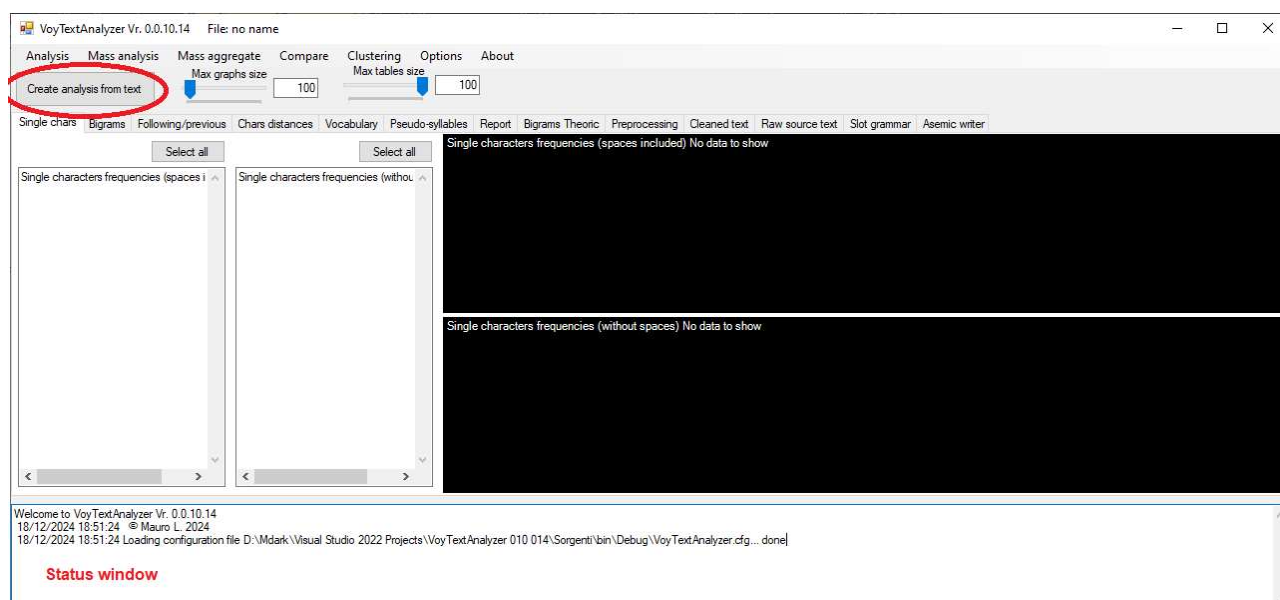VoyTextAnalyzer adds the slot grammars functions to a previous program, TextAnalyzer, which is unpublished and was meant as a statistical tool for texts in any language[1]. All the original statistical functions of TextAnalyzer are available in VoyTextAnalyzer but are not described here. A manual is available: "Manuale Text Analyzer", unfortunately only in Italian at the moment.

**Installation**

No installer is provided. If you are not interested in the source files just copy VoyTextAnalyzer.exe in any directory and run it by double-clicking on it. The first time VoyTextAnlyzer is run it will create a configuraton file in the same directory (it's just an XML file with .cfg extension).
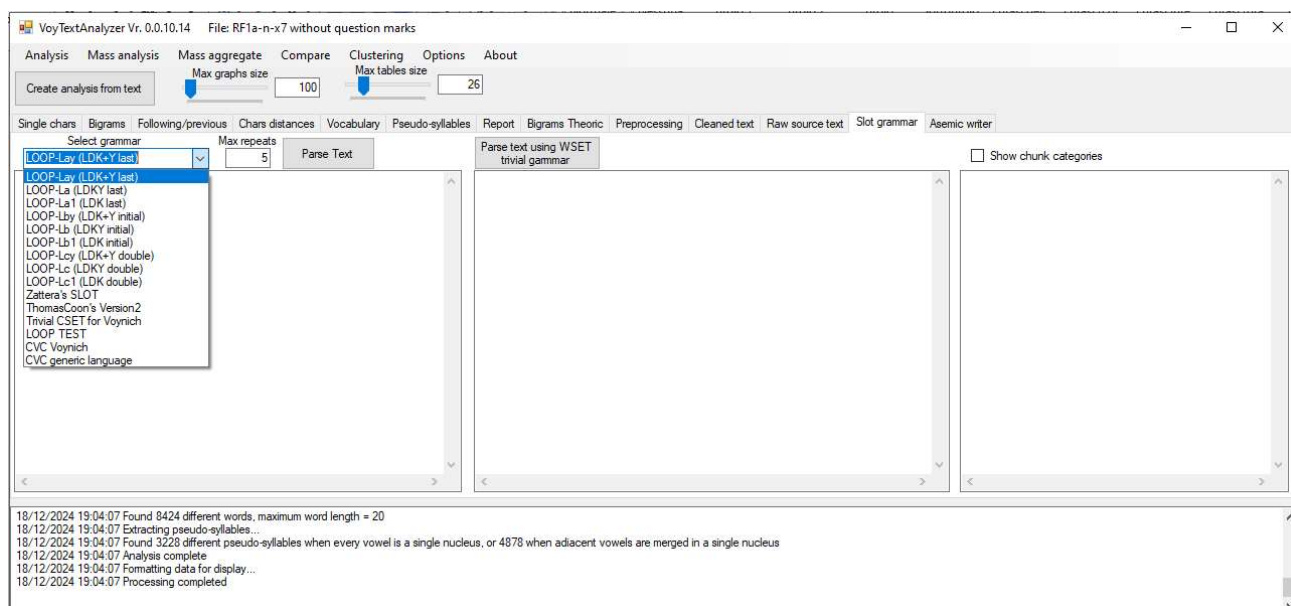
**Use**



