

# Measuring the Information Content of Architectural Plans

John S. Gero and Julie R. Jupp

Email: {john, jupp\_j} @arch.usyd.edu.au

Key Centre of Design Computing and Cognition

University of Sydney NSW 2006 Australia

**Summary:** (English) This paper describes and develops a preliminary approach to the measurement of the information content of two-dimensional design drawings. We utilise a general method for extracting information from an encoded string of symbols as a canonical representation of architectural plans. The information content of each drawing or set of drawings is determined by measuring its entropy. We present two classes of qualitative representation of shape and space. The first uses a qualitative representation of the outline of shapes in the drawing. The second uses a qualitative representation of the spaces described in the drawing. We describe the preliminary implementation of the method to a time-evolution of two formally described design styles, Romanesque and Gothic cathedral plans.

**Summary:** (Portuguese) Este documento descreve e desenvolve uma abordagem preliminar de medida do conteúdo informativo de desenhos bidimensionais do projeto. Nos utilizamos um método geral para extração da informação derivado de uma série de símbolos codificados como uma canônica representação do plano arquitetônico. O conteúdo informativo de cada desenho, ou série de desenhos, é determinado pela medida de sua entropia. Nos apresentamos duas classes de representação qualitativa do espaço e da forma. A primeira usa uma representação qualitativa de esboços das formas no desenho. A segunda usa uma representação qualitativa dos espaços descritos no desenho. Nos descrevemos a implementação preliminar do método em relação a período-evolução de dois estilos formalmente descritos, projetos de catedrais do Romanesco e do Gótico.

## 1. Introduction

This paper develops an approach to a formal computational measure of design drawings. These ideas are explored within a representation of 2D architectural plan drawings. In an earlier paper Gero and Park (1997), presented an approach to representing shape features using qualitative reasoning. These descriptions represented the outline of two-dimensional drawings. Those ideas were extended directly into an information theoretic framework for a measure of the complexity and similarity of architectural design drawings in Gero and Kazakov (2001). This choice of representational formalism of shape and associated ontology was made because of its ability to map onto feature space that in turn intuitively relates to concepts of shape complexity and similarity. This implementation presented a comparison of the design corpus of two architects, Alvar Aalto and Louis Kahn using partial representation of architectural drawings (only the outline of the plan). We describe an extended schema for the purpose of calculating the information content of a more complete plan.

## 2. Symbolic Representation of Design Drawings

Gero and Park (1997) developed an efficient description of shapes using symbols enabling a representation of qualities rather than a numerically based description of quantities within an encoding schema known as Q-codes. This simple and effective process of symbolic mapping is useful in modelling design variables and attributes, and transforming possible numeric value ranges into small sets of discrete and finite symbols (Mantyla; 1988, Gero and Park 1997). This type of qualitatively representation of shape features as a symbol provides a better framework for computer-aided tools for human spatial reasoning (Engenolfer and Shariff 1998) as it deals with classes of shape features rather than simply the instance of a shape. Space, as well as its organization, is an underlying category of human cognition. Its ability to structure activities and relationships with the external world and many of our reasoning capabilities is fundamental to a formal system of representation. The qualitative modelling of two-dimensional architectural plans should therefore consider space as well as shape as its two fundamental primitives.

### 2.2 Two Class Encoding Schema

We extend the Q-codes schema to include symbolic representations of information derived from spatial characteristics with relational values by introducing new symbols while maintaining the feature-based approach and the equivalent analogy with language. The extension provides representations of the organization of spaces. This offers an extended schema for the computation of an architectural plan's information content in order to link the modelling of shape and space with measurement.

#### 2.2.1 Q-code Schema

Gero and Park's schema uses a set of landmarks placed on the outline of an architectural drawing. Landmarks are defined as points that are considered as distinguished by the coding procedure and the direction of the outline has a

discontinuity on a coarse “macroscopic” scale. The general characteristics of these strings are symbol(s) plus their sign value(s) and maintain a counter clockwise direction of scanning that can commence at any landmark point. Within the symbolic representation qualitative values in the number range of values are described as a ‘landmark’ set with values in the range  $\{-, 0 \}$ . The set of intervals are then described as  $\{-, 0\}$ ;  $[0,0]$ ;  $(0, \_)$  which correspond to the qualitative set Q with the sign values  $\{+, 0, -\}$  (Weither, 1994, Gero and Park, 1997). The four types of shape attributes developed within the original schema are as follows:

