools are available for removing stop words such as the Natural Language Toolkit (NLTK) [95].

Secondly, Since the NLP models are case-sensitive, they need to be consistent. For example, lowercasing letters throughout the cor-

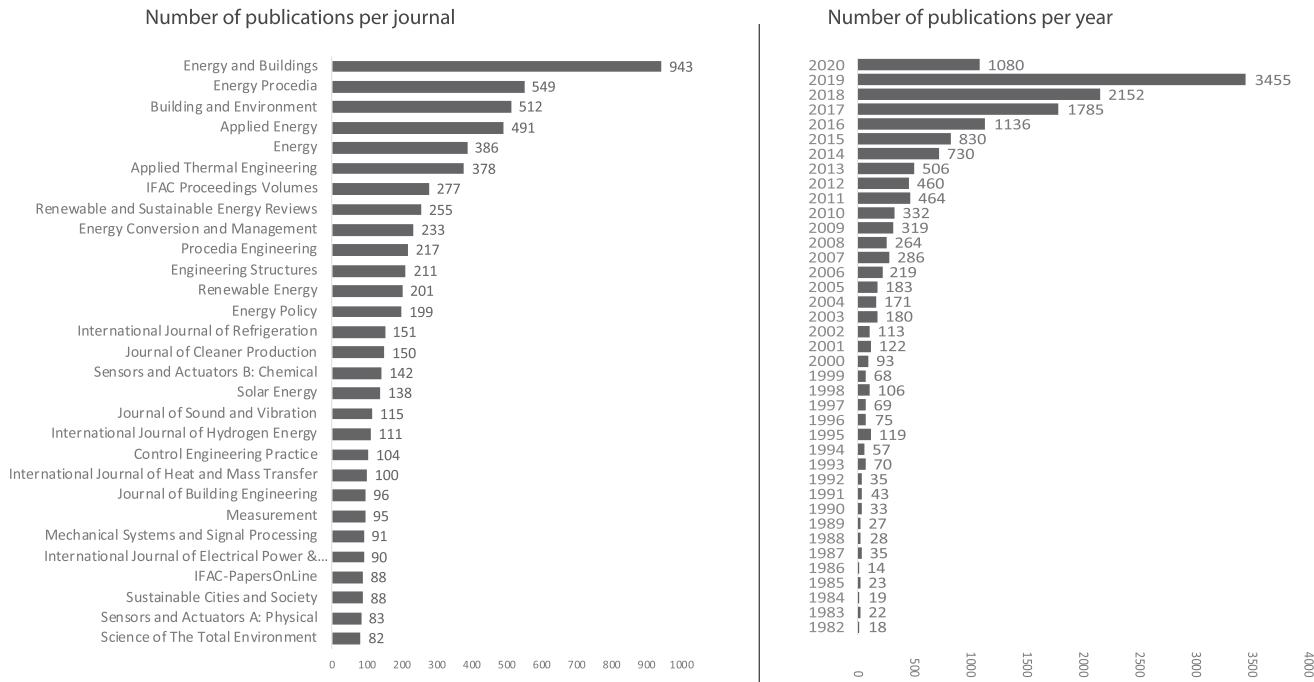


Fig. 3. Number of collected papers per journal and per year.

pus. Also, converting regular plural nouns into singular ones by removing "s". There are many other sets of tools for making the text consistent called text stemming and lemmatization. However, the current study will only use two of these methods, which have resulted in a better accuracy (Fig. 4). Firstly, the common root of different words was used. Secondly, compound words were converted into a single word with "_" separating them. The compound words, however, were extracted from each article's keyword section. We included only the keywords section as it is known to contain the main important acronyms and definitions. After making the full text consistent, it is now ready to be prepared for the NLP text mining process.

2.4. NLP text mining using Word2Vec

Word2Vec is a word embeddings algorithm that is used to extract the semantic similarities between different words in a text [100,101]. This similarity is indicated by assigning each word in the text to a multi-dimensional vector. Then the Euclidean distance between each word can be calculated using the cosine of the angle between these vectors: $\text{sim}(A, B) = \cos(\theta) = \frac{AB}{\|A\|\|B\|}$. The closer the words to each other, the more similar they are likely to be. The Word2Vec training process aims to predict a word (known as the central word) from the context within which this word falls (context words)[100,54]. This central word is initially masked, then the algorithm tries to predict it from a window of n words before and after (in our model, we used a window of 20 words). This window was decided based on hyperparameters fine tuning. After reaching a reasonable accuracy in predicting each word in the corpus, the training stops. Then, the hidden layer is extracted as an embedding vector. Deciding the dimension of the hidden layer (embedding

vector) is a best-practice-driven process and is subject to hyperparameter fine-tuning. In our case, we assigned a vector of 300 dimensions to each word which was proven to give the highest accuracy for our model. The architecture of the word2vec model is illustrated in Fig. 5.

2.5. Extracting the relationship between categories

In word2Vec, two words are similar if they frequently appear in similar contexts. For example, if the word architect and the word early_design_phase are frequently appearing among similar words, then these two words will be assigned to relatively near vectors. Concurrently, two, or more, words can be added or subtracted from each other by adding/subtracting their corresponding vectors. For example, adding artist + engineer results in a vector that is closest to the word architect. This metric is used to extract the relationship between words from different categories in two main steps.

Firstly, there can be many words that refer to the same term. In this case, the words are mapped to that term. For example, the word early_design_phase is found to have many other synonyms that are similar to it such as: 'early_design_stage', conceptual_design, early_design_development, early_design, concept_design, early_design_stages, and others. These words can differ slightly in using "_" rather than "-" or using the word stage rather than phase. Another example is the use of acronyms that refer to the same term such as gbrs and green_building_rating_system were easily captured using the word2vec similarity metric. Thus, the similarity metric is used to extract these similar words which makes it easier to implement the following step i.e. extracting the relationships.

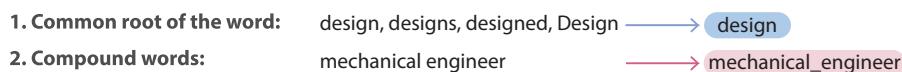


Fig. 4. Each word can have one or more similar synonyms which are mapped to the original word. The figure shows two types of text stemming and lemmatization.

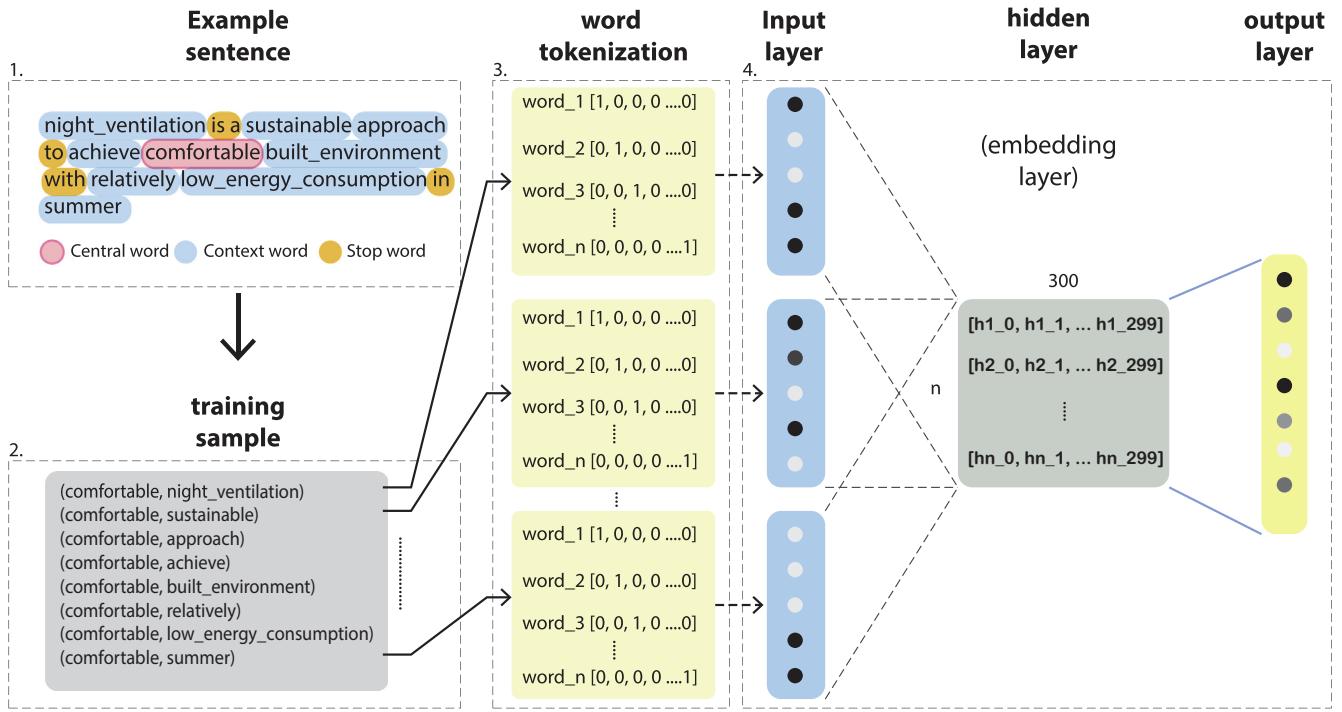


Fig. 5. Word2vec architecture. This architecture is used to extract the embedding matrix (the hidden layer) which is the vector representation of each word in the latent space.

