



# Keyword-based patent citation prediction via information theory

Farshad Madani <sup>a</sup>, Martin Zwick<sup>b</sup> and Tugrul Daim<sup>a</sup>

<sup>a</sup>Department of Engineering and Technology Management, Maseeh College of Engineering and Computer Science, Portland State University, Portland, OR, USA; <sup>b</sup>Systems Science PhD Program, Portland State University, Portland, OR, USA

## ABSTRACT

Patent citation shows how a technology impacts other inventions, so the number of patent citations (backward citations) is used in many technology prediction studies. Current prediction methods use patent citations, but since it may take a long time till a patent is cited by other inventors, identifying impactful patents based on their citations is not an effective way. The prediction method offered in this article predicts patent citations based on the content of patents. In this research, Reconstructability Analysis (RA), which is based on information theory and graph theory, is applied to predict patent citations based on keywords extracted from the abstracts of selected patents. After applying three classes of RA (variable-based analysis without and with loops and state-based analysis), nine specific IV states of a predicting model are extracted. These states involve the four keywords of “chamber”, “hous”, “main”, and “return”. Lastly, the abstracts of the patents are examined to identify the technical subjects relevant to smart building technologies for which these keywords are proxies.

## ARTICLE HISTORY

Received 22 March 2018  
Accepted 28 August 2018

## KEYWORDS

Patent mining; patent citation analysis; patent citation prediction; information theory; reconstructability analysis; OCCAM

## 1. Introduction

Due to the rapid growth in the number of patents filed around the world, patent databases have become tremendous technological resources. To elicit technical knowledge from patents and use it in managerial and engineering decisions, three major types of patent analysis have been developed: bibliometrics analysis, citation analysis, and keyword-based analysis. These methods are based upon on the three major parts of patent documents: metadata, main body, and citations. In bibliometrics analysis, the relation between inventors, firms, research institutes, countries, etc. are analyzed. In citation analysis, citations are considered as proxies for technological impacts, which means that when a patent cites another patent, the citing patent is impacted by the cited patent. In keyword-based patent analysis, keywords are employed as proxy for patent content, extracted by text mining tools, and analyzed by machine learning methods, specifically by clustering analysis. This paper is an example of this third type of analysis which is keyword-based.

Patent citation analysis is applied for various purposes, because patent citations can represent technological changes. Assessing future technological impacts (Lee et al. 2012), monitoring technological trends (Lee, Jeon, and Park 2011), predicting emerging technologies (Érdi et al. 2012), and exploring technology diffusion (Chang, Lai, and Chang 2009) are examples of such patent citation applications. But citation analysis inherently suffers from a drawback, namely that it is not able to rely on technical content of patents (Yoon and Park 2004). However, there is *correlation* between the number of citations and the technical richness of patents' contents; thus, we exploit the link between patent citations and patent content by using keyword-based analyses to make citation predictions. Predicting citations based on patent content may allow us to discover technologies that will be impactful in the future. Therefore, the research question of this paper is that if there are a set of keywords in the abstract of patents that the frequency number of the keywords can be applied to predict patent citation.

Reconstructability Analysis (RA) is chosen as the main methodology for citation prediction since RA is very well suited for predictive modeling that is exploratory as opposed to confirmatory. Other methodologies, such as network analysis and cluster analysis, are unable to predict citations. In this work, RA is specifically applied to smart building technologies to predict which patents will be cited in the future.

## 2. Literature review

Patents are good representatives for technological events. When a patent cites another patent, this means the citing patent is impacted by the cited patent (Lee et al. 2012). Patent citations are deployed to study many technological events such as linkage between science and technology (Li et al. 2014; Kousha and Thelwall 2017), knowledge flow and diffusion (Yoon and Lee 2008), patent valuation (Harhoff et al. 1999; Hu, Rousseau, and Chen 2012), stock market valuation (Hall, Jaffe, and Trajtenberg 2005), technological convergence between industries (Karvonen and Kässi 2013), emerging research areas (Kajikawa and Takeda 2009), emerging technologies (Breitzman and Thomas 2015; Kim and Bae 2017), the future impact of current technologies (Lee et al. 2012) and technology diffusion (Chang, Lai, and Chang 2009). Different methodologies like stochastic analysis (Lee et al. 2012), cluster analysis (Chang, Lai, and Chang 2009; Liu and Shih 2010; Érdi et al. 2012; Breitzman and Thomas 2015), and network analysis (Érdi et al. 2012) are applied to patent citation analysis, but patent citation analysis is subject to some drawbacks (Yoon and Park 2004). First, patent citation analysis only discovers individual relations between two patents. There might be semantically relation between two patents without any citations between them, Therefore, patent citation analysis does not identify the overall relationships among all patents. Basically, citations cannot cover the richness of potential information, so they limit the scope of analysis. Second, citation analysis is not able to consider semantic relationships between patents; it may even produce superficial or misleading indices.

To remedy the above-mentioned problem, patent researchers have developed keyword-based methods by applying text mining tools because keywords are good representatives of the content of patents (Madani and Weber 2015). The majority of keyword based methods are developed for the purpose of technology prediction (Yoon and Park 2004; Yoon and Park 2007; Jun 2014; Kim, Park, and Jang 2015) or other related topics such as technology monitoring (Lee, Jeon, and Park 2011), technology discovery (Lee, Yoon, and Park 2009),

and road mapping (Lee et al. 2008). The dominant methods deployed in the keyword based methods are cluster analysis methods and network analysis. Clustering methods such as *k*-mean (Kim, Park, and Jang 2015), principle component analysis (Lee, Yoon, and Park 2009) are utilized to group or to map patents based on their semantic similarity. Network analysis analyzes the relation between patents based on the intense of their similarity.

There is, however, a huge gap between patent citation analysis and keyword based patent analysis. As introduced above, many patent citation analysis methods and many keyword-based patent analysis methods have been developed, but no research has been done to bridge these two types of patent analysis. Bridging these two types of analysis would allow us to predict patent citation, which would improve significantly the accuracy of patent citation based methods (Yan et al. 2011).

None of the existing methodologies – network analysis or cluster analysis – is able to predict patent citations based on the content of patents (keywords). Network analysis studies the structure of relationships between entities like citations, and cluster analysis groups a set of similar objects, like keywords. To improve patent citation prediction, we need to apply a method that predicts patent citations based on the frequency number of keywords. Such a capability is provided by Reconstructability Analysis (RA) which yields models that predict citations based on keyword frequencies. RA is introduced in the next section.

### 3. Methodology

#### 3.1. Patents extraction

Patents extraction provides the data set for analysis, so the more accurate the extraction, the more accurate analysis we will have. Since we are prediction patent citations based on patent contents, we extract patents based on keywords representing smart building technologies. It is very important to have the correct keywords that address smart building technologies. In this research, an initial dictionary of keywords is developed based on the literature and experts' judgment, and then it is expanded by applying Google Adword capabilities. To identify the keywords which address smart building patents, the concept of smart building is divided into three main categories: 1) energy, 2) efficiency, and 3) building (Table 1). The first column titled “energy” represents all forms of energy consumption in buildings. To cover all possible wordings in the patents, the asterisk symbol (\*) is used as a wildcard. For instance, heat\* covers all possible variations such as heat, heating, heater, etc. The second column addresses all possible words representing the concept of “smart”, such

**Table 1.** “Smart Building” keywords used in the query shown in Figure 1.

Energy	Efficiency	Building
Heat*	Sav*	Home
Light*	Optimi*	House
Cool*	Manag*	Floor
Ventil*	Reduc*	
Refrig*	Control*	
Pump*	Smart	
HVAC	automat*	
	Sustainab*	
	Intelligent	

ABST((energy OR heat\* OR light\* OR Cool\* OR Ventil\* OR Refrig\* OR Pump\* OR HVAC) AND (efficien\* OR Sav\* OR optimi\* OR Manag\* OR Reduc\* OR control\* OR smart OR automat\* OR Sustainab\* OR intelligent) AND (building OR home OR house OR floor))

**Figure 1.** The query used to search in USPTO database through “LexisNexis Academic Universe”.

as saving, managing, reducing, etc. The last column signifies all possible words approximating “building”. The key words are extracted through several rounds by applying the literature (Kim, Stumpf, and Kim 2011) and Google AdWords.

