

## Neural Network Based Chemical Structure Indexing

S. D. D. V. Rughooputh\* and H. C. S. Rughooputh

University of Mauritius, Reduit, Mauritius

Received July 17, 2000

Searches on chemical databases are presently dominated by the text-based content of a paper which can be indexed into a keyword searchable form. Such traditional searches can prove to be very time-consuming and discouraging to the less frequent scientist. We report a simple chemical indexing based on the molecular structure alone. The method used is based on a one-to-one correspondence between the chemical structure presented as an image to a neural network and the corresponding binary output. The method is direct and less cumbersome (compared with traditional methods) and proves to be robust, elegant, and very versatile.

### 1. INTRODUCTION

There exist several ways to search for chemical information from the Internet based World-Wide-Web system using a web browser such as Netscape or Microsoft Explorer or Mosaic. Examples of WWW-based chemical search server/engines include Chemical Abstracts Services, ChemExper Chemical Directory, ChemFinder WebServer (CambridgeSoft), NIST database, ChemIDplus (Specialized Information Services), Hazardous Substances Databank Structures (HSDB) (Specialized Information Services), NCI-3D (Specialized Information Services), and general-purpose databases of WWW contents (such as Yahoo and Alta Vista). WebServer capacities range from several thousands to several millions with databases varying from general to specialized chemicals such as liquid crystals, pesticides, polycyclic aromatic hydrocarbon, drugs, environmental pollutants, potential toxins, etc. Most of these databases are freely accessible to researchers from academic and industrial laboratories.

One can search various electronic databases for chemicals by their chemical names (some accept wildcards and/or typographic variations in names), CA index, IUPAC names, common names, trade names or synonyms, molecular formula, molecular weight, Chemical Abstracts Service (CAS) Registry Numbers (CASRNs are unique identifiers for chemical compounds with standard format being xxxxxx-xx-x), catalog number, chemical characteristics, 2D chemical structures and substructures, and molecular descriptors. The databases will identify the type of search requested and provide the hits accordingly. Today's chemical databases are more versatile with faster processing and can also correct for obvious errors as well as invalid CAS RNs.

Although sparse for the time being, some of the chemical databases also provide additional information such as 2D/3D chemical structures (as Windows metafiles or molfiles) and useful references. Helper applications or Viewers are normally needed to display chemical structural records of the compounds. Web Browsers cannot read these without a helper application and the appropriate plug-ins. To display a structure, a structure-drawing program or WWW viewer

must be used. Examples of software for chemical structures/viewers include ChemDraw (xxx), Chem3D, ChemOffice and ChemOffice Pro from CambridgeSoft, ISIS/Draw (MDL Information Systems, Inc.), Wetlab (Molecular Simulations, Inc.), ChemWeb (Softshell International, Ltd.), Accord Internet Viewer (Synopsis Scientific Systems), and Rasmol viewer. It should be noted that most of the databases are not user-friendly requiring hours of training. Also, one normally finds that not all structures in the database currently have chemical formulas/molecular weights assigned.

Quite often one would like to identify a particular compound or a related compound (simple or complex) directly from its chemical structure alone without a priori knowledge of the CAS number, molecular formula, chemical functionality details, and so on. Recognition of chemical structures can normally be a slow process requiring (in most cases) electronic submissions to the server.

Recent studies of the visual cortex of the cat highlight the role of temporal processing using synchronous oscillations for object identification.<sup>1,2</sup> A few oscillatory neural network models have been developed using related strategies for image processing. However, most of these approaches either use linear or phase-oscillator models and as such ignore dendritic processing. In this paper, the original neural network model of Eckhorn<sup>1</sup> is modified according to the proposal of Johnson<sup>2</sup> for recognition of molecular structures. Chemical structures are presented as standard drawings to the neural network. The output of the network will be binary barcode-like features with a 1:1 correspondence with the<sup>1–3</sup> corresponding input images. We demonstrate the robustness of the recognition process of such structures using some 100 chemical compounds with diverse complexities. Assuming a standard entry format, our method proves to be simple and robust for compound recognition. Because of the 1:1 correspondence between the input image and the generated binary barcode, the latter can be used as a replacement for the CAS RN.