Extracting the relationships comes after creating a list of all the words and their synonyms from each of the four categories. We used a method called **n-gram** to extract these relationships [87]. The n-gram model searches for the similarity between two words by sampling n samples from contiguous sequence of their synonyms. For example, the objective is to extract the similarity between the two words W_a and W_b such that W_a is "energy_consumption" which has other synonyms such as W_{a_1} ("building_energy_use") and W_{a_2} ("energy_consumption_data"); and the word W_b "energy_benchmarking". The 1-gram model will look for the the similarity by taking one word at time, while the 2-gram model will look for the similarity by taking pairs of words at time. At the end, the total similarity between the two words is given by the average of all the similarities:

$$\bar{S}(W_a, W_b) = \frac{\sum_{n=1}^{\max(\text{len}(W_a), \text{len}(W_b))} n - \text{gram}(W_a, W_b)}{\max(\text{len}(W_a), \text{len}(W_b))}$$

Where $\bar{S}(W_a, W_b)$ is the average similarity between two lists of words W_a and W_b and their synonyms $W_a = [w_{a_1} \dots w_{a_n}]$, $W_b = [w_{b_1} \dots w_{b_n}]$. n is the is defined by the maximum number of synonyms of the two words. If $n = 1$, then it is called unigram; if $n = 2$, it is called digram; if $n > 2$ it is referred to as n-gram. The n-gram is obtained by the cosine similarity between the two word lists W_a and W_b as follows:

$$n - \text{gram}(W_a, W_b) = \text{Sim}\left(\sum_{i=1}^n W_{a_i}, \sum_{j=1}^n W_{b_j}\right)$$

An n-gram similarity is a number within the range $[-1.0, 1.0]$. If the two words are identical (e.g. w_a is the same as w_b), their similarity = 1.0, if they are perfectly semantically opposite, their similarity will be -1.0 theoretically. However, 0.0 means that there is no semantic similarity between the two words. These numbers are converted into triplets $\{W_a, W_b, \bar{S}(W_a, W_b)\}$ which is then con-

verted into a directed weighted graph. The results will be explained in the following Section 3. For example, the n-gram similarity between the word "fault detection and diagnosis" which has the synonyms : ["fault_detection_and_diagnosis", "fdd", "fault_detection"] and the word "neural_network" which has the synonyms ["neural_networks", "deep_learning", "cnn"] will follow 3-gram similarity which results in a value of 0.42 in this case which is relatively high.

3. Results

The methodology outlined a process of using text-mining and NLP methods to extract and process various concepts from a large corpus of research publications related to the convergence of data science and building performance. This section focuses on the detailed visualization of the aspects of drawing relationships between these categories. The key output of this work lies in the ability to quantify in relative terms the strength of relationships between the words found in the various categories being studied: the data sources, energy efficiency applications and life cycle phases of the built environment versus the data science techniques available to researchers. Fig. 6 shows the framework of this process starting with the definition of the categories and selection of words to the visualization of similarity of words and clustering of words into concepts.

3.1. Vector representation and relationships of extracted words

This first method of visualizing and drawing relationships comes in the form of a scatter plot that illustrates the various words extracted from the corpus and the directional nature and magnitude of their differences according to the vector model. Fig. 7 illustrates this situation by showing the embedding vector of words projected into a two-dimensional space. The keywords are categorized according to the four dimensions of the analysis:

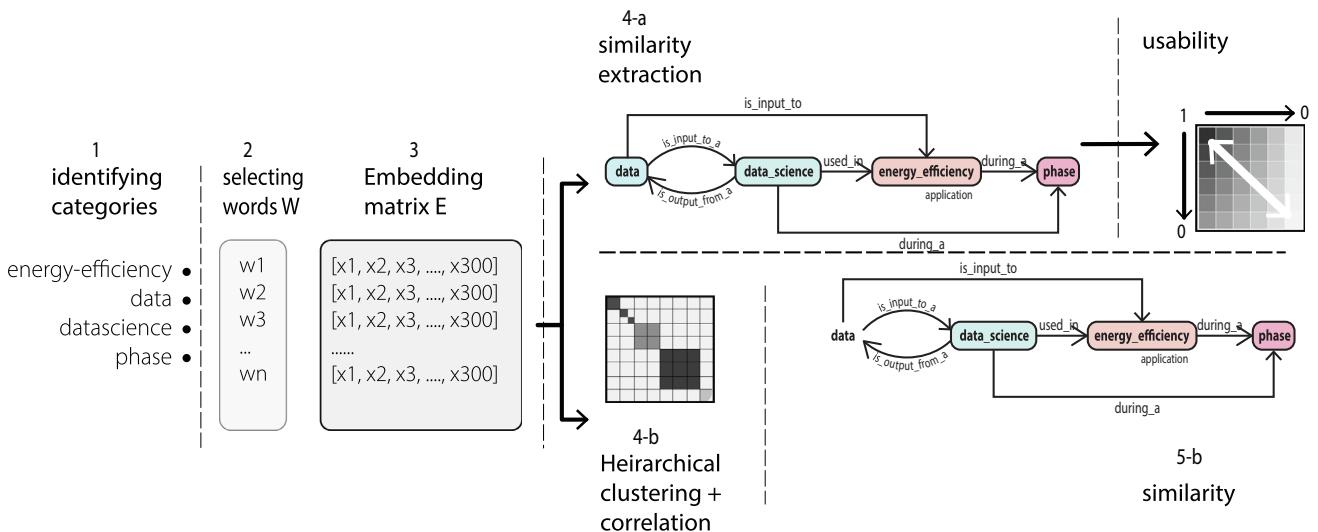


Fig. 6. Overview of the ways of showcasing the results of the text-mining and NLP process. 1) Identifying 4 categories and 2) assigning the corresponding words under each category. After that, 3) the embedding vector of each word is extracted. Then, there are two main approaches: a) is the usability relation extraction (Section 3.2) including 4-a) graph relation extraction using only the similarity metric, and 5-a) sorting the words based on its usability; and b) is the clustering of concepts (Section 3.3) including 4-b) unsupervised hierarchical clustering of words based on the embedding vector of each word (from step 3) and then 5-b) the graph relations between categories based on the clustered data.

data, data science, energy efficiency and life-cycle phase. The various words are clustered according to their relationship with each other in the vector model. The scatter plot shows how the words most closely associated with various life cycle phases of buildings can be extracted as a pattern of points from the lower left to the upper right portion of the diagram.

3.2. Usability-based similarity relation extraction

The next method to visualize relationships was the comparison of several word categories against each other to show the correlations between various concepts. These visualizations are used to illustrate the ranking of lowest to highest correlations of various data and data science concepts in both the energy efficiency applications in buildings and when those techniques are generally utilized.

3.2.1. Data sources used in building energy efficiency applications

The first comparison in this process was to show the relationship between words referring to data sources with selected energy efficiency applications. Fig. 8 shows a heat map of the various data source words extracted from the literature and their relationship strength with words extracted that related to energy efficiency applications from the life-cycle phase of the building. The horizontal axis (energy efficiency applications) is grouped according to the life cycle phases of buildings and the vertical axis (data sources) is sorted according to the average strength of relation for each data source as compared to the applications.

It can be observed that data are used mainly during the operation and maintenance and the design phases of the building lifecycle. However, data are underutilized in the commissioning phase. On the one hand, there are some energy efficiency applications that use data most frequently, such as passive design, demand-controlled ventilation, model predictive controls (MPC), fault detection and diagnosis, and retrofit analysis. On the other hand, there are other energy efficiency applications that do not use data frequently such as Measurement and verification (M&V), operation and maintenance (O&M), HVAC optimization, parametric design, and district energy systems.

