

International Journal on Document Analysis and Recognition
Generation of Search-able PDF of the Chemical Equations segmented from Document Images.
 --Manuscript Draft--

Manuscript Number:	
Full Title:	Generation of Search-able PDF of the Chemical Equations segmented from Document Images.
Article Type:	Original Paper
Corresponding Author:	Sekhar Mandal, Ph.D Indian Institute of Engineering Science and Technology Howrah, West Bengal INDIA
Corresponding Author Secondary Information:	
Corresponding Author's Institution:	Indian Institute of Engineering Science and Technology
Corresponding Author's Secondary Institution:	
First Author:	Sekhar Mandal, Ph.D
First Author Secondary Information:	
Order of Authors:	Sekhar Mandal, Ph.D Prerana Jana, B. Tech Anubhab Majumdar, B. Tech Bhabatosh Chanda, Ph. D
Order of Authors Secondary Information:	
Funding Information:	
Abstract:	PDF format of scanned document images is not searchable. OCR tries to remedy this adversity by converting document images into editable and searchable data, but it has its own limitations in presence of equations - both mathematical and chemical. OCR system for mathematical equation is already a major research area and has provided successful result. However, chemical equation segmentation has been a less ventured road. In this paper, we present a novel method for automated generation of searchable PDF format of segmented chemical equations from scanned document images by performing chemical symbol recognition and auto-correction of OCR output. We use existing OCR system, pattern recognition technique, contextual data analysis and a standard L A TEX package to generate the chemical equation in searchable PDF format. The effectiveness of the proposed method is verified through exhaustive testing on 234 document images.
Suggested Reviewers:	

[Click here to view linked References](#)**IJDAR manuscript No.**

(will be inserted by the editor)

1
2
3
4
5

Generation of Searchable PDF of the Chemical Equations segmented from Document Images

10 Prerana Jana · Anubhab Majumdar · Sekhar Mandal · Bhabatosh Chanda
11
12
13
14
15
16
1718 Received: date / Accepted: date
1920 **Abstract** PDF format of scanned document images is
21 not searchable. OCR tries to remedy this adversity by
22 converting document images into editable and search-
23 able data, but it has its own limitations in presence
24 of equations - both mathematical and chemical. OCR
25 system for mathematical equation is already a major
26 research area and has provided successful result. How-
27 ever, chemical equation segmentation has been a less
28 ventured road. In this paper, we present a novel method
29 for automated generation of searchable PDF format of
30 segmented chemical equations from scanned document
31 images by performing chemical symbol recognition and
32 auto-correction of OCR output. We use existing OCR
33 system, pattern recognition technique, contextual data
34 analysis and a standard L^AT_EX package to generate the
35 chemical equation in searchable PDF format. The ef-
36 fectiveness of the proposed method is verified through
37 exhaustive testing on 234 document images.
38
3940 **Keywords** Chemical equations, mathematical sym-
41 bols, morphological operation.
4243
44 P. Jana
45 Indian Institute of Engineering Science and Technology,
46 Shibpur, India
47 E-mail: prerana.jana@gmail.com48 A. Majumdar
49 Indian Institute of Engineering Science and Technology,
50 Shibpur, India
51 E-mail: anubhabmajumdar93@gmail.com52 S. Mandal
53 Indian Institute of Engineering Science and Technology,
54 Shibpur, India
55 E-mail: sekhar@cs.iiests.ac.in56 B. Chanda
57 Indian Statistical Institute, Kolkata, India
58 E-mail: chanda@isical.ac.in

1 Introduction

Text keywords are used for retrieving documents from WWW using search engines like Google. A large number of documents are being digitized today for the purpose of archival, transmission and browsing. The existing OCR systems show high accuracy in interpreting text portions, but fail to process other components like graphics, half-tones, chemical and mathematical equations properly.

A few studies [19], [20], [21] are directed toward math-symbol or math equation recognition assuming that the math-zones are already marked. A number of work has been done over the past decade to detect and extract the mathematical equations present in heterogeneous document images.

Fatemian et al. [4] proposed a scheme which utilised character size, font information etc. to identify the connected components. Two bags, namely *text* and *math* are defined. The *text* bag is used to keep all letters and numbers; whereas the *math* bag collects punctuation, special symbols, Roman digits, italic letters, lines and dots. Objects in the *math* bag are then grouped together according to their spatial proximity. Grouping of items in *text* bag is redefined next followed by review and correction to move isolated items to their proper destinations. Math component segmentation is done in [5] through physical and logical segmentation using spatial characteristics of the math zone as well as identifying some math-symbols. The document is then segmented to characters, words, lines and blocks by physical segmentation. The logical segmentation process that follows consists of two steps; first the displayed math is detected by identifying their usual center position and in the next step in-line maths is detected by identifying special symbols.

Kacem et. al. [6] extracted the equations using fuzzy logic by detecting mathematical operators like '+', '-', etc. Their method was tested on a dataset consisting of 300 expressions and the success rate is about 93%. As some of the operators like '+', '-' , '(' and ')' do appear in chemical equations as well, it leads to the miss-classification of chemical equations as mathematical equations reducing the success rate. A similar method has been proposed in [7] to segment the mathematical expression in printed documents. The statistical approach taken by Garain [9] on the corpus of 400 pages to differentiate normal text lines and lines containing equations/expressions is on the basis of their white spacings which are usually larger in math-equation than the normal text. However, the chemical equations in the documents bear the same property. Jin et. al.[11] proposed a similar method to extract displayed formulas using Parzen classifier.

Drake and Baird [12] came up with a graphical approach; similarly Guo et. al.[13] developed a Gaussian mixture model to describe spatial relationships between sub-components of a math expression. Another method to check text style (regular, italic, bold) at the character level has been proposed in [10]. Garain [8] proposed a method to segment the displayed and embedded mathematical formulas from the documents using a bunch of features. The method is tested on a dataset of 200 images containing 1163 embedded and 1039 displayed expressions and the success rate is 88.3% and 97.2% respectively for embedded and displayed expressions. A method proposed by Chu and Liu [14] used features based on centroid fluctuation information on non-homogeneous regions to detect displayed and embedded formulas.

In a nutshell, in all the above methods emphasis is given only in mathematical equation. In eventual segmentation/classification, the chemical equations would automatically be included as a part of mathematical (or other) equations thereby reducing the success rate of the segmentation and effectiveness of the subsequent classification, if any.

There are some methods that are used to reconstruct chemical formula from scanned image. They have used chemical datasets. Algorri et al. [23, 24] proposed a system that reconstructs chemical molecules from scanned document. They have used connected component analysis and their own vectorisation algorithm for character recognition. Connected components that are not recognised by the OCR engine, are used to produce a graph of vectors. A rule based approach reconstructs the formula from the vector graph and the character information. ChemReader [25] starts with connected component. Alphanumerics are recognised using the GOCR

open source OCR tool. Graphical components are identified using Hough transforms, corner detection and other bespoke algorithms.

In this paper, we propose a fully automated segmentation and detection technique of chemical equations present in heterogeneous document images followed by generation of L^AT_EX file of the segmented equations.

The paper is organized as follows. Proposed method of segmentation, detection of chemical equation and its conversion to PDF form is presented in section 2. Section 3 presents experimental results. We conclude the paper in section 4.

2 Proposed Work

The proposed method consists of four distinct steps and they are as follows: (i) Locating of displayed equations; (ii) Extracting of chemical equations; (iii) Refining OCR output; and (iv) Converting the extracted chemical equation into search-able PDF format.

2.1 Locating displayed equations

A skew free heterogeneous binary image is the input to the proposed algorithm. The tables and graphics are extracted out from the heterogeneous documents using the techniques presented in [16] and [17] respectively. The rest of the document contains normal text lines and displayed equations if any. Our main motivation is to classify chemical and non-chemical equations, we restrict our study to displayed equations (DE) only, as embedded chemical equations are not frequent present in document.

Major sub-steps for locating of displayed equations (DE) are - (i) Text line segmentation; (ii) Blob formation using morphological tools; and (iii) DE zone extraction. The details of the aforesaid steps are described in the following subsections.

2.1.1 Text line segmentation

To detect DE zones, text lines have to be segmented first from which the operators are identified to determine whether a text line is a displayed equation or not. We have taken the horizontal projection profile of the document page to segment the text line. A part of a document page and its horizontal projection profile are shown in Fig. 1.

