

## 21:14 Counting words with a state machine.

by Robert Graham

In this paper we implement `wc`, the classic Unix word count program, using an asynchronous state machine parser. We implement this twice: first a simplified version supporting ASCII, then a more complete program supporting Unicode UTF-8 encoding. We implement this algorithm in both C and JavaScript. Even the latter is significantly faster than the standard versions of `wc`, such as the GNU Coreutils `wc` that comes with Linux.

### Introduction

A parser is software that translates external data into some internal data structures.

At university, they teach you abstract and formal parsers, often in a class that builds a compiler. However, little of that theory is used elsewhere in your coursework. In your networking class, the code they teach uses concrete and ad-hoc parsers, discarding everything you learned in your parser class. While parser theory they teach you is useful, even academics struggle to use it in practice.

In this paper, we do a mix of theory and practice. On one hand, we look at abstract theory of state-machines and deterministic/non-deterministic finite automata. On the other hand, we build the state-machine by hand, banging the bytes together.

The reason this concept is important is demonstrated by Nginx replacing Apache as the dominant web server on the Internet. Apache parses input the legacy way they taught you in networking class. The newer Nginx parses input using a state-machine. This parsing is more scalable, allowing much higher loads on the web server.<sup>39</sup>

In this paper, we demonstrate state-machines by re-implementing the classic Unix command-line program `wc`. Over the last year, it has been popular for proponents of various languages to re-implement `wc` in order to show that their favorite language can compete with C in performance. In this case, we do this to demonstrate our favorite algorithm is better than existing algorithms, implementing it in two different languages.

<sup>39</sup>This overstates the importance of just the parsing. Nginx scales better than Apache for a lot of reasons. However, these reasons are all interconnected: if you write an asynchronous server, then state-machine parsers are a much better way of parsing the requests.

These re-implementations are usually incomplete, only parsing ASCII. In this paper, we do a more complete version, correctly parsing UTF-8.

The intent of this article isn't that you should go and parse everything with state-machines. It puts a burden on future programmers trying to read the code, most of whom are unfamiliar with the technique. On the other hand, when performance and scalability are needed, state-machines are a good choice. You probably wouldn't want to use them for `wc` in the real world, as the program doesn't need to be especially fast. We choose `wc` in this paper only because it's a popular benchmark target, the simple thing that more complex endeavors are compared against.

### What is WC?

This command-line utility has been part of Unix since time began on the first of January, 1970. As defined in the POSIX standard, it counts the number of lines, words and characters, when the corresponding flag of `-l`, `-w`, and `-c` is set. If no parameters are set, then the default is all three, `-lwc`.

```
2 $ echo "basic input/output" | wc
   1      2     19
```

We see here that the program has reported one line, two words, and 19 characters. Words are counted by the number of strings of non-spaces separated by spaces. Thus, this example is only two words, not three.

Modern character encodings can use multiple bytes per character, such as UTF-8 or various character sets for Chinese, Japanese, and Korean. In such cases, the `-m` parameter replaces the `-c` parameter, counting the number of multi-byte characters instead of the number of bytes. As we see in these two examples, changing from `-c` to `-m` changes the character count:

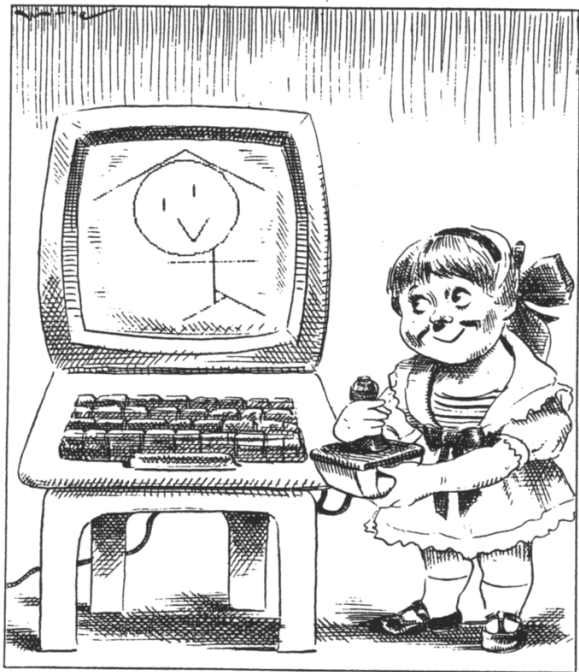
```
$ echo わたしは にほんごがすこししか はなせません | wc -lwc
   1      3     67
$ echo わたしは にほんごがすこししか はなせません | wc -lwm
   1      3     23
```

## How do they implement WC?

There are many versions of the program, such as GNU's Coreutils for Linux, BusyBox, macOS, FreeBSD, OpenBSD, QNX and SunOS. Most implementations count words by counting the number of times a space is followed by a non-space. Think of it as an edge-triggered condition, going from space to non-space. As we'll soon see, this can also be treated as a state-machine with two states.