Fig. 8 also shows that data sources also varies in their utilization. Some of these data are frequently used such as energy consumption data, building envelope, energy conservation measures (ECM), occupant behaviour, cost analysis, and calibrated models. Nonetheless, other data are underutilized related to HVAC design, weather and thermal comfort such as inlet/outlet temperature, condenser fan power, and mass flow rate; dew point, noise level, mean radiant temperature, and dry-bulb temperature; and clothing insulation, thermal sensation, and thermal comfort indices.

3.2.2. Data science techniques that utilize the various data sources from the built environment

The next comparison similarly uses the words related to data sources, but instead compares them to various data science techniques selected for this analysis. Fig. 9 outlines the relationship between the various data science techniques versus the data sources created in the built environment. This time both axes are sorted according to the average strength of relation for both the data science techniques (horizontal axis from right to left) and data sources (vertical axis from top to bottom).

This relationship is dominated by energy simulation, optimization, regression, and validation. However, the figure shows that there is abundant room for further data use in generative Adversarial Networks (GANs), dimensionality reduction, segmentation, and anomaly detection. There is another pattern that can be observed for applications such as factor analysis, reinforcement learning, and multi-objective optimization. These data-science applications are used frequently but with no significant relation to data sources. These relations have various observations from the data-sources perspective.

From the data source perspective, a different order from the previous heatmap can be observed. While energy consumption data is still dominating the use in data-science applications, historical data, real-time data, thermal comfort, and schedules are the highest frequently used data sources for different data-science applications. On the other side of the spectrum, HVAC design elements such as condenser fan power, inlet/outlet temperature, CAV, and fan power; as well as passive design strategies such as thermal mass are under-used in data science applications.

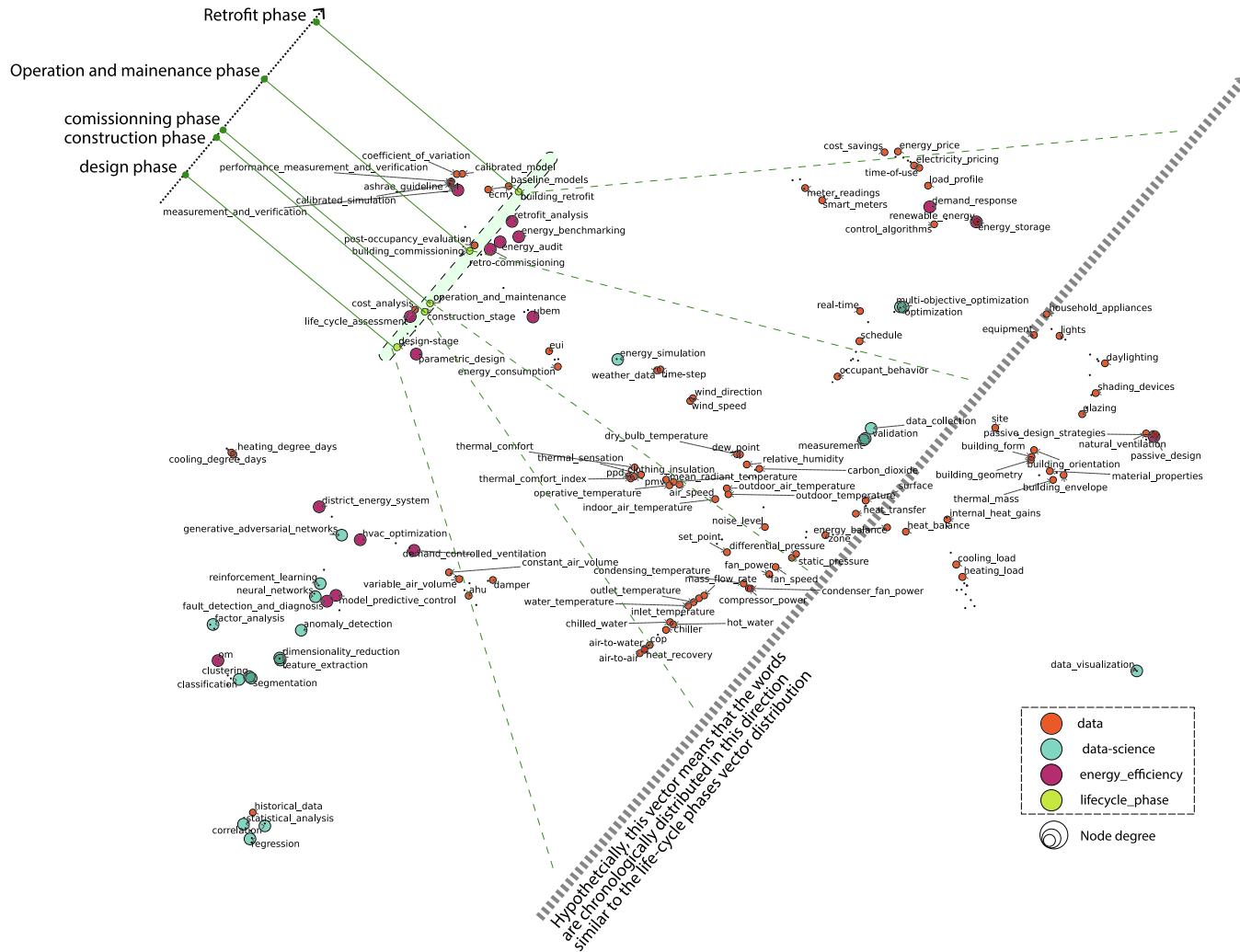


Fig. 7. The vector representation of the words from each category. These words are located based on their embedding vector. The embedding vector of each word is dimensionally reduced from 300 dimensions to 2 dimensions for the sake of visualization. The euclidean distance between words indicates the semantic similarity between these words. The degree of a specific node refers to the number of nodes connected to that specific node. The x and the y axis here represent the components of a 2D euclidean space.

3.3. Clustering of concepts

The next visualization method utilizes hierarchical clustering instead of sorting the words from strongest to weakest relation. Clustering allows for words with similarities within each category to be grouped and observed. Hierarchical Agglomerative Clustering (HAC) was used for this process using Ward's method. This algorithm is applied to the embedding vector of words in each category to group similar words together based on the euclidean distance between words in the vector space. This grouping is visualized in the form of a tree called a dendrogram (Figs. 10 and 11).

3.3.1. Hierarchical agglomerative clustering of concepts

The HAC has been applied for words in each distinct category using the Ward's method [111]. On the one hand, Fig. 10 shows the HAC of energy_efficiency category (on the left) and the HAC of the data_science category (on the right). The energy_efficiency category has been clustered into three groups. These groups are likely to be grouped based on the life-cycle phase, namely, Operation and maintenance phase, design phase, and the commissioning phase. However, the data_science category has been clustered into five different groups/subgroups. These are, Machine Learning (**ML**), Deep Learning (**DL**), Data pre/post-processing (**PP**), Optimization (**OP**), and Statistical methods (**St**). On the other hand, Fig. 12 shows

the HAC of the **data** category. In this figure, the **data** category is clustered into nine groups of keywords based on their similarities:

- 1. Passive systems (PS)** which includes data that are used in passive design such as building geometry, orientation, glazing, materials, shading devices, and natural ventilation.
- 2. Heat recovery ventilation data (HR)** such as air-to-air, air-to-water, and heat recovery.
- 3. Building Energy Modelling data (BEM)** including heat/energy balance, zones, surfaces, and heating/cooling loads.
- 4. Measurement and verification data (M&V)** including energy conservation measures (ECM), baseline model, ASHRAE guideline 14, calibrated simulation, and post-occupancy evaluation (POE).
- 5. Energy consumption related data (EC)**. This includes the energy price, cost saving, time of use, energy use intensity (EUI), heat gains from appliances and equipment, load profile, smart meters, and others.
- 6. HVAC related data** including **HVAC-AF** airflow data and **HVAC-T** temperature related data. HVAC airflow data include AHU, VAV, CAV, fan speed and power, etc. However, HVAC temperature related data include inlet/outlet temperature, set-point, chiller water, and hot water.