- (i) Angle measured at the node, A-code: Each node is coded as  $\{A+\}$ ,  $\{A0\}$  or  $\{A-\}$  correspondingly, depending on which of the intervals the angle between two segments adjacent to this landmark belongs to;
- (ii) Relative length of line segments, L-code: Each segment is coded as  $\{L+\}$  if its length is longer than the length of the previous segment,  $\{L0\}$  if they have the same length and  $\{L-\}$  if it is shorter than the previous;
- (iii) Angle measured at a node for two tangents, C-code: This code is a generalization of the A-code to curvilinear line segments at a node with values of  $\{C-\}$ ,  $\{C0\}$  or  $\{C+\}$ .
- (iv) Relative curvature of a line segment, K-code: describes the curvature of a curvilinear line segment with values of  $\{K-\}$  “convex”,  $\{K0\}$  “straight” and  $\{K+\}$  “concave”.

A sequence of Q-codes forms a word to represent a shape pattern of significance such as an architectural plan. The Q-codes are directly employed by Gero and Kazakov (2001), and used to calculate the complexity and to compare the similarity of architectural drawings. In this way the problem has been reduced to estimating complexity and similarity for symbol strings and corpuses made of symbol strings. Since cognitively humans recognise and identify space not only through complex forms, (by registering their characteristic features as in the existing schema), but also their configurations (Treisman and Gelade, 1980), the existing representational schema requires further description. In order to utilise this schema completely within an information theoretic model a need exists to extend the description of landmark points beyond a simple representation of those that exist on the boundary of the drawing towards one that is capable of describing internal nodes and the spaces they define.

### 2.3 Extended Q-code Schema

Additional to the four existing codes a fifth symbol and value is added to the representation of internal nodes within the design drawing. This extends the properties of shape features in two-dimensions and is described as follows:

Relative area of spatial region, R-code; Each area is coded as  $\{R+\}$  if it is larger than the area of the previous region,  $\{R0\}$  if they have the same area, and  $\{R-\}$  if it smaller than the previous one. A 2D drawing containing spaces defined as regions has relational values that are influenced by a nodes’ location, regions that lies on the outline (u) of a plan identify a special case. In the instance of nodes occurring on the boundary a special condition in the R-code schema creates a fourth value where  $R_u$  defines this relation. The new code is established by the arcs of the graph of the spatial topology for each space (or “walls” separating the spaces) and they collapse onto the adjacent node; measured in a counter clockwise direction, Figure 1.

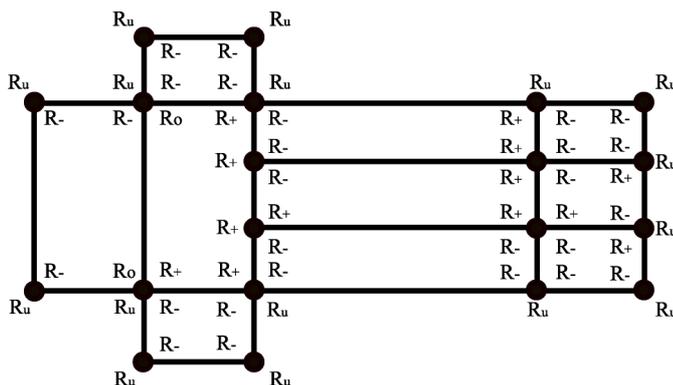


Figure 1.

Example plan encoded with R-codes that describes the comparisons of the area of the two adjacent areas ( $a_i$ ) with the values of “bigger”  $\{R+\}$ , “equal”  $\{R0\}$ , “smaller”  $\{R-\}$  and “unbounded”  $\{R_u\}$

The new code is defined by the following symbols and values:

$$R_- \equiv a_{i-1} < a_i$$

$$R_+ \equiv a_{i-1} > a_i$$

$$R_0 \equiv a_{i-1} = a_i$$

$$R_u \equiv a_{i-1} \gg a_i \quad \text{where: } a \text{ equals the area.}$$

The extended Q-code string describes in more detail the physicality of shapes and spaces as a sequence of symbols that is assumed to denote its pictorial characteristics. The resulting symbol string is understood as a circular structure, i.e., the last symbol is followed by the first.

