



I chose the New York Times crossword as the focus of the project because it is the best written and edited crossword, by general consensus among avid solvers. Consequently, the best crossword solving computer program should be able to solve the New York Times puzzle.

```
## MONDAY SEPTEMBER 5 2015 ##
## ACROSS
0 0 5 the top of a wave
0 6 4 heed a red light
0 11 4 tanginess
1 0 4 do-it-yourselfer's book genre
1 6 4 nurse duty with a hammer
1 11 4 part of the eye
2 0 8 Chris who sang "Wicked Game", 1991
2 6 4 early of rising Sun Records
2 11 4 word repeated before "pants on fire"
3 0 13 shoutline series named after an old fiction genre
3 6 4 proverbial mouse
4 7 4 "when all I failed, read the instructions"
4 12 3 young-sounding wildebeest
5 0 4 spyden's hair
5 5 4 -cola
5 10 5 cousins of ostriches
6 0 6 early of rising Sun Records
6 4 4 cheese off
```

Part of the program's input

(\*[rɪŋg ˈlɪt, ˈlɑːd])  
 (\*[rɪŋg ˈlɪquɪd, ˈkɒnɒl])  
 (\*[rɪŋg ˈlɪquɪd, ˈoʊl])  
 (\*[rɪŋg ˈmɛdɪəm, ˈkɒnɒl])  
 (\*[rɪŋg ˈmɛdɪəm, ˈdɛɛpɑːt])  
 (\*[rɪŋg ˈmɛdɪəm, ˈlɑːd])  
 (\*[rɪŋg ˈmɛss, ˈspɑːtɜːr])  
 (\*[rɪŋg ˈnɛd, ˈhɑːz])  
 (\*[rɪŋg ˈpæn ˈbɑːtɜːr, ˈspɑːtɜːr])  
 (\*[rɪŋg ˈpæn ˈbɪʃəp, ˈspɑːtɜːr])  
 (\*[rɪŋg ˈpæn ˈsaʊnd, ˈsɛs])  
 (\*[rɪŋg ˈpæn ˈsprɪ, ˈpæn])  
 (\*[rɪŋg ˈpæn, ˈgɪdɪtɛɪ])  
 (\*[rɪŋg ˈpæn, ˈspɪɪdɜːr])  
 (\*[rɪŋg ˈpæn, ˈspɪɪdɜːs])  
 (\*[rɪŋg ˈpæn ˈsɛs])  
 (\*[rɪŋg ˈpæn ˈfɔːr, ˈpæn])  
 (\*[rɪŋg ˈbʊtɜːr ˈsaʊnd, ˈsɛs])  
 (\*[rɪŋg ˈpæn ˈkɔːstɪŋ, ˈtɛfʊbɪk])  
 (\*[trɔːlɔːr ˈfɪll, ˈdɛɛpɑːt])  
 (\*[sɜːr ɪs ˈnɛt ɑːr ɪs, ˈɪk])  
 (\*[sɜːr ɪs ˈɔːr, ˈsɔːr])  
 (\*[tɪ... (former military base near Monterey, cal, ˈɔːr)])  
 (\*[tɪ... former californian military base, near monterey, ˈɔːr])  
 (\*[tɪ... (former military base near monterey, cal, ˈɔːr)])

A screenshot of the data used in solving the puzzles

Finding the right data sources was one of the most important aspects of this project, since my approach relies so heavily on searching through data. My main source is a corpus of close to five million clue-answer pairs, taken from years of previous crosswords. I wrote a web scraper to gather this data from crossword solving sites, which I then combined with an existing data set.

I performed a significant amount of preprocessing on this set, starting with sorting it alphabetically by row. Because of the large size of the collection, and because of the nature of crossword puzzles (certain clues appearing over and over), this resulted in many duplicate clues. I removed all entries that had both an identical clue and an identical answer. The next important step was partitioning the large set into six smaller ones by answer length. I did this to improve performance on my slower performance search algorithms. Other preprocessing included translating encoded characters to their corresponding plaintext, matching single and double quotes, and removing answers that were too short (under three letters) or too long (over 21 letters).

In addition to the main data set of clues and answers, I used what amounts to a crossword dictionary that I also scraped from the internet, and a list of all Wikipedia article titles. The Wikipedia data is useful because, while the dictionary only contains single word answers, Wikipedia titles provide multiple word answers, acronyms, geographic and historical information, popular culture, etc. I performed similar preprocessing on the Wikipedia titles, including partitioning them by length to increase performance.