To extract the patents information from the LexisNexis database, the query shown in Figure 1 is designed based upon field tags used by the database. To gather smart building patents, the query parameters were set for US patents issued between 1990 and 2013. All keywords mentioned in Table 1, were deployed to look in titles and abstracts of patents. After applying the query, 2483 patents were recognized. In this case, “abstract”, “application number”, “filing date”, and “cited patents” are represented by ABST, APPL-NO, FILED-DATE, and REF-CITED tags, respectively, to design the query shown in Figure 1. The output was in text file format.

### 3.2. Patents data preprocessing

To do Reconstructability Analysis, we need to elicit two main groups of numbers from the patents data set. The two groups are: (1) the numbers of citations of all the patents (the DV), and (2) the frequencies of all the keywords (the IVs). So the extracted patents data, which are in text format, need to be organized in two separate databases. The first database is provided to figure out the number of citations between the extracted patents. There might be patents citing the extracted patents, but they are not considered in the data set because they are not basically relevant to smart technologies, and consequently their keywords are not considered for the citation prediction. The second database is for text mining purposes to figure out the frequency numbers. The first database is a  $2483 \times 2483$  matrix reflecting relations between the patents. If a cell contains “1”, it means the patent located in the column is cited by the patent located in the row and if the cell shows “0”, it means there is no citation. The data fields of the second database are “patent number” which comprises the patent numbers of the extracted patents, and 235 fields that contain the frequency numbers of keywords in the abstracts.

Several tasks including syntax tagging, word stemming, and stop-word elimination are required to extract the keywords. In syntax tagging, words or terms are distinguished in the sentences, then their suffixes are removed via a stemming procedure. To omit words such as *the*, *is*, *at*, *which*, and *on* from the reserved words, a list of stop-words, given from RANKS.NL Website (“Stopwords,” “Ranks” Company 2014), are deployed.

To extract the most repeated keywords in the patents, we used Weka (Hall et al. 2009), written in Java and developed at the University of Wikato in New Zealand. We applied Weka’s StringToWordVector filter which removed numbers and stop words. We also removed general keywords, such as “Winter”, “mold”, and “dark,” that are not relevant

to the technologies we are trying to predict. At the end of this filtering process, 235 keywords remained. With these keywords, we created a database including patent number, abstract, frequency number of the keywords, and also citation number of the patents. The citation numbers come from the matrix mentioned earlier.

With this database, we apply a two-step preprocess to the data set to make it ready for RA. First, we create a spreadsheet containing the IVs, which are the frequency numbers of the keywords extracted from text mining, and the DV, which is the number of citations of the patents. Then the frequencies of the IVs are binned, because, as mentioned before, only nominal variables or discretized quantitative variables converted to nominal variables are usable in RA. After binning, the binned data including IVs and DV are organized in a text file in OCCAM format (Fusion, Willet, and Zwick 2012).

The DV is binned to two bins, because in the majority of cases (patents) citation number is zero. In this case, 92.9% of the patent citations are zero. Therefore, where  $DV = 1$  the citation number is zero, meaning the patent is not cited, and  $DV = 2$  means the citation number is more than zero, i.e. the patent is cited at least once.

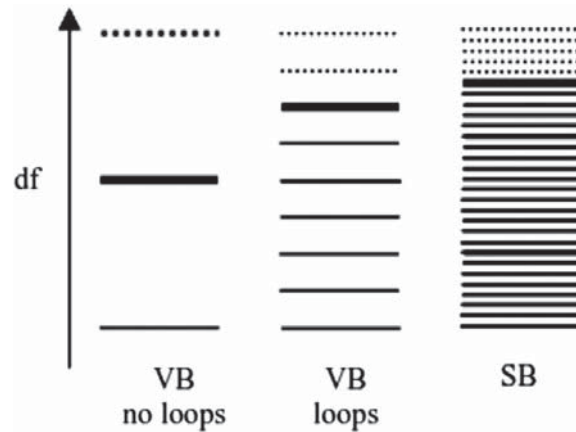
### 3.3. Reconstructability analysis

Reconstructability Analysis (RA), introduced first by Ashby (1964), is a method based on information theory and graph theory. RA has similarities with log-linear methods applied in social sciences (Bishop, Feinberg, and Holland 1978) and can be considered a machine learning technique such as those widely used in computer science (Perkowski et al. 1997). RA, like log-linear and machine learning methods, is applied in many different fields, including time-series analysis, classification, decomposition, compression, pattern recognition, prediction, control, and decision analysis (Zwick 2004a).

Basically, RA decomposes a probability or frequency distribution into component distributions (Klir 1985) by applying statistical multivariate analysis similar to log-linear methods (Zwick 2011) and logistics regression. RA also overlaps mathematically with Bayesian networks. Despite these similarities, RA, log-linear, and Bayesian Networks methodologies each has its own unique capabilities (Zwick 2011). For example, RA can analyze set-theoretic relations and arbitrary functions, is capable of state-based analysis (Zwick 2004c; Johnson 2005), and has a Fourier version (Zwick 2004b). RA and log-linear, but not Bayesian Networks, can utilize models with loops in both directed and neutral systems. (In directed systems, IVs and DVs are distinguished, and IVs are used to predict DVs. In neutral systems, IVs and DVs are not distinguished. Neutral systems simply try to find the relations between all the variables.)

Three classes of RA models are applied in this research: (1) variable-based (VB) models without loops, (2) variable-based models with loops, and (3) state-based (SB) models (which typically have loops). These models give us coarse, refined, and ultra-refined models, schematically depicted in Figure 2 (Zwick 2011). The bold lines in the figure show how complex a model might be accepted in each of the classes. In VB loopless models, the indicated bold line might be the most complex model that could be statistically acceptable. The dotted lines above the bold line indicate models that are statistically insignificant. Because VB models with loops make finer discriminations, a more complex such model might be statistically significant. Since SB models make even finer discriminations, they may allow a still more complex model to be accepted as statistically significant.





**Figure 2.** Degree of refinement of RA models (Zwick 2011).

In loopless models, there is only a single “predicting component”. For example, in the model ABCD:ABZ, A and B predict Z. (The first component, ABCD, allows for associations between the IVs.) Models with loops have more than one predicting component. For example, in ABCD:AZ:BZ model, two variables, A and B, separately predict Z, and these predictions are fused in the model. State-based models like ABCD:A<sub>1</sub>B<sub>2</sub>Z:B<sub>1</sub>C<sub>3</sub>Z predict Z with more specific information about states of IVs. The states are indicated by subscripts.

A web-based program developed at Portland State University (Portland, Oregon) was used to perform the RA analysis. The program is named “Organizational Complexity Computation and Modeling” (OCCAM), which is also an allusion to Occam’s (sometimes spelled Ockham’s) Razor, the principle of parsimony, important in modeling. The first precursor to OCCAM was written by Zwick and Hosseini (Hosseini and Zwick 1986). A detailed overview of RA (Zwick 2004a) OCCAM manual (Fusion, Willet, and Zwick 2012), and a review of OCCAM architecture (Willet and Zwick 2004) are available at the RA web site, <http://www.pdx.edu/sysc/research-discrete-multivariate-modeling>.

Three basic functions of OCCAM were employed for this data patent analysis:

- *Loopless analysis*, as the first step in RA, is used to discover what keywords (IVs) are individually the most predictive of patent citation (DV).
- *With-loops analysis* find three best models for patent citation predicting. These best models involve multiple and/or interacting IVs are based on BIC (Bayesian Information Criterion), AIC (Akaike Information Criterion), and Incremental p-values criteria. The AIC and BIC best model criteria select models based on a linear tradeoff of the information captured in a model and its complexity, with the BIC criterion penalizing models more heavily for complexity. The Incremental-p model picks the highest information model whose difference from the reference (independence) model is statistically significant, where there is also a path from the reference to the model, where every incremental step of the path is statistically significant. In all of the analyses, significance level is chosen as 5%. Actually, models are chosen based on  $\Delta$ AIC and  $\Delta$ BIC, which are *differences* of AIC and BIC of the models from values for the reference (the independence model); this results in good models having higher positive values of these differences. (Normally, when BIC or AIC values are given, they’re given as values for

particular models, not as differences from a reference model, and then good models are those that have lower or negative values.) These models with loops are the key outcomes of this research.

