

Recognition and Grouping of Handwritten Text in Diagrams and Equations

Michael Shilman, Paul Viola, and Kumar Chellapilla
Microsoft Research, Redmond, Washington (USA)
{shilman,viola,kumarc}@microsoft.com

Abstract

We present a framework for grouping and recognition of characters and symbols in online free-form ink expressions. The approach is completely spatial; it does not require any ordering on the strokes. It also does not place any constraints on the layout of the symbols. Initially each of the strokes on the page is linked in a proximity graph. A discriminative recognizer is used to classify connected sub-graphs as either making up one of the known symbols or perhaps as an invalid combination of strokes (e.g. including strokes from two different symbols). This recognizer operates on the rendered image of the strokes plus stroke features such as curvature and endpoints. A small subset of very efficient image features is selected, yielding an extremely fast recognizer. Dynamic programming over connected subsets of the proximity graph is used to simultaneously find the optimal grouping and recognition of all the strokes on the page. Experiments demonstrate that the system can achieve 94% grouping/recognition accuracy on a test dataset containing symbols from 25 writers held out from the training process.

Keywords: symbol recognition, handwriting, segmentation, mathematics recognition

1 Introduction

Handwritten text recognition is a maturing technology that has spawned many software products. In these systems the user writes words in a structured fashion, either along a line or in an “input region”. The recognition system can then process the entire line of text using dynamic programming to find the optimal recognition *and* grouping of the strokes. When freed from the rigid “input region” requirement, users frequently generate free form handwritten notes which include handwritten text, diagrams, and annotation. These notes require significant initial processing in order to group the strokes into “lines” of text which can then be passed to the recognizer (see for example [10]). The grouping process is inherently difficult, and the best perfor-

mance is achieved for simple paragraph structures in which there are a number of longer lines physically separated from drawing and annotations. While the grouping process could be integrated with the recognition process, the complexity of connected cursive recognition favors the two step process in which grouping precedes recognition. While it is likely that an integration of grouping and recognition would yield better results, this remains an open problem.

There are however a number of ink recognition problems which provide few constraints on the high-level layout of the page (Figure 1). One example is mathematical equation recognition, which incorporates many types of geometric layouts and symbols. Other examples include chemical structures, editing marks, musical notes, and so on. These scenarios are particularly important to pen computing because they exploit the flexibility of a pen to quickly express spatial arrangements, which is something that is currently difficult using a mouse and keyboard alone.

Therefore, we pose the problem of a system that performs integrated grouping and recognition of symbols over a page of handwritten ink. The system should not constrain writing order, because it is common to add extra strokes to correct characters after the fact. It should not make strict assumptions about the layout of the page. It should also scale to large pages of ink such as freeform notes, which can contain thousands of strokes in some cases.

Layout and timing-insensitive character recognition and grouping is not an easy problem. Symbol recognition is a well-known problem, for which many methods have been proposed [3]. The handwriting recognition community has developed countless techniques for optimizing grouping and recognition over a fixed spatial or temporal order, and for recognizing isolated characters [8, 12]. The closest related systems are those that deal with the processing of mathematical expressions [2, 4, 6, 7, 11]. Unlike some of these systems, we are trying to solve the problem in a way that does not require time ordering of strokes, does not require a linear organization of strokes on the page, and deals in a principled fashion with symbols that contain multiple strokes, some of which can be interpreted in isolation as another symbol.

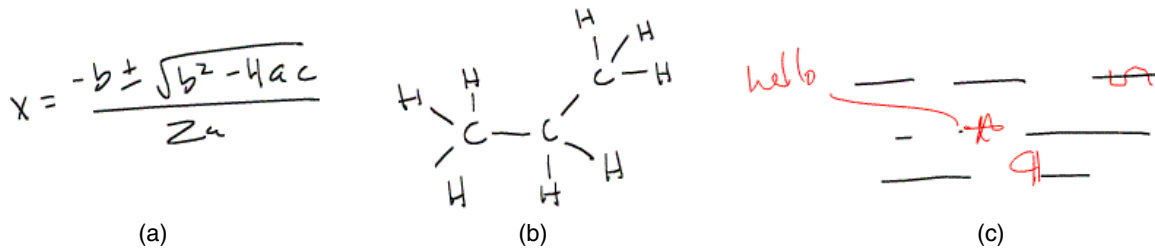


Figure 1. The pen is a particularly useful input device when layout is unconstrained, as is the case in (a) mathematics, (b) chemistry, and (c) document annotations.