### 2. PCNN ALGORITHM AND IMPLEMENTATION

The PCNN,<sup>3</sup> Figure 1, has received much attention from the image processing community. The PCNN algorithm is an iterative algorithm whereby a binary output image is

\* Corresponding author phone: (230)454-1481; fax: (230)465-6928; e-mail: sdr@uom.ac.mu.

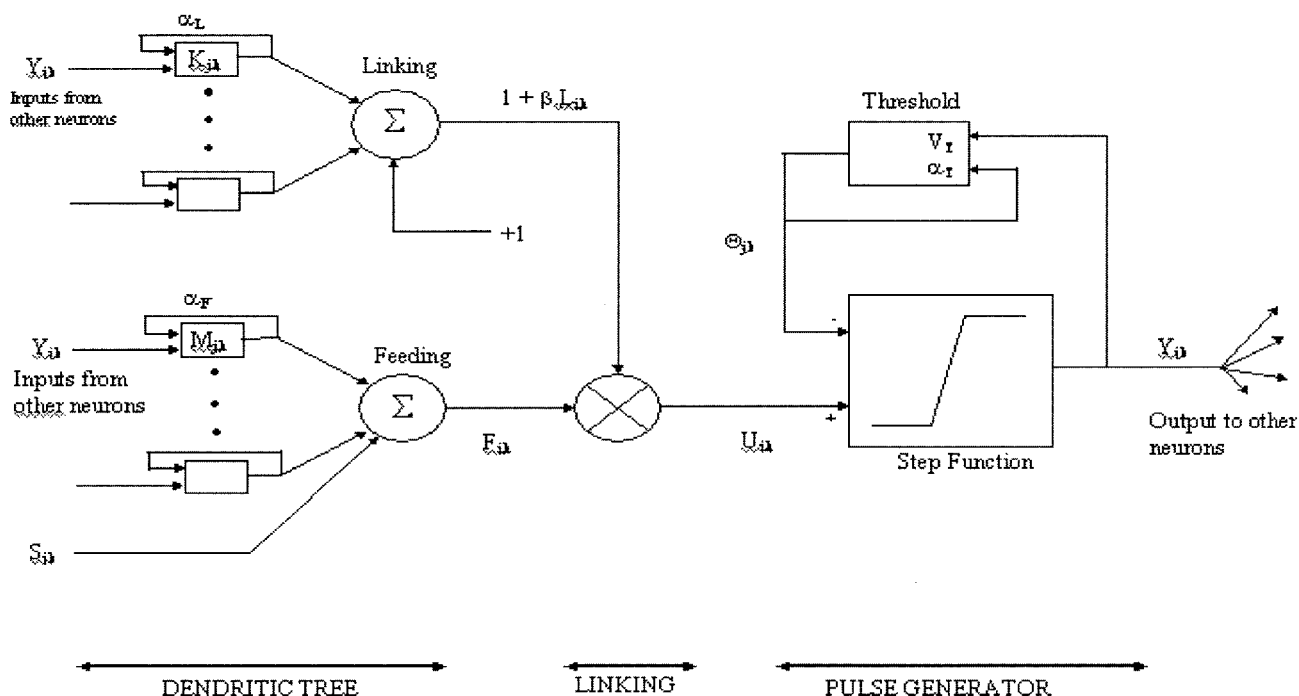


Figure 1. PCNN neuronal model.

produced for each iteration cycle. The output image at each iteration typically represents some segments or edges of the input image.

The network is comprised of two-dimensional integrate-and-fire neurons with one neuron for each input pixel. Each neuron receives input signals from a feeding synapse and a linking synapse. The PCNN neuronal model consists of four basic elements: the feeding element  $F$ , the linking element  $L$ , the internal activation element  $U$ , the output element  $O$  and  $Y$  is the output. The input unit (dendritic tree) includes the image pixel input (real value) as well as feedback (binary) inputs from surrounding neurons. The linking and feeding branches merge together to compute the internal activation  $U$ . The integrated signals from the linking synapse plus an offset term of "1" are multiplied with the integrating signals from the feeding synapse to produce a membrane voltage  $U$ . The pulse generator compares a threshold value  $\Theta$  with  $U$ . If  $U > \Theta$ , then the PCNN output becomes equal to 1 and is reset to a maximum value  $V_\Theta$ , else the output is 0 and  $\Theta$  is decreased exponentially by a time constant  $\alpha_\Theta$ .