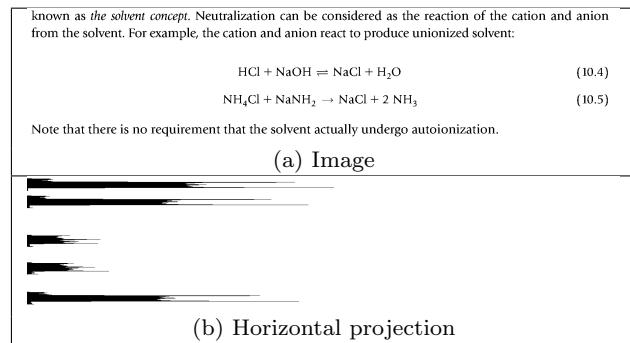


Fig. 1: A document page and its horizontal projection profile.

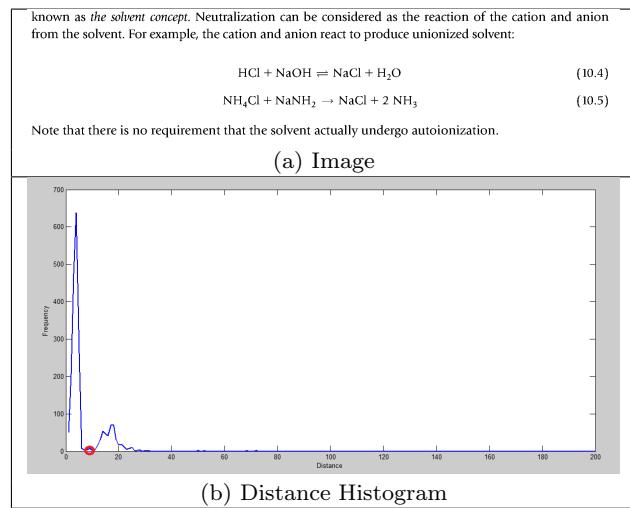


Fig. 2: A document page and its distance histogram.

2.1.2 Creation of word blobs

This is done by coalescing the characters in a word using morphological closing operation. Such character coalescing process depends on the accuracy in detecting the normal character gap and the gap between the consecutive connected components in that text line. The component analysis is done first.

The mathematical formulation for blob formation is as follows. Consider a binary image $I_{P \times Q}$, which consists of connected components $C_k (k = 1, 2, \dots, M)$, as defined in literature. Let C_a and C_b be the two consecutive connected components in a text line. L_b is the left most x-coordinate of C_b and R_a is the right most x-coordinate of C_a .

Let F be a function which ensures that the two connected components lie in the same text line. $T(C_k)$ and $B(C_k)$ return the topmost y-coordinate and the bottommost y-coordinate of C_k . Then F may be represented as

$$F(C_a, C_b) = \begin{cases} 1 & \text{if } (T(C_a) \leq B(C_b) \text{ AND } B(C_a) \geq T(C_b)) \\ 0 & \text{otherwise} \end{cases}$$

Cluster formation requires information on inter-word gap. The histogram H , in Fig. 2 shows the distribution of gaps or distance between two consecutive characters. The distance function, D , obtained from the histogram H , is defined for computing the horizontal distance between any two consecutive connected components C_a and C_b , as shown below

$$D(C_a, C_b) = L(C_b) - R(C_a)$$

where $b = \min_x \{L(C_x) - R(C_a)\}$ such that $(F(C_a, C_x) = 1 \text{ AND } L(C_x) > R(C_a))$.

The distance histogram is a multi-modal histogram. The first peak corresponds to the character gaps. An example image and corresponding distance histogram are shown in Fig. 2(a) and (b) respectively.

Our intention is to find out character gaps in running texts of a document page so that we can combine the consecutive characters into a single blob. Hence, we consider the upper boundary (l) of the first hump as the length of structuring element. Morphological close operation with a line structuring element of length (l) is carried out to form the word blobs. The result of blob formation is shown in Fig. 3 for the input image in Fig 2(a).

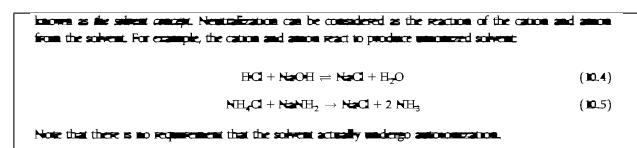


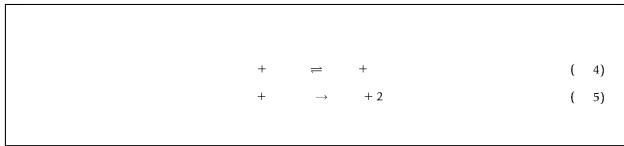
Fig. 3: Blob formation of the image shown in Fig. 2(a).

2.1.3 DE zone extraction

In a mathematical or chemical equation, one or more operators are present. These operators signal us the presence of displayed equations in a document. So, we have identified the common *operators* used in chemical and mathematical equations to segment the displayed equations from a text document image. We have considered the set of *operators* (+, -, →, ←, ↔, ↗, ↘) which are commonly used both in chemical equations as well as mathematical equations to fulfil our aim to identify displayed zones containing chemical or other equations.

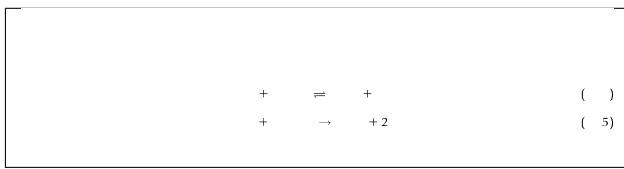
After blob formation, small component like dots of i and j are eliminated on the basis of area. The region (R_c) corresponding to each blob is cropped using its bounding box information from the original image and the number of connected component(s) present in R_c is

1 counted. If the number of components is more than 1,
 2 then that blob is not an operator and is removed from
 3 the blob image. The result of this operation is shown in
 4 Fig. 4 for the image in Fig. 3.
 5



13 Fig. 4: Single components extracted from the word blobs shown
 14 in Fig. 3

17 The remaining components in the blob image are
 18 operators along with some alphanumerics like *a*, *A*, (,),
 19 etc. The logical AND operation is performed between
 20 the blob image and the original image. The Euler num-
 21 ber of the operators that we have considered is 1 and
 22 based on this feature some of the alphanumerics are dis-
 23 carded and the resultant image is denoted by I_s (Fig. 5).
 24



32 Fig. 5: After removal of alphanumerals based on Euler number
 33 from the image shown in Fig. 4

Operators	horizontal projection	vertical projection
+		
—		
→		

43 Fig. 6: Operators and their horizontal and vertical projection
 44 profiles

47 Our next task is identify operator from I_s and for
 48 this, we have used neural network based classifier. For
 49 this purpose, we employ the following feature set.
 50

- 51 • Aspect ratio: (f_a) of each component
- 52 • Density:

$$55 f_d = \frac{\#pixels_o}{\#pixels_b},$$

58 where $\#pixels_o$ denotes the number of object pixels
 59 and $\#pixels_b$ denotes area of the bounding box.

- The horizontal and vertical projection profiles (see Fig. 6) of each component is obtained.

- (i) Spike for horizontal projection profile

$$f_{sh} = \frac{\#pixels_{on}}{w}$$

where $\#pixels_{on}$ denotes the number of on-pixels at
 the middle of the projection profile and w denotes
 the width of the profile.

- (ii) Spike for vertical projection profile

$$f_{sv} = \frac{\#pixels_{on}}{h}$$

where $\#pixels_{on}$ denotes the number of on-pixels at
 the middle of the projection profile and h denotes
 the height of the profile.

- Ratio of

$$f_{dr} = \frac{\#pixels_{on}}{\#pixels_{off}},$$

where $\#pixels_{on}$ denotes the number of object pixels
 and $\#pixels_{off}$ denotes the number of background pixels along the diagonal of each component. The ratio is determined for both the right(rd) and left(l) diagonals.

- A binary variable (f_{open}): $f_{open} \in \{1, 0\}$.

The value of f_{open} is obtained using morphological opening operation with a line like structuring element (SE) of length $\frac{w}{2}$, where w is the width of the component. If the output of the opening operation has a single component, f_{open} is set to 1, else f_{open} is 0.

- Number of end points (f_{ep}): f_{ep} is number of end points of a connected component. This is determined using thinning operation. After thinning each connected component in I_s , if a pixel has a single 8-connected neighbor, then that pixel represents an end point.