2 Optimized Recognition and Grouping

In this paper, we present an efficient, purely spatial approach to simultaneously group and recognize handwritten symbols on a page. Our approach is an optimization over a large space of possible groupings in which each grouping is evaluated by a recognizer (Figure 2). This is in contrast to approaches where grouping and recognition are carried out as separate steps (e.g. systems with a separate layout analysis step).

In this approach the recognizer carries the burden of distinguishing good groupings from bad groupings and also must assign correct labels to good groupings. This sort of recognizer must evaluate quickly in order to process the large number of possible stroke groupings for a page of ink in a reasonable time.

Given such a recognizer, there are several benefits to this factoring of the problem. Improving the accuracy or performance of the system is simply a function of improving the accuracy or performance of the recognizer. Introducing new features to the system, such as rotation- or scale-invariance is simply a matter of changing the recognizer, rather than changing both the recognizer and the layout analysis. Perhaps most significantly, it enables our system to be nearly entirely learned from examples rather than relying on hand-coded heuristics. This last point bears repeating: ours is a monolithic system which once developed, requires no hand-constructed geometric features. All thresholds and parameters are learned automatically from a training set of examples.

Our system operates in the following manner. As a pre-processing step, it first builds a neighborhood graph of the ink in which nodes correspond to strokes, and edges are added when strokes are in close proximity to one another. Given this graph, we iterate efficiently over connected sets of nodes in the graph using dynamic programming and fast hashing on collections of nodes. For each set of nodes of up to size K , we perform a discriminative recognition on the set. This allows us to incorporate non-local information that rules out spurious answers that might result from

a generative model. We use dynamic programming to optimize over the space of possible explanations. The resulting system achieves high accuracy rates without any language model, places no stroke ordering requirements on the user, and places no constraints on the way in which symbols must be laid out on the page.

2.1 Optimization

Given a page of ink, we wish to minimize a global cost function:

$$C(\{V_i\}) = \Phi(R(V_0), R(V_1), \dots, R(V_n)) \quad (1)$$

In Equation 1, each V_i is a subset of the vertices which form a partition of the page, R is the best recognition result for that set of vertices, the function Φ is a combination cost (such as *sum*, *max*, or *average*), and C represents the overall cost of a particular grouping $\{V_i\}$.

To implement this optimization efficiently, we need a way to iterate over valid sets V_i (*graph iteration*), an efficient and accurate symbol recognizer R (*recognition cost*), a cost function to combine the cost of two subgraphs Φ (*combination cost*), and a way to reuse computation (*dynamic programming*).

2.2 Graph Iteration

In order to constrain the set of possible groupings, we say that a grouping is only valid if the strokes in that grouping are in close proximity to one another.

Thus, from a page of ink we construct a neighborhood graph $G = (V, E)$ in which the vertices V correspond to strokes, and edges E correspond to neighbor relationships between strokes, as shown in Figure 2b. We use the terms *strokes* and *vertices* interchangeably.

In our system, vertices are neighbors if the minimum distance between the convex hulls of their strokes is less than a threshold. However, we expect that any reasonable proximity measure would generate similar recognition results