Module 1:  
\_\_\_ Beta Kappa (3)' ➡ ['PHI']

Module 2:  
“\_\_\_ Beta Kappa (3)' ➡ ['PHI', 'TAU', 'ETA', 'RHO', 'CHI', 'PSI']

Module 3:  
“P?I’ ➡ ['PPI', 'PLI', 'PSI', 'PTI', 'PCI', 'PAI', 'PHI', 'PRI', 'PII', 'PEI', 'PJI', 'PKI', 'PNI']

An example of how the three modules would run on a single clue. The first generates the fewest candidates, and the third the most, but the reliability of the answers decreases in the same direction.

The first step in solving every puzzle is generating answer lists for each of the puzzle's unencoded clues. This is the most important step, since none of the rest is possible without reliable answer candidates. It is also what most of my work programming was devoted to. Treating this step primarily as a data analysis problem, I developed three search 'modules', each taking a slightly different approach to generating potential answers.

The first module is based on binary search, but incorporates Levenshtein Distance—a metric for determining the similarity of two strings—instead of exact matching, the usual criterion. Using binary search, I find the location the desired clue would occupy in the corpus. This gives me a small range of neighboring clues, whose edit distances to the desired clue is then computed. If any of these distances falls within a certain threshold, then their corresponding answer is added to the clue's candidate list.

The theory behind this approach is that across thousands of crossword certain words, and therefore certain clues and answers—or their near variants—show up in common. These words, used for their advantageous vowel patterns or letter pairings, are commonly considered to constitute ‘crosswordese’, and we can expect their clues to appear with only minor changes in our data set. For early week puzzles, which rely more on common knowledge and simple phrases, and less on wordplay, ambiguity, and long answers like puzzles later in the week tend to do, this holds especially true.

This module is extremely fast, but has limited scope. The majority of clues in most puzzles will not match closely enough with any element from the data set to be considered reliable. What's more, this technique accounts for moderate wording changes, but not when the first word of the clue is changed, as the data is sorted alphabetically.

The results of the first module are considered very reliable, so the second module is only applied to clues that had no matching answers from the first. Because performance is not an issue with binary search, the first module is run on the entire data set; the fuzzy search, however, runs in order linear time on the *number of words* in the data set. To compensate for this enormous hit, this module is run only on the set of clue-answer pairs whose answers are the same length as the answers needed for the desired clue. This was not done for the binary search module, because to search the partitioned data, each file must be passed around in memory. Since performance of the algorithm itself was fast enough, I chose to avoid this.

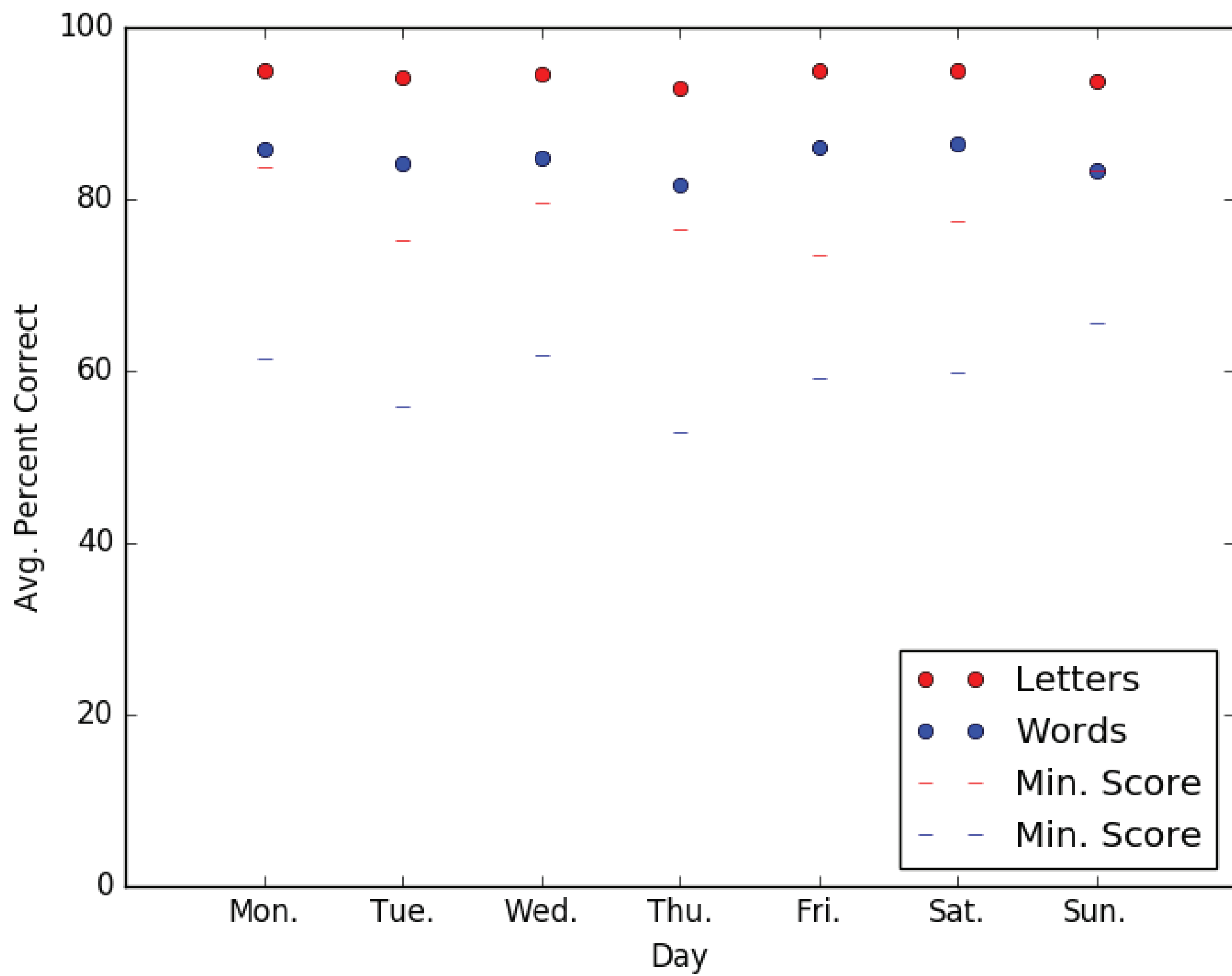
To perform the fuzzy search I identify the clue's key words by eliminating those perceived not to be important to the clue's meaning (a, and, of, to, the, etc.). This remaining set is then compared to the key words of each clue in the partitioned data using Python's regex library. If more than half of the key words , an answer is considered possible.

