Worksheet Activity: Text	Compression
--------------------------	-------------

	Name:
<b>Directions</b> : Compress the verses of the tongue twister shown below, us you've learned. Use the Pitter Patter example as your guid	
<b>Notes:</b> A space (indicated by _) counts as a character. So do punc The grids provided for the dictionary and compressed text a	
	Dictionary
A_tutor_who_tooted_the_flute_	
Tried_to_tutor_two_tooters_to_toot	1
Said_the_two_to_their_tutor,_	2
·	3
"ls_it_harder_to_toot_	5
Or_to_tutor_two_tooters_to_toot?"	6
149 bytes	7
1.0 23,000	8
How to count:	9
If there's a character in a square it counts as a byte.  Index numbers on the dictionary also count as bytes.	
Notice that once you go past the 10th dictionary entry	(B)
you need to invent a single character symbol as a dictionary index.	Bytes in Dictionary
Compressed Text:	
Bytes in Compressed Text	•
Bytes in To	ext Bytes in Dictionary Total
	<b>+</b>
	Percent Compressed
	(1-(total/149))*100:

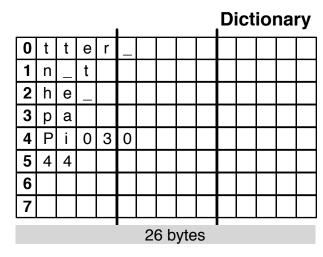
## Text Compression:: Pitter Patter example

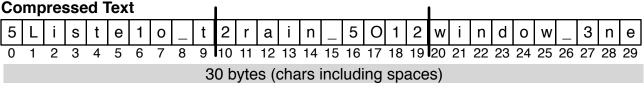
Here is the "Pitter Patter" poem example we did in class. The original poem, dictionary, and compressed text is shown.

Recall that to arrive at this final dictionary, we had to go through several revisions looking for patterns and realizing new ones.

Also note that the dictionary is self-referencing, with later entries referencing earlier ones.

Original Poem
Pitter\_Patter\_
Pitter\_Patter\_
Listen\_to\_the\_rain\_
Pitter\_Patter\_
Pitter\_Patter\_
On\_the\_window\_pane





**Total** = 30+26 = 56 bytes =  $\sim 54\%$  compression

Note: There is actually one more 2-character pattern that occurs in the text - "in" The pattern "in" occurs twice and it was not included in the dictionary because it would actually increase the total number of bytes in the compressed text rather than reduce it. Do you see why?

## How it might actually be stored:

Recall that the ASCII code for english text requires only 7 bits since there are only 128 ASCII characters. But a byte has 8-bits. Therefore each byte carries with it one unused bit. We can use that last bit to indicate whether the byte should be read as a normal ASCII character or as a dictionary entry. Additionally, the dictionary can use the 8th bit to denote when the next dictionary entry is starting.

## Is this type of compression actually used?

Yes. It's called Lempel-Ziv compression or "LZ" compression for short. Lempel and Ziv are Isreali computer scientists who invented it in the late 70s. An MIT proffessor, Terry Welch, improved LZ compression slightly and so the current version in use is called LZW compression. It is the compression scheme used in ".zip" files. Of course the real Limpel-Zev compression works by looking for patterns in the 0s and 1s since a computer program that might do the compression would be ignorant of language. LZ compression works very well for any kind of file in which there are a lot of patterns. Thus, it works very well for compressing plain text files since languages tend to have a lot of patterns.