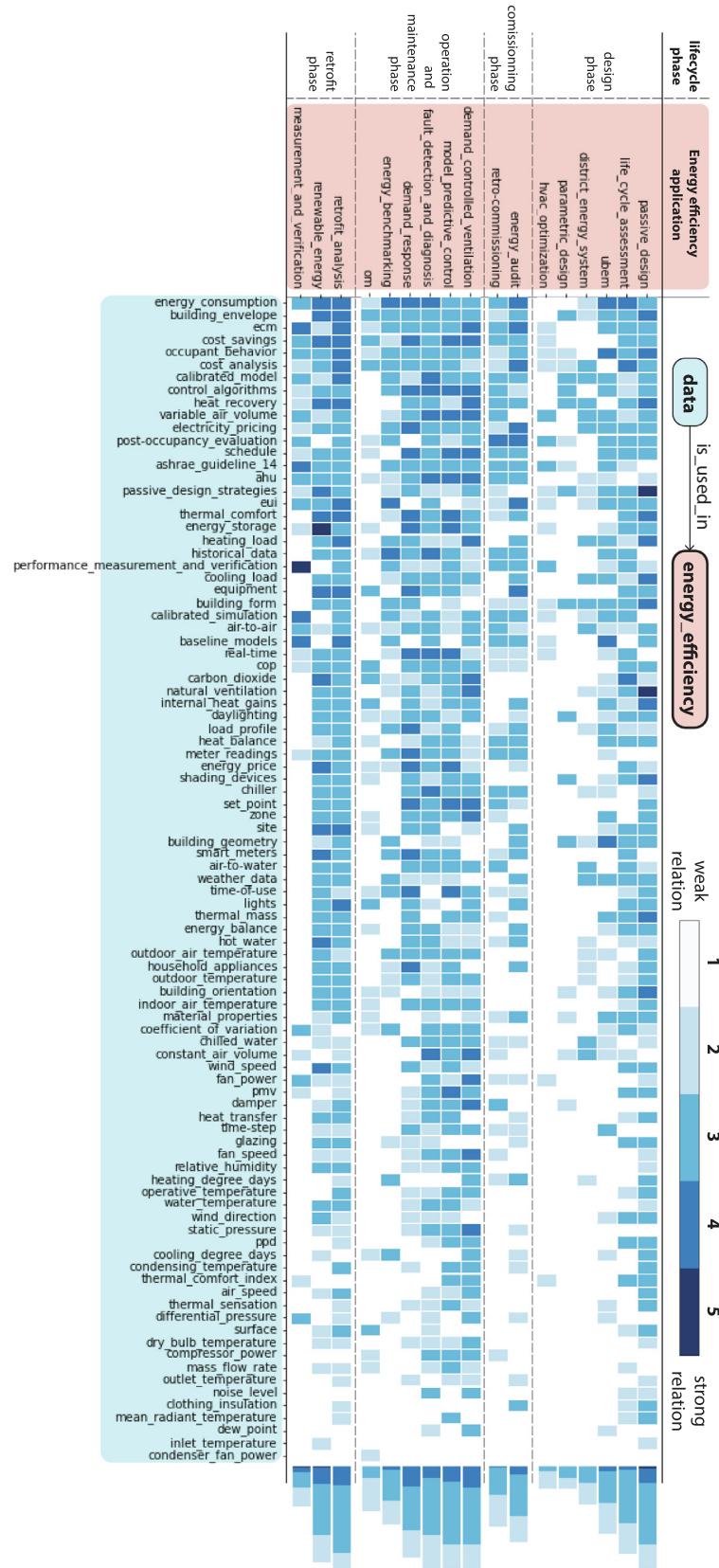


Fig. 8. The relation between data points and different energy-efficiency applications. The energy efficiency applications (shown on the Y-axis with red highlight) are chronologically grouped based on the building life-cycle phases. The data-points (shown on the X-axis blue highlight) are sorted based on their appearance frequency. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

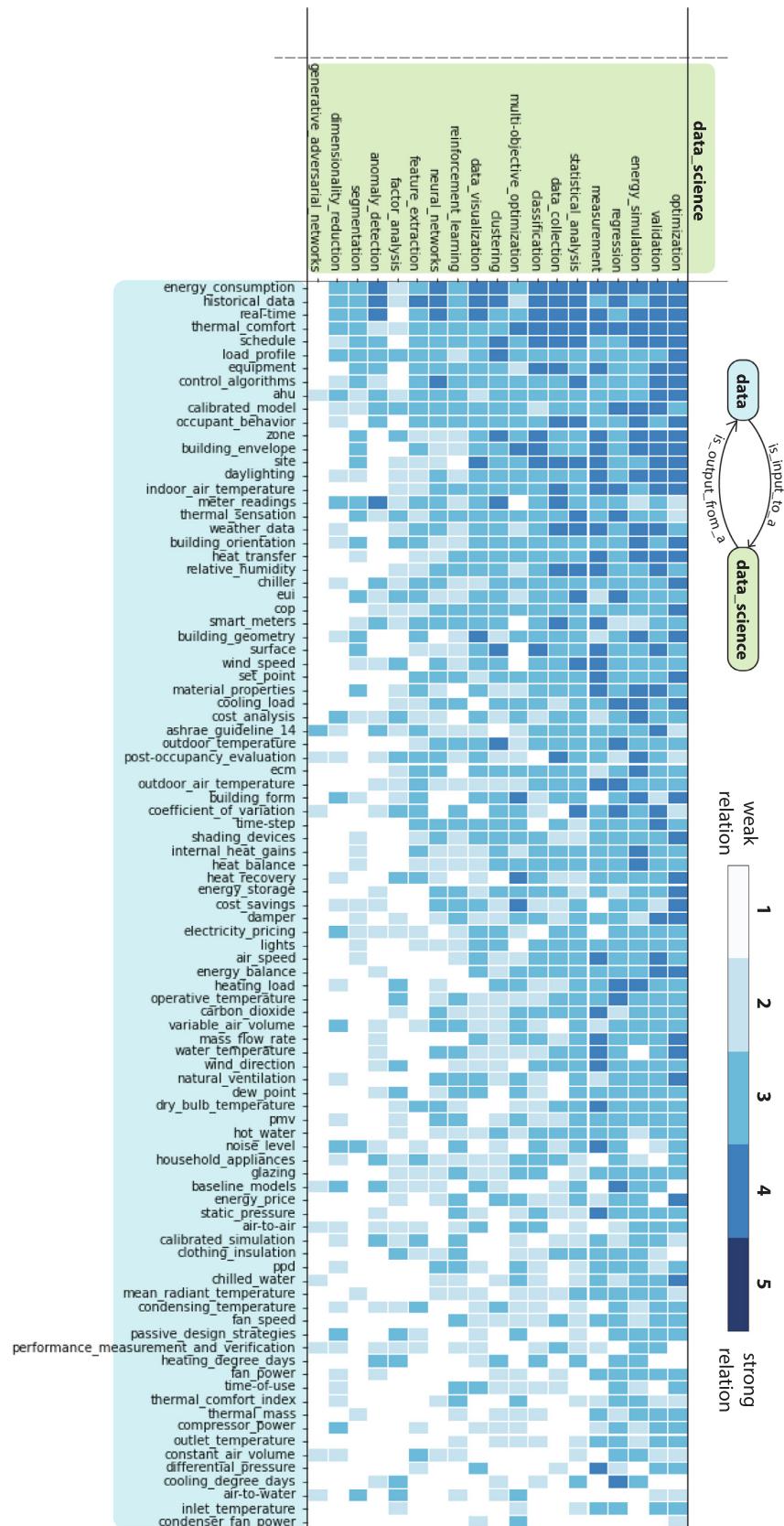


Fig. 9. The relation between data points (shown on the X-axis with blue highlight) and data-science algorithms (shown on the Y-axis with green highlight). Both of them are sorted based on the sum of the similarity per each row/column. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

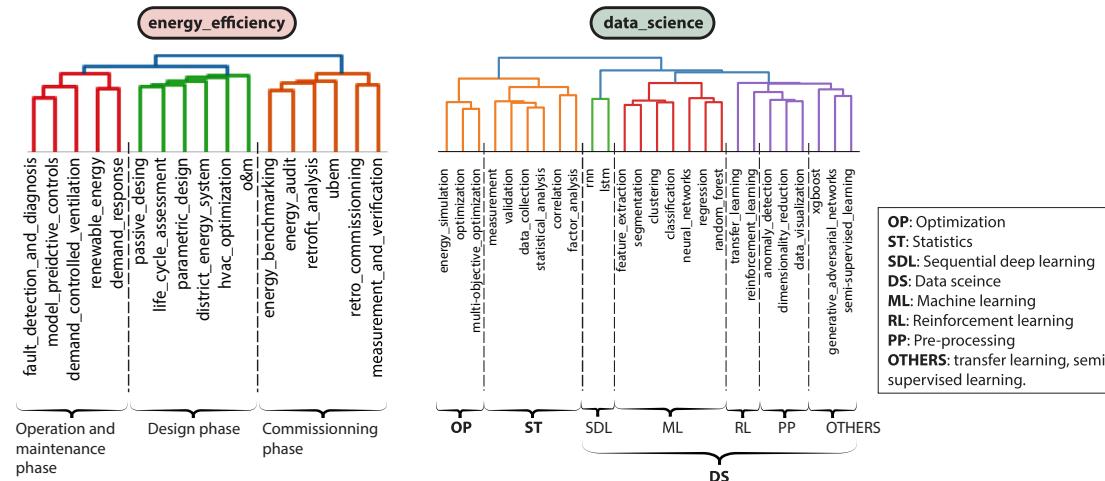


Fig. 10. The hierarchical agglomerative clustering (HAC) of the `energy_efficiency` and `data_science` categories using Ward's method.