Now, $[f_a, f_d, f_{sh}, f_{sv}, f_{rd}^{trd}, f_{ld}^{lrd}, f_{open}, f_{ep}]$ is the feature vector for classification of operators from I_s .

We classify all single components in I_s into following 4 classes:

1. Arrows ($\rightarrow, \leftarrow, \leftrightarrow, \leftarrow, \rightarrow$)
2. Minus (-)
3. Plus (+)
4. Others (, (,), etc)

The classification is done using a two-layer feed-forward network, with sigmoid hidden and softmax output neurons and having 100 hidden layers (Fig. 7). The network is trained with scaled conjugate gradient back propagation.

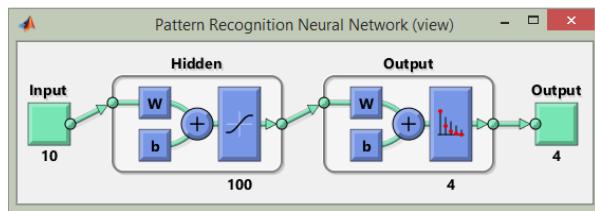


Fig. 7: The neural network used to classify single characters from Fig. 5

Table 1: Results of classifier for identification of *operators*

	1	2	3	4
1	319	0	0	2
2	0	120	0	0
3	0	0	2267	3
4	0	0	0	335

We have taken a set of 7046 samples from our image dataset. This set consists of aforesaid operators and other symbol/character. Out of 7046 images of aforesaid operators and single characters; 1000 plus, 1000 minus, 1000 arrows and 1000 other single characters are taken for the entire training set. Total training set consists of 4000 images.

The remaining 3046 samples are used for testing the classifier. The accuracy of the network is depicted by the confusion matrix of the test dataset in Table 1. The reason for high number of ‘+’ sign in the dataset is because it is the most frequently encountered *operator* as compared to the other operators or single characters. Fig. 8 shows the *operators* identified from Fig. 1(a).

We also need to identify the direction of the arrowhead for its correct representation. The direction of arrowhead is identified by measuring the height of the arrow elements near its two ends.

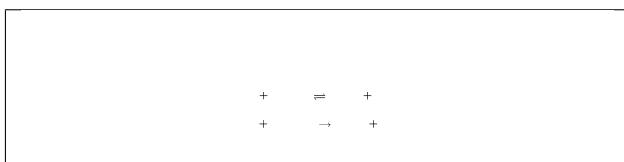


Fig. 8: Extracted operators from the image shown in Fig. 1(a)

To detect ‘=’ or ‘ \rightleftharpoons ’ one extra step is required. For each ‘-’ and ‘ \rightarrow ’, a rectangular window of size $l \times l/2$ is placed below the symbol to check if there is ‘-’ and ‘ \leftarrow ’ respectively within the window; if present, they are considered to form either ‘=’ or ‘ \rightleftharpoons ’ sign. l be the length of the symbol. The upper boundary of the window coincides with the lower boundary of the bounding box of each aforesaid symbol.

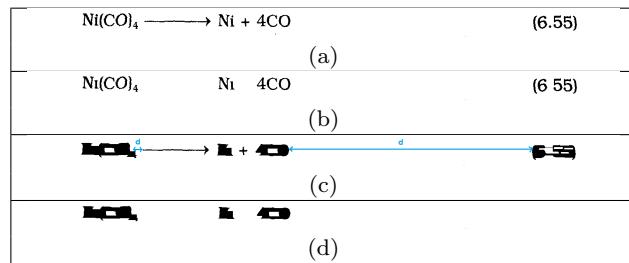


Fig. 9: The output of closing operation on portion of an image (a) a part of an image; (b) same part without operators; (c) result of closing operation; (d) after equation number removal.

Initially, all the text lines consisting of at least one *operator* are considered candidate displayed equations (CDE). The *operators* are separated from CDE. The upper boundary (u_v) of the second hump of distance histogram (Fig. 2(b)) is obtained which represents the word gaps in the text line. For each CDE zone morphological closing operation with a line structure element of length u_v is carried out. If the distance between two neighbouring components is less than u_v , it means they belong to a same word and are merged by closing operation. The output of closing is shown in Fig. 9.

Equation numbers are common in the displayed equation zones. These numbers have to be removed because for each CDE we have counted the number of *operators* and other corresponding components in the output of closing operation. If the number of components $\leq 2 \times$ number of *operators*, then the CDE is considered displayed equation; otherwise some embedded formulae/equations may exist in the line. To eliminate the equation number from the output of closing operation, the *operators* are moved to the output of closing operation and the component analysis is done. From both ends, distance (d) (Fig. 9(c)) between the first two consecutive components is measured and if $d > 5 \times u_v$, then the first component from the end is considered the equation number and is removed. (Fig. 9(d)).

2.2 Extracting of chemical equations

Each segment of a displayed equation is divided into three zones; namely upper zone, middle zone and lower zone (Fig. 10) along vertical direction. To identify the three zones of a DE zone, uppermost and lowermost co-ordinates of each connected component below the same DE zone are also obtained. The median of uppermost coordinate, and median of lowermost co-ordinate of such components in DE zone are computed. A virtual horizontal line, called the baseline, separates the middle zone and lower zone of DE zone. Similarly, the

median of uppermost co-ordinate of the components in the DE zone generates a horizontal line, called top line, which separates the middle and upper zones of the DE zone.

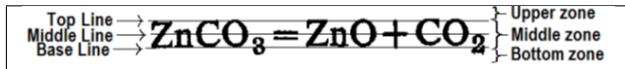


Fig. 10: Different zones of DE equation

The subscripts in a DE zone belong to lower-half of the middle zone and lower zone whereas the superscripts belong to upper zone and upper-half of the middle zone. Based on the location of the components in a DE zone we have detected the subscripts and superscripts and are separated from the DE zone. The operators are also separated from DE zone.

Now, each displayed equation is an input to an in-built OCR of MATLAB R2014a. The OCR returns each DE zone as a text string. We made a dictionary out of all the elements in the periodic table. An important observation is that an element always starts with a capital letter. Using this property, an element can be expressed by a regular expression $[A-Z][a-z]^*$. It means an element's symbol starts with an upper case letter and may or may not have one or more lower case letters (for example H , He , Uut etc). We have designed a parser to extract the sub-string matching the regular expression mentioned above with the following grammar:

```

start → capital.follow
follow → small.follow | ∈
capital → A|B|...|X|Y|Z
small → a|b|...|x|y|z

```

The working principle of the parser is depicted in Fig. 11. Each of the substring returned by the parser is matched against the aforesaid dictionary and if it is a positive match then that substring is considered as a symbol of the chemical element. Let us consider, the number of substrings extracted from the OCR output by the parser is n and the number of positive matches the aforementioned dictionary is m . If $m:n$ ratio is more than a threshold value β then this DE is considered a Chemical Equation. This threshold (β) is set to 0.7 experimentally by running the proposed algorithm on dataset containing 1390 displayed equations. The reason for the ratio not being 1 are: (i) Limitations of OCR, and (ii) Touching and broken characters.

The \uparrow and \downarrow are frequently used in chemical equation to represent the state of compounds and are thus important to detect. For each identified chemical equation, blob formation is done as discussed in Sec. 2.1.2. After

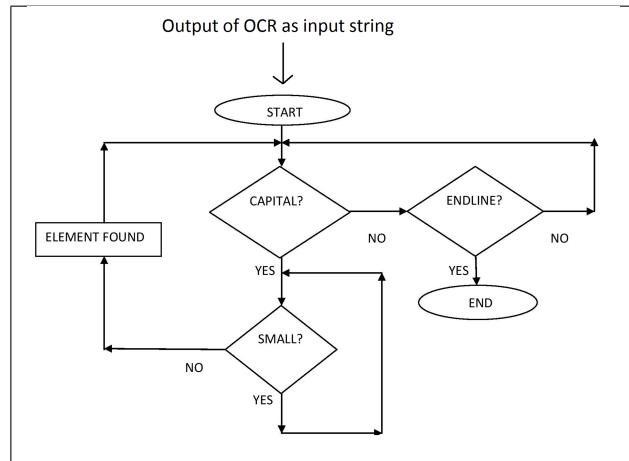


Fig. 11: Working flow chart of the parser

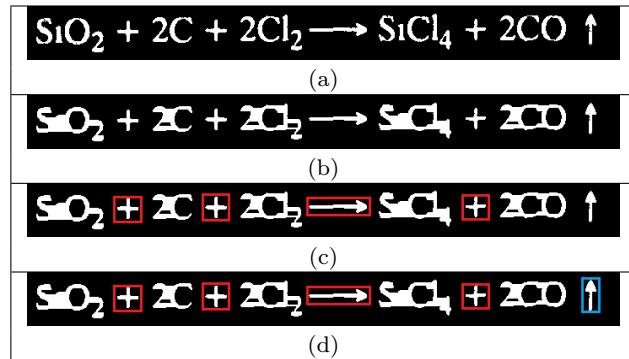


Fig. 12: Up and down arrow detection (a) A chemical equation; (b) Image after blob formation ; (c) Operators are marked in red; (d) Detected arrow in blue.

blob formation, we apply component analysis method. Then operators are marked (Fig. 12(c)) in the blob image, as they are already identified. For each component in blob image (which is not an operator), we check its immediate left component (C_l) and right component (C_r). If none of C_l or C_r is an operator, the component under consideration is an up arrow or a down arrow. Fig. 12(d) shows the detected up arrow in blue. The identification of the arrowhead is done in the same way as discussed in operator detection.