Important: before you can use slot grammars, you must first load and 'analyze' the source text. Use the 'Create analysis from text' button to open a file dialog and load the desired source text file. The file must be in .txt format. VoyTextAnalyzer will automatically remove all 'non-word characters' (according to Regex definition) and redundant spaces. The loading time is very fast (a couple of seconds with a text the size of the Voynich manuscripts).

Once the source text has been loaded you can proceed to the 'Slot grammar' tab.

---

[1] And indeed the slot grammars functions can be applied to any language, not only to Voynichese. The pre-defined grammar "CVC generic language" is a rough example of a grammar usable on any language.

The drop-down menu allows to choose which of the pre-defined grammars to use. Notice you cannot add or modify grammars without recompiling the software, this will be explained later.

Every grammar has a default number of repetitions, which can be changed with the 'Max repeats' box, from 1 up to a maximum determined by the grammar.

Use the button 'Parse text' to parse all the word types into chunks and calculate the Huffmann codes and $N_{bits}$. The Trivial WSET grammar cannot be selected in the drop-down menu: use the "Parse text using WSET trivial grammar" button instead. The elaboration takes about 10 seconds (most of which is writing the output to the text boxes)[2].

---

[2] Notice the elaboration is not inside a System.Thread (sorry), this mean it could timeout and crash the program. This may be corrected in a future version (the rest of TextAnalyzer functions are threaded), however the processing time is rather fast and a text such as the VMS should not cause any problems (no timeouts should occur when writing to the textboxes).

You can click on any textbox, press Ctrl-A to select everything inside and then Ctrl-C to paste all the data in any Windows program, for instance Excel (the data may seem weirdly aligned in the textboxes but they will be neatly divided in columns once copied to Excel).

*Leftmost textbox: Report*

It shows the data which have been used in the , particularly the coverage along word types and the total number of bits required ($N_{bits}$).

$N_{chunktypes}$ is the total number of unique chunks found. See 'Asemic writer' for $N_{chunktokens}$ and $N_{transitions}$.

More informations follow in the report (use the scroll bars or copy and paste).

*Center textbox: all the text word types as divided in chunks by the grammar*

Self-explicative. There are more columns on the right side: the list of all the word types not found by the grammar and the list of the word types discarded because they include 'rare' characters.

*Rightmost textbox: chunks dictionary*

Self-explicative. Clicking on 'Show chunk categories' will change the visualization to one where chunks have been divided in 'categories' (not used in the paper [1] and not further developed).

**Asemic writer tab**

All the data in this tab have not been used in the [1]  and may be the subject of a future paper.

The first two textboxes are written by the 'Parse text' function of the previous tab. You need instead to click the "Write asemic" button to generate the random text in the last box.

*Leftmost textbox: chunkified grammar*

Briefly, the chunkified grammar is generated from the word types divided in chunks (see previous 'Slot grammar' tab, central textbox). It's a derived slot grammar [chunkified grammar = CHUNK(original grammar)] where the n-th slot is built by taking all the chunks which appear in the n-th position of a word type. For example: 'daiin' = [d] + [aiin], the [d] goes in the first slot of the chunkified grammar, the [aiin] goes in the second.

The total number of chunks used in the slots of the chunkified grammar is the $N_{chunktokens}$ value displayed in the report.

*Center textbox: Markov chain transitions table*

It's computed from the chunkified slot grammar *and* the frequency with which each word token appears in the text. It gives the probability of transitioning, say, form START to a [d] as 1st chunk, then from [d] in the 1st chunk to [aiin] as 2nd chunk and so on.