- *State-based analysis* helps to more deeply scrutinize the IVs selected in prior steps. There might exist some interactions between the IVs selected in prior steps that are undetectable in variable-based analysis. State-based models can use these subtle interactions for more accurate predictions.

#### 4. RA results

As mentioned before, RA is designed for nominal variables or for discretized quantitative variables converted to nominal variables. To make quantitative variables usable, it is necessary to bin them. Table 2 shows the binning of the four IVs that analysis (discussed below) revealed to be predictive. The four keywords shown in Table 2 are the best predictors in the variable-based analysis with loops (see section 4.2 and Table 4). For instance,  $Db = 2$  means “Chamber” is frequented at least one time, or  $Dg = 3$  means “hous” is frequented between 2 and 15. Also,  $Z$ , which is the dependent variable (DV), is binned to 1 and 2. If  $Z = 2$ , it means the patent is cited; otherwise it is not cited.

To bin the number of frequency of the keywords, a binning software program is applied. The binning program tries to create a set of bins that are equally sampled. Since for all keywords, most cases have no citations of the key words, a uniform distribution is impossible to achieve, but the binning program attempts a reasonable compromise. The first bin is necessarily always for zero frequency of keywords (IVs) or citations (DV). The number of additional bins then depends on how many cases are spread over 1 frequency/citation, 2 frequencies/citations, etc. For keywords *chamber*, *main*, and *return*, the number of cases where there are 1 or more is small (167, 137, and 84 compared to over 2000 for no frequency), so it makes sense only to add one extra bin for these keywords. For *hous*, however, there are 410 cases where this keyword is frequented once, and 267 where it is frequented more than once, so it makes sense to assign a separate bin for each of these situations. Other keywords are treated similarly.

**Table 2.** Relation between the bin numbers and the frequency of occurrence of the IVs (the keywords).

IV	Keyword	Values			
Db	chamber	Bin	1	2	
		Frequency	0	1–12	
		#Patents	2316	167	
Dg	hous	Bin	1	2	3
		Frequency	0	1	2–15
		#Patents	1806	410	267
Iz	main	Bin	1	2	
		Frequency	0	1–11	
		#Patents	2346	137	
Ja	return	Bin	1	2	
		Frequency	0	1–11	
		#Patents	2399	84	

#### 4.1. Loopless variable-based analysis

The result of loopless analysis is tabulated in Appendix 2, where uncertainty reduction progressively increases as models go from the bottom (the independence model) to the top (the most complex model considered). Uncertainty reduction is the primary information theoretic measure of the goodness of a model; it is analogous to %variance explained, but because of the logarithm term in the Shannon expression of uncertainty even small uncertainty reductions (for example, even 8%) can represent large effect sizes. Since uncertainty reduction is specific to information theoretic analyses, a more general measure of goodness of models, namely %correct in prediction, is also offered in this paper. The independence model, IV:Z, by definition, has no uncertainty reduction in the DV since no IVs are associated with it, i.e. so  $\% \Delta H(DV)$  is always zero in level 0. Also, %Correct(Data) starts from 94.7%, which means that a default prediction (the prediction without any knowledge of the presence or absence of any keywords) of no citation ( $Z = 1$ ) will be 94.7% correct. This is the default prediction, because that percentage of the patents is not cited by other patents in the sample. This means that we are struggling to predict very infrequent occurrence of  $Z = 2$  in this data set.

Among the IVs, three single IVs, namely Ac, Iz, and Ja, give the most uncertainty reduction. Of these three IVs, Ja, which is “return”, is a much stronger predictor than Iz (“main”) and Ac (“light”). As shown in Table 3, it has much higher values of both uncertainty reduction and %Correct. Note that Iz and Ac do not improve %Correct over its default value, but their predictive effect is captured in their non-zero uncertainty reduction.

As shown in Appendix 2, IV:IzJaZ is the best BIC model in the loopless analysis. Iz and Ja together predict Z with 83.33% uncertainty reduction and  $\Delta DF = 3$ . Also, IV:DbIzJaZ and IV:DgIzJaZ are the best AIC and incremental-p models, respectively.

#### 4.2. Variable-based analysis with loops

Allowing loops delivers more powerful variable-based models, as shown in Table 4. Since BIC is more conservative than AIC and Incremental  $p$ -value, the BIC model is the most reliable result of this study. Not surprisingly, the first and second single IV predictors are in the BIC model. In addition, two more IVs, Db and Dg, are added in the BIC model as new predictors. The AIC and Incremental  $p$ -values models are identical, and they include most of the same predicting variables. In comparison to BIC, they keep Dg, Iz, and Ja and drop

**Table 3.** Best single IVs prediction patent citations.

IV	$\% \Delta H(DV)$	$p$	% correct
Ja	54	.00	98.1
Iz	15	.00	94.7
Ac	3.3	.00	94.7

**Table 4.** Best variable-based models selected based on three criteria.

Criterion	Model	$\% \Delta H(DV)$	% correct	$\Delta DF$	$p$
BIC	IV:DbZ:DgZ:Iz:JaZ	86.6	98.8	5	0
AIC	IV:AmZ:AsZ:DgZ:EbZ:FdZ:Iz:JaZ	90.9	99.3	15	0
Incremental p	IV:AmZ:AsZ:DgZ:EbZ:FdZ:Iz:JaZ	90.9	99.3	15	0



Db. Also, the AIC and the p-value models have more uncertainty reduction and prediction correctness, but the delta degrees of freedom,  $\Delta DF$ , has increased greatly from 5 to 15, which means these two models are 3X more complex. All variable-based models generated by OCCAM are available in Appendix 3.

RA also gives specific information about any given model. The “Do Fit” option in OCCAM examines the given model in detail, and states exactly what DV values it predicts for all possible values of the predicting IVs (see Appendix 4). Since we are interested in predicting whether a patent will be cited or not, we look at column “Z = 2” for the model probabilities. If the probability of Z = 2 is considerably larger than its marginal probability of 5.27%, patent citation is more likely, assuming that the difference between the model probabilities and the independent model margins are statistically significant. This difference is examined by a Chi-squared test to see if its p-value is less than 0.05. We also consider only IV states where the calculated conditional probability distribution is different from uniform, and this difference is statistically significant with the same p-value cutoff. Eight combinations of the IVs that have passed both of the above criteria are shown in Table 5. For example, in the first row, having Db = 1, Dg = 1, Iz = 1, and Ja = 2, the model 100% predicts Z = 2. Or for the second row, having Db = 1, Dg = 1, Iz = 2, and Ja = 1, the model predict Z = 2 for 22.13%. The column named “rule” is the DV (Z) prediction result. The prediction rule indicates which DV state one should predict for any particular IV state. The rule is determined from the conditional probability distribution for the model, namely  $q(DV|IV)$ . Specifically, the rule indicates the *most probable* DV state, the DV state whose  $q(DV|IV)$  is maximum, for the particular IV. In the present case, if rule = 2, it means that one should predict Z = 2, i.e. that the patents containing the IVs (the keywords) will be cited. If rule = 1, it means that one should predict Z = 1, i.e. that the patents will not be cited. So, for example, in Table 5, for IV state (Db, Dg, Iz, Ja) = (1,1,2,1), the probability of Z = 1 is .779, while the probability of Z = 2 is .221, so one predicts Z = 1. In addition to Table 5, OCCAM generates individual tables for each of the predicting IVs in the model. These tables, shown in Appendix 4, reveal that none of the other three IVs alone ever predict Z = 2, but Ja predicts Z = 2 100% when Ja = 2. This shows that this analysis method can pick up predictive interaction effects involving several IVs even though the IVs are not individually predictive.