An example of chemical equation segmentation from a sample document image is given in Fig. 13. Fig. 13(c) shows the output in PDF without any correction.

2.3 Refinement of OCR output

Due to the limitations of OCR, the output of OCR is not fully correct in the paradigm of chemical equation/formula. Out of 3406 chemical compounds in our dataset, the accuracy of OCR conversion is only 43.13%.

Table 3: Part of the Error Hash Map

Erroneous OCR Output	Possible Input Set
8	g 3 a
3	g
S	g s
0	O
'E	3
s	3
w	3
3.	a
21	a
E1	a
8.	a
1	l i
I	l i
l1	n u
1'1	n
I1	n
11	n u
X1	n
ll	n
Cl	q
Q	q
1-1	H
1-1	H
l-1	H
l-1	H
7	2
4	2
Z	2
z	2
1	I i
U	u
ll	u
1l	u
'S	5
A	4
2	Z
C	e c

occurrence must be 1 or greater and * indicates the occurrence is 0 times or greater. 0 and 1 are excluded from the first digit as number of molecules or atoms cannot be 0 and if the number is 1, the numeric coefficient is not mentioned by default. Matched coefficients are stored in S_{coeff} .

2.3.2 State separation

The four physical states of a chemical compound - solid, gaseous, liquid, and aqueous are denoted by '(s)', '(g)', '(l)', and '(aq)', respectively. To detect the physical state of the compound, regular expression $[()A-Za-z0-9]^+$ is used and the checking starts from the end of $S_{chemical}$. The matched substring, S is extracted from $S_{chemical}$ and Algorithm 1 is run. As mentioned earlier, OCR output for each character is stored in a cell of S .

In this algorithm, S (after removing first and last character - opening and closing brackets) and H are taken as inputs and all possible Combinations of OCR output is produced by *GetAllCombinations* (See Fig. 15). For each cell element in S , corresponding values from hash map, H is assigned to a set, *InputSet* (See Line 3 in Algorithm 1). This set contains all possible inputs to the OCR system. Now, the key itself is added with its corresponding values in the hash table to make the *InputSet* if the length of key is 1. For example, '8' is added to the *InputSet* as its length is 1 (Fig. 15). On the contrary, in the second *InputSet*, 'Cl' is not included as its length is 2 (Fig. 15).

Each cell element of S gives one *InputSet*. Now, cartesian product of all the *InputSet* is taken to give us all possible Combinations and only one of the combinations is correct under proper chemical context. These Combinations are compared with 's', 'g', 'l' and 'aq'. If no match is found, the substring extracted from $S_{chemical}$ is concluded as a radical (Fig. 16, the compound contains a radical having the same regular expression mentioned earlier) , not a state; else, we separate the state from the compound and store it in S_{state} (after adding '(' and ')' at the start and end of S).

Algorithm 1 Get All Combinations from the Error Hash Map

```

1: procedure GETALLCOMBINATIONS( $S, H$ )
2:   for all element(s)  $\in S$  do
3:      $InputSet \leftarrow H.Get(element)$ 
4:     if  $InputSet$  is NULL then
5:       RETURN                                 $\triangleright$  Not in Error Map
6:     else
7:       if length(element) = 1 then
8:          $InputSet = \cup [element]$ 
       $\triangleright$  Element might be correct output but still in the error
      list for other inputs
9:       else
10:        Ignore
       $\triangleright$  Input is one character, output length > 1 means error
11:       end if
12:     end if
13:   end for
14:   Combinations  $\leftarrow$  CartesianProduct( $InputSet$ )
15:   Return Combinations
16: end procedure

```

2.3.3 Refinement of the formula unit

After extracting coefficient and state, only the formula unit is left in $S_{chemical}$. The algorithm for auto correction of each formula unit is done in two steps using Algorithm 1 and Algorithm 2. First, Algorithm 1 is performed on the formula unit to get all possible combina-

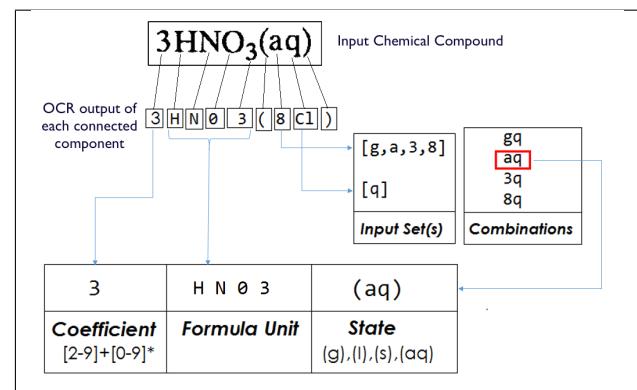


Fig. 16: Example of chemical compound having a radical in the end.

tions. Next, the output of Algorithm 1(*Combinations*) is taken as the input of Algorithm 2. This algorithm is used to match the *Combinations* against a nearly exhaustive list of all molecules, chemical compounds, radicals and atoms namely *ChemList* collected from Wikipedia.¹

Three cases may arise as follows–

(i) Exactly one match –

Fig. 17(b) shows the output of Algorithm 1. This is matched against *ChemList* using Algorithm 2 and algorithm finds one exact match as indicated by the red rectangle. This match is considered as the *Corrected* formula unit.

(ii) No match –

Longest common substring(s) (LCS) between *Combinations* and *ChemList* is computed and the formula unit in *Chemlist* having the longest common substring with *Combinations* is considered as *SubMatch*. There can be multiple such *SubMatch*. If there is only one, then the corresponding formula unit in *ChemList* is considered as the *Corrected* formula unit; else the *SubMatches* having the same length as that of the *Combinations* are considered as *PossibleFormulaUnits*. Fig. 18 (a) is a sample chemical compound. The OCR converted string is ‘N 21 B 1’. Algorithm 1 returns ‘NaBI’ and ‘NaBl’ as the two combinations. None of them match with any chemical compound in *ChemList*. Hence, LCS is computed between these two possible combinations and

¹ http://en.wikipedia.org/wiki/Dictionary_of_chemical_formulas

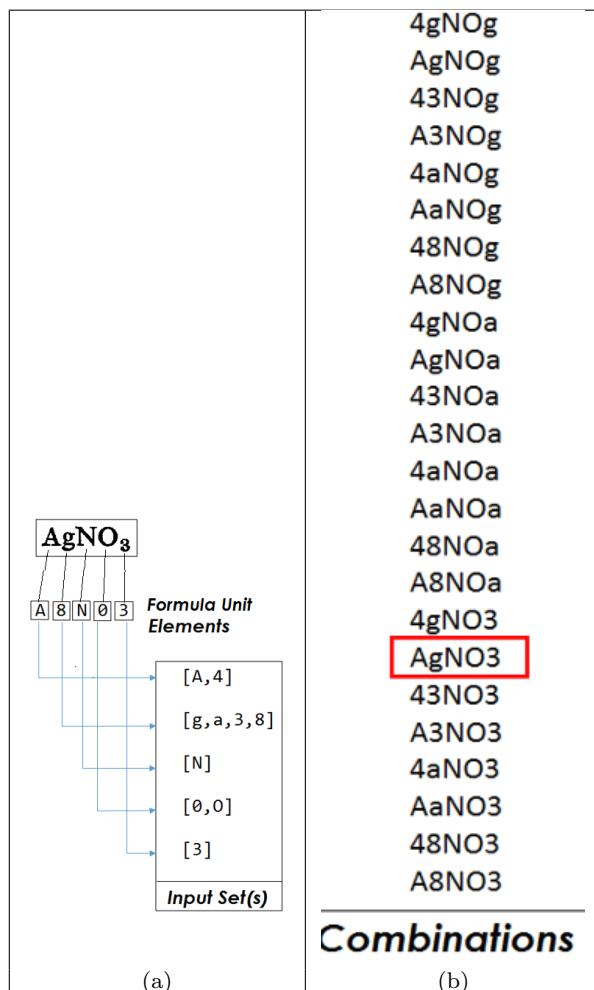


Fig. 17: Example of auto correction of a formula unit.