The total number of transitions is the $N_{transitions}$ value displayed in the report.

*Rightmost textbox: asemic random text*

The asemic text is automatically generated with as many words as the original text to make comparisons easier.

VoyTextAnalyzer generates a very good pseudo-Voynich text, indistingushable from the original on many statistical properties (data no shown, you can check yourself by copying the text and using your favourite analysis

tool, or the built-in statistical function of VoyTextAnalyzer). The main statistical difference is the asemic writer generates less hapax legomena and rare words than found in the original Voynich. This may be the subject of a future paper.

## Adding a grammar using source files

VoyTextAnalyzer has no grammar editor. Grammars can be added/modified only in source files. The environment in which VoyTextAnalyzer has been developed is Microsoft Visulat Studio 2022 free edition (.NET framework 4.5.2).

All the grammars are in ParseGrammarsLibrary.cs



```
                                                          Reference name
12          public static ParseVoynich.ParsingGrammar LOOP_LDK_final()
13          {
14
15              ParseVoynich.ParsingGrammar grammar = new ParseVoynich.ParsingGrammar();
16
17              grammar.name = "LOOP-La1 (LDK last)";          You can write anything here,
18              grammar.notes = "QCSHY at start, LDK at the end";   displayed in the Report
19
20              grammar.default_loop_repeats = 5;   Self-explicative
21              grammar.max_loop_repeats = 8;
22
23              // Take care! Overlapping strings must be ordered by length to avoid queer chunkifications. Ie. [iii, ii, i], not [i, ii, iii]   Heed this remark!
24              grammar.slots_table = new string[6][];
25                                                                                              This number must be exactly the
26              grammar.slots_table[0] = new string[] { "q", "ch", "sh", "y" };                  highest number [n] in the slots table
27              grammar.slots_table[1] = new string[] { "eee", "ee", "e" };                      plus one, or  the program will crash
28              grammar.slots_table[2] = new string[] { "o" };          These numbers must start from zero and be   in initializatio, or generate weird results
29              grammar.slots_table[3] = new string[] { "a" };          strictly progressive, or the program will crash
30              grammar.slots_table[4] = new string[] { "iii", "ii", "i" };   in initialization, or generate weird results
31              grammar.slots_table[5] = new string[] { "l", "d", "k", "r", "s", "t", "p", "f", "cth", "ckh", "cph", "cfh", "n", "m" };
32
33
34              // Word types containing 'rare' characters which do not appear in the grammar cannot be found. The following list is used to esclude them from the counts
35              grammar.rare_characters = new char[] { 'c', 'h', 'g', 'x', 'v', 'z', 'b', 'j', 'u' };
36              // But, some rare characters could be included in some of the slot glyphs (this is typical with Voynic manuscript).
37              //   For instance: let's say 'tchodypodar' is not found because it needs too many chunks. It contains the rare characters 'c' and 'h', but we need to avoid it to be
38              //     discarded, because 'c' and 'h' are in the 'ch' group, which the grammar can generate. The following list defines the groups
39              //       wherein rare characters are not considered 'rare'.        This two arrays are used to manage the 'rare'
40              grammar.rare_chars_groups = new string[] { "ch", "sh", "ckh", "cth", "cph", "cfh" };   characters not used by the grammar. See
41                                                                                  the comments. You can just copy them if
42                                                                                  you want to treat the VMS the standard way
43              return grammar;
44          }
```

After creating a new grammar, it must be added in the function ParseGrammarsLibrary.GenGrammarsList() for it to appear in the drop-down menu:

```
564    ⊟      public static List<ParseVoynich.ParsingGrammar> GenGrammarsList()
565           {
566               List<ParseVoynich.ParsingGrammar> grammars_list = new List<ParseVoynich.ParsingGrammar>();
567
568               ParseVoynich.ParsingGrammar new_grammar = new ParseVoynich.ParsingGrammar();
569
570
571               new_grammar = LOOP_LDK_Y_final();
572               grammars_list.Add(new_grammar);
573
574               new_grammar = LOOP_LDKY_final();
575               grammars_list.Add(new_grammar);
576

                         etc. etc.

621               return grammars_list;
622           }
```

## Short description of the source files

I will only describe the classes used by the 'Slot grammars'/'Asemic writer' tabs, not the whole program. All the names in the sources are in English and as self-explicative as possible, but comments are in Italian (I'm sorry).

The launch routines of the functions from the buttons and the data visualization in the tabs are all in Form1.