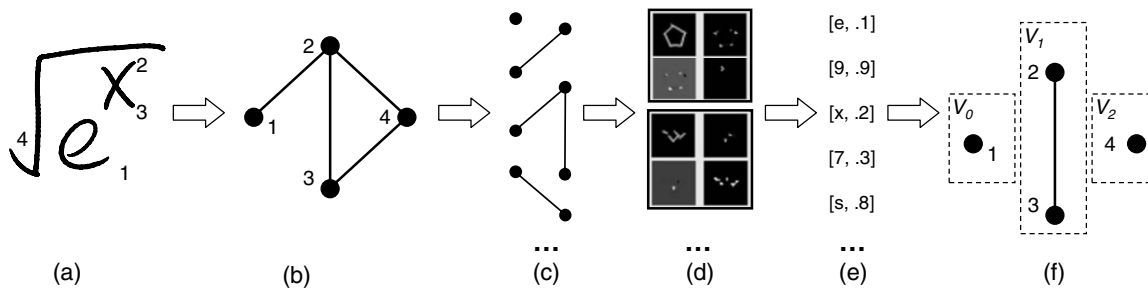


Figure 2. An overview of our approach. (a) A user sketch containing several strokes. (b) A neighborhood graph of the strokes in the sketch. (c) Connected subsets of the neighborhood graph of up to a fixed size K . (d) Rendered images of the subsets that are passed to an AdaBoost recognizer. (e) Results from the recognizer include a symbol hypothesis and a score. (f) An optimization partitions the graph to jointly maximize the recognizer scores.

as long as the neighborhood graph contains edges between strokes in the same symbol.

Given this neighborhood graph, we wish to enumerate all connected subsets of the nodes V_i in V where $|V_i| \leq K$. Each subset V_i becomes a symbol candidate for the recognizer.

To our knowledge, there is no efficient way to enumerate subsets of up to size K without duplicating subsets. We iterate by first enumerating all subsets of size 1. We then expand each subset by all of the edges on its horizon, eliminate duplicates, expand again, and so on, up through size K . This eliminates the propagation of duplicates through each round.

The subsets V_i that are generated for the graph in Figure 2b are $\{1\}, \{2\}, \{3\}, \{4\}, \{1, 2\}, \{2, 3\}, \{2, 4\}, \{3, 4\}, \{1, 2, 3\}, \{1, 2, 4\}, \{2, 3, 4\}, \{1, 2, 3, 4\}$.

2.3 Recognition Cost

The second implementation detail of the optimization process is the recognition cost R (see Section 3 for a detailed discussion). The simple requirements are that R should return relatively low costs for subsets of the graph V_i that correspond to symbols (such as in Figure 3a). Similarly, R should return relatively high costs for subsets of the graph that do not correspond to symbols (such as in Figure 3b,c,d).

In fact, this is not as easy as it sounds. Many of the subsets that are passed to the recognizer are *invalid*, either containing strokes from multiple characters or do not contain all the strokes of a multi-stroke symbol. We call such subgraphs garbage. While some of the garbage doesn't look like any symbol in the training set, some invalid examples are indistinguishable from training samples without the use of context. For example a single stroke of an X can be easily interpreted in isolation as a back-slash (Figure 3c).

Therefore we also pass the context $X(V_i, E)$ into the recognizer to help it spot garbage. We define the context to be the set of nodes in $V - V_i$ that are connected to V_i in E , and show an example in Figure 3d.

2.4 Combination Cost

The third implementation detail of the optimization is the combination cost, $\Phi(c_1, c_2)$. The combination cost is a function of the costs of the two subsets of the graph. We considered several alternative costs:

- **Sum.** $\Phi(c_1, c_2) = c_1 + c_2 + \varepsilon$. The sum of the costs makes intuitive sense: if the costs are negative log likelihoods then the sum corresponds to a product of probabilities. The ε penalty can be used to control over/under grouping (higher values of ε force grouping into fewer symbols).
- **Max.** $\Phi(c_1, c_2) = \text{Max}(c_1, c_2)$. This function penalizes the worst hypothesis in the set.
- **Average.** $\Phi(c_1, c_2) = (c_1 + \omega c_2)/(1 + \omega)$. This function averages the scores across all of the symbols in the hypothesis. ω is a weight corresponding to the number of symbols in the best interpretation for $V - V_i$.

2.5 Dynamic Programming

Finally, Because the function we wish to optimize cleanly partitions the graph into a combination of $R(V_i)$ and $C(V - V_i)$, we are able to use dynamic programming to avoid redundant computation. In other words, if we have already computed C for a subset of strokes in the neighborhood graph, we can reuse the result by looking it up in a hash table. We hash on sets of strokes by XOR'ing stroke ID's.