## 2.4 Additional Schema

The second representational class builds upon the inclusion of regions within the first class to establish relations within the drawing that describe information concerning the connectivity of spaces. This method commences with the standard method of graph theoretic representation of spaces. An initial graph is generated by the connection of spaces to describe information about their topological relationships. This is the standard dual of the graph describing the physical shape of the drawing but with the addition of the spaces making it a semantic graph. The spatial morphology of the architectural plan can then be articulated in relation to constraints placed upon each graph.

**Definition 1 (Space)** Let a two-dimensional space be the symbol  $\mathbf{S}$  and be defined by the position of the space relative to the boundary according to the following:

- $S^e$  is the external or unbounded space and is unique space,
- $S_i$  is an internal or bounded space, where “i” denotes the space

**Definition 2 (Vertex)** Let  $v$  be the symbol  $v_1, v_2, \dots, v_n$ , and denote a vertex of a spatial graph according to the following:

- $v_i^e$  denotes a vertex that occurs at a boundary, ie it occurs at the boundary of at least one internal space and the external space and
- $v_i$  denotes a vertex that occurs inside the plan more generally  $v_i$  refers to the vertex at the boundary of i+1 spaces

**Definition 3 (Dual)** Let  $d$  be the symbol  $d_1, d_2, \dots, d_n$ , and denote the vertices of the dual of the initial spatial graph according to the following:

- $d_i^e$  denotes vertices of the type  $v_1, \dots, v_{i-1}$  and  $v_{i+1}, \dots, v_n$
- $d_i$  denotes one or more vertices of the type  $v_i, \dots, v_{i-1}$

**Definition 4 (External Dual Node)** Let  $ds^e$  be the instance of the external or unbounded dual that defines the unique space of the dual’s graph.

**Definition 5 (Dual of Dual)** Let  $dd$  be the symbol  $dd_1, dd_2, \dots, dd_n$ , and denote the vertices of the dual of the dual graph according to the following:

- $dd_i$  denotes one or more duals of the type  $d_n, \dots, d_{i-1}$

The nature of spatial connectivity within the initial graph is represented at a *vertex* ( $v$ ), within the first abstraction of the graph at the dual, and within the second at the dual of the dual. Using these definitions these examples follow for vertex, dual and the dual of dual.

This abstraction enables the explicit representation of information relating to various levels of spatial connectivity, previously implicit within the two-dimensional drawing. This reasons over class spatial knowledge to describe information contained within architectural drawings based on deduction. The order of abstraction is illustrated in Figure 2. The qualitative representation of spatial connectivity is an auxiliary class to the parent Q-code language.

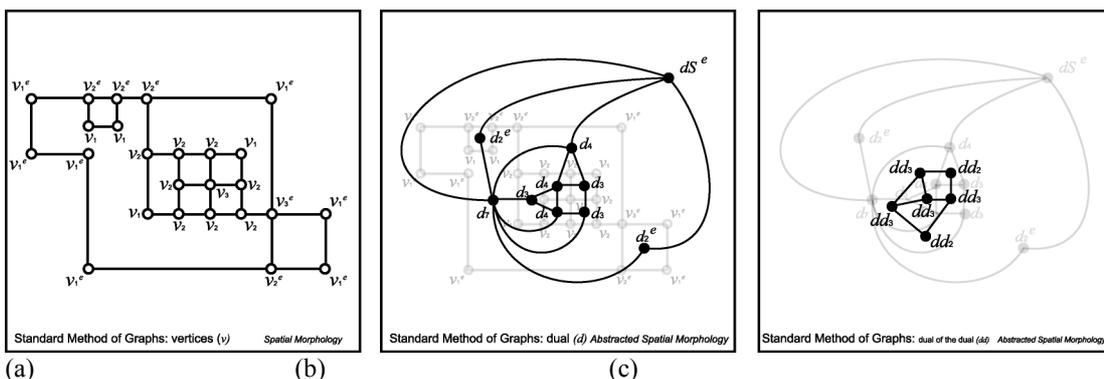


Figure 2.