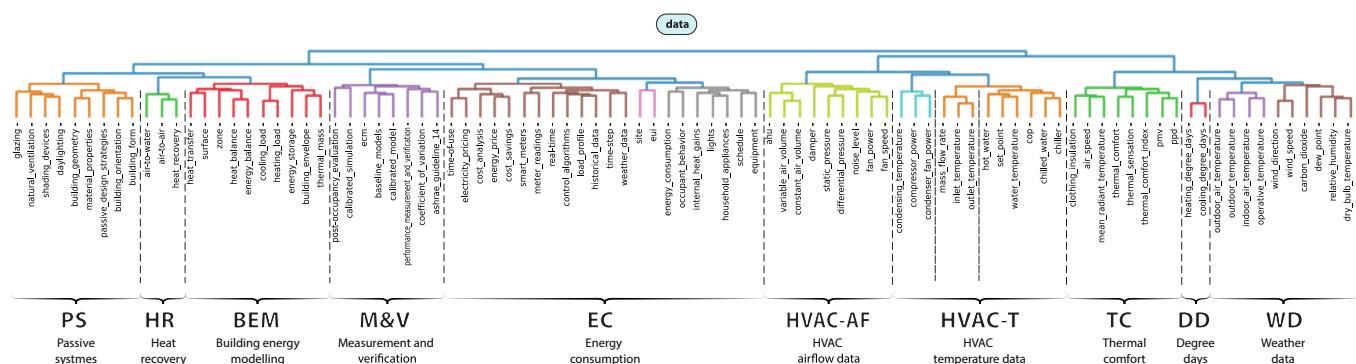


Fig. 11. The hierarchical agglomerative clustering (HAC) of the `data` keywords.

7. **Thermal comfort data (TC)** including the clothing insulation, thermal sensation, Predicted Mean Vote (PMV) and Percentage of People Dissatisfied (PPD), and Mean Radiant Temperature (MRT).
8. **Design degree days (DD)**, including heating degree days, and cooling degree days.
9. **Weather data (WD)** such as wind speed and direction, CO₂, dew point, temperature and relative humidity. From Fig. 12, it can be observed that the weather data (WD) has the highest correlation with all the other data groups. This insight could be attributed to the generic nature of weather data because it is used as input for most of the applications that use the rest of the groups except M&V. M&V on the other hand shows the lowest correlation with all the other data groups (with an average of 0.21). This conclusion might pertain to the unique nature of M&V that requires the evaluation of energy conservation measures (ECM) on building performance (in the operation phase) versus a baseline model (in the design phase). Such overlap between the operation and the design phases applications is rare. For example, Thermal comfort applications are used either in the design phase (i.e., Sizing, BEM, HVAC) or in the operation phase (i.e.: operation controls, Post-occupancy evaluation). Fig. 12 Also shows that Thermal Comfort (TC) and Energy consumption (EC) constitute the median of the categories with average correlations 0.67 and 0.58 respectively. These two categories falls in the area between the essential inputs/output data of energy-efficiency applications (Group1)

and the fine-tuner data (Group2). Any energy efficiency application is meant to balance between these two categories, i.e. to trade-off between Thermal comfort and energy consumption. Finally, Group 2 consists of categories of data that have high potential in different energy efficiency and data science applications but are not fully matured.

3.3.2. Data use across other categories

The final visualization in this section focuses on converging the three categories of data, energy_efficiency, and data_science with the clustering techniques (Fig. 13). Data science applications such as optimization (OP), machine learning (ML), statistical methods, and sequential deep learning (SDL such as RNN and LSTM) have relatively high relation with data. Those methods require a large amount of data from various resources for more accurate results [21]. However, it is can be seen that pre and post-processing methods (PP), Reinforcement Learning and transfer learning (RL), and other emerging models such as GANs and XGBoost do not have a strong relationship with many energy efficiency applications, especially those that do not change rapidly over the building lifecycle such as passive systems and building energy modeling. These types of data are usually generated during the early design phase of the building. Fig. 15 also confirms this claim as it shows very few data-science applications utilization during the design phase of the building compared to other phases.

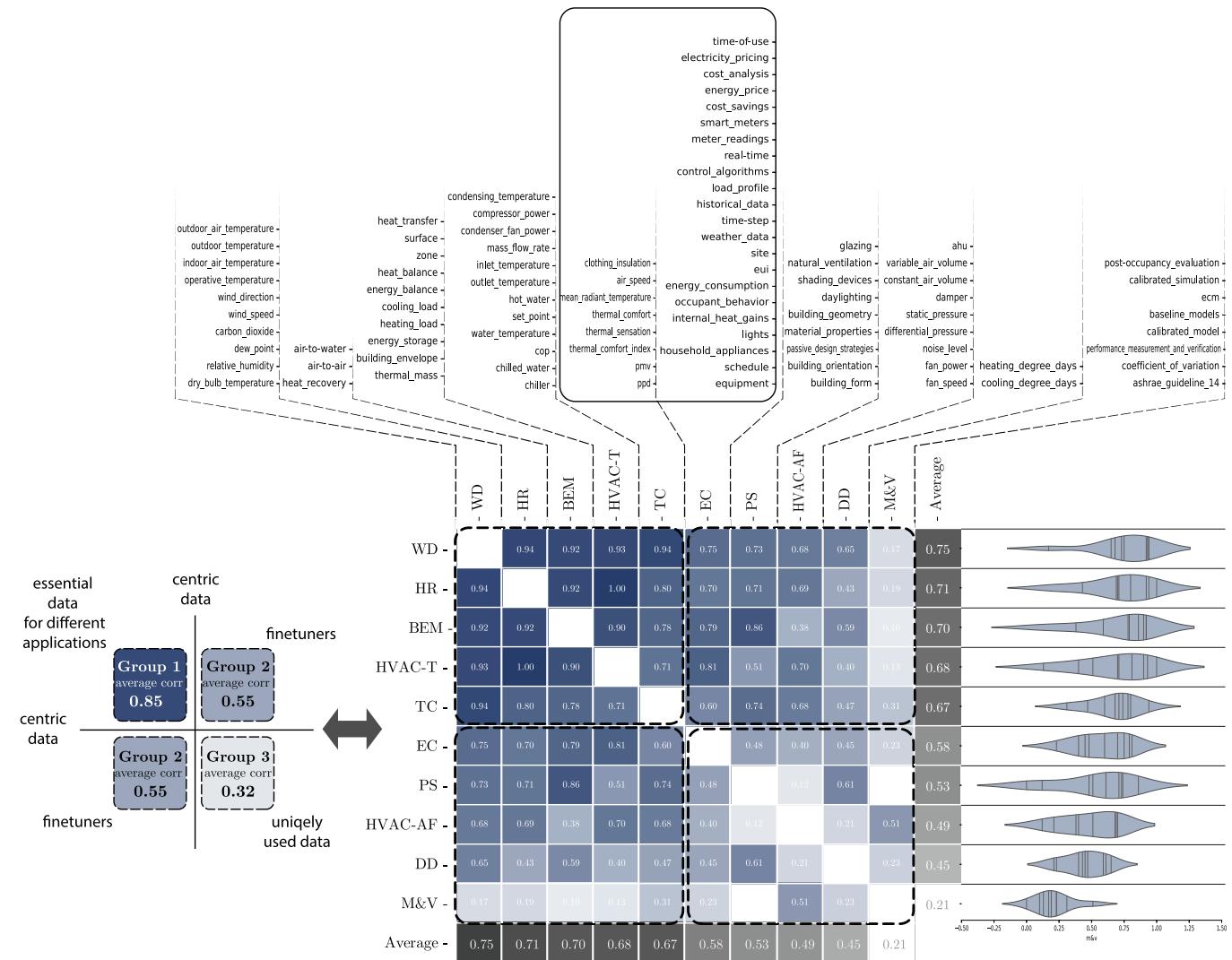


Fig. 12. The average correlation of the `data` keywords. The detailed breakdown of each individual keyword can be shown at <https://doi.org/10.6084/m9.figshare.13989653.v1>.