**Table 5.** Variable based analysis for the BIC model (IV:DbZ:DgZ:Iz:JaZ) (Probabilities  $p$  and  $q$  are shown as %).

IVs				Data			Model		rule
				freq	obs. p(DV IV)		calc. q(DV IV)		
					Z = 1	Z = 2	Z = 1	Z = 2	
Db	Dg	Iz	Ja						
1	1	1	2	56	0	100	0	100	2
1	1	2	1	81	77.8	22.2	77.9	22.1	1
1	1	2	2	4	0	100	0	100	2
1	2	1	2	4	0	100	0	100	2
1	2	2	1	14	85.7	14.3	85.56	13.4	1
1	3	1	2	8	0	100	0	100	2
1	3	2	1	13	23.1	76.9	22.1	77.9	2
2	1	1	2	5	0	100	0	100	2
2	3	2	1	7	0	100	3.8	96.2	2
–	–	–	–	2483	94.7	5.3	94.7	5.3	1

### 4.3. State-based analysis

In the state-based analysis done here, only the IVs (Db, Dg, Iz, and Ja) are used; the remaining IVs are ignored. The result of state-based analysis is given in Appendix 5. The best model for all three criteria (BIC, AIC, and incremental-p) is IV:Db1Z:Dg3Iz2Ja1Z:Iz1Ja1Z:Ja2Z:Z. This state-based model has four specific states added in the order listed to the independence model DbDgIzJa:Z. As shown in Table 6, this model is slightly better in terms of  $\Delta DF$ ,  $\% \Delta H(DV)$ , and  $\Delta BIC$  than the (BIC) best variable-based model with loops, namely IV:DbZ:DgZ:Iz:JaZ, previously given in Table 4.

The conditional probability distribution of the state-based model is shown in Table 7, which also provides other detail omitted in Table 5.  $p$ -rule is the  $p$ -value for testing the prediction rule, i.e. the deviation of the calculated conditional probability distribution  $q(DV|IV)$  from equal likelihood of  $Z = 1$  or  $2$ , and  $p$ -margin is the  $p$ -value for testing its deviation from the marginal  $p(DV)$  distribution, which is (94.7%, 5.3%). The table lists all the IV states for which  $p$ -rule  $\leq .05$ . Six IV states (not shaded) have a lower probability than the margin of being cited ( $Z = 2$ ) and nine IV states (shaded) have a higher probability of being cited. Components of the state-based model, IV:Db1Z: Dg3 Iz2 Ja1 Z: Iz1 Ja1 Z: Ja2 Z: Z, point to important IV states: (Dg = 3, Iz = 2, Ja = 1) and (Ja = 2) states always predict citation, while (Iz = 1, Ja = 1) states never predict citation.

### 4.4. Prediction rules

As pointed out earlier, each of the three classes of RA analysis has a different degree of refinement. Variable-based (VB) analysis without loops, VB with loops, and state-based

**Table 6.** Comparing the models resulted from variable-based analysis and state-based analysis

Model		$\Delta DF$	$p$ -value	$\% \Delta H(DV)$	$\Delta AIC$	$\Delta BIC$	Inc. $p$ -value	$\%C$ (Data)
State-based	IV:Db1Z:Dg3Iz2Ja1Z:Iz1Ja1Z:Ja2Z:Z	4	0	87.3	887.0	863.8	0.00	98.8
Variable-based	IV:DbZ:DgZ:Iz:JaZ	5	0	86.6	878.8	849.7	0	98.8

**Table 7.** State-based BIC model IV:Db1Z:Dg3Iz2Ja1Z:Iz1Ja1Z:Ja2Z:Z (Probabilities  $p$  &  $q$  are shown as %).

				Data					Model				
				obs. p(DV IV)			calc. q(DV IV)		rule	#correct	%correct	p-rule	p-margin
IV				freq	Z = 1	Z = 2	Z = 1	Z = 2					
Db	Dg	Iz	Ja										
1	1	1	1	1583	100	0	100	0	1	1583	100	0.00	0.00
1	1	1	2	56	0	100	0	100	2	56	100	0.00	0.00
1	1	2	1	81	77.8	22.2	79.4		1	63	77.7	0.00	0.00
1	1	2	2	4	0	100	0	100	2	4	100	0.05	0.00
1	2	1	1	343	100	0	100	0	1	343	100	0.00	0.00
1	2	1	2	4	0	100	0	100	2	4	100	0.05	0.00
1	2	2	1	14	85.7	14.3	79.4		1	12	85.7	0.03	0.01
1	3	1	1	208	100	0	100	0	1	208	100	0.00	0.00
1	3	1	2	8	0	100	0	100	2	8	100	0.00	0.00
1	3	2	1	13	23.1	76.9	21.2	78.8	2	10	76.9	0.04	0.00
2	1	1	1	65	100	0	100	0	1	65	100	0.00	0.06
2	1	1	2	5	0	100	0	100	2	5	100	0.03	0.00
2	2	1	1	42	100	0	100	0	1	42	100	0.00	0.13
2	3	1	1	27	100	0	100	0	1	27	100	0.00	0.22
2	3	2	1	7	0	100	3.4	96.6	2	7	100	0.01	0.00
-	-	-	-	2483	94.7	5.3	94.7	5.3	1	2454	98.8	-	-

All of the IV states are statistically significant ( $p$ -rule  $\leq .05$ ).

**Table 8.** Patent citation predictions: the probabilities of being cited for the 9 keyword combinations all exceed the marginal probability (on the last line) of 5.3%.

#	Keyword frequency				Freq.	Probability (%)	
	Chamber (Db)	Hous (Dg)	Main (Iz)	Return (Ja)		Not cited	Cited
1	0	0	0	> 0	56	0	100
2	0	0	> 0	0	81	79.3	20.6
3	0	0	> 0	> 0	4	0	100
4	0	1	0	> 0	4	0	100
5	0	1	> 0	0	14	79.3	20.6
6	0	> 1	0	> 0	8	0	100
7	0	> 1	> 0	0	13	21.2	78.7
8	> 0	0	0	> 0	5	0	100
9	> 0	> 1	> 0	0	7	3.4	96.5
–	–	–	–	–	2483	94.7	5.3

All of the IV states are statistically significant ( $p\text{-rule} \leq .05$ ).

(SB) analyses allow coarse, refined, and ultra-refined modeling respectively (Kramer et al. 2012). In this case, four variables with the largest effects in DV prediction, for VB without loops analysis, are Db, Dg, Iz, and Ja. Furthermore, more detailed models made of the four variables are recognized in VB with loops and SB analyses. Db, Dg, Iz, and Ja represent “chamber”, “hous”, “main”, and “return” keywords correspondingly. While both Table 5 and Table 7 contain detailed information about Db, Dg, Iz, and Ja as predictors of Z, Table 7 is based on the more refined state-based model, so we use Table 7 as the basis of the prediction. Table 8, based on Table 7, shows the citation prediction of the IV states whose  $p$ -values are less than 0.05.

To display which combination of keyword frequencies predicts the patent citations, we need to convert the bin numbers of the four IVs in Table 7 into their corresponding frequency numbers or ranges; this is done in Table 8. The relation between the bin numbers and the frequency numbers are given above in Table 2. So, for example bin number 1 for Db, the upper left number in Table 7, corresponds to frequency 0, the upper left number in Table 8. Table 8, shows how we can predict the patent citations via the keywords frequencies. Specifically, it shows nine predictions from keywords frequencies of whether patents with these frequencies would be cited. Note that two of the IV states ( $\{Db = 1, Dg = 1, Iz = 2, Ja = 1\}$  and  $\{Db = 1, Dg = 2, Iz = 2, Ja = 1\}$ ) in Table 7 which have a prediction rule of 1 (which predict not being cited), their  $q(DV = 2|IV)$  are 20.6%, almost 4 times bigger than 5.276%, which is the default probability of being cited in the data.