The input stimulus  $S$  (the pixel intensity) is received by the feeding element, and the internal activation element combines the feeding element with the linking element. The value of internal activation element is compared to a dynamic threshold which gradually decrease at iteration. The internal activation element accumulates the signals until it surpasses the dynamic threshold and then fires the output element and the dynamic threshold increases simultaneously strongly. The output  $Y$  of the output neuron is then iteratively fed back to the element with a delay of one iteration. The interconnections  $M$  and  $W$  refer to the Gaussian weight functions with the distance as the argument. There are three potentials  $V$  and decay constants associated with  $F$ ,  $L$ , and dynamic threshold  $T$ .

If a digital image is applied as input to such a network, the network will group image pixels based on spatial proximity and brightness similarity. Segment extraction

occurs since groups of neurons in a similar state tend to pulse in unison. A two-dimensional texture image can be mapped into a one-dimensional output function for e.g. time-series (iteration number). Spatial averaging of neural responses gives a temporal output signal that is shown to be valid as a coding mechanism. The neuronal dynamics can be implemented by iterating the following equations:<sup>3</sup>

$$F_{ij}[n] = \exp(-\alpha_F)F_{ij}[n-1] + S_{ij} + V_F \sum_{kl} M_{ijkl} Y_{kl}[n-1] \quad (1)$$

$$L_{ij}[n] = \exp(-\alpha_L)L_{ij}[n-1] + V_L \sum_{kl} W_{ijkl} Y_{kl}[n-1] \quad (2)$$

$$U_{ij}[n] = F_{ij}[n] (1 + \beta L_{ij}[n]) \quad (3)$$

$$Y_{ij}[n] = \begin{cases} 1 & \text{if } U_{ij}[n] > \Theta_{ij}[n-1] \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

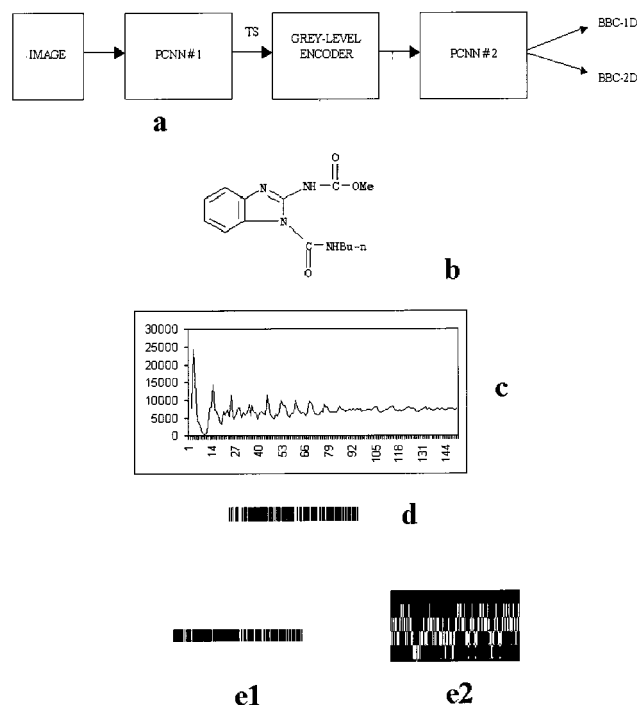
$$\Theta_{ij}[n] = \exp(-\alpha_\Theta)\Theta_{ij}[n-1] + V_\Theta Y_{ij}[n-1] \quad (5)$$

The time signal  $G[n]$ , as computed by eq 6, is the number of "on pixels" in each iteration.

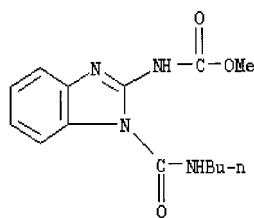
$$G[n] = \sum_{ij} Y_{ij}[n] \quad (6)$$

This time signal has been shown to be invariant with regard to changes in rotation, scale, shift, or skew of the input image.<sup>3</sup> Furthermore, it has been shown<sup>3</sup> that there is a 1:1 correspondence between images and the time signatures (icons). It is opportune therefore to exploit these properties for the purpose of recognition of chemical structures.