4. Discussion

The high level quantification of relationships between the data sources, data science techniques and various applications in the built environment provide the foundation for several key takeaways. This section expands upon the analysis of the results to provide high-level insights that can be used to guide future research. Each insight includes the discussion of a representative publication that illustrates the momentum or gap for that particular point.

Fig. 14 shows a comparison of the various data science techniques as compared to the energy efficiency applications for buildings as well as the life cycle phases. The following subsections outline key takeaways for the research community to consider.

4.1. What are the most common data analysis techniques?

The top five data science-related techniques found in the left side of Fig. 14 are intuitively those related to the traditional building energy domain techniques of simulation, optimization, neural networks, reinforcement learning, and statistical analysis.

The literature for the application of energy simulation and optimization for building energy efficiency applications is the most

voluminous due to the major efforts for decades of open-source simulation projects like EnergyPlus [36] and optimization engines such as BEopt [30], GenOpt [156], and jEplus [165]. More recently, to ease the application of machine learning and statistical analysis to building simulation, there has been development on interfacing with open-source programming languages such as Python [139] and R [75]. As illustrated in Figs. 14 and 15, building simulation, also known as building performance simulation or building energy modeling, plays a vital role throughout the building's lifecycle (passive and parametric design, M&V, FDD, LCA, energy audit, and retrofit analysis). The evident developments in the discipline of building performance simulation are supported by the rapid growth of the International Building Performance Simulation Association (IBPSA) over the last two decades and research efforts under the International Energy Agency's Energy in Buildings and Communities (IEA-EBC) program. To aid applications of building performance simulation, a wide variety of tools have been developed with more than 200 software tools and programs listed on the Building Energy Software Tools directory [71]. Crawley et al. provides an overview into the capabilities of twenty major building performance simulation programs [35]. Despite its long standing history and developments, challenges remain leading to opportunities in research and development. Hong, Langevin and Sun lists

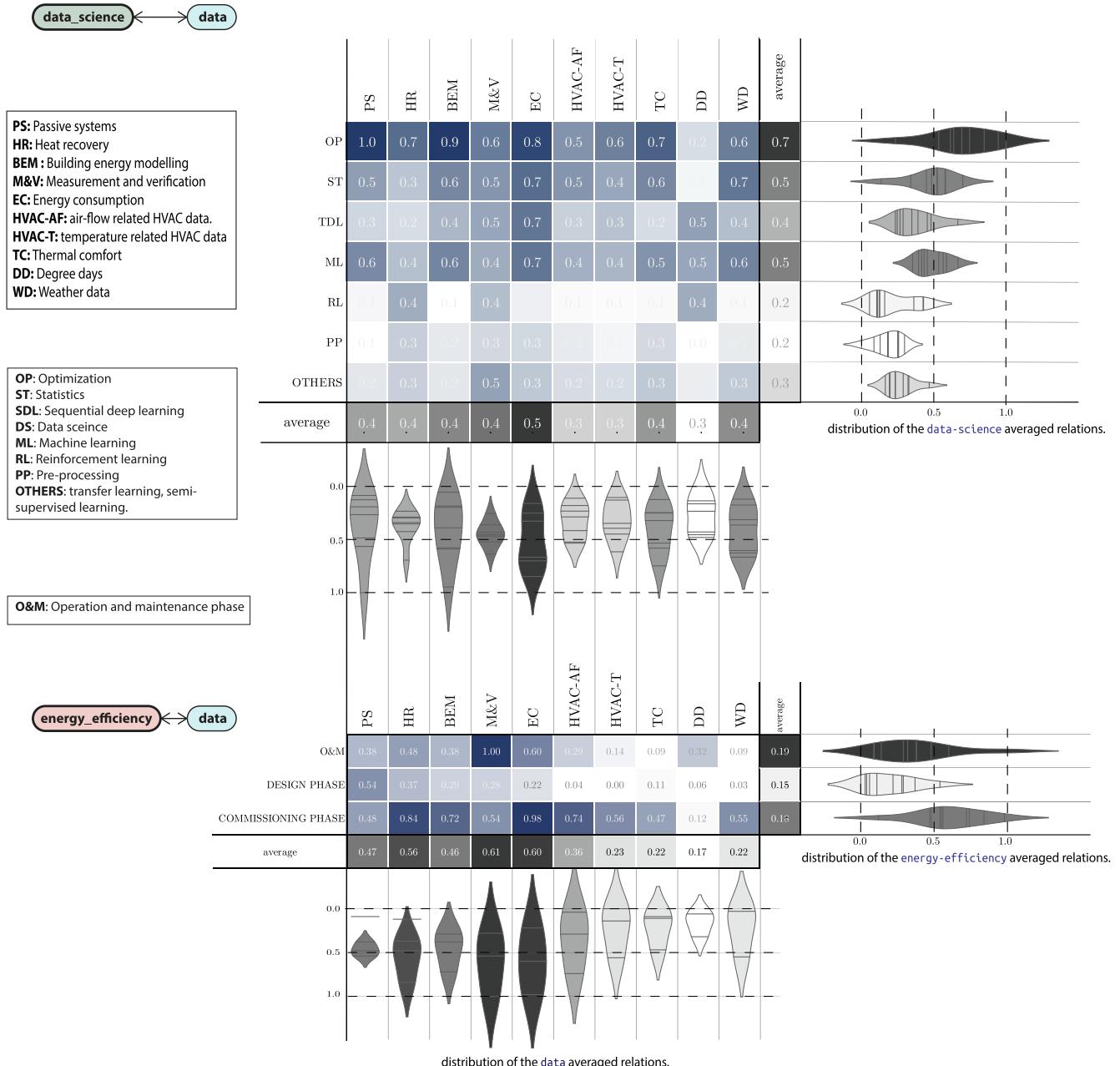


Fig. 13. The heatmap on the top shows the average relations between each pair of elements from the **data_science** and the **data** categories. The heatmap on the bottom shows the average relations between each pair of elements from the **energy_efficiency** and the **data** categories. A detailed version of the relations between every pair of keywords can be viewed at [2]..

the ten BPS challenges [67]. Table 1 lists each of these ten challenges. Additionally, we include the relevant publications and existing open-source repositories and data-sources that form the foundation in addressing these challenges.

4.2. What are the most explored building energy efficiency applications for data science?

From the perspective of applications using data-driven methods, automated fault detection and diagnosis followed by retrofit analysis, model predictive control, demand response, and energy benchmarking emerges as the most popular.

Fault detection and diagnosis (AFDD) is a field that has been growing rapidly since the early 1990s as a means of finding and fixing problems in building systems that result in energy waste and inefficiency. Katipamula and Brambley found the field to be maturing as early as 2006 [82]. Although matured, there have been recent developments in AFDD as a result of advancements in Artificial Intelligence techniques [168] and anomaly detection [121,122]. A challenge in the AFDD of building energy systems lies in that it is a class-imbalanced classification problem (i.e., there are few or no faulty training data). A Generative Adversarial Networks (GANs) integrated AFDD framework that generates artificial faulty samples in an adversarial way provides an innovative way to augment the training dataset, and have been shown to outperform traditional air handling unit [162] and chiller [161] AFDD methods.