This module surmounts the first's problem of limited scope, but can overcompensate, producing too many false positives on short clues. If our sought clue has only one or two key words, for instance, then a very long clue in the data set has a good chance of also containing these key words, without actually being related in meaning.

The final module discards the clue, focusing only on the length of the answer and the letters already known from crossing answers. These patterns are then compared with my crossword dictionary and list of Wikipedia article titles. Depending on the length of the answer and the portion of known letters this module can overcompensate even more than the fuzzy search, since most of the letters need to be known to significantly narrow the possibilities. This becomes an even bigger problem with short answers, when we consider the vast number of three and four letter acronyms with Wikipedia pages. Thus, even 'P?' generates nearly every possible letter combination, even with two of the three letters filled in, as in the example above. This module is also quite slow, and is likewise run on data partitioned by answer length.

Having generated a list of candidate answers for each clue, the next step is to seek the configuration that will result in the fewest conflicting word intersections. This is generally referred to as constraint satisfaction. In a constraint satisfaction problem, or CSP, where the solution space is relatively small and well understood, we can use techniques that essentially start with the assumption that one answer or configuration is correct, and then continue solving the problem under that assumption until a conflict is discovered, at which point the algorithm will backtrack. This works well for, say, sudoku, because each square must contain one of the numbers 0 through 9, and at some point we are guaranteed to find a solution. This is not the case with crossword puzzles, at least not with my program. In order for this to work, I would need to have a guarantee beforehand that each candidate list contains the correct answer, which is not the case.

Instead, Regis Filbin solves the puzzle the way a human solver would, that is, ranking the clues in terms of confidence in their answers, and, having filled some of them in, culling candidates of other clues based on clues that have to be in certain positions. This allows for a greater degree of elimination among clues that need it the most (those with the longest answer lists), but it also means that an early mistake can have far-reaching consequences later in the puzzle. Even without relying on what I call the sudoku method, this program could incorporate some element of backtracking, namely by making multiple passes of the puzzle, identifying where perceived conflicts are, and then searching the original answer list to find a more suitable answer. I have not done this, mainly because I did not have time to expand the project so much, and also because I get good results without it.