ChemList. Six compounds with LCS length 3 is found as shown in Fig. 18(b). Since the number of *SubMatches* is six, the *SubMatch* having the closest length as that of *Combination* (i.e. 4) is considered as *PossibleFormulaUnit*. In this case, it is ‘NaBr’.

(iii) More than one exact match –

In Fig. 19, ‘u’ of ‘Cu’ in left hand side of the equation is ‘11’ as the output of OCR. Among all possible combinations returned by Algorithm 1, ‘Cu’ and ‘Cn’ both match with *ChemList*. Hence, more than one exact match are found and both are considered as *PossibleFormulaUnit*.

The above steps are precisely mentioned in Algorithm 2. This algorithm returns *Corrected* and *PossibleFormulaUnits* upon which context analysis is done and is discussed in the next section.

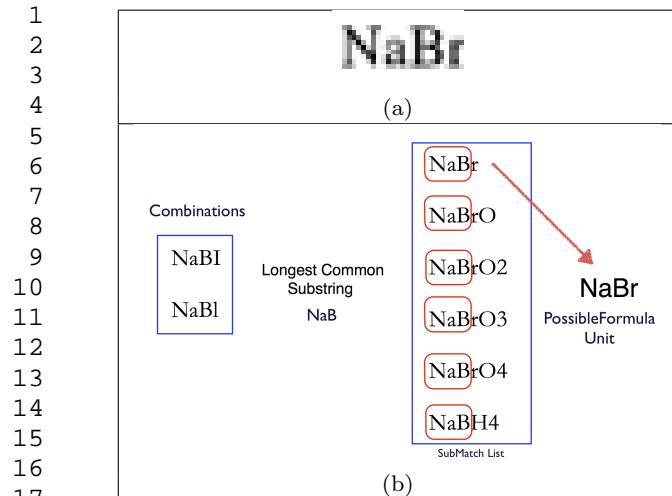


Fig. 18: (a) Sample Chemical Compound; (b) LCS Match.

Algorithm 2 Find Match between ChemList and Combinations derived from Algorithm 1

```

23: procedure FINDMATCH(ChemList,Combinations)
24:   for all Combinations do
25:     match with ChemList
26:   end for
27:   if #(Match Found) = 1 then
28:     Corrected  $\leftarrow$  Match
29:     Return Corrected
30:   else if #(Match Found) = 0 then
31:     SubMatch  $\leftarrow$  LCS match(Combinations,ChemList)
32:     SubMatch  $\leftarrow$  SubMatch
33:     Return SubMatch
34:   else
35:     PossibleFormulaUnit(s)  $\leftarrow$  SubMatches
36:     having closest length as Combination
37:     Return PossibleFormulaUnit(s)
38:   end if
39:   else
40:     PossibleFormulaUnits  $\leftarrow$  Matches
41:     Return PossibleFormulaUnit(s)
42:   end if
43: end procedure

```

2.3.4 Auto correction of the entire equation using Context Table

Here, we have *Corrected* and *PossibleFormulaUnit(s)* and try to find out the *FinalEquation* in the context of the equation itself. If a chemical equation does not have any *PossibleFormulaUnit*, context analysis is not required. The process exits after performing Line 2 of Algorithm 3.

Algorithm 3 takes all *Corrected* and *PossibleFormulaUnits* and returns the *FinalEquation* by forming the Context Table. As the universe is a

Chemical Equation	$\text{Cu} + 2\text{HNO}_3 + \text{O} \rightarrow \text{Cu}(\text{NO}_3)_2 + \text{H}_2\text{O}$				
OCR output	I1		I1		
Auto Corrected Formula Unit	Cu	2HNO3	O	Cu(NO3)2	H2O
Context Table	Reactants			Products	
	Case 1	Case 2		Cu	
	Cu ✓	Cn ✗		H	
	H	H		N	
	O	O		O	

Fig. 19: Formation of context table.

closed system, all chemical equations have the same periodic elements in the left hand side, called *Reactants* as that in the right hand side, called *Products*. All the periodic elements follow the regular expression [A-Z][a-z]*. So, for each *PossibleFormulaUnit*, the set of periodic elements in the *Reactants*, P_R and in the *Products*, P_P are computed and stored in the *ContextTable*. When the set difference of P_R and P_P in the table is empty, that *PossibleFormulaUnit* is considered as *Corrected* (Fig. 19 Case 1). In the Case 1 of P_P , the empty set condition satisfies. Hence, ‘Cu’ will be the *Corrected* formula unit, not ‘Cn’. But if the above condition comes true for multiple possibilities, we cannot decide which of the possible formula units are actually in the original equation. This is considered an *ERROR* case. Finally, S_{coeff} and S_{state} (if any) are added with their corresponding *Corrected* formula unit after the context analysis and this results in *FinalCompounds* for each equation.

Now according to the stoichiometry of the chemical reaction, pre-recognised operators along with *FinalCompounds* are concatenated together. This gives us the final auto-corrected chemical equation.

2.4 Generation of chemical equation in search-able PDF format

The final auto-corrected chemical equation is then converted to LATEX using the format specified by mhchem² package which provides commands or typesetting chemical molecular formulae and equations. This produces the searchable PDF format.

3 Experimental Result

We have implemented our algorithm in MATLAB 8.3.0.532 (R2014a) in a PC (Intel(R) Core(TM) i5-3337U CPU

² <ftp://www.ctan.org/tex-archive/macros/latex/contrib/mhchem/mhchem.pdf>

1 **Algorithm 3** Auto-Correction of the entire equation
 2 using chemical context table
 3

```

4 1: procedure GETFINALEQN(PossibleFormulaUnit,Corrected)
5   2:   Include all Corrected units in FinalCompound
6   3:   Count  $\leftarrow$  0
7   4:   for every PossibleFormulaUnit do
8     5:     compute(PR)
9       6:        $\triangleright$  PR : Set of periodic elements in Reactants
10      7:     compute(PP)
11       8:        $\triangleright$  PP : Set of periodic elements in Products
12      9:     if PR - PP =  $\emptyset$  then
13        10:       Corrected  $\leftarrow$  PossibleFormulaUnit
14        11:       Count  $\leftarrow$  Count + 1
15      12:     end if
16    13:   end for
17  14:   if Count = 1 then
18    15:     FinalCompound  $\leftarrow$  [Scoeff, Corrected, Sstate]
19  16:   else if Count = 0 then
20    17:     Break
21  18:   else
22    19:     Multiple Corrected compounds
23  20:    $\triangleright$  ERROR
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

```

@ 1.80GHz running Windows 8). The proposed method has been tested on a dataset consisting of 234 document images. Out of 234 pages 50 are taken from ICDAR 2013 Math-zone segmentation datasets and other document pages are scanned from different Mathematics and Chemistry books. The summary of the experimental results is shown in Table 4. Out of 3406 chemical formula in the test dataset, 114 formula were partially corrected and 52 formula could not be corrected at all. The overall accuracy of complete refinement is 95.12%. This is measured by (#Completely Corrected formula / #Total number of formula). Due to the longest common substring match and then performing context analysis of the entire equation, there is a very small window of zero correction. Zero correction is the case when there have been no correction to the OCR output by the auto-correction algorithm. For example, OCR output of ‘Mg’- *I*³ could not be corrected by our auto-correction algorithm as this erroneous conversion was not in the error hash map.

With our dataset, zero correction rate is 0.01% (It is computed as #Zero Correction Compound / #Total Compounds). These results are quite encouraging.