In the present work, we use two PCNNs to preprocess these images for our purpose as shown in Figure 2a. Time signatures of chemical structures are obtained by presenting the chemical drawings (Figure 2b) to the first PCNN (PCNN#1). As we are dealing with images, the chemical structures are drawn following a standard format. This



**Figure 2.** a. Block diagram of barcoding technique; b. image of chemical structure; c. time signature of b; d. grey level barcode of c; e1. 1D binary barcode; and e2. 2D binary barcode.



**Figure 3.** Chemical structure to be identified.

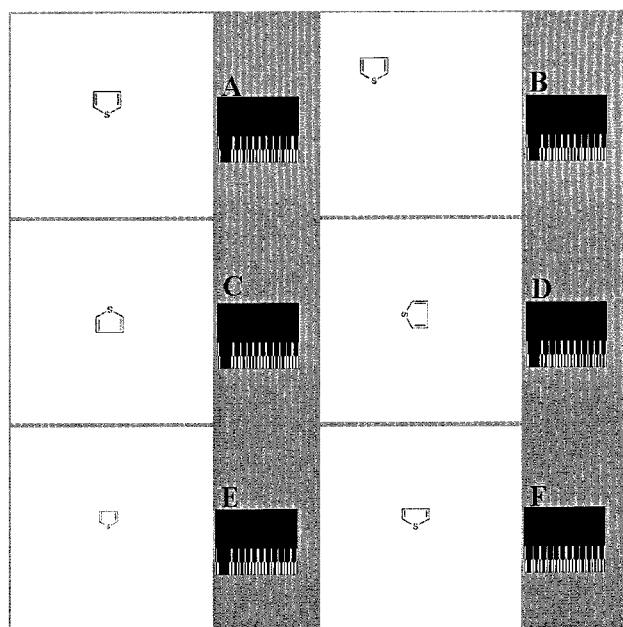
implies that all the drawings must satisfy a set of simple rules, viz. specifications for the representation of bonds (line length and thickness), atomic/molecular symbols (font, size, thickness), and ring dimensions. This requirement for drawing chemical structures can be met using a simple software (far less sophisticated than existing ones). The time signature, for example in Figure 2c, as such is not convenient for computer use. We further convert the time signal to an image by using a gray level encoder. The 8-bit gray level image of the time signal (Figure 2d) is presented to a second PCNN (PCNN#2) to produce (by image segmentation) a 1D binary barcode (Figure 2e1) for the first iteration cycle. Subsequent iterations will likewise produce corresponding 1D barcodes that can be grouped together as a 2D barcode (Figure 2e2). As previously mentioned, there is a one-to-one correspondence between these barcoded PCNN#2 output and the corresponding converted time-signal input image (Figure 2d). Thus, each chemical structure will produce a unique barcode (1D or 2D).

### 3. CURRENT SEARCHING TECHNIQUE

The traditional method of retrieving the chemical information on a compound such as in Figure 3, in the event the name is either not familiar or cannot be generated, is used when a molecular formula is first obtained from the chemical structure. The elements therein are then arranged in the so-