Of the nine prediction rules (rows) in Table 8, seven are combinations of the keywords that predict that patents will be cited. These seven can be grouped to summary rules I & II which are:

**I: return > 0**

**predicts  $p_{\text{citation}} = 100\%$ , #patents = 77**

This encompasses individual rules 1, 3, 4, 6 and 8 in the above table. The second summary rule has primary and secondary components which encompass individual rules 7 and 9, respectively, in the table. This summary rule is

For Return = 0,

**II: hous > 1 & Main > 0**

**predicts  $P_{\text{citation}} = 78.7\%$ , #patents = 13**

**Ila: hous > 1 & Main > 0 & Chamber > 0**

**predicts  $P_{\text{citation}} = 96.5\%$ , #patents = 7**

## 5. Discussion

To see how these four keywords enhance patent citation prediction, one must dig into the patents to understand to what the keywords are referring in the abstracts. To do this, rules I and IIa are considered as the basis of the discussion (see section 4.4). Rule I expresses the fact that if a smart building patent contains “return” or its derivatives in its abstract, the patent will definitely be cited (100% probability). Rule IIa indicates that a smart building patent will be cited with 96.5% probability if it contains “main” and “chamber” and at least two derivatives of “house” in its abstract.

“Return” has two main roles in the abstracts of smart building patents. First, “return” occurs in the abstract of smart building patents when something such as air, or water is circulating in a HVAC system and also occurs when someone such as a user or an occupant is going back to a building. For example, “return air” refers to how “return air” may be used as a cooling source in HVAC systems in US20080265046; “return water” acts as a part of pool heater in US5560216, an automatic washer in US5241843, a water heater system in US20080265046, and a steam heating system in US20080223947. In addition, “occupants return” and “user return” are mentioned in the abstract of patents related to communicating control systems (such as US20080217419, US8386082 and US20110125329) where the system manages the energy consumption based the presence or absence of a user or an occupant. Second, “return” is used in the abstract of smart building patents when a main part of the invention functions to return something. “A return plenum pressure controller” and “a return fan control system” in an HVAC system (US8326464, US20130096722, and US20100057258), “return line connecting an output of the chiller to an input of the cooling tower” in a climate control system (US20100201125), “return ducts” in a dehumidifier

**Table 9.** Predicted inventions that will be cited.

Prediction Rules	Key Phrases	Inventions
<i>Rule I:</i> (return > 0)	<i>return users, return occupants</i> <i>return fan control system</i> <i>return line connecting an output of the chiller to an input of the cooling tower</i> <i>return duct</i> <i>non-return valve</i>	Communication Control Systems HVAC Systems Climate Control Systems  Dehumidifier Systems Water Supply Systems Power Converters
<i>Rule IIa:</i> (main > 0, hous > 1, chamber > 0)	A casing <i>houses</i> semiconductor modules, constituting a <i>main</i> circuit for power conversion; . . . Within the casing, a cooling <i>chamber</i> including a coolant passage is formed, and a <i>chamber</i> wall of the cooling chamber is formed with a thermally conductive material. At least the semiconductor modules, are <i>housed</i> inside the cooling <i>chamber</i> , and at least the capacitor and the control circuit are disposed outside the cooling <i>chamber</i> . An automatic biomass fuel burning heating device and method comprising a burn <i>chamber</i> having . . . A burn <i>chamber</i> having a feed auger opening in a burn <i>chamber</i> wall. The control maintain supply air into the brooder <i>house</i> at a constant set temperature by varying the volume rate of air flow based on air inlet temperature and the temperature in the burn <i>chamber</i> .	Waste litter heater

system (US20060086112), and “non-return valve” in a water supply system (US7066197) are examples of the main parts that function to return something.

“Main”, “chamber”, and derivatives of “house” in the smart building patents provide inventions that generate a form of energy, more specifically heat, or convert heat to another form of energy. A group of patents including US7978468, US7969735, US20110235270, and US20100188814 explain “a power converter that converts input power to a specific type of power and outputs the power resulting from the conversion”. In this technology, a semiconductor module is housed inside a cooling chamber with the peripheral wall thereof constituted of a thermally conductive material to reduce the effect of heat from the semiconductor module on the other components. In addition, patent US20060236906 describes another invention that provides an automatic waste burning heating device. This technology is applied in poultry industry to burn poultry litter to generate and consume heat in buildings for different purposes.

In summary, it is expected that smart building inventions related to communication control systems, HVAC systems, climate control systems, dehumidifier systems, and water supply systems will be cited, according to rule I. In addition, smart building inventions related to power converters and waste litter heater will be cited, according to rule IIa. The summary of predictions and their related rules are illustrated in Table 9.

## 6. Conclusion

Patent citation is used as a proxy for technological impact studies, but is subject to the limitations of each patent’s contents (Yoon and Park 2004). The importance of patent citation prediction is revealed when we see 94% of smart building patents of this study are not cited. In this research, a keyword-based method is developed to predict patent citations. The keywords are extracted by applying Weka, a text mining software program developed by the University of Waikato (New Zealand) (Hall et al. 2009). Keywords are analyzed by Reconstructability Analysis (RA) (Ashby 1964) to discover keyword patterns in promising patent citations. This method enables us to both predict patent citations as proxies of technological impacts and to find out which aspects of technologies cause the impacts by interpreting the associated keywords.

Three different classes of RA searches are applied: (1) variable-based models without loops, (2) variable-based models with loops, and (3) state-based models. These models give us coarse, refined, and ultra-refined models, as shown in Figure 2 (Zwick 2011). As a result of RA analysis, four keywords, including “chamber”, “house” and its derivatives, “main”, and “return”, emerged as the keywords whose frequency in a patent related to smart building technologies which were the most predictive indicators of likelihood to be cited by other patents. The specific combinations of frequency of the keywords are summarized in Table 8. According to the keywords investigated in the abstracts, some smart building inventions, shown in Table 9, are predicted to be cited.

Practitioners can take the advantage of our method to not only explore in patent databases to find those patents leveraging technological changes in their industry like smart buildings, but also dig into the patents more efficiently by considering the keywords recognized by RA to identify the technologies. In keyword-based studies, like ours, having a comprehensive thesaurus or dictionary of the research area is ideal for the researchers, but professional thesaurus or dictionaries are rarely available for

specific emerging technologies, such as smart building technologies. Semantic analysis can be employed as a remedy in the future studies to compensate for the lack of professional thesaurus or dictionaries. WordNet, a lexical database made by Princeton University (Princeton University 2010) enables researchers to analyze their corpus semantically, e.g. patents, to extract all possible keywords related to the main concept of their research.

The use of Reconstructability Analysis in this research makes several theoretical contributions. First, the analytical procedures described in this paper show how data can be analyzed in an exploratory mode, eliminating the need to hypothesize and then test specific predictive relations whose form is explicitly specified. Second, these procedures allow one to detect interaction effects involving multiple predictors, even when the individual IVs do not have significant predictive effects. Third, the predictive models are conceptually transparent, being simple conditional probability distributions. This differentiates these methods from other data mining techniques, such as neural networks, which are “black boxes,” where predictive models are not as easy to interpret and are not strongly associated with statistical measures. Fourth, this work outlines a hierarchical procedure where exploratory analysis is done at different degrees of refinement; this allows flexibility in the analysis under varying conditions of sample size and computational capabilities. (RA requires sample sizes that are larger than those required for standard linear regression analyses.)

Ideally, these methods should be applied in conjunction with cross-validation techniques, i.e. subjecting models obtained from training data to separate test data; this has not been done in this study. Future research should include cross-validation assessment of predictive models. It should expand % correct calculations to include the explicit analysis of sensitivity and specificity, i.e. the subdivision of incorrect predictions into false negatives and false positives. Also, more specifically, as an expansion of this particular study, analysis should consider the predictive keywords found in models favored by the less conservative model selection criteria (AIC and IncrP) that are not found in models favored by the preferred model selection criterion (BIC). As shown in Table 4, these keywords include variables Am, As, Eb, and Fd, which correspond to keywords control, dev, head, and flow. In addition, although state-based calculations take considerable computer time, state-based runs can be done with additional predictors, i.e. 5 or 6, as opposed to the 4 reported in this paper.