Consider the sample image (Fig. 13(a)) and its corresponding segmented displayed chemical equations are shown in Fig. 13(b). Fig. 13(c) shows the direct OCR output where ‘i’ has been wrongly identified as ‘I’, ‘T’ and ‘1’ (for *Si* in all the lines of Fig. 13(c)). Similarly ‘O’ results in ‘0’ (line 1,2,3). ‘S’ sometimes is detected as ‘5’ (line 5). Our auto correction algorithm remedies these

Table 4: Summary of Experimental Results

#Total Images	234
#Total DEs	1390
Operator recognition	99.8%
DE segmentation accuracy	98.63%
Chemical DE Classification Accuracy	98.83%
#Total Chemical Operands	3406
Complete refinement accuracy	95.12%
Zero Auto correction rate	0.01%

issues. Fig. 20 demonstrates the effect of our auto correction algorithm. This algorithm is targeted towards chemical equation with linear representation. Organic bonds cannot be detected in this system.

Some sample experimental results are shown in Fig. 23, Fig. 24, Fig. 25 and Fig. 26. More results are shown in <https://sites.google.com/site/chemeqndb/home>.

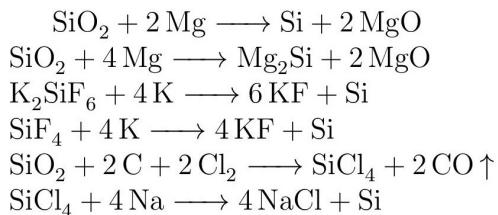


Fig. 20: Auto corrected output of Fig. 13(c).

Next, we try to analyse the sources of some of the errors and shortcomings of our algorithm which have negative effect on the performance figures for auto correction.

Case 1 : Chemical equations sometimes contain some texts such as ‘and’, ‘or’ etc between two chemical compounds (See Fig. 21(a)). Sometimes two chemical equations are conjuncted by these words in the same line. If these words are not in chemical context and OCR does not convert them correctly, our autocorrection algorithm cannot match them against *ChemList*, hence the error occurs. But OCR conversion has a high accuracy rate for such type of texts. Therefore, this is not a severe error.

Case 2 : When the chemical compound is written in formats such as (*Na₂SiO₃*)_n (Fig. 21(a)), only *Na₂SiO₃* is detected based on *ChemList* and LCS matching. This is considered as a partial autocorrection case.

Case 3 : Some equations have conditions (pressure, temperature) written over the arrows (See Fig. 21(b)). In this work, we only concentrated on chemical compounds in the equation. This does not effect the auto-

correction accuracy rate as most of the time they get segmented in separate text lines; else we ignore the over arrow conditions beforehand.

Case 4 : Fractions in the numeric coefficients (Fig. 21(c)) are not dealt with in our autocorrection algorithm as they are not very common in chemical equations.

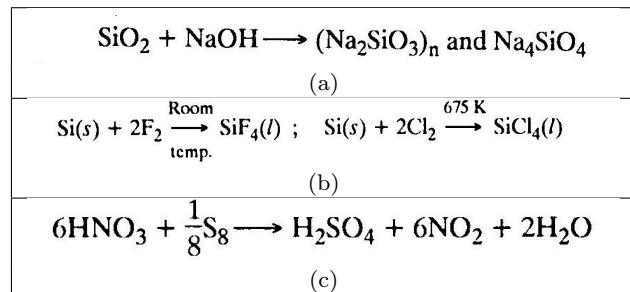


Fig. 21: Sample error cases (a) presence of non-chemical words in segmented chemical equation; (b) presence of over arrow conditions; (c) fractional coefficient.

Case 5 : Here, in the equation shown in Fig. 22, both the ‘g’s in reactant and product side have been converted to ‘S’ by the OCR which results in multiple auto-corrected formula units on both side. At this point, we reach Step 17 of Algorithm 3 where context table formation cannot conclude which one is the final corrected compound. However, normally the probability of occurrence of such situation is extremely rare, so no further steps are taken to rectify this.

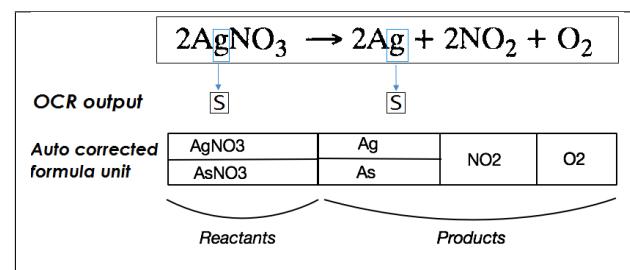


Fig. 22: Error case of multiple *Corrected* compounds

4 Conclusion

We have presented an automated chemical equation segmentation and chemical context based auto correction system that is able to provide the exact searchable format of linear chemical equations in any document image. The experimental results demonstrate the efficiency of our proposed method. One of the drawbacks of our system is the time complexity as the search space

in the *ChemList* is quite big and is growing over time due to discovery of new compounds. The search method can be improved and made more efficient. Since our proposed method is novel, we have not concentrated on making the system time efficient yet but more on the accuracy of the auto correction. This work leads to several research avenues. Chemical context horizon can be widened. Auto correction on non-linear or bond structure representations of chemical equations could be ventured in.

References

- D. Blostein and A. Grabavac. “Recognition of Mathematical Notation”, *Handbook of Character Recognition and document Image Analysis*, 557–582, 1997.
- K-F. Chan and D-Y. Yeung. “Mathematical Expression Recognition: A Survey”. *IJDAR*, Vol. 3, no, 1, 3–15, 2000.
- U. Garain and B. B. Chaudhuri. *On OCR of Printed mathematical Expressions*. “Digital Document Processing”, Ed. B. B. Chaudhuri, *Advances in pattern Recognition*, 235–259 , 2007.
- R. Fateman, T. Tokuyasu, B. Berman, and N. Mitchell. “Optical character recognition and parsing of typeset mathematics”, *Visual Commun. And Image Representation*, Vol 7, no 1, 2–15, 1996.
- J. Y. Toumit, S. Garcia-Salicetti, and H. Emptoz. “A hierarchical and recursive model of mathematical expressions for automatic reading of mathematical documents”. In Proc. of ICDAR, 116–122, 1999.
- A. Kacem, A. Beliad and M. Ben Ahmed. “Automated Extraction of printed mathematical formulas using fuzzy logic and propagation of context”, *IJDAR*, vol.4 no. 2, 97–108, 2001.
- M. Suzuki, F. Tamari, R. Fukuda, S. Uchida and T. Kanahori. “INFTY - An Integrated OCR system for Mathematical Documents”, Proc. of ACM Symposium on Document Engineering, 95–104, 2003.
- Utpal Garain. “Identification of Mathematical Expressions in Document Images”, Proc. of ICDAR, 1340–1344, 2009.
- Utpal Garain. “Recognition of Printed Handwritten Mathematical Expressions”, Ph.D Thesis, ISI, Kolkata, India, 2005.
- B. B. Chaudhuri and U. Garain. “Extraction of type atyle based meta-information from Imaged documents”, *IJDAR*, vol. 3 no. 3, 138–149, 2001.
- J. Jin, X. Han and Q. Wang. “Mathematical formulas extraction”, Proc. of ICDAR, 1138–1141, 2003.
- D. M. Drake and H. S. Baird. “Distinguishing mathematical notation from english text using computational geometry”, Proc. of ICDAR, 1270–1274, 2005.
- Y.-S. Guo, L. Huang and C.-P. Liu. “A new approach for understanding of structure of printed mathematical expressions”, Proc. of ICMLC, 2633–2638, 2007.
- We-Te Chu and Fan liu. “Mathematical formula detection from heterogeneous document Images”, Proc. of CTAAI, 140–146, 2013.
- S. P. Chowdhury, S. Mandal, A. K. Das, and B. Chanda. “Segmentation of Text and Graphics from Document Images”, in Proc. of ICDAR, 619–623, 2007.
- S. Mandal, S. P. Chowdhury, A. K. Das, and B. Chanda, *A simple and effective table detection system from Document Images*, IJDAR, Vol. 8(2), 172182, 2006.