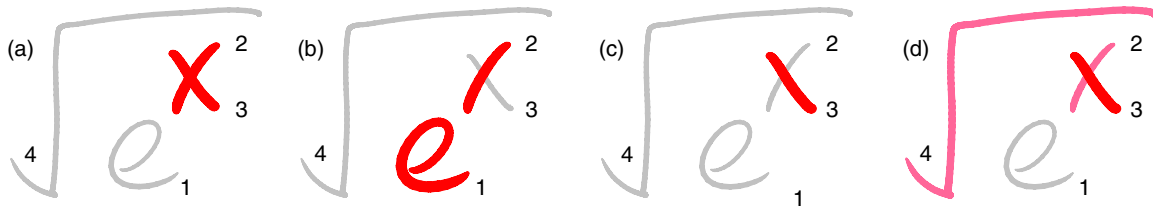


Figure 3. A set of inputs to the recognizer. (a) A full symbol that is passed as a candidate to the recognizer. (b) Overgrouped garbage that is passed as a candidate to the recognizer. (c) Garbage that is ambiguous with a back-slash when passed to the recognizer without context. (d) Neighborhood context that makes the stroke in (c) unambiguously garbage.

3 AdaBoost Symbol Recognizer

The recognizer utilized in the dynamic programming system described above is based on a novel application of AdaBoost [5].

The basic framework used is most closely related to the work of Viola and Jones, who constructed a real-time face detection system using a boosted collection of simple and efficient features [13]. We chose this approach both because of its speed and because it is easily extensible to include additional feature information.

We have generalized Viola-Jones in two ways: our classification feature information. We have generalized Viola-Jones in two ways: our classification problem is *multi-class*, and we have added additional input features to the image map. These additional features are computed directly from the on-line stroke information and include curvature, orientation, and end-point information.

The input to the recognition system is a collection of images. The two principle images are the *candidate image* and the *context image*. The candidate image is quite conventional, the strokes of the current candidate sub-graph are rendered into an image which is 29x29 pixels. The geometry of the strokes is normalized so that they fit within the central 18x18 pixel region of the image. Strokes are rendered in black on white with anti-aliasing. The context image is rendered from the strokes which are connected to some candidate stroke in the proximity graph.

3.1 Additional Feature Images

Each of the principle images are augmented with additional stroke feature images. This is much like the earlier work on AMAP [1]. The first additional image records the curvature at each point along each stroke. The angle between the tangents is a signed quantity that depends on the direction of the stroke, which is undesirable. The absolute value of this angle provides direction invariant curvature information.

Two additional feature images measure orientation of the stroke. Orientation is a difficult issue in image processing, since it is naturally embedded on a circle (and hence 2π is identical to 0). We have chosen to represent orientation in terms of the normal vector (perpendicular vector) to the stroke (which is measured from the same nearby points used to measure curvature). The two components of the normal are represented as two images the **normalX** image, and the **normalY** image (by convention the normal has a positive dot product with the previous tangent).

The final additional feature image contains only the end-points of the strokes, rather than the entire stroke. This measure can be useful in distinguishing two characters which have much ink in common, but have a different start and end point, or a different number of strokes (for example '8' and '3').

3.2 The Viola-Jones Filters

A very large set of simple linear functions are computed from the input images define above. The form of these linear functions was proposed by Viola and Jones, who call them 'rectangle filters'. Each can be evaluated extremely rapidly at any scale (see Figure 5). The filters measure the differences between region averages at various scales, orientations, and aspect ratios. The rigid form of these features arises from the fact that each can be computed extremely rapidly, in 6 or fewer add/multiplies.

For these experiments a set of one and two rectangle filters were constructed combinatorially. A set of filters of varying location, size, aspect ratio, and location are generated. The set is not exhaustive; some effort is made to minimize overlap between the filters, resulting in 5280 filters. Such a large set is clearly overcomplete in that requires only 841 linear filters to reconstruct the original 29 by 29 image. Nevertheless this overcomplete basis is very useful for learning. Each filter can be evaluated for each of the 10 feature images, yielding a set of 52,800 filter values for each training example. Clearly some approach for selecting a critical subset of these will improve performance.