Parsing words is easy, the hard part is character-sets. We could hard-code ASCII values into our program, such as `0x20` for space and `0x0A` for newline, but this wouldn't work for non-ASCII systems. IBM mainframes that use the EBCDIC character-set will represent a space using `0x40`.

Thus, instead of using hard-coded values these programs use the standard `isspace()` function to test if a character is a space. Recently, many people have re-implemented `wc` in their favorite language to show that they can be just as fast as C. In fact, most of the processing time is spent in the `isspace()` function, so all they really proved is that hard-coded constants like `0x20` in other languages are faster than `isspace()` function calls in C.



<sup>40</sup>The 'r' in `mbrtowc()` means "re-entrant." If parsing at the end of a fragment, it saves state before resuming at the start of the next fragment.

The problem is worse for multi-byte character-sets like UTF-8. The program must first parse multiple bytes into wide characters using functions like `mbtowc()` (or `mbrtowc()`),<sup>40</sup> then test if they are a space with `iswspace()`. Re-implementations often do only ASCII. This paper includes two re-implementations, the first for ASCII, the second for UTF-8.

## How do we implement it?

Our first version supporting ASCII is shown in Figure 17. In the GitHub project accompanying this article, the program is `wc2o.c`, where the 'o' stands for "obfuscated C version." This program is pretty darn opaque when trying to figure out how it counts words. On the other hand, it exposes the idea of state-machine parsing.

Line 5 declares the state-machine table consisting of four states. Each state is a row of three transitions. (Table 1)

Line 7 declares a table that will translate bytes. All 256 ASCII values translate into one of three possible values: `word(0)`, `space(1)`, and `newline(2)`. Specifically, the character `0x0A` or '`\n`' translates to `newline(2)`, and the characters '`\b\t\m\v\f`' translate to `space(1)`. All other values translate to `word(0)`. The reason we include this translation step is that so that the state-machine on line 5 is  $4 \times 3$  states rather than  $4 \times 256$  states. In our final version, we don't do this translation, and just have large state-machines instead.

Line 15 loops getting the next byte of input, one byte at a time. Calling `getchar()` here for every character is potentially expensive, but we aren't benchmarking this program, just showing the algorithm. In our final version, we read input a buffer at a time instead of a byte at a time.

Line 16 does the state transition, in other words, it parses the input. We translate the byte into one of the three column values, `word(0)`, `space(1)`, or `newline(2)`. We look up that in the current row, then set the next row according to the transition. Thus how in the `was-space(0)` state, if we receive a `non-space(0)` character, we transition to `new-word(2)` state.

Line 17 processes what we parsed. In our case, the processing is trivialized to just counting the number of times we visit each state.

```

#include <stdio.h>
2 int main(void)
{
4     static const unsigned char table[4][3] = {
        {2,0,1}, {2,0,1}, {3,0,1}, {3,0,1}
6     };
    static const unsigned char column[256] = {
8         0,0,0,0,0,0,0,0,0,1,2,1,1,1,0,0,0,
        0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,1,0,
10    };

12    int state = 0;
    int c;

14    while ((c = getchar()) != EOF) {
16        state = table[state][column[c]];
        counts[state]++;
18    }
    printf("%lu %lu %lu\n", counts[1], counts[2],
20        counts[0] + counts[1] + counts[2] + counts[3]);
    return 0;
22 }

```

Figure 17: wc2o.c, an obfuscated word counter for ASCII.

	word(0)	space(1)	newline(2)
was-space(0)	new-word(2)	was-space(0)	new-line(1)
new-line(1)	new-word(2)	was-space(0)	new-line(1)
new-word(2)	was-word(3)	was-space(0)	new-line(1)
was-word(3)	was-word(3)	was-space(0)	new-line(1)

Table 1: Simple word-count state machine.