- a) depicts representations of vertices, types  $v_1, v_1^e, \dots, v_n, v_n^e$ . This first case represents spatial configurations already explicit within the drawing. b) depicts representations of the dual, types  $d_1, d_1^e, \dots, d_n, d_n^e$ . The second case makes explicit

the implicit connectivity of in the original graph representation, by describing properties of the spatial relations and identifying constraints on the vertex. This assists in the reasoning of local spaces as not all the spatial information can be made explicit within the initial graph of spaces. c) depicts the dual of the dual, types  $dd_1, dd_2, \dots \dots dd_n$ . The third case makes explicit the implicit connectivity of in the dual graph representation, by describing properties of the spatial relations of the dual.

### 3. Applying Information Theory the Extended Q-codes

Once these drawings are encoded in this canonical form we can translate the description into a design semantic to discover comparisons by calculating the design drawing's information content. Approaching this from classic information theory requires the calculation of the entropy of the string of codes, as a measure of its complexity. The approach to the measurement of these symbol strings' information content is based on a data-compression technique.

#### 3.1 Similarity and Complexity Measurements for Drawings

A measure of information content has previously been studied within linguistics using data-compression techniques. Benedetto, Caglioti and Loreto, (2002) have utilised this method. For an encoded design drawing the entropy of a string of characters is defined here as the length (in bits) of the smallest code, which produces as output the string. A zipper algorithm takes the file of the design drawing's encoding and transforms it into the shortest possible file. This is not the best way to encode the file but is an excellent approximation of it that enables this investigation to measure a large body of design drawings. Using a common compression algorithm, the Lempel and Ziv algorithm (LZ77) (Lempel and Ziv, 1977), we can compute for individual shapes and their organization of space those that belong to different groups. This has the potential to show how individual drawings can be used to track changes over time, providing a powerful tool to measure the entropy of a sequence by zipping it. There does however exist several ways to measure the relative entropy and a deeper understanding of these definitions can be seen within the work of Wyner, (1995).

#### 3.2 Formal Measures of Design Drawings

The overall application of this method is aimed at investigating the automatic recognition of a design's 'style' in which a given design drawing, represented as a symbol string can be classified via its relative entropy. An example of this encoding procedure for a cathedral plan is illustrated in Figure 3. Our hypothesis maintains that for any single string the method will recognise the design within a group of other designs using this method of encoding.

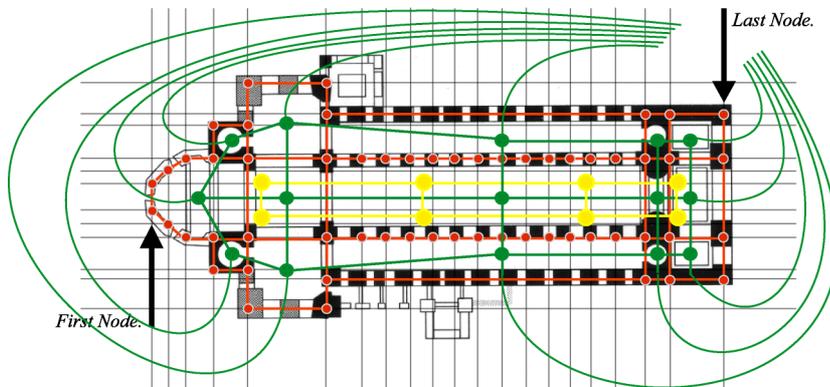


Figure 3.

Two-class representation of a Romanesque cathedral plan (Speyer, 1030AD). The encoding is carried out on the basis of numbering nodes as a grid laid over the plan in matrix form, beginning with the corner to which the bold arrow is pointed.  $\{L_0, A_+, R_0, L_0, A_-, R_-, L_0, A_+, R_0, L_0, A_-, R_-, \dots \dots L_+, A_+, R_0, L_+, A_-, R_-\}$ , and  $\{(v_1^e, v_1^e, v_1^e, v_1^e, v_1^e, v_1^e, v_1^e, v_2^e, v_2^e, v_1^e, v_1^e, v_2^e, v_3^e, \dots \dots d_4, d_4^e, d_3^e, d_4^e, d_3^e) > [dd_2, dd_2, dd_3, dd_3, dd_3, dd_3, dd_2, dd_2]\}$ .

##### 3.2.1 Measuring Time Evolution of Designs

The initial application demonstrated here illustrates the time evolution of Romanesque to Gothic architecture as exemplified by their floor plans. We have chosen for our tests Romanesque and Gothic cathedrals due to their distinct chronology, the amount of historical documentation and the large corpus of plans that enables an evaluation of our method.