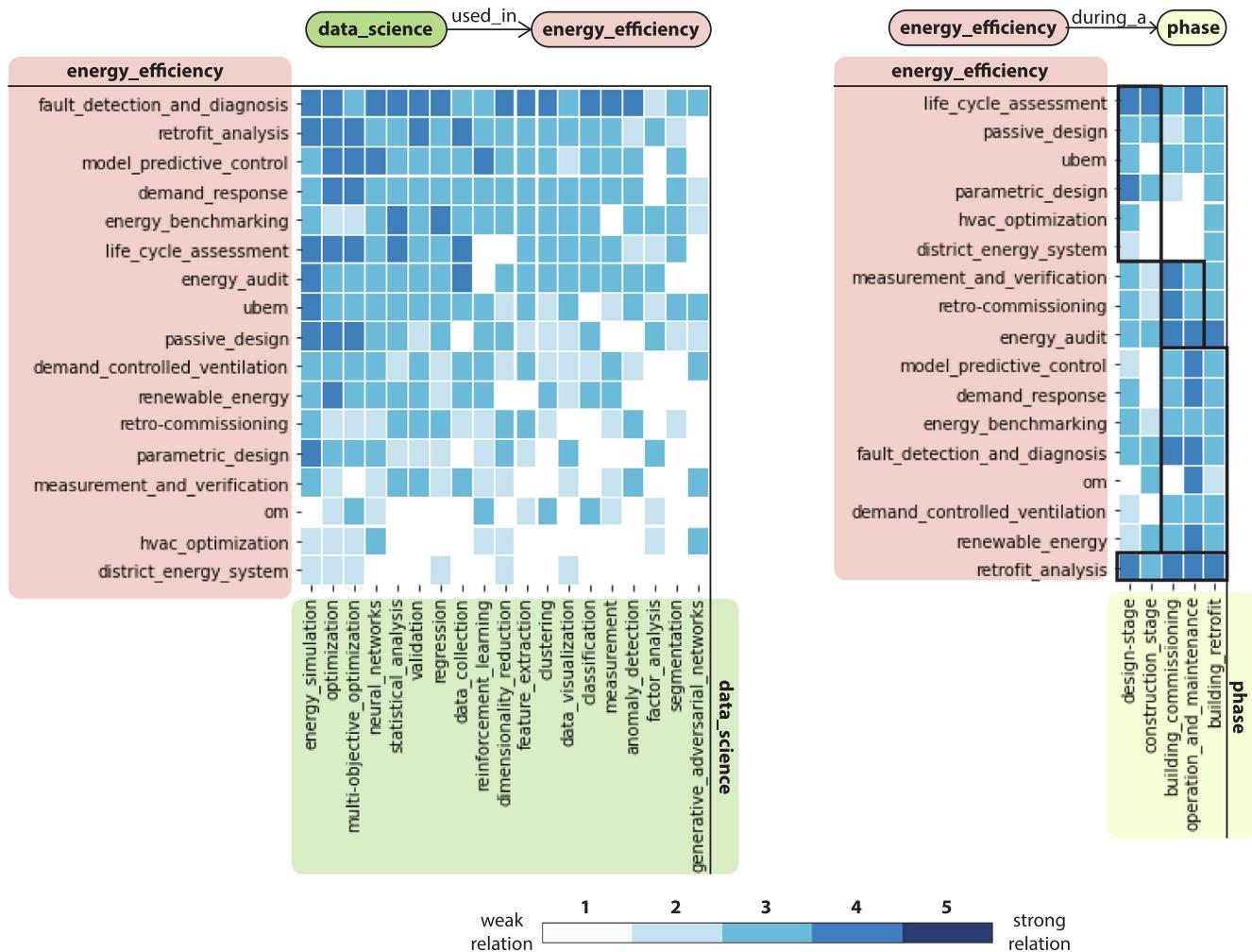


Fig. 14. The figure on the left shows the relation between data-science algorithms and energy-efficiency applications sorted based on usability. On the right, the relation between energy-efficiency applications and life-cycle phases are illustrated. The heatmaps' axes are all sorted based on strength correlation except life-cycle phase which is in a chronological order.

However, amongst the data science techniques listed in Fig. 14, GANs remain the least applied across various energy efficiency applications investigated. This is not surprising since it is a relatively new machine learning technique that might be beginning to emerge. GANs has also been applied on thermal comfort for generating balanced dataset [123,124]. Additionally, GANs recently have shown promising results in semi-real-time simulation of urban solar radiation simulation as well as urban wind simulation using Pix2Pix [32].

Retrofit analysis emerged as another top application across most techniques due to the influence of studies showing the large potential of upgrading the building stock [10]. As shown in Fig. 14, retrofit analysis often involves the use of physics-based simulation models, data collection and validation. The aim of retrofit analysis is to better understand the impacts of various factors on the retrofit of an existing building. However, buildings are made up of continuously changing sub-systems dynamically interacting with one another [62]. Since during a retrofit analysis training data of different scenarios is often not available, it is not surprising that retrofit analysis typically involves physics-based modeling that describes the complex dynamic interactions in buildings by a set of mathematical equations. Data collection followed by model calibration is often carried out to ensure the model's validity and thus credibility for the subsequent retrofit analysis [63]. Since building

operation and characteristics may change over time, continuous model calibration and data assimilation methodologies have also been proposed to ensure the simulation model remains reasonably representative of the actual physical building system [28,155]. pModel predictive control (MPC) has gained traction in the last two decades through numerous case study-based implementations [5]. Data science techniques are essentially used to obtain the predictive model and to solve the receding horizon control problem. Simulation (white-box), data-driven (black-box), and hybrid (gray-box) are the three main categories of controller models [7]. Neural network models are becoming more popular due to their stronger modeling capability [6]. In addition to model identification, optimization techniques are also used to determine the optimal control actions in the coming horizon. Typical algorithms include gradient-free methods such as GA and PSO, and gradient-based methods such as NLP and MILP [43]. Over the past few years, reinforcement learning (RL) is becoming a major competitor of MPC with its advantages of lower requirements on the predictive model and better adaptability [166]. However, it also comes with other problems such as higher data requirements and lower interpretability. The comparison between these two categories of optimal control approaches will be a major research topic in the future. Another important topic for further exploration is to reduce the implementation cost and promote the application of optimal con-