**NEW FROM XITEX**


## \$95 MORSE TRANSCIVER

**SEND:**

- 1 to 150 WPM (set from terminal)
- 32 character FIFO buffer with editing
- Auto Space on word boundaries
- Grid/Cathode key output
- LED Readout for WPM and Buffer space remaining

**SERIAL INTERFACE:**

- ASCII (110, 300, 600, 1200) or Baudot (45, 50, 57, 74) compatible
- Simplex Hi V Loop or T'L electrical interface
- Interfaces directly with the XITEX® SCT-100 Video Terminal Board; Teletypes® Models 15, 28, 33, etc.; or the equivalent



**MRS-100 CONFIGURATIONS:**

- \$95 Partial Kit (includes Microcomputer components and circuit boards; less box and analog components)
- \$225 Complete Kit (includes box, power supply, and all other components)
- \$295 Assembled and tested unit (as shown)

Overseas Orders and dealer inquiries welcome

**COPY:**

- 1 to 150 WPM with Auto-Sync.
- Continuously computes and displays Copy WPM
- 80 HZ Bandpass filter
- Re-keyed Sidetone Osc. with on-board speaker
- Fully compensating to copy any 'fist style'

See your local dealer or contact XITEX® direct.

MC/Visa accepted

**XITEX CORP**  
13628 Neutron • P. O. Box 402110  
Dallas, Texas 75240 • (214) 386-3859



Scalar Unicode Value	First Byte	Second	Third	Fourth
00000000 00000000 0xxxxxxx	0xxxxxxx			
00000000 00000yyy yyxxxxxx	110yyyyy	10xxxxxx		
00000000 zzzzyyyy yyxxxxxx	1110zzzz	10yyyyyy	10xxxxxx	
000uuuuu zzzzyyyy yyxxxxxx	11110uuu	10uuzzzz	10yyyyyy	10xxxxxx

Table 2: UTF-8 Bit Distribution, Unicode 6.0

As we receive the bytes of an HTTP request, we enter the `method` state the first time we receive a non-space character. We remain in the `method` state until we receive a space character, at which point we transition to the `space1` state.

In C, we might process each byte of input with a function like the following `switch/case` logic:

```

1 int http_parse(int state, unsigned char c,
2               ...) {
3     switch (state) {
4         ...
5     case METHOD: /*GET, POST, HEAD, ...*/
6         if (c == '\n') {
7             ...
8             return EOL;
9         } else if (isspace(c)) {
10            ...
11            return SPACE1;
12        } else {
13            ...
14            return state; /* no change in state */
15        }
16        ...
17    }
18 }


```


Most major web servers that aren't Apache use this method. Nginx calls this state `sw_method`, which you can see in the open-source online.<sup>42</sup>

You can test on a live network whether a web server is parsing requests using a state-machine. Send a request to the server consisting of GET, followed by five billion spaces and only then the rest of the request. If the server acts like Apache buffering a complete header, then it'll run out of buffer space. If instead the server acts like Nginx and parses input with a state machine, it'll happily keep reading spaces as long as it's in that state. If the connection terminates prematurely, it'll be because of a timeout instead of running out of buffers. (It takes a while to send five gigabytes.)

<sup>42</sup>See near line 159 of `ngx_http_parse.c`.

This example uses a `switch/case` block of code to handle the transitions. In our state-machines for counting words, we use a lookup table instead. A third choice is to use a mixture. The program `masscan`, for example, does a lot of parsing of such protocols like FTP, SMTP, X.509, and so. It uses a mixture of `switch` statements and lookup tables.






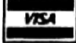
**(ALMOST) FREE CLUES:**

If you've ever been stuck in an adventure game, you need **The Book of Adventure Games** by Kim Schuette. This fantastic book contains complete legible (typeset) maps, magnificent illustrations, and all the hints you need to complete 77 of the all-time most popular adventure games including **Zork I, II, III, Deadline, Starcross, Witness, Planetfall, Enchanter, Sorcerer, Infidel, Suspended (with map), Wizardry, Knight of Diamonds, Legacy of Lylgamyn, All Scott Adams, All Sierra On-Line** including **Time Zone, Ultima I, II, III** and many more! Best of all, the book doesn't spoil your fun! At about 25¢ an adventure, it's the biggest bargain around. So stop getting ripped off by \$10 cluebooks and call:

**1-(800)-821-5226 Ext. 500**  
24 hrs. a day, 7 days a week  
or write:  
Witt's End  
42 Morehouse Rd., Dept. 6  
Easton, CT 06612

*Free UPS shipping. Add \$3.00 for C.O.D. APO's FPO's o.k. Add \$5.00 for foreign shipping. No charge for credit cards. We accept Visa: Mastercard, Personal Check (allow 2 weeks to clear), Certified Check or money order.*


All Trademarks are acknowledged.