The first Romanesque cathedrals were designed around late 700 AD in Germany. The Romanesque period lasted up until 1200AD and cathedrals of this era can be mostly found in Germany, France and Italy. The first Gothic cathedral was constructed around 1190 AD and was built in Germany, with the 'style' spreading throughout France, Italy and the Netherlands by the late 1200's. There are clearly two distinct time periods of Romanesque and Gothic design, as well as a division within the Gothic era, which occurred around 1350AD and is described historically in terms of Early and Late periods of Gothic design. The overall sample for this demonstration includes 8 plans of Romanesque cathedrals and an additional 8 Gothic cathedrals. The measure obtained is essentially the number of

cumulative distinct Q-words in the Q-code shape and space description. The idea being that the most complex drawings are those ones whose description cannot be compressed. At the least we expect to find that the proposed technique is able to differentiate between Romanesque and Gothic cathedral plans. This is based on the assumption that there have been changes in the complexity of plan drawings over the transition of Romanesque and Gothic cathedral design. Each plan was encoded as two separate sentences, using both schemas and labelled with the year the building was completed. Each Q-code sentence includes three parts based on the relative angle, length and region for the first sentence. For the additional schema each encoding also includes three parts based on the original graph,  $(v)$ , the abstract graph  $\langle d \rangle$ , and the second abstraction  $[dd]$ .

The method of application for this compression algorithm is simple and can be understood if we consider  $A$  and  $B$  as two floor plans each representing a design period, in this instance we have chosen Romanesque and Gothic architectural designs. We take a long Romanesque string  $A$  and we append to it a short Gothic string,  $b$ . The zipper begins reading the file starting from the Romanesque string, after time it is able to encode optimally the Romanesque file. When the Gothic string begins, the zipper starts encoding it in a way, which is optimal for the Romanesque. i.e., it finds most of the matches in the Romanesque string. The first part of the Gothic file is encoded with the Romanesque code. After a while the zipper “learns” Gothic, i.e., it tends progressively to find most of the matches in the Gothic string with respect to the Romanesque one, and changes its rules. The relative entropy  $E_A$  per character between  $A$  and  $B$  will be estimated by:

$$E_{AB} = (\Delta_{Ab} - \Delta_{Bb}) / |b|, \quad (1)$$

where:  $|b|$  is the number of characters of the sequence  $b$  and  $\Delta_{Bb} |b| = (l_{B+b} - l_B |b|)$  is an estimate of the entropy of the source B.

Therefore if the length of the Gothic file is “small enough”, i.e., if most of the matches occur in the Romanesque string, the calculation given from (1) returns its relative entropy.

### 3.2.3 Results from Comparing Romanesque and Gothic Cathedrals

The results of these experiments on time evolution are outlined in Figures 4(a) and 4(b). Both graphs refer to the similarity between each cathedral plan for each class of encoding. From each graph we can see that the complexities vary significantly between the two eras and confirms the discriminatory power of the approach. Examining the results for Romanesque plans in Figure 4(a) we can see that the sample of plans produced between 789 and 1200 return similar values. The significant increase in complexity at 1200AD (first cathedral of the Gothic era) reflects the transition that occurred during this time between the two different design periods. The representation of the cathedrals’ spatial morphology is presented in Figure 4(b). Examining this graph we can see there is a smaller difference in the complexity of these symbol strings for each era, confirming that there is similar spatial morphologies belonging to each design period. A comparable increase in the complexity values occurs at the same time as in Figure 4(a), indicating that the spatial morphology within Gothic cathedrals became more complex. Overall these results imply that the complexity increased over the two periods of cathedral design both in terms of the classes of shape features and their spatial organization.