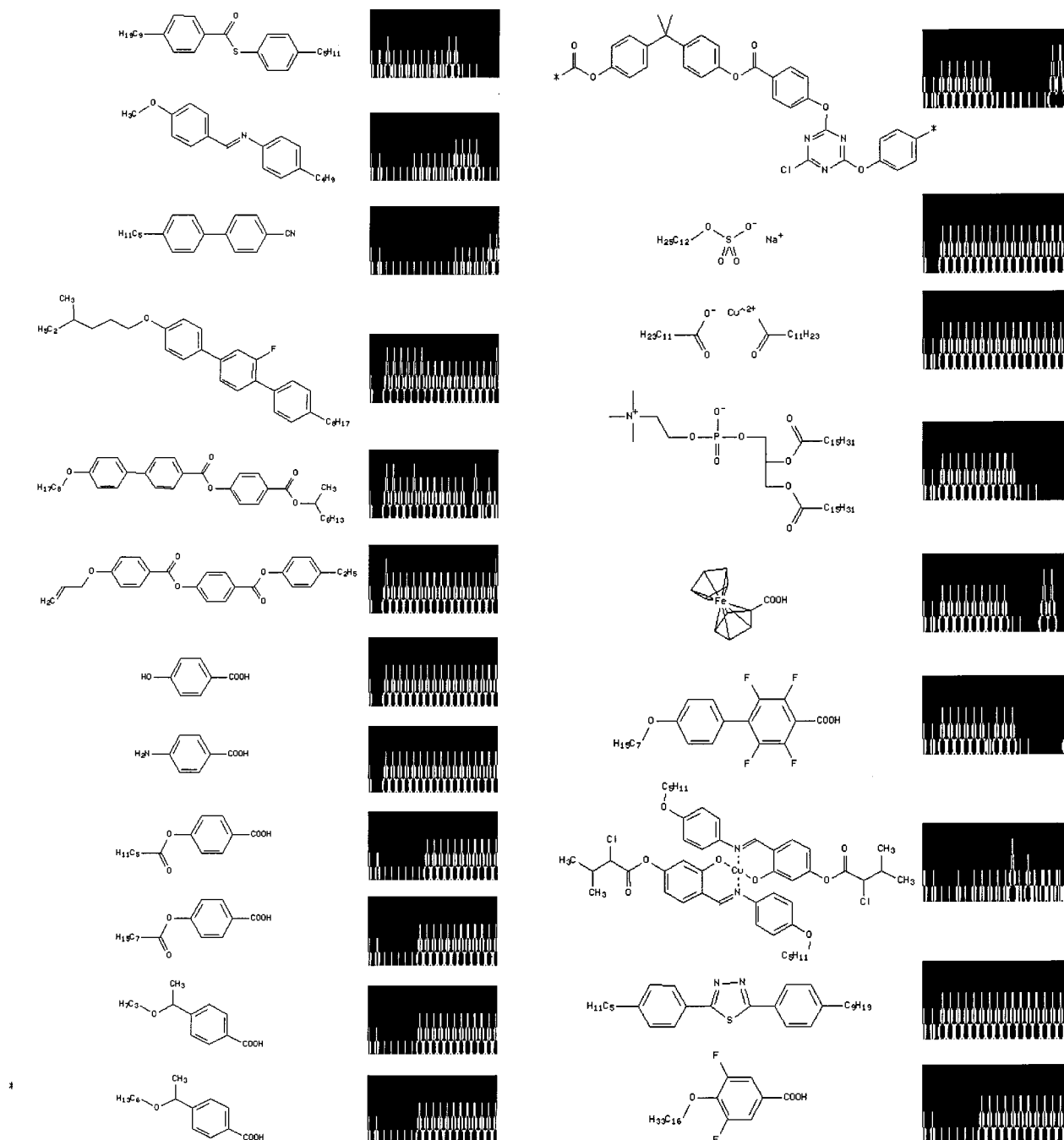
**Table 1.** Traditional Searching Technique (CAS Search)

Molecular Formula Search	DISPLAY SCAN option
=> FILE REGISTRY	=> D RN IN STR 1-10
=> E C14H18N4O3/MF	L2 ANSWER 1 OF 10 COPYRIGHT 1994 ACS
E1 1 C14H18N4O2TL/MF	RN 125705-88-6 REGISTRY
E2 1 C14H18N4O2ZR/MF	IN ***Carbamic acid, [5-[(3-methyl-1-oxobutyl)amino]-1H-benzimidazol-2-yl]-, methyl ester (9CI)***
E3 248 -> C14H18N4O3/MF	
E4 C14H18N4O3.(C2H4O)NC15H24	
O/MF	
E5 1 C14H18N4O3.(CH2O)X/MF	
E6 1 C14H18N4O3.1/2CL4PT.H/MF	
E7 1 C14H18N4O3.1/2H2O4S/MF	
E8 1 C14H18N4O3.2C2H6O3S/MF	
E9 1 C14H18N4O3.2C7H3IN2O3/MF	
E10 1 C14H18N4O3.2C7H6O2/MF	
E11 3 C14H18N4O3.2CLH/MF	
E12 3 C14H18N4O3.BRH/MF	
	L2 ANSWER 10 OF 10 COPYRIGHT 1994 ACS
	RN 17804-35-2 REGISTRY
	IN ***Carbamic acid, [1-[(butylamino)carbonyl]-1H-benzimidazol-2-yl]-, methyl ester (9CI)***
=> S E3 AND BUTYL	
248 C14H18N4O3/MF	
526169 BUTYL	
L1 28 C14H18N4O3/MF AND BUTYL	
=> S L1 AND AMINO AND METHYL	
1927517 AMINO	
6665678 METHYL	
L2 10 L1 AND AMINO AND METHYL	