- 1 17. S. P. Chowdhury, S. Mandal, A. K. Das, and B. Chanda,
2 *Segmentation of Text and Graphics from Document Im-*
3 *ages*, In Proc. of ICDAR, pp. 619623,2007.
4 18. R. C. Gonzalez and R. Wood, *Digital Image Processing*,
5 Addison-Wesley, 1992.
6 19. D. Blostein and A. Grabavec, *Recognition of Mathemat-*
7 *ical Notation*, Handbook of Character Recognition and
8 document Image Analysis, 577–582, 1997.
9 20. K-F. Chan and D-Y. Yeung, *Mathematical Expression*
10 *Recognition: A Survey*, IJDAR, Vol. 3, no: 1, 315, 2000.
11 21. U. Garain and B. B. Chaudhuri *An OCR of Printed*
12 *mathematical Expressions*, *Digital Document Processing*,
13 Ed: B. B. Chaudhuri, Advances in pattern Recognition,
235259 , 2007.
14 22. A. Fujiyoshi, M. Suzuki, S. Uchid, *Grammatical Verifi-*
15 *cation for Mathematical Formula Recognition Based on*
16 *Context-Free Tree Grammar*, Mathematics in Computer
17 Science, 279–298, 2010.
18 23. M. E. Algorri, M. Zimmermann, C. M. Friedrich, S. Akle,
19 and M. Hofmann-Apitius, *Reconstruction of chemical*
20 *molecules from images*, In Proc. 29th Annual International
21 IEEE Conference on Engineering in Medicine and Biology
Society, 4609–4612, 2007.
22 24. M. E. Algorri, M. Zimmermann, and M. Hofmann-
23 Apitius, *Automatic recognition of chemical images*, In
24 pro. Eighth Mexican International Conference on Current
Trends in Computer Science, 41–46, 2007.
25 25. J. Park, G. R. Rosania, K. A. Shedden, M. Nguyen,
26 N. Lyu, and K. Saitou, *Automated extraction of chemical*
27 *structure information from digital raster images*, Chemistry
28 Central journal, vol. 3(1), 2009.
29 26. A. K. Jain, J. Mao, K. M. Mohiuddin *Artificial Neural*
30 *Network: A Tutorial*, Computer (Volume:29 , Issue: 3),
Mar 1996.

31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

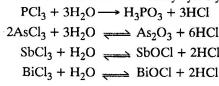
Chemistry of Non-metallic Elements and Their Compounds

785

3. Formation of Halides. All the elements of nitrogen family form trihalides such as NCl_3 , PCl_3 , AsCl_3 , SbCl_3 and BiCl_3 . The trihalides of nitrogen are not stable. However, NF_3 is stable because of the small size of fluorine.

Trihalides have pyramidal structure like ammonia. Central atom in trihalides is sp^3 hybridised and one of the hybrid orbitals is occupied by a lone pair.

Except nitrogen trihalides (which cannot form co-ordinate bonds with water due to the absence of d -orbitals), trihalides of all other elements are easily hydrolysed by water as :

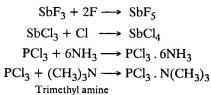


Phosphorus trifluoride, PF_3 is stable and does not undergo hydrolysis because of stronger P–F bonds. P–F bonds are stronger than P–O bonds. Therefore, P–O bonds are not formed by breaking P–F bonds. These halides are predominantly covalent. However, their ionic character increases as we move down the group. BiCl_3 is quite ionic.

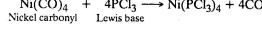
NF_3 has low dipole moment as compared to ammonia. NF_3 and NH_3 both have pyramidal geometry with a lone pair of electrons on nitrogen. The directions of dipoles in case of three N–F bonds are towards F because F is more electronegative than N and oppose the effect of unshared electron pair on nitrogen. The directions of dipoles in the three N–H bonds are towards N and add to the effect of the electron pair (Fig. 7.45).

NF_3 does not act as a Lewis base like NH_3 . Fluorine due to its higher electronegativity, withdraws electrons from nitrogen towards itself. Thus, N-atom is devoid of its electron donor character and hence does not act as a Lewis base.

Trifluorides and trichlorides of phosphorus and antimony also act as Lewis acids (i.e. tendency to accept lone pair of electrons). It is because of the tendency of these elements to accept lone pair of electrons in their vacant d -orbitals. For example :



PCl_3 also acts as Lewis base (electron donor nature) because P-atom is able to donate its lone pair of electrons to the vacant d -orbitals of other elements like nickel. For example :



Nickel carbonyl Lewis base

All members of nitrogen family except nitrogen and bismuth form pentahalides especially pentafluorides. Nitrogen does not form pentahalides due to the non-availability of the d -orbitals. Bismuth does not form a pentahalide because of the reluctance of 6s electrons of bismuth to participate in bond formation (inert pair effect).

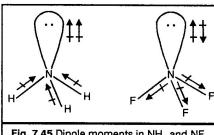
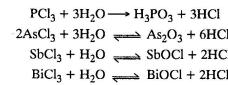
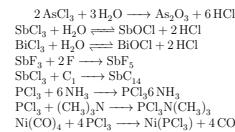
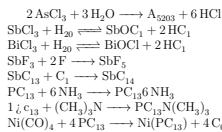


Fig. 7.45 Dipole moments in NH_3 and NF_3 .



(a) Input Image

(b) Segmented Displayed Chemical Equations



(c) Direct OCR output

(d) Auto-corrected OCR output

Fig. 23: Experimental result: Sample 1.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Chemistry of Non-metallic Elements and Their Compounds 7/83

state is exhibited when only the three p -electrons of the valence shell take part in bonding while the s -electrons do not take part in bonding due to inert pair effect. The stability of compounds with + 3 state increases on moving down the group.

Nitrogen also exhibits oxidation states of + 1, + 2 and + 4 in compounds such as N_2O , NO and NO_2 respectively.

7.34.4 General Trends in Chemical Reactivity

1. Reactivity. Nitrogen is a diatomic gas. Each atom completes its octet by sharing three pairs of electrons. The ($\text{N}=\text{N}$) triple bond is very strong bond (distance 1.094\AA). This is indicated by its very high dissociation energy (900 kJ/mole). The collision between nitrogen molecules at room temperature are unable to break the strong triple bond. Hence, nitrogen is inert at room temperature. Only at very high temperature, sufficient energy is supplied by collisions to break the triple bond. Thus, at higher temperatures, nitrogen molecule gets dissociated and shows some chemical activity. Since nitrogen is an inert element, its presence in air dilutes the effect of oxygen. If nitrogen was not present in air, burning process could not have been controlled. Presence of nitrogen in air also helps in checking the process of rusting.

Unlike nitrogen, phosphorus is a reactive element because in P_4 molecule, each phosphorus atom is linked to other phosphorus atoms through single covalent bonds. The breaking of P-P bond requires much less energy. Hence, phosphorus is quite reactive element.

2. Formation of hydrides. All members of the nitrogen family form hydrides having the general formula MH_3 . The typical hydrides of these elements are :

NH_3	PH_3	AsH_3	SbH_3	BiH_3
Ammonia, ₂	Phosphine	Arsine	Sibine	Bismuthine

These hydrides can be obtained by the action of water or dilute mineral acids with compounds like nitrides, phosphides, etc.

$\text{Mg}_3\text{N}_2 + 6\text{H}_2\text{O} \longrightarrow 3\text{Mg}(\text{OH})_2 + 2\text{NH}_3$	$\text{Mg}_3\text{N}_2 + 6\text{H}_2\text{O} \longrightarrow 3\text{Mg}(\text{OH})_2 + 2\text{NH}_3$
Magnesium nitride	Magnesium nitride
$\text{Ca}_3\text{P}_2 + 6\text{H}_2\text{O} \longrightarrow 3\text{Ca}(\text{OH})_2 + 2\text{PH}_3$	$\text{Ca}_3\text{P}_2 + 6\text{H}_2\text{O} \longrightarrow 3\text{Ca}(\text{OH})_2 + 2\text{PH}_3$
Calcium phosphide	Phosphine
$\text{Zn}_3\text{As}_2 + 6\text{HCl} \longrightarrow 3\text{ZnCl}_2 + 2\text{AsH}_3$	$\text{Zn}_3\text{As}_2 + 6\text{HCl} \longrightarrow 3\text{ZnCl}_2 + 2\text{AsH}_3$
Zinc arsenide	Arsine
$\text{Mg}_3\text{Sb}_2 + 6\text{HCl} \longrightarrow 3\text{MgCl}_2 + 2\text{SbH}_3$	$\text{Mg}_3\text{Sb}_2 + 6\text{HCl} \longrightarrow 3\text{MgCl}_2 + 2\text{SbH}_3$
Magnesium stibide	Sibine
$\text{Mg}_3\text{Bi}_2 + 6\text{HCl} \longrightarrow 3\text{MgCl}_2 + 2\text{BiH}_3$	$\text{Mg}_3\text{Bi}_2 + 6\text{HCl} \longrightarrow 3\text{MgCl}_2 + 2\text{BiH}_3$
Magnesium bismuthide	Bismuthine

Nitrogen also forms hydrides like N_2H_4 (Hydrazine) and N_2H_2 (hydrazoic acid).