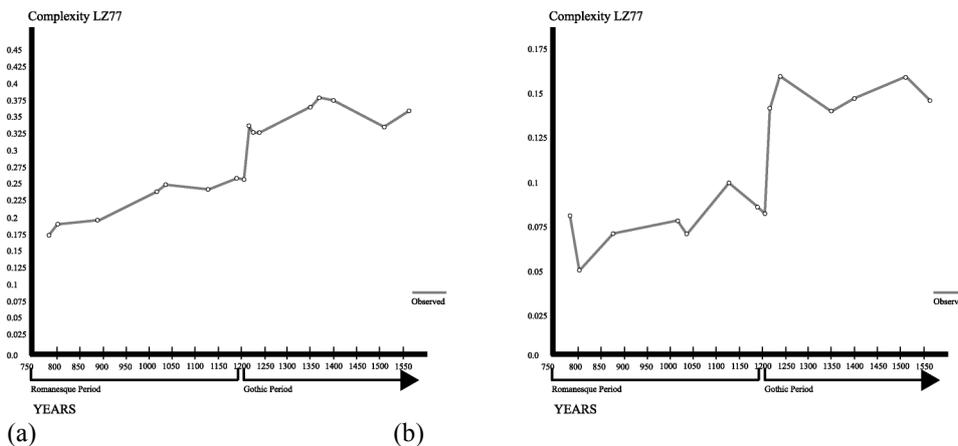


Figure 4

Preliminary results of complexity over time using compression algorithm LZ77, for Romanesque and Gothic cathedrals: 798AD to 1568 AD; (a) Shape Features (extended schema) and (b) Spatial Morphology (additional schema)

## 4. Discussion

This research is ongoing, however, these preliminary results are convincing with respect to the variations in the complexity values for both the original strings of Romanesque and Gothic drawing plans and their appended portions. These findings feature some interesting results for a model of time evolution and offer a framework for its

future application within design recognition and classification models. These results also provide the potential of a meaningful design semantic from this translated description.

This translation of design drawings into a qualitative representation plays an important role in this method of measurement. The development of this two-part encoding schema affords an information theoretic framework of measurement a more robust representation. The new schema demonstrates the descriptive power of symbol strings as a type of language for specialized qualitative sequences. The information content of architectural drawings presented within this encoding schema provides the basis for measuring complexity for comparing and categorizing architectural plans automatically and lays the foundation for the development of digital representations of style. If we can digitally represent style we will be able to formally determine the derivational structures of various architectural styles.

### Acknowledgements

The research described here has been supported by a grant from the Australian Research Council and an Australian Postgraduate Award. Computing resources have been provided by the Key Centre of Design Computing and Cognition.

### References

- Benedetto, D, Caglioti, E and Loreto, V: 2001, Language trees and zipping, La Sapienza, University of Rome, Italy. [http://arxiv.org/PS\\_cache/cond-mat/pdf/0108/0108530.pdf](http://arxiv.org/PS_cache/cond-mat/pdf/0108/0108530.pdf)
- Egenhofer, MJ and Shariff, A R: 1998, Metric details for natural-language spatial relations, *ACM Transactions on Information Systems* 16, (4), 295-321.
- Gero, JS and Damski, J: 1999, Feature-based qualitative modelling of objects, in G Augenbroe and C Eastman (eds.), *Computers in Building*, Kluwer, Boston, pp. 309-320.
- Gero, JS and Kazakov, V: 2001, Entropic similarity and complexity measures for architectural drawings, in JS Gero B Tversky and T Purcell (eds.), *Visual and Spatial Reasoning in Design II*, Key Centre of Design Computing and Cognition, University of Sydney, Sydney.
- Gero, JS and Park, S-H: 1997, Computable feature-based qualitative modelling of shape and space, in R Junge (ed.), *CAAD Futures 1997*, Kluwer, Dordrecht, pp. 821-830.
- Ziv, J and Lempel, A: 1978, Compression of individual sequences via variable-rate coding, *IEEE Transactions on Information Theory* 24: 530-536.
- Mantyla, M: 1988, *An introduction to solid modelling*, Computer Science Press, Rockville.
- Shah, J: 1991, Assessment of feature technology, *CAD* 23, (5) 331-343.
- Treisman AM and Gelade, G: 1980, A feature-integration theory of attention, *Cognitive Psychology* 14: 97-136.
- Wyner, AD: 1995, Typical sequences and all that: Entropy pattern matching and data compression, *IEEE Information Theory Society Newsletter*.

This is a copy of the paper: Gero, JS and Jupp, J (2002) Measuring the information content of architectural plans, in PL Hippolyte and E Miralles (eds), *SIGraDi Caracas 2002*, Ediciones Universidad Central de Venezuela, Caracas, pp. 155-158.