**Figure 4.** Effects of translation, rotation, scaling, and drawing format on generated barcodes.

called Hill order. For carbon-containing compounds, the Hill order means that carbons are listed first, hydrogens are listed next, and then all other elements are listed in alphabetical order. For compounds that do not contain carbon, the elements are arranged in alphabetical order and the number of each element is indicated. A molecular formula search for this compound is then carried out. Obviously, there may be many chemical structures that may satisfy the molecular formula provided. Thus, there may be a need for further refinement. Table 1 shows a typical chemical search result using the chemical structure in Figure 3. In this example, a search using the molecular formula shows that there are many possible candidates. To reduce the number of answers from the molecular formula search, the latter is combined with name fragments in the Basic Index of the Registry File. The



**Figure 5.** Examples of chemical structures and their corresponding 2D-binary barcodes.

process is continued until a smaller set is obtained. Then the display option can be used to identify the RN and the name of the chemical (RN: 125705-88-6, index name: carbamic acid, [5-[(3-methyl-1-oxobutyl) amino]-1H-benzimidazol-2-yl]-,methyl ester).

The traditional method can thus be very time-consuming and tedious to the casual user. Searches on chemical databases are presently dominated by the text-based content of a paper which can be indexed into a keyword searchable form. Such traditional searches can prove to be very time-consuming and discouraging to the less frequent scientist.

#### 4. RESULTS

As described in section 2, the PCNN provides a 1:1 correspondence between images and their time signatures (icons), and this property is exploited for molecular structure

recognition. In Figure 4 we illustrate the invariance of the technique to translation (A,B) rotation (A, C, D) scale (A,E) using a thiophene molecule. The uniqueness property is sensitive to the drawing specification format as shown in (A,F) where the ring size and the bond lengths have been altered. Thus, it is important that all chemical structures are drawn according to a standard entry format. In Figure 5 we show the molecular structures of a number of diverse chemical structures and their corresponding binary barcodes. The uniqueness of the chemical structures and the binary barcodes suggest that this technique can be easily exploited for direct structure recognition of chemicals. The binary barcodes generated from chemicals forms a database. Comparing its binary sequences to those in a database can then search the barcode of a chemical structure. Because of the

1:1 correspondence, the hit from the database is thus a one-step procedure.

## 5. CONCLUSIONS

We have reported a simple chemical indexing based on the molecular structure. The method used is based on a 1:1 correspondence between the chemical structure presented as an image to a neural network and the corresponding binary output. We make use here of Pulse-Coupled Neural Networks (PCNNs) to produce binary barcodes of images of chemical structures. A number of parameters are adjustable to suit individual applications that renders this barcoding technique secure, versatile, and robust. The method is direct and less cumbersome (compared with traditional methods) and proves to be robust, elegant, and very versatile as demonstrated in this paper. We propose that the binary outputs can be used

as a replacement for CAS RNs. Our work is now being extended to the substructure scale.

## REFERENCES AND NOTES

- (1) Eckhorn, R.; Reitboeck, H. J.; Arndt, M.; Dicke, P. Feature linking via synchronization among distributed assemblies. *Neural Computation* **1990**, 2, 293–307.
- (2) Johnson, J. J. Pulse-Coupled Neural Nets: Translation, rotation, scale, distortion, and intensity signal invariances for images. *Appl. Optics* **1994**, 33(26), 6239–6253.
- (3) Lindblad, T.; Kinser, J. M. *Image Processing using Pulse-Coupled Neural Networks*; Springer-Verlag: London, 1998.
- (4) Rughooputh, H. C. S.; Rughooputh, S. D. D. V. Hybrid neural network technique for identification of Narcotics and Explosives. In *Proceedings of the 7th International Conference on Image Processing and its Applications (IPA'99)*; 1999; pp 784–788.
- (5) Rughooputh, H. C. S.; Rughooputh, S. D. D. V. Forensic application of a novel hybrid neural network. In *Proceeding of the International Joint Conference on Neural Networks (IJCNN'99)*; 1999.

CI000394D