3.3 AdaBoost Feature Selection and Learning

The above sections describe a processing pipeline for training data: a rendering process for candidate and context, a set of additional feature images, and set of rectangle filters. The machine learning problem is to generate a classifier for this data which correctly determines the correct symbol of the candidate strokes, or possibly that the set of strokes is garbage. We use AdaBoost to learn a classifier which selects a small set of rectangle filters and combines them.

For these experiments the “weak learner” is a classifier which computes a single rectangle filter and applies a threshold (this is a type of decision tree known as a decision stump). In each round of boosting the single best stump is selected, and then the examples are reweighted. We use the multi-class variant of confidence rated boosting algorithm proposed by Schapire and Singer [9].

After N rounds, the final classifier contains N weak classifiers. Since each weak classifier depends on a single rectangle filter only N filters need to be evaluated. Excellent performance is achieved with between 75 and 200 filters. On a training set of 3800 examples from 25 writers, 0 training errors is observed with 165 weak classifiers. On a test set of 3800 examples from a different set of 25 writers 96% of the characters were classified correctly.

4 Evaluation

To evaluate our approach, we ran tests on a corpus of automatically-generated mathematical expressions. We collected a modest set of handwritten characters, digits, and mathematical operators from 50 users with 5 examples per class. Of these examples, we synthesized short expressions containing digits and operators with a generative grammar (Figure 6). Our generated expressions are intentionally



Figure 4. The left most images in each row are the candidate and context images rendered at 29x29 pixels. The remainder of each row shows various feature images computed from these principle images. From left to right: stroke curvature, stroke normal X, stroke normal Y, and endpoints.

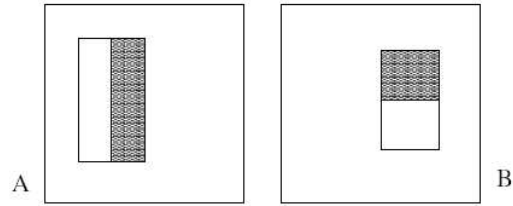


Figure 5. Example rectangle filters shown relative to the enclosing classification window. The sum of the pixels which lie within the white rectangles are subtracted from the sum of pixels in the grey rectangles. Rectangle filters which contain two rectangles are shown in (A) and (B).

dense, in order to make the grouping problem more interesting. Also it is worth noting that although each of our test examples is horizontally-oriented, our technique applies independent of the layout. We have manually applied the technique to examples with more interesting layouts and show that it works in practice, although our test data does not reflect this condition.

We separated the generated expressions into training and test data, such that 25 users’ data made up the training set and the other 25 users made up the test set. This split ensures that we are testing the generalization of the recognizer across different populations.

We applied the above system to the test data with three different combination cost functions: *sum*, *max*, and *avg*, as described in Section 2.4. For *sum* we varied the value of ϵ to see its effect on the overall accuracy. For all of these approaches we measured the total number of symbols in the test data, the total number of false positives and false negatives in the results. A false negative occurs any time there is a group of strokes with a specific symbol label in the test data, and that exact group/label does not occur in the test

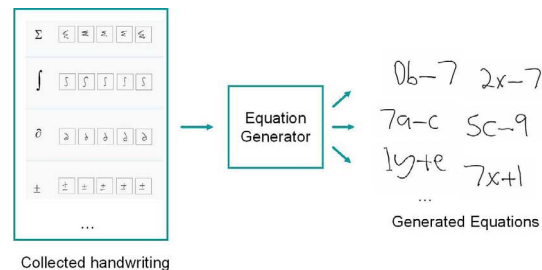


Figure 6. Collected truth data is rendered into two sets of mathematical expressions, which serve as training and test data, respectively.

data. A false positive is the converse. Our best results are 94% accuracy for grouping and recognition for the *avg* combination cost. The full results are shown in Table 1.

Cost	False Pos	False Neg	Total
Sum ($\epsilon = -.2$)	299	243	3840
Sum ($\epsilon = -.25$)	308	248	3840
Max	267	243	3840
Avg	225	202	3840

Table 1. Results on 3840 symbols in the context of generated mathematical expressions.