H	O	T	E	A	P	S	G	M
C	O	N	V	E	M	E	D	O
B	O	T	O	R	I	W	O	N
O	S	T	H	E	C	D	I	E
A	T	B	A	A	P	A	A	E
P	A	R	A	D	O	T	A	D
T	P	E	N	C	A	L	E	D
E	N	A	S	E	V	E	N	E
G	E	A	B	F	A	L	E	I
R	E	A	R	T	R	F	E	R
P	H	E	R	O	N	G	A	I
I	E	S	E	T	A	T	O	N

[illegible]

S	O	B	E	D	E	A	S	T	H	M	A
H	E	A	R	E	U	S		A	T	H	R
I	N	T	E	A	R	S		H	O	U	L
R	O	S	A					S	H	O	L
			A	B	E	R	I	N	T	H	
	S	H	E	E	T	E				S	C
	T	S	E	E	T	E		Y	A	M	M
	S	O	A	M	I					A	R
	S	A	R	R	I	V	A	L		A	R
T	D	S			N	E	W	B	I	O	S
		K	I		N	G	W	M	I	N	O
A	P	R	O	N		O	D	S		O	P
	L	O	A	D	E			J	E	T	B
	S	E	C	A	P	E		J	E	R	A
	S	E	E	K	T	O		A	N	S	T

### Three examples of rebus puzzles

Many crosswords have themes, related longer answers that appear throughout the grid and are generally tied together with a single "theme clue". These themes almost always involve wordplay, intentional misspelling, or other humor; famously difficult problems for computers. Regis Fillbin's natural language processing algorithms will fail to correctly interpret such wordplay.

Rich text, non-alphanumeric characters, and so-called "rebus" puzzles also exceed the program's capabilities. These puzzles will have a special character, multiple characters, or even an entire word in a single square, generally requiring some level of human intelligence to interpret them correctly. Other puzzles will have a shape described in the black squares of the grid itself.

Clearly, there is room to expand Regis Filbin's capabilities in these areas, but some remain essentially impossible. More realistically, I could make large improvements in speed. Regis Filbin is already faster (I think) than any human solver. Still, computers can usually perform tasks thousands or millions of times faster than people, so some improvement is definitely possible. Specifically, the slowness of the program comes from the size of my data. Based on my tests, if I could cut the data in half, I would see roughly a 25 percent speed up in solving times, which would bring many puzzles under one minute. Obviously there is also potential for improving Regis Filbin's accuracy. With more time, I would do this by implementing the backtracking described in the section on filling the puzzle.

## Puzzles Solved by My Program

F	I	J	I	L	B	R	O	A	D	F	L	O	G
O	M	O	O	L	R	O	A	D	D	R	I	C	O
O	P	E	N	S	E	S	A	M	E	O	N	T	O
				S	L	A	I	N		A	N	D	E
D	I	O		O	R	E	G	O	N	S	T	A	T
E	N	E	R	G	E				M	U	S	S	
M	E	N	U			G	E	E	N	A		A	M
O	P	E	R	A	T	I	N	G	S		S	T	E
S	T	S		V	O	L	G	A		N	O	D	E
			A	R	O				W	O	O	L	E
O	R	I	G	I	N	A	L	S	I	N		L	A
V	O	C	A	L			I	N	N	E	S		
E	D	I	T		O	C	E	A	N	S	P	R	A
R	E	N	E		A	F	I	R	E		S	K	E
T	O	G	S		T	O	N	E		A	F	R	O

Monda

H	U	N	T		G	L	E	E		A	T	A	D
O	R	E	O		E	A	R	T	H		N	A	M
B	A	R	R	Y	W	H	I	T	E		S	K	I
O	L	D	N	A	G		C	A	N	T	E	E	N
				W	A	H				R	E	L	A
R	O	B	I	N		W	I	L	L	I	A	M	S
E	D	U	C	E		D	O	O	K		T	A	B
A	I	L	E	D		E	G	O		T	I	A	N
M	E	L			T	H	I	S		R	O	N	A
			M	A	U	R			C	E	S	E	N
C	R	A	D				R	U	B				
		A	D	T		A	S	K		P	L	A	T
A	N	K	A		T		T	H	E	B	E	E	E
I	T	E	M		S	A	P	O	R		O	R	E
R	O	T	S			G	I	B	B		G	I	N

Tuesday

W	U	E	L	A	T	H	O	R	A	M	E	P	S	O
U	E	E	B	A	T	H	O	R	A	M	I	N	G	O
S	T	E	A	M	R	E	U	O	A	L	I	G	O	R
S	O	M	B	R	E	T	O	M	A	L	L	O	W	E
		A	S	E				B	R	E	A			
S	C	R	A	P	P	A	P	E	R		T	V	A	D
T	R	E		E	M	I	R	S		L	I	E	T	O
R	O	L	F	E		R	O	T		O	N	N	E	D
I	W	I	L	L		W	O	O	E	R		S	I	G
P	E	T	A		W	A	F	F	L	E	C	O	N	E
			T	H	E					E	L	O		
A	B	A	S	E	D		F	L	O	V	E	R	O	N
B	U	X	O	M										
B	R	E	D	A		A	M	I	N	A	D	E	E	
A	L	L	A	N		D	A	D	S		L	E	A	R