**ELIZABETH GRANT**  
HIGH CLASS MILLINERY

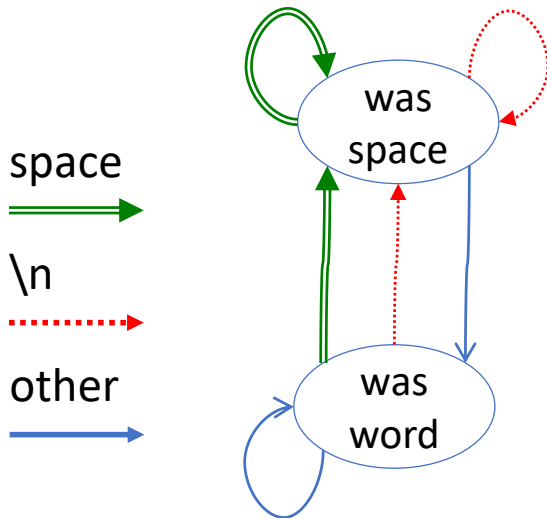
12 WEST STREET (Over Bigelow-Kennard's)

Smart Tailored and Dress Hats. Made of fine materials and of the best workmanship. Exclusive styles. No two hats alike. Courteous attention whether you buy or not.  
PRICES, SIX DOLLARS AND UP.

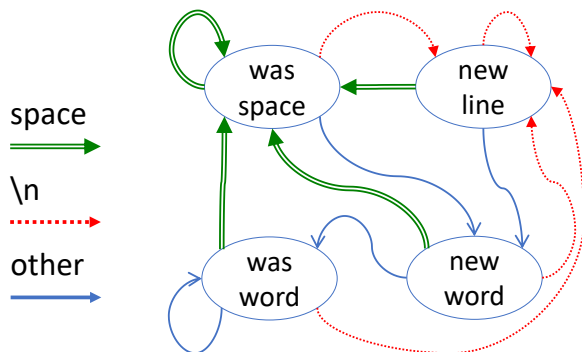


## How can we construct a state-machine for word-counting?

Most implementations of `wc` effectively use a machine with two states which can be represented with the following diagram. Note that they aren't designed explicitly as a state-machine, but that's effectively how the code works.



In the `wc2o.c` program, we changed this to a machine with four states. This is the table with three types of transitions and four states:



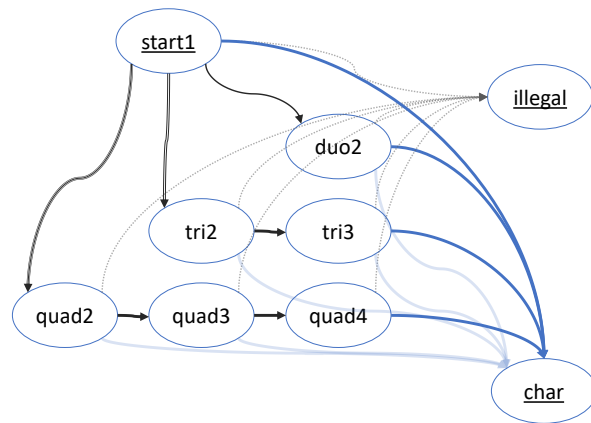
We did this in order to be overly clever in how we were going to process the data.

Remember, there's two things going on here. One step is parsing the input. The next step is processing the results. Thus, we first need to parse out things like words, characters, and newlines. Then we need to process this information, which for word-counting, is done by counting each time we enter a state. We've cleverly collapsed the processing into a simple operation.