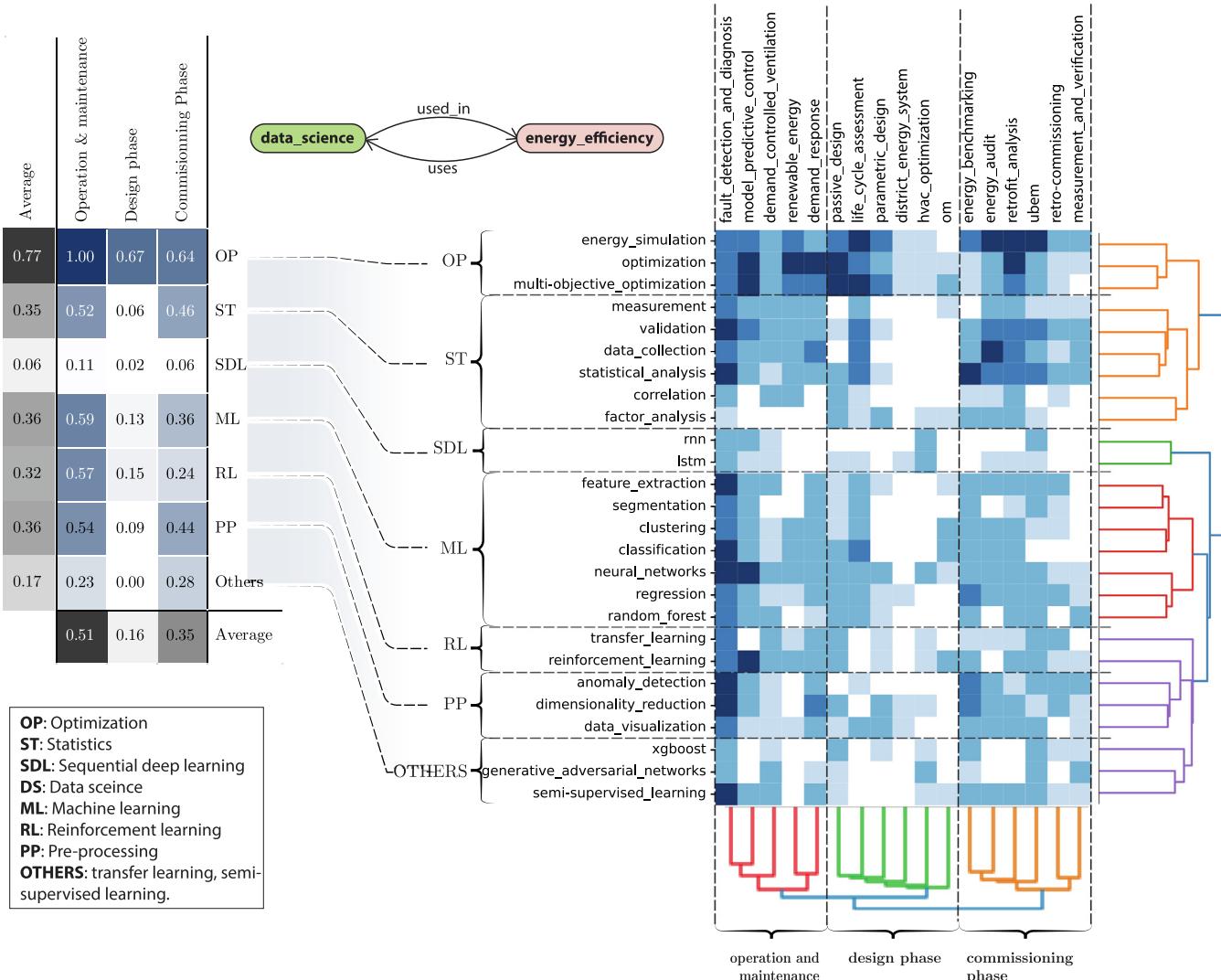


Fig. 15. The relationships between `energy_efficiency` applications and data science techniques: the HAC (on the right hand side) between energy efficiency (X-Axis) and data science (Y-Axis). The heatmap on the left hand side is a summarized version of the relation map by taking the average of each cluster.

trol. As a potential approach, transfer learning can be integrated with both MPC [25] and RL [159].

Demand response (DR), or demand-side management, is to reduce the energy cost by controlling the end-use customers' energy consumption with respect to energy prices. With the increasing penetration of renewable energy, the availability of renewable sources is another important factor to consider [129]. While reducing or shifting the electricity load, buildings still need to cater to the occupants' necessary needs such as lighting, office equipment, and environment conditioning. Thus, with no surprise, optimization is the most critical data science technique used for demand response [80]. Many other techniques are also involved in the process of designing DR programs (grid side) and helping the customers react to the programs (demand side). For example, energy simulation is a useful tool to design and evaluate the DR strategies [29]. Also, many clustering algorithms are applied to extract the typical load profiles to better understand the end users, estimate the participants' potential, and help decide the scheduling schemes [93]. Besides, at the city or grid level, the volume and variety of data generated when applying DR are both enormous. Therefore, techniques of data collection and dimension reduction are also essential in real implementation [76].

Building energy benchmarking is a concept which also comes up in the top five applications of data science. This field has grown based on the success of energy labeling schemes and city-wide data disclosures. Recent work in this area focuses on updating the modelling techniques [12] and even redefining the way buildings are categorized for benchmarking [118,164]. The large increase in open data sets available has created opportunities to target specific strategies to cities based on their specific needs [163] and using a combination physics-based and data-driven methods [135]. Data-driven improvements have been suggested related to generalizability [104] and interpretability [105].

4.3. What are the emerging application areas in which there are gaps?

On the other axis, it can be seen which energy efficiency techniques have the lowest relation to the data science concepts, indicating the gaps and opportunities for novelty. District energy systems shows up as the weakest, likely due to the only recent focus on the simulation and modelling of such techniques in the domain. Johansson et al. [78,79] looked at district energy systems and raised up some practical limitations such as the availability and quality of sensors. Also, district energy prediction is dependent

Table 1

Challenges of building performance simulation, reviews or key publications on the topic, and corresponding open-source repositories and data-sources.

Challenge	Relevant review(s) / publication(s)	Code	Open Data
1. Addressing the building performance gap	Type and definition [39]; Causes [148]; Credibility gap [16]	ObepME [81]; WinProGen (Occupant-behaviour gap) [20]	[72,108,20]
2. Modeling human-building interactions	Occupant modeling methodology in BPS [160] Challenges and opportunities [116]	Buildings.Occupants [154]	Occupant behavior [70]
3. Model calibration	Calibration methods and techniques [34,126,45]; Sensitivity analysis [146]	Bayesian calibration [27,125]; Optimization [22]	OpenStudio Calibration examples [113]
4. Modeling operation, controls and retrofits	Retrofit toolkits [89,90]; Model based commissioning [152]	Crowd-sourced ML for buildings [106]	Large, open meter data [107]
5. Modeling operational faults	Energy performance optimization [15]	Openstudio Fault Models Gem [26]	Open fault detection data [55]
6. Zero-net-energy and grid-responsive buildings	Gaps and needs [13,86]; Grid-responsive buildings [114]	BeOpt [31]	NZEB occupant behaviour [85]; Watts per person[117]
7. Urban-scale building energy modeling	Modeling methodology and workflows [127]; Challenges and future opportunities [66]	PyCity [140,144,1]	City buildings dataset [24] SynCity [135], NYC-UBEM [133]
8. Evaluating the energy-saving potential of building technologies at national or regional scales	E3 [98], Building stock energy prediction [92]	INTERDYME [9], PortableDyme [57], Scout [115]	Open data use for city-wide benchmarking [134]
9. Modeling energy efficient technology adoption	[44,73,53]	N.A.	Air-Conditioning Heating and Refrigeration Institute (AHRI) open data [61]
10. Integrated modeling and simulation	Progress, prospects, and requirements [33]	IFC[19], GBXML[42], OpenFOAM[74], EnergyPlus[36], Co-simulation e.g. obFMU[68]	Physics-based and data-driven modelling for NYC [135]

on both outdoor weather as well as control and social behaviour of consumers [58].

HVAC optimization and om (operations and maintenance) are contemporary topics, but only have small overlap with some of the more recent innovative data science techniques emerging. On the one hand, HVAC optimization has been reviewed by Selamat et al. [141] in three areas: HVAC operational parameters optimization, HVAC control system optimization, and building design optimization. His survey concluded that predictive optimization has more potential energy consumption reduction compared to conventional methods. Not only on the HVAC system scale, but also optimization should be done on the building design and building thermal dynamics. Other implementations of the data science for HVAC control in om have been conducted in recent years [59,17,46]. These issues include user security regarding data collection and storage, the lack of standardized data exchange schemes, and the lack of personnel with proper data science and domain knowledge.

Parametric design is also seen to be under-utilizing data science applications although it has gained much momentum in the last two decades [64]. This momentum is attributed to the advancement in Computer-Aided design CAD software as well as the emergence of user-friendly programming languages such as visual programming languages (VPL) [56]. Visual programming tools such as Revit Dynamo, and Rhino-Grasshopper has enabled end-use programmers to use data science algorithms in the design process. For example, Machine learning tools such as ANT [3] Lunchbox and OWL [84]; Optimization and multi-objective optimization such as OPOSSUM [157], octopus, Galapagos, and Optimus [37]; Energy Modelling such as Ladybug tools, [136], BuildFit [4]; Data visualization and deep learning using Gh_CPython [97]. These tools have grasped the attention of a large body of researchers and end-use programmers recently and may have a great potential for converging data science into the design process.

5. Conclusion

This paper outlined the text mining analysis of approximately 30,000 publications found in the top journals in the built environment analytics domain. This process aims to review the data science methods used in different building energy efficiency applications by mining large corpus of structured text from ELSEVIER journals. This process discovered high-level trends and potential gaps in the literature. Some data science methods have been extensively used in energy efficiency applications such as optimization, neural networks, statistical analysis, and energy simulation. However, there is still room for more opportunities of using other algorithms such as anomaly detection, factor analysis, segmentation, and GANs. Additionally, data-science methods are observed to be under-utilized during the commissioning and design phases of the building while saturated during the operation and maintenance phase. This could be attributed to the availability of ground-truth data during these lifecycle phases. Furthermore, different data sources are used frequently such as energy consumption-related data and BEM-related data. While other data sources are underutilized such as thermal-comfort related data, as well as HVAC-optimization related data. These results are extracted using a model based on the Word2Vec similarity metric. The results from this metric have shown consistency with the previous studies. Thus, researchers in this domain should utilize these results to determine which avenues are saturated, and therefore, will require much more effort to differentiate their work, and those which are emerging and have unexplored potential.

Having said that, we acknowledge some limitations related to this method. Firstly, each paper was treated equally over the text mining procedure, neglecting the effect of those seminal studies. Correspondingly, one future work is to introduce bibliometrics as a feature to distinguish the paper's significance. Also, this method results in a non-directed relationship graph. This means that word such as "occupant" can appear as a data (e.g. number of occupants)