## Abbreviations

Abbreviation	Keyword	Description
RA	Reconstructability Analysis	The methodology used to predict patent citations based on keywords extracted from abstract of patents.
IV	Independent Variable	IVs are the inputs (keywords) of the prediction model.
DV	Dependent Variable	DV is the output of the model, predicted by the IVs.
VB	Variable-based	Three different classes of RA models are applied in this research: 1) variable-based (VB) model without loops, 2) variable-based (VB) model with loops, and 3) state-based (SB) models (which typically have loops).
SB	State-based	



Abbreviation	Keyword	Description
BIC	Bayesian Information Criterion	BIC and AIC are generally used to define the best models. These criteria trade off uncertainty reduction and model simplicity in different ways (Kramer et al. 2012). BIC, the more conservative criterion, is favored in this study. A third criterion, Incremental-p, is also used; it picks the highest information model whose difference from independence is statistically significant and for which a path from independence exists where each increment of complexity is also statistically significant,
AIC	Akaike Information Criterion	
DF	Degree of Freedom	The number of parameters in a model, the measure of its complexity.
H	Uncertainty	Information theoretic measure of variable unpredictability. The reduction of uncertainty of a DV achieved by knowing the IVs is roughly analogous to the %variance explained, except small values of uncertainty reduction can indicate big effect sizes.

## Notes on contributors

**Farshad Madani** recently received his PhD in technology management from Portland State University, Portland, Oregon. He received his BSc and MSc in Industrial Engineering from Sharif University of Technology, Tehran, Iran. His current research includes technology intelligence, patent mining, and data mining. Farshad received the Dr Dryden memorial award for his excellent academic performance in spring 2017.

**Martin Zwick** was awarded his Ph.D. in Biophysics at MIT in 1968, and joined the Biophysics Department faculty of the University of Chicago in 1969. Initially working in crystallography and macromolecular structure, his interests shifted to systems theory and methodology, the field now known as the study of chaos, complexity, and complex adaptive systems. Since 1976 he has been teaching and doing research in the Systems Science Program at Portland State University. During the years 1984-1989 he was director of the program; in 2009, he was awarded the Branford Price Millar award for faculty excellence. His main research areas are information theoretic modeling, machine learning, theoretical biology, game theory, and systems theory and philosophy. Scientifically, his focus is on applying systems theory and methodology to the natural and social sciences, most recently to biomedical data analysis, the evolution of cooperation, and sustainability. Philosophically, his focus is on how systems ideas relate to classical and contemporary philosophy and how they can help us understand and address societal problems.

**Dr. Tugrul Daim** is a Professor and Director of the Technology Management doctoral program at Portland State University. He is the Editor-in-Chief of International Journal of Innovation and Technology Management and North American Editor of Technological Forecasting and Social Change. He received his BS in Mechanical Engineering from Bogazici University in Turkey, MS in Mechanical Engineering from Lehigh University in Pennsylvania, MS in Engineering Management from Portland State University, and Ph.D. in Systems Science: Engineering Management from Portland State University in Portland Oregon.

## Disclosure statement

No potential conflict of interest was reported by the authors.

## ORCID

Farshad Madani  <http://orcid.org/0000-0002-4451-9678>

## References

- Ashby, W. R. 1964. *Constraint Analysis of Many-dimensional Relations*. General Systems Yearbook. Fort Belvoir, VA: Defense Technical Information Center.
- Bishop, Y., S. Feinberg, and P. Holland. 1978. *Discrete Multivariate Analysis*. Cambridge, MA: MIT Press.
- Breitzman, A., and P. Thomas. 2015. "The Emerging Clusters Model: A Tool for Identifying Emerging Technologies Across Multiple Patent Systems." *Research Policy* 44 (1): 195–205.
- Chang, S.-B., K.-K. Lai, and S.-M. Chang. 2009. "Exploring Technology Diffusion and Classification of Business Methods: Using the Patent Citation Network." *Technological Forecasting and Social Change* 76 (1): 107–117.
- Érdi, P., K. Makovi, Z. Somogyvári, K. Strandburg, J. Tobochnik, P. Volf, and L. Zalányi. 2012. "Prediction of Emerging Technologies Based on Analysis of the US Patent Citation Network." *Scientometrics* 95 (1): 225–242.
- Fusion, J., K. Willet, and M. Zwick. 2012. *OCCAM: A Reconstructability Analysis Program*. Portland State University. <http://www.pdx.edu/sysc/research-discrete-multivariate-modeling>.
- Hall, M., E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten. 2009. "The WEKA Data Mining Software: An Update." *ACM SIGKDD Explorations Newsletter* 11 (1): 10–18.
- Hall, B. H., A. Jaffe, and M. Trajtenberg. 2005. "Market Value and Patent Citations." *The Rand Journal of Economics* 36 (1): 16–38.
- Harhoff, D., F. Narin, F. M. Scherer, and K. Vopel. 1999. "Citation Frequency and the Value of Patented Inventions." *Review of Economics and Statistics* 81 (3): 511–515.
- Hosseini, J., and M. Zwick. 1986. *Segment Congruence Analysis Via Information Theory*, G62–G77. Philadelphia, PA: International Society for General Systems Research.
- Hu, X., R. Rousseau, and J. Chen. 2012. "A New Approach for Measuring the Value of Patents Based on Structural Indicators for Ego Patent Citation Networks." *Journal of the American Society for Information Science and Technology* 63 (9): 1834–1842.
- Johnson, J. 2005. *State-based Systems Modeling: Theory, Implementation, and Application*. Portland, OR: Portland State University.
- Jun, S. 2014. "A Technology Forecasting Method Using Text Mining and Visual Apriori Algorithm." *Applied Mathematics & Information Sciences* 8 (1L): 35–40.
- Kajikawa, Y., and Y. Takeda. 2009. "Citation Network Analysis of Organic LEDs." *Technological Forecasting and Social Change* 76 (8): 1115–1123.
- Karvonen, M., and T. Kässi. 2013. "Patent Citations as a Tool for Analysing the Early Stages of Convergence." *Technological Forecasting and Social Change* 80 (6): 1094–1107.
- Kim, G., and J. Bae. 2017. "A Novel Approach to Forecast Promising Technology Through Patent Analysis." *Technological Forecasting and Social Change* 117: 228–237.
- Kim, G. J., S. S. Park, and D. S. Jang. 2015. "Technology Forecasting Using Topic-based Patent Analysis." *Journal of Scientific & Industrial Research* 74 (5): 265–270.
- Kim, H., A. Stumpf, and W. Kim. 2011. "Analysis of an Energy Efficient Building Design Through Data Mining Approach." *Automation in Construction* 20 (1): 37–43.
- Klir, G. J. 1985. *The Architecture of Systems Problem Solving*. New York: Plenum Press.
- Kousha, K., and M. Thelwall. 2017. "Patent Citation Analysis with Google." *Journal of the Association for Information Science and Technology* 68 (1): 48–61.
- Kramer, P., S. K. Westaway, M. Zwick, and S. Shervais. 2012. "Reconstructability Analysis of Genetic Loci Associated with Alzheimer disease." 6th Int. Conf. Soft Comput. Intell. Syst. 13th Int. Symp. Adv. Intell. Syst., 2104–2110.
- Lee, C., Y. Cho, H. Seol, and Y. Park. 2012. "A Stochastic Patent Citation Analysis Approach to Assessing Future Technological Impacts." *Technological Forecasting and Social Change* 79 (1): 16–29.
- Lee, C., J. Jeon, and Y. Park. 2011. "Monitoring Trends of Technological Changes Based on the Dynamic Patent Lattice: A Modified Formal Concept Analysis Approach." *Technological Forecasting and Social Change* 78 (4): 690–702.