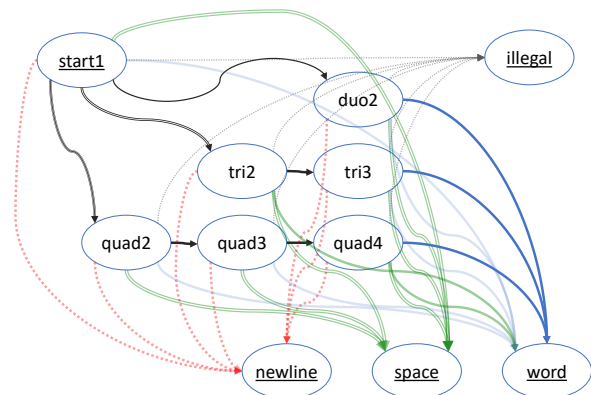
SAM COUPE AND SPECTRUM MAGAZINE!  
 PROGRAMS, UTILITIES, INFO, IDEAS! NEWS, REVIEWS AND HOMEGROWN SOFTWARE MONTHLY SINCE 1987!  
**"OUTLET"**  
 GRAPHICS, AND HELP PAGES, SERIOUS SOFTWARE  
 SPECIAL OFFER! Latest issue £2.50 to newcomers on:-  
 +3, DISCIPLE/+D, MICRODRIVE, OPUS, TAPE, SAM DISC  
 CHEZRON SOFTWARE, 605 LOUGHBOROUGH RD., BIRSTALL, LEICESTER LE4 4NJ

The lesson here isn't that parsers can completely trivialize processing as we've done here, but instead that we often add artificial states to benefit later processing.

Now let's talk about UTF-8. Using the original table as a guide, we might construct a state-machine for parsing 1-byte sequences, 2-byte sequences called a "duo," 3-byte sequences "tri," and 4-byte sequences "quad."



However, our needs are simpler. We don't need to parse out the code point and test with `iswspace()` but can instead include that functionality within the state-machine itself, where the output is one of four values: word, space, newline, or illegal.



There are, in fact, more states than just this. Instead of a simple path for 3-byte characters, we must add additional states that recognize 3-byte characters that result in spaces. This creates a table of roughly thirty states that's too complex to draw here.

Instead, here are snippets of the code that take an existing table and adds states for characters like U+1680 Ogham Space Mark. It clones existing states that follow the same path, but at the end marks the character as a `space` instead of a `word`:

```

/* clone existing states */
2 memcpy(table[TRI2_E1], table[TRI2],
        sizeof(table[0]));
4 memcpy(table[TRI3_E1_9a], table[TRI3],
        sizeof(table[0]));
6 /* link in new states */
table[0][0xE1] = TRI2_E1;
8 table[TRI2_E1][0x9a] = TRI3_E1_9a;
table[TRI3_E1_9a][0x80] = SPACE;

```

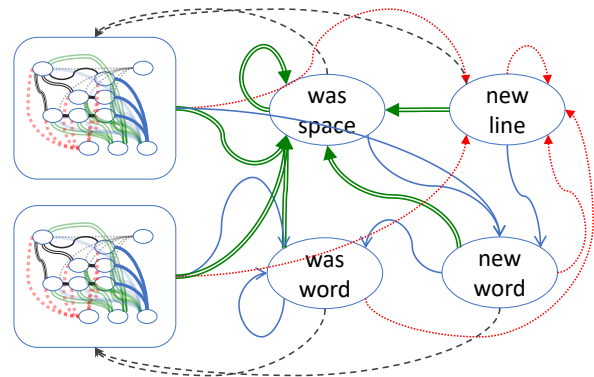
What this code is doing is exactly what generic regex code would do. All we are doing here is creating manually what regex libraries would do based upon expressions. What we are doing here is manual optimization for concepts that exist abstractly.

Now let's combine our UTF-8 state-machine parser with our word-count state-machine parser. There's two ways of doing this. The obvious way is to feed the output of one as input to the other. The other way is to combine the two into a single state-machine.

This is a multiplicative process. That means replicating one state-machine for every state in the other state-machine. Again, let's talk regex theory. There are two ways of representing such a thing. One way increases computation, what we call an NFA or non-deterministic finite automata, which is what would happen if we fed the output of one as the input to the next. The other way keeps computation the same but increases the size of the table. This is a DFA or deterministic finite automata. As you build complex regexes, you cause either computation to explode or memory to explode. In this case, we've chosen DFA, so memory explodes.