5 Conclusion

This paper presents an integrated grouping and recognition system of on-line freeform ink. Grouping is a requirement for recognition in such tasks because each symbol may have a number of strokes. Simple heuristics that group intersecting strokes may work in some cases. In domains which include multi-stroke symbols such as '=' (equals) or ' π ' (pi), these heuristics fail. Conversely, it is not uncommon to see strokes from different characters come very close to or intersect each other.

This integrated system first constructs a proximity graph which links pairs of strokes if they are sufficiently close together. The system then enumerates all possible connected subgraphs looking for those that represent valid characters. The notion of proximity is defined so that strokes from the same symbol are always connected. This definition of proximity will necessarily link strokes from neighboring symbols as well. These connected subgraphs are not interpretable as a valid symbol, and will be discarded as garbage. Note, a garbage subgraph can also arise if a symbol is undergrouped: e.g. only one of the strokes in a multi-stroke character is included. A fast recognizer based on AdaBoost is trained to recognize all symbol classes as well as a unique class called garbage, which includes subgraphs of strokes that are not valid. In order to address the undergrouping problem, the recognizer operates both on the current candidate strokes as well as the context of the surround strokes.

Dynamic programming is used to search for the minimum cost decomposition of the initial proximity graph into connected subgraphs, each of which can be interpreted as a valid symbol. The set of all possible connected subgraphs is efficiently enumerated using an incremental hashing scheme which grows subgraphs one node at a time and efficiently removes duplicates.

The recognizer is trained on symbols which come from 25 writers. The final system achieves a 94% simultaneous

grouping and recognition rate on test data from 25 different users which was not used during training.

6 References

1. Y. Bengio and Y. LeCun, "Word Normalization for On-line Handwritten Word Recognition," in Proc. of the International Conference on Pattern Recognition, (IAPR, ed.), (Jerusalem), pp. 409-413, Oct. 1994. 5 pages.
2. D. Blostein and A. Grbavec, "Recognition of Mathematical Notation," in Handbook of Character Recognition and Document Image Analysis, Eds. H. Bunke and P. Wang, World Scientific c, 1997, pp. 557-582.
3. A. Chhabra, "Graphic Symbol Recognition: An Overview." In Proceedings of Second International Workshop on Graphics Recognition, Nancy (France), pages 244-252, August 1997.
4. K. Chan and D. Yeung, "Mathematical Expression Recognition: A Survey," Int'l J. Document Analysis and Recognition, vol. 3, no. 1, pp. 3-15, Aug. 2000
5. Y. Freund, R. Schapire. "Experiments with a New Boosting Algorithm." ICML 1996: 148-156
6. N. Matsakis, "Recognition of Handwritten Mathematical Expressions", Massachusetts Institute of Technology, Cambridge, MA", May 1999
7. E.G. Miller and P.A. Viola, "Ambiguity and Constraint in Mathematical Expression Recognition," Proc. 15th Nat'l Conf. Artificial Intelligence, pp. 784-791, 1998
8. R. Plamondon, S. Srihari. "On-Line and Off-Line Handwriting Recognition: A Comprehensive Survey." IEEE Transactions on Pattern Analysis and Machine Intelligence 22(1): 63-84 (2000)
9. R. Schapire, Y. Singer. "Improved Boosting Algorithms using Confidence-Rated Predictions." COLT 1998: 80-91
10. M. Shilman, Z. Wei, S. Raghupathy, P. Simard, D. Jones: "Discerning Structure from Freeform Handwritten Notes." ICDAR 2003: 60-65
11. S. Smithies, K. Novins, and J. Arvo, "A Handwriting-Based Equation Editor," *Graphics Interface '99*, June 1999.
12. C. Tappert, C. Suen, T. Wakahara. "The State of the Art in Online Handwriting Recognition." IEEE Transactions on Pattern Analysis and Machine Intelligence 12(8): 787-808 (1990)
13. P. Viola, M. Jones: Robust Real-Time Face Detection. ICCV 2001: 747