Hydrazine (N_2H_4) is prepared by oxidising NH_3 with sodium hypochlorite (NaOCl)

$$2\text{NH}_3 + \text{NaOCl} \longrightarrow \text{N}_2\text{H}_4 + \text{NaCl} + \text{H}_2\text{O}$$

It is a strong reductant. Hydrazine as well as its derivatives are used as rocket fuels.

Structure. In all the hydrides having general formula MH_3 , the central atom is sp^3 hybridized and forms four sp^3 hybrid orbitals. One of the four sp^3 hybrid orbitals contains a lone pair of electrons and the remaining three sp^3 hybrid orbitals are used to make three M-H bonds. Thus MH_3 molecule has a pyramidal geometry. The structures of NH_3 and PH_3 molecules are shown in Fig. 7.44.

(a) More repulsion, bond angle is large
(b) Less repulsion, bond angle is less.

Fig. 7.44 Structure of NH_3 and PH_3

(a) Input Image

(b) Segmented Displayed Chemical Equations

$\text{Mg}_3\text{N}_2 + 6\text{H}_2\text{O} \longrightarrow 3\text{Mg}(\text{OH})_2 + \text{ZNH}_3$
 $\text{Ca}_3\text{P}_2 + 6\text{H}_2\text{O} \longrightarrow 3\text{Ca}(\text{OH})_2 + \text{ZPH}_3$
 $\text{Zn}_3\text{As}_2 + 6\text{HCl} \longrightarrow 3\text{ZnCl}_2 + \text{ZASH}_3$
 $\text{Mg}_3\text{Sb}_2 + 6\text{HCl} \longrightarrow 3\text{MgCl}_2 + 2\text{SbH}_3$
 $\text{Mg}_3\text{Bi}_2 + 6\text{HCl} \longrightarrow 3\text{MgCl}_2 + 2\text{BiH}_3$
 $\text{ZNH}_3 + \text{NaOCl} \longrightarrow \text{N}_2\text{H}_4 + \text{NaCl} + \text{H}_2\text{O}$

$\text{Mg}_3\text{N}_2 + 6\text{H}_2\text{O} \longrightarrow 3\text{Mg}(\text{OH})_2 + 2\text{NH}_3$
 $\text{Ca}_3\text{P}_2 + 6\text{H}_2\text{O} \longrightarrow 3\text{Ca}(\text{OH})_2 + 2\text{PH}_3$
 $\text{Zn}_3\text{As}_2 + 6\text{HCl} \longrightarrow 3\text{ZnCl}_2 + 2\text{AsH}_3$
 $\text{Mg}_3\text{Sb}_2 + 6\text{HCl} \longrightarrow 3\text{MgCl}_2 + 2\text{SbH}_3$
 $\text{Mg}_3\text{Bi}_2 + 6\text{HCl} \longrightarrow 3\text{MgCl}_2 + 2\text{BiH}_3$
 $2\text{NH}_3 + \text{NaOCl} \longrightarrow \text{N}_2\text{H}_4 + \text{NaCl} + \text{H}_2\text{O}$

(c) Direct OCR output

(d) Auto-corrected OCR output

Fig. 24: Experimental result: Sample 2.

<p>Chemistry of Non-metallic Elements and Their Compounds</p> <p>7.30. CARBON-NITROGEN COMPOUNDS*</p> <p>An important compound containing carbon as well as nitrogen is calcium cyanamide, CaCN_2. It is obtained by heating CaC_2 with nitrogen at 1373 K.</p> $\text{CaC}_2(s) + \text{N}_2(g) \longrightarrow \text{CaCN}_2(s) + \text{C}(s)$ <p>Mixture of CaCN and C is used as a fertilizer under the name nitrolim. It is also used to manufacture melamine plastics. Calcium cyanamide is the starting material for the manufacture of sodium cyanide which is obtained by fusing calcium cyanamide with C and Na_2CO_3.</p> $\text{CaCN}_2 + \text{C} + \text{Na}_2\text{CO}_3 \longrightarrow \text{CaCO}_3 + 2\text{NaCN}$ <p>Sodium cyanide is used for the extraction of silver and gold from their ores. On treatment with strong acids, sodium cyanide liberates HCN which is a colourless gas and behaves as a weak acid in aqueous solution ($\text{pK}_a = 9.0$).</p> <p>On a large scale HCN is obtained by heating ammonia with methane at a high temperature.</p> $\text{CH}_4(s) + \text{NH}_3(s) \xrightarrow[\text{1500 K}]{\text{Pt. Catalyst}} \text{HCN}(g) + 3\text{H}_2(g)$ <p>Cyanides and HCN are extremely poisonous and their ingestion or inhalation may prove fatal. HCN is used in the manufacture of methyl methacrylate polymers and adiponitrile, which is an intermediate for nylon.</p> <p>Two other compounds containing carbon and nitrogen are cyanogen, $(\text{CN})_2$ and cyanamide, H_2NCN. Cyanogen has perfect resemblance to halogens (X_2) and is referred to as a pseudohalogen. Cyanogen can be obtained by the oxidation of HCN by O_2 using a silver catalyst or by the oxidation of CN^- by Cu^{2+}.</p> $4\text{HCN} + \text{O}_2 \xrightarrow{\text{Ag}} 2(\text{CN})_2 + 2\text{H}_2\text{O}$ $4\text{CN}^- + 2\text{Cu}^{2+} \longrightarrow 2\text{Cu}(\text{CN}) + (\text{CN})_2$ <p>Cyanogen is a poisonous gas like HCN. It has linear structure and disproportionates in basic solution to cyanide and cyanato ions.</p> $(\text{CN})_2 + 2\text{OH}^- \longrightarrow \text{CN}^- + \text{OCN}^- + \text{H}_2\text{O}$ <p>Ca CN on treatment with water gives cyanamide which is a solid having m.p. 318 K.</p> $\text{CaCN} + \text{H}_2\text{O} \longrightarrow \text{CaO} + \text{H}_2\text{NCN}$ <p>7.31. SILICON</p> <p>It is the second member of group 14. Silicon appears just below carbon in the periodic table. Its atomic number is 14 and therefore, it has the electronic configuration $1s^2 2s^2 2p^2 3s^2 3p^2$. Silicon is expected to give characteristics similar to that of carbon since the two have similar electronic configuration ($ns^2 np^2$). This is true in certain cases. For example, silicon forms compounds such as SiH_4 and SiCl_4 which are covalent compounds and have tetrahedral geometry just like CH_4 and CCl_4. However, carbon and silicon differ in most of their characteristics. For example,</p> <ul style="list-style-type: none"> (i) CO_2 is a gas while SiO_2 is a solid. (ii) Melting point of carbon (3773 K) is much higher than that of silicon (1700 K). (iii) CCl_4 is not hydrolysed by water while SiCl_4 is hydrolysed. <p>*For Entrance Examinations.</p> <p>(a) Input Image</p> <p>(b) Segmented Displayed Chemical Equations</p> <p>(c) Direct OCR output</p> <p>(d) Auto-corrected OCR output</p>	<p>$\text{CaC}_2(s) + \text{N}_2(g) \longrightarrow \text{CaCN}_2(s) + \text{C}(s)$</p> <p>$\text{CaCN}_2 + \text{C} + \text{Na}_2\text{CO}_3 \longrightarrow \text{CaCO}_3 + 2\text{NaCN}$</p> <p>$\text{CH}_4(g) + \text{NH}_3(g) \longrightarrow \text{HCN}(g) + 3\text{H}_2(g)$</p> <p>$4\text{HCN} + \text{O}_2 \longrightarrow 2(\text{CN})_2 + 2\text{H}_2\text{O}$</p> <p>$4\text{CN}^- + \text{Cu}^{2+} \longrightarrow 2\text{Cu}(\text{CN}) + (\text{CN})_2$</p> <p>$(\text{CN})_2 + 2\text{OH}^- \longrightarrow \text{CN}^- + \text{OCN}^- + \text{H}_2\text{O}$</p> <p>$\text{CaCN} + \text{H}_2\text{O} \longrightarrow \text{CaO} + \text{H}_2\text{NCN}$</p> <p>1</p> <p>1</p>
--	--

Fig. 25: Experimental result: Sample 3.