Thus, where one state-machine needs 35 states and the other just four, that means the combination may needs as many as  $4 \times 35 = 140$  states. However, we are going to do a small trick. The `was-space` and `new-line` states are clones of each other, as are `was-word` and `new-word`. Thus, we

only need to double rather than quadruple the UTF-8 state-machine. This produces something that may be represented like:



## The Final Code

The final code is in `wc2.c`. It's a few hundred lines so is not included in this article but is instead available on GitHub.<sup>43</sup>

The complicated part that takes hundreds of lines is where it builds that state-machine table. This results in a table roughly with 70 states (rows), and 256 columns, where each column represents the transition that will happen when a byte of input is received.

Once we've built the table, we simply process chunks of input analogous to the following. The actual code looks slightly different, with the inner loop separated into a `parse_chunk()` function.

```

1 unsigned counts[MAX_STATE];
  //Get the next chunk of input.
3 length = fread(buf, 1, sizeof(buf), fp);
  //For all bytes in that chunk,
5 for (i=0; i<length; i++) {
  //Get the next byte.
7   c = buf[i];
  //Do the state transition.
9   state = table[state][c];
  //Do the counting.
11  counts[state]++;
  }
13 //Report the results.
word_count = counts[NEW_WORD];
15 line_count = counts[LINE_COUNT];
char_count = counts[0] + counts[1]
17           + counts[2] + counts[3];

```

<sup>43</sup>[git clone https://github.com/robertdavidgraham/wc2](https://github.com/robertdavidgraham/wc2) || `unzip pocorgtfo21.pdf wc2.zip`







sized blocks, like 64k; you can't choose the size of blocks depending upon the parsed contents. When data is received, it is dispatched to the appropriate copy of the state for parsing each file. We just need an 8-bit integer for every file to hold all the per file parser state.

This would be a silly thing to do with files but is an important thing for networking services. Apache is broken and can't scale beyond 10,000 concurrent TCP connections because it struggles with 10,000 threads in the system. All its major competitors use a single thread (or single thread per CPU core) and handle things asynchronously.

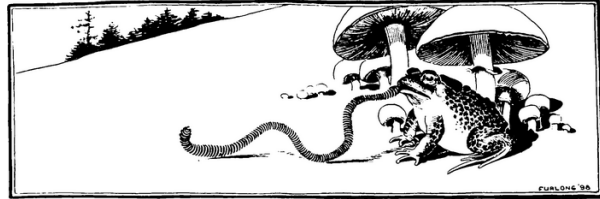
Again, the parser can't influence data reception. The network stack simply receives packets from the other side, whatever size of packets those might be. The parser must handle the case where it receives exactly as much data as it was expecting, or too much data, or not enough data.

## Conclusion

There have been many posts over the last year of people implementing `wc` in their favorite language, such as Haskell or Python. In this paper, instead of a different programming language, we've chosen a fundamentally different algorithm, that of an asynchronous state-machine parser. We've implemented the same algorithm in both C and JavaScript, to show that the speed is property of the algorithm instead of the language.

Instead of a simplified problem of just handling ASCII, we've demonstrated the algorithm using the difficult problem of UTF-8 encodings. Given something bizarre like Ogham text, we still produce the same answer as compliant `wc` programs. While only the UTF-8 encoding is implemented, the concept extends to any character-set, including the CJK (Chinese, Japanese, Korean) multi-byte character-sets.

Such state-machine parsers are costly in terms of code maintainability: most programmers are unfamiliar with them. However, they have clear advantages for writing scalable, secure code for modern Internet applications.



## PRO-LOG IS PROVIDING LOCAL STOCKING FOR THEIR STD MICROPROCESSOR BOARD PRODUCTS

Effective  
September 21, 1981

Pro-Log will begin a local distribution plan for their microprocessor board product line. Off the shelf delivery will be available through their stocking representative network in selected locations across the United States. The initial stocking locations will include:

<b>Northern California</b>	Manco 3350 Scott Blvd., Bldg. 55 Santa Clara, CA 95051
<b>Southern California</b>	Advanced Digital Group 15904 Strathern Sr., Suite 21 Van Nuys, CA 91406
<b>New York Area</b>	Tecnimat 500 Grand Avenue Englewood, NJ 07631
<b>New England States</b>	Martindale Associates 212 Main Street North Reading, MA 01864

Watch for additional stocking locations in the very near future!

108560A 11K 9/81