- Lee, S., S. Lee, H. Seol, and Y. Park. 2008. "Using Patent Information for Designing new Product and Technology: Keyword Based Technology Roadmapping." *R&D Management* 38 (2): 169–188.
- Lee, S., B. Yoon, and Y. Park. 2009. "An Approach to Discovering New Technology Opportunities: Keyword-based Patent Map Approach." *Technovation* 29: 481–497.
- Li, R., T. Chambers, Y. Ding, G. Zhang, and L. Meng. 2014. "Patent Citation Analysis: Calculating Science Linkage Based on Citing Motivation." *Journal of the Association for Information Science and Technology* 65 (5): 1007–1017.
- Liu, D.-R., and M.-J. Shih. 2010. "Hybrid-patent Classification Based on Patent-network Analysis." *Journal of the American Society for Information Science and Technology* 62 (2): 246–256.
- Madani, F., and C. Weber. 2015. "The Evolution of Patent Mining: Applying Bibliometrics Analysis and Keyword Network Analysis." *World Patent Information* 46: 32–48.
- Perkowski, M., M. Marek-Sadowska, L. Jozwiak, T. Luba, S. Grygiel, M. Nowicka, R. Malvi, Z. Wang, and J. Zhang. 1997. "Decomposition of Many-valued Relations," *ISMVL*, 13–18.
- Princeton University. 2010. "About WordNet." WordNet. Princeton University. <http://wordnet.princeton.edu>.
- "Stopwords," "Ranks" Company. 2014. <http://www.ranks.nl/stopwords>.
- Willet, K., and M. Zwick. 2004. "A Software Architecture for Reconstructability Analysis." *Kybernetes* 33: 997–1008.
- Yan, R., J. Tang, X. Liu, D. Shan, and X. Li. 2011. "Citation Count Prediction: Learning to Estimate Future Citations for Literature." Proceedings of the 20th ACM International Conference on Information and Knowledge Management – CIKM '11, Glasgow, Scotland, UK, 1247.
- Yoon, B., and S. Lee. 2008. "Patent Analysis for Technology Forecasting: Sector-specific Applications." 2008 IEEE International Engineering Management Conference, Estoril, Portugal, 1–5.
- Yoon, B., and Y. Park. 2004. "A Text-mining-based Patent Network: Analytical Tool for High-technology Trend." *The Journal of High Technology Management Research* 15: 37–50.
- Yoon, B., and Y. Park. 2007. "Development of New Technology Forecasting Algorithm: Hybrid Approach for Morphology Analysis and Conjoint Analysis of Patent Information." *IEEE Transactions on Engineering Management* 54 (3): 588–599.
- Zwick, M. 2004a. "An Overview of Reconstructability Analysis." *Kybernetes* 33 (5): 877–905.
- Zwick, M. 2004b. "Reconstructability Analysis with Fourier Transforms." *Kybernetes* 33: 877–905.
- Zwick, M. 2004c. "State-based Reconstructability Analysis." *Kybernetes* 33: 1041–1052.
- Zwick, M. 2011. "Reconstructability Analysis of Epistasis." *Annals of Human Genetics* 75 (1): 157–171.

## Appendices

### Appendix 1. The extracted keywords from the patents

IV	Keyword	IV	Keyword	IV	Keyword	IV	Keyword	IV	Keyword	IV	Keyword
aa	power	bp	panel	de	threshold	eu	Condim	gj	sampl	ia	follower
ab	energ	bq	valv	dg	hous	ev	Scor	gk	condenser	ib	gener
ac	light	br	reflect	dh	coil	ew	Beam	gl	wir	ic	assemb
ad	ga	bs	switch	di	room	ex	Test	gm	model	id	wast
ae	member	bt	van	dj	dur	ey	Cool	gn	dat	ie	enclosur
af	bal	bu	cit	dk	select	ez	Condit	go	inform	if	hom
ag	golf	bv	electron	dl	timer	fa	Funct	gp	st	ig	load
ah	hydr	bw	fir	dm	door	fb	Cover	gq	ba	ih	steam
ai	finger	bx	thermostat	dn	lock	fc	Cast	gr	saving	ii	3-d
aj	lift	by	transfer	do	fuel	fd	Flow	gs	stat	ij	mot
ak	food	bz	port	dp	liquid	fe	Ha	gt	rot	ik	picture
al	cel	ca	carbur	dq	treatm	ff	Hav	gu	imag	il	cabinet
am	control	cb	commun	dr	semiconduc	fg	Invent	gv	hot	im	fig
an	il	cc	transformer	ds	structur	fh	Melt	gw	crop	in	exchang
ao	circuit	cd	inst	dt	firebox	fi	Metal	gy	exerc	io	loc
ap	wiper	ce	park	du	detect	fj	Molt	gz	provid	ip	point
aq	system	cf	wind	dv	upper	fk	NoZI	ha	interfac	iq	motor
ar	aud	cg	air	dw	air-m	fl	Par	hb	sh	ir	rf
as	dev	ch	shel	dx	posit	fm	Pat	hc	lower	is	stick
at	pump	ci	cabl	dy	mach	fn	Pool	hd	diffuser	it	light-em
au	channel	cj	sol	dz	sid	fo	Prov	he	composit	iu	guest
av	endotrach	ck	weather	ea	el	fp	Strip	hf	layer	iv	chain
aw	tub	cl	em	eb	head	fq	Substr	hg	outer	iw	user
ax	compon	cm	ar	ec	print	fr	Surface	hh	protect	ix	appl
ay	opt	cn	alarm	ed	fireplac	fs	U	hi	period	iy	dish
az	hvac	co	brush	ee	chair	ft	uniform	hj	led	iz	main
ba	pressur	cp	water	ef	sen	fu	weir	hk	acoust	ja	return
bb	electr	cq	launch	eg	temper	fv	step	hl	reson	Z*	Citation
bc	heat	cr	vacuum	eh	oper	fw	filter	hn	refriger		
bd	network	cs	ozon	ei	concentr	fx	vibr	ho	se		
be	hydraul	ct	sc	ej	displ	fy	clamp	hp	barrel		
bf	sign	cu	therm	ek	tank	fz	box	hq	ccfl		
bg	unit	cv	oil	el	greenh	ga	build	hr	bod		
bh	modl	cw	miner	em	bask	gb	lin	hs	oil-st		
bi	level	cx	subst	en	remov	gc	sewer	ht	tissu		
bj	setpoint	cy	brak	eo	row	gd	mater	hu	mirror		
bk	compr	cz	dist	ep	cut	ge	veloc	hv	clean		
bl	sect	da	fluid	eq	cl	gf	infrar	hw	fixtur		
bm	extern	db	chamber	er	roof	gg	zon	hx	receiv		
bn	wal	dc	dril	es	combust	gh	forml	hy	sub-chamber		
bo	laser	dd	spac	et	Floor	gi	siz	hz	bottom		

\*Z is the dependent variable (DV).

## Appendix 2. OCCAM results – loopless variable-based analysis

ID	Model	Level	$\Delta DF$	p-value	% $\Delta H(DV)$	$\Delta AIC$	$\Delta BIC$	%C(Data)
13	IV:AqBglzJaZ	4	95	0.00	88.92	722.10	169.46	98.5
12	IV:BaDglzJaZ	4	35	0.00	88.38	836.60	633.00	98.9
11	IV:AmDglzJaZ	4	71	0.00	88.35	764.27	351.25	98.9
10*	IV:DglzJaZ	3	11	0.00	85.75	857.58	793.59	98.6
9	IV:BglzJaZ	3	15	0.00	85.34	845.41	758.15	98.1
8*	IV:DbzJaZ	3	7	0.00	84.98	857.75	817.03	98.5
7*	IV:lzJaZ	2	3	0.00	83.33	848.76	831.31	98.1
6*	IV:DbJaZ	2	3	0.00	58.46	593.67	576.22	98.1
5*	IV:DgJaZ	2	5	0.00	57.11	575.77	546.69	98.1
4*	IV:JaZ	1	1	0.00	54.89	561.03	555.21	98.1
3*	IV:lZ	1	1	0.00	15.46	156.56	150.75	94.7
2*	IV:AcZ	1	4	0.00	3.32	26.08	2.81	94.7
1*	IV:Z	0	0	1.00	0.00	0.00	0.00	94.7

\*Indicates those models whose difference from their lower level progenitor is statistically significant

Best Model(s): by  $\Delta BIC$  is 7\*, by  $\Delta AIC$  is 8\*, and by Information, with all Inc. p-value < 0.05 is 10\*.