Wednesday

P	O	P	S		A	T	E	I	T		K	I	S	S
S	N	O	W		V	O	L	V			I	R	A	Q
S	H	A	L	F	N	E	L	S	O	N		S	E	A
A	H	L		F	E	N	D	E	R		A	M	A	R
W	O	O	D	O	U	T				A	M	E	L	I
P	I	E	C	E	O	F	W	R	I	T				I
		L	O	O		R	A	I				Z	E	E
A	S	H	E			P	O	X					N	E
A	S	P	A		B	R	O			P	S			
P	A	R	T	I	A	L	E		C	L	I	P	S	E
I	M	D	O	N	E			H	A	N	D		X	A
R	A	C	K	S		B	O	O		I	D		F	A
R	A	L	O			E	E	M	I		F	I	N	A
T	O	R	N		C	I	S	C	O				O	R
T	E	T	E		I	N	K	E	R		H	I		I

Thursday

K	A	T	I	C	S	O	M	F		A	S	P	E	N	S
L	A	T	O	N	E	F	A	R		R	I	B	E	T	A
P	A	K	E	T	T		R	E		E	X	A	C	T	
O	K	I	E	s		D	E	I	T		S	T	A	W	N
W	E	E	L	A					T	O	V		P	O	L
					G	A	I	T	S				D	I	R
			G	O	D	C	H	A	R	L	O	T	T	E	
	Z	I	P	A	L	E	E	D	O	O	A	H			
F	A	M	I	L	E	R	E	U	N	I	O	N			
A	C	M	E				A		A	U	D	I	T		
	T	E	E			S	C	I			N	E	W	C	A
B	F	F			N	O	N	E			R	H	O	E	
A	R	I	S	E			V	W	B	E	E	T	T	L	E
C	O	V	E	R	T		A	A	A	R	E	A	T	E	d
K	N	E	A	S			N	B	A	S	T	A	R	S	

Friday

B	I	G	M	A	C		S	A	D	S	O	N	G	S
O	G	L	A	C					O	U	T	T	H	E
D	O	A	J	I	G		F	T	S	U	M	T	T	E
Y	T	D		T	E	N	T		O		N	I	T	T
C	O	S				S	E	T	H	S		C	L	E
A	R	O	A	R		L	O	A	T	H		E	L	S
S	I	M	I	A	N		P	R	O	A	M			
T	O	E		S		O	O	P		P	O	L	E	M
			S	U	M		A	C		D	O	S	I	D
O	R	S		L	A	P	A	Z		S	A	L	E	M
N	A	T	L		N	A	P	E	S		S	L	E	A
S	C	R	A	M		B	I	D	E	T		P	L	O
T	E	A	C	A	K	E	S		T	O	T	O	I	V
A	M	N	E	S		I	A	C		B	R	O	N	Z
R	E	D	S	C	A	R	E		Y	O	N			

Saturday

X	W	H	A	T	N	A	S	A	S	W	A	B	U	T	E
X	P	H	I	L	O	N	O	D	I	N	O	N	O	B	E
X	S	E	E	P	A	R	E	S	E	S	E	E	E	E	E
X	O	L	S	T	S	E	S	O	R	P	E	C	I	A	L
J	M	C	K	E	I	N	C	O	K	E	R	O	C	K	I
B	K	V	I	N	O	S	E	S	E	S	W	A	T	H	E
S	P	A	T	A	T	E	S	E	S	S	A	S	H	A	N
S	P	A	T	A	T	E	S	E	S	S	A	S	H	A	N
N	E	W	S	T	A	N	D	L	E	D	A	G	G	R	E
A	R	E	A	S	R	A	P	I	V	E	D	A	G	E	O
G	E	O	S	E	S	N	O	N	O	N	O	N	O	N	O
E	V	E	N	I	N	G	E	R	S	E	R	I	S	T	H
E	V	E	N	I	N	G	E	R	S	E	R	I	S	T	H
L	E	A	P	E	R	M	S	K	A	L	L	O	U	R	C
M	R	S	A	P	E	R	M	S	K	A	L	L	O	U	R
T	M	I	N	C	A	G	I	N	E	S	S	M	L	Z	E
T	M	I	N	C	A	G	I	N	E	S	S	M	L	Z	E
C	O	M	E	T	E	Z	E	A	J	A	O	A	S	I	S
M	A	N	K	U	R	P	E	S	E	S	E	I	N	P	A
M	A	N	K	U	R	P	E	S	E	S	E	I	N	P	A
P	S	I	K	E	E	L	S	E	E	R	A	P	S	E	

Sunday