## Appendix 3. OCCAM results – variable-based analysis with loops

ID	MODEL	Level	$\Delta DF$	p-value	% $\Delta H(DV)$	$\Delta AIC$	$\Delta BIC$	Inc. p-value	%C(Data)
22*	IV:AmZ:AsZ:DgZ:EbZ:FdZ:lZ:JaZ	7	15	0.00	90.9	902.5	815.2	0.00	99.3
21*	IV:AmZ:AsZ:CbZ:DgZ:EbZ:lZ:JaZ	7	15	0.00	90.9	902.2	814.9	0.00	99.2
20*	IV:AmZ:AsZ:BgZ:CbZ:DgZ:lZ:JaZ	7	17	0.00	90.8	897.5	798.6	0.00	99.3
19*	IV:AmZ:AsZ:DgZ:EbZ:lZ:JaZ	6	13	0.00	89.9	895.7	820.1	0.00	99.2
18*	IV:AmZ:AsZ:DcZ:DgZ:lZ:JaZ	6	13	0.00	89.6	893.0	817.4	0.00	99.1
17*	IV:AmZ:AsZ:CbZ:DgZ:lZ:JaZ	6	14	0.00	89.6	890.8	809.4	0.00	99.1
16*	IV:AmZ:AsZ:DgZ:lZ:JaZ	5	12	0.00	88.5	883.4	813.6	0.00	99.0
15*	IV:AmZ:CbZ:DgZ:lZ:JaZ	5	11	0.00	88.2	882.4	818.4	0.00	98.9
14*	IV:AmZ:DgZ:EbZ:lZ:JaZ	5	10	0.00	88.1	883.9	825.7	0.00	99.1
13*	IV:AmZ:DgZ:lZ:JaZ	4	9	0.00	86.9	873.6	821.2	0.00	98.9
12*	IV:BgZ:DgZ:lZ:JaZ	4	7	0.00	86.7	875.3	834.6	0.00	98.7
11*	IV:DbZ:DgZ:lZ:JaZ	4	5	0.00	86.6	878.8	849.7	0.00	98.8
10*	IV:DgZ:lZ:JaZ	3	4	0.00	85.3	867.4	844.1	0.00	98.7
9*	IV:BgZ:lZ:JaZ	3	5	0.00	84.9	861.1	832.0	0.00	98.2
8*	IV:DbZ:lZ:JaZ	3	3	0.00	84.6	861.8	844.3	0.00	98.6
7*	IV:lZ:JaZ	2	2	0.00	83.0	847.1	835.4	0.00	98.1
6*	IV:DbZ:JaZ	2	2	0.00	58.4	595.3	583.7	0.00	98.1
5*	IV:DgZ:JaZ	2	3	0.00	57.1	579.6	562.1	0.00	98.1
4*	IV:JaZ	1	1	0.00	54.9	561.0	555.2	0.00	98.1
3*	IV:lZ	1	1	0.00	15.5	156.6	150.7	0.00	94.7
2*	IV:AcZ	1	4	0.00	3.3	26.1	2.8	0.00	94.7
1*	IV:Z	0	0	1.00	0.0	0.0	0.0	0.00	94.7

\*Indicates those models whose difference from their lower level progenitor is statistically significant.

Best Model(s): by  $\Delta BIC$  is 11\*, by both  $\Delta AIC$  and Incremental-p is 22\*.

#### Appendix 4. OCCAM results – variable-based analysis with loops: the distributions corresponding to the four separate components of IV:DbZ:DgZ:Iz:JaZ Model

IV	Data			rule	#correct	%correct
	freq	obs. p(DV IV)				
		Z = 1	Z = 2			
Db						
1	2316	95.50	4.49	1	2212	95.50
2	167	83.83	16.16	1	140	83.83
	2483	94.72	5.27	1	2352	94.72

IV	Data			rule	#correct	%correct
	freq	obs. p(DV IV)				
		Z = 1	Z = 2			
Dg						
1	1806	94.96	5.03	1	1715	94.96
2	410	97.31	2.68	1	399	97.31
3	267	89.13	10.86	1	238	89.13
	2483	94.72	5.27	1	2352	94.72

IV	Data			rule	#correct	%correct
	freq	obs. p(DV IV)				
		Z = 1	Z = 2			
Iz						
1	2346	96.675	3.325	1	2268	96.675
2	137	61.314	38.686	1	84	61.314
	2483	94.724	5.276	1	2352	94.724

IV	Data			rule	#correct	%correct
	freq	obs. p(DV IV)				
		Z = 1	Z = 2			
Ja	2399	98.04	1.95	1	2352	98.04
2	84	0	100	2	84	100
	2483	94.72	5.27	1	2436	98.10



### Appendix 5. OCCAM results – state-based analysis for Db, Dg, Iz, and Ja variables

ID	MODEL	Level	ΔDF	P-value	%ΔH(DV)	ΔAIC	ΔBIC	Inc. p-value	Prog.	%C(Data)
16	IV:Db1Dg1Iz1Ja2Z:Db1Z:Dg3Iz2Ja1Z:Iz1Ja1Z:Ja2Z:Z	5	5	0	87.25	885.01	855.93	0.91	13	98.8
15	IV:Db1Dg1Iz1Ja2Z:Db1Z:Dg3Iz2Z:Iz1Ja1Z:Ja2Z:Z	5	5	0	87.25	885.01	855.92	0.89	12	98.8
14	IV:Db1Dg1Iz1Ja2Z:Db1Z:Dg3Ja1Z:Iz1Ja1Z:Ja2Z:Z	5	5	0	87.23	884.83	855.74	0.93	11	98.8
13*	IV:Db1Z:Dg3Iz2Ja1Z:Iz1Ja1Z:Ja2Z:Z	4	4	0	87.25	887.03	863.76	0.00	9	98.8
12*	IV:Db1Z:Dg3Iz2Z:Iz1Ja1Z:Ja2Z:Z	4	4	0	87.25	887.02	863.76	0.00	10	98.8
11*	IV:Db1Z:Dg3Ja1Z:Iz1Ja1Z:Ja2Z:Z	4	4	0	87.23	886.83	863.56	0.00	8	98.8
10*	IV:Dg3Iz2Z:Iz1Ja1Z:Ja2Z:Z	3	3	0	85.82	874.32	856.87	0	7	98.6
9*	IV:Dg3Iz2Ja1Z:Iz1Ja1Z:Ja2Z:Z	3	3	0	85.82	874.32	856.87	0	7	98.6
8*	IV:Dg3Ja1Z:Iz1Ja1Z:Ja2Z:Z	3	3	0	85.76	873.70	856.25	0	7	98.6
7*	IV:Iz1Ja1Z:Ja2Z:Z	2	2	0	83.42	851.71	840.08	0	2	98.1
6*	IV:Db1Iz2Ja1Z:Iz1Ja1Z:Z	2	2	0	82.80	845.30	833.67	0	4	98.5
5*	IV:Iz1Z:Ja2Z:Z	2	2	0	82.79	845.27	833.64	0	2	98.1
4*	IV:Iz1Ja1Z:Z	1	1	0	71.83	734.83	729.01	0	1	96.6
3*	IV:Db1Iz1Ja1Z:Z	1	1	0	55.03	562.46	556.64	0	1	94.7
2*	IV:Ja2Z:Z	1	1	0	54.88	560.92	555.10	0	1	98.1
1*	IV:Z	0	0	1	0	0	855.93	0	0	94.7

\*Indicates those models whose difference from their lower level progenitor (Prog.) is statistically significant  
 Best Model(s): by ΔBIC is 13\*, by ΔAIC is 13\*, and by Information, with all Inc. p-value < .05 are 13\*, 12\*, and 11\*.