Files used by the Slot grammars processing:

| ParseGrammarsLibrary.cs | Definition of the grammars, see above |
| --- | --- |
| ParseVoynich.cs | ParseVoynich.parse_words<br>   1) Parses the word types into chunks<br>ParseVoynich.get_chunks_data<br>   2) Creates of the chunks dictionary. |
| ParseHuffmann.cs | ParseHuffmann.get_Huffmann_codes<br>   3) Adds the codes to the chunk dictionary |

Files used by the Asemic writer processing:

| PackGrammarAndWrite.cs | PackGrammarAndWrite.calculate_chunkified_grammar<br>   1) Compilation of the chunkified grammar<br>   2) Compilation of the transitions table<br>PackGrammarAndWrite.Markov_asemic_writer<br>Asemic writer |
| --- | --- |

## Input class for the slot grammars processing

The dictionary of all the word tokens found in the source text is used as the input of the slot grammars processing. It is computed by the original TextAnalyzer functions when 'Creating an analysis from text'.

```
public List<EValueOcc> vocabulary_words_distribution = new List<EValueOcc>();
```

An EvalueOcc is a basic class, formed by a string + a number.

```
public class EValueOcc
{
    public string element;

    public long value;
}
```

In this case **element** stores the word type, while **value** stores the number of times the word type appears in the text (occurrences).

`vocabulary_words_distribution` example:

| element | value |
|---------|-------|
| daiin   | 834   |
| ol      | 577   |
| aiin    | 529   |
| …       | …     |

**Slot grammar processing classes (ParseVoynich.cs and ParseHuffmann.cs)**

The main class is `ParsingResults`. It stores the word types parsed into chunks and the chunks dictionary. It is computed from `vocabulary_words_distribution`.

```
public class ParsingResults
{
    public ParsingGrammar grammar;

    public int total_wordtypes_found;                            // Nhits
    public int total_wordstypes_not_found_but_with_rare_characters;
    public float total_grammar_wordtypes;                        // Ngrammarsize

    public int total_number_of_chunktypes;     // Nchunktypes

    public List<WordParseData> data = new List<WordParseData>();   // Chunkified word types

    public Dictionary<string, Chunk> chunks_dictionary = new Dictionary<string, Chunk>();


    Two more dictionaries here for chunk 'categories'
    Not further developed and not documented here

}
```

| `public ParsingGrammar grammar;` | The slot grammar used for parsing (pre-loaded by choosing a grammar from the drop-down menu). |
|---|---|
| `public List<WordParseData> data` | The list of all word types divided in chunks (provided there was a 'hit') |

```
public Dictionary<string, Chunk> chunks_dictionary
```
The chunks dictionary

```csharp
public class WordParseData
{
    public string word;
    public long occurrences;

    public bool word_found;
    public bool word_not_found_but_has_rare_characters;

    public int number_of_words_chunks;

    public List<List<EValue>> parsed_word = new List<List<EValue>>(); // Contains the chunks which compose a word type, divided slot by slot
}
```

Example of `parsed_word` (An EValue is like an EvalueOcc, but uses an 'int' rather than a 'long' numeric type. In parsed_word `value` stores the number of the slot which scored a 'hit'. It's used for the chunk 'categories' and is not further detailed).

The word type is 'daiin', the grammar is LOOP-Lay

| "d", 5 | | |
|--------|---------|--------|
| "a", 3 | "ii", 4 | "n", 5 |

```
public Dictionary<string, Chunk> chunks_dictionary
```

The key of the dictionary is the chunk string (for instance: "aiin").

```csharp
public class Chunk
{
    public string chunk;
    public long number_of_times_used;

    public string Huffman_code = "";
}
```

`Chunk` stores the chunk string (same as the dictionary key), the number of times the chunk is used in the text and the assigned Huffman code. Huffmann codes are added to the pre-compiled dictionary by `ParseHuffmann.get_Huffmann_codes`

**Calculation of N$_{bits}$**

It's calculated at visualization time, in Form1. `display_parsing_report`

```csharp
// Nbits metric: length of the Huffmann-compressed text (text + chunks dictionary)
int charset_size = text_analyzer.analysis_results.monograms_distribution_excluding_spaces.Count - Voynich_parsing_results.grammar.rare_characters.Length;
double bits_per_character = Math.Log(charset_size, 2);

long Huffmann_compressed_text_total_required_bits = 0;
double Huffmann_dictionary_total_required_bits = 0;
foreach (ParseVoynich.Chunk chunk in Voynich_parsing_results.chunks_dictionary.Values)
{
    Huffmann_compressed_text_total_required_bits += chunk.number_of_times_used * chunk.Huffman_code.Length;

    Huffmann_dictionary_total_required_bits += chunk.chunk.Length * bits_per_character + chunk.Huffman_code.Length;
}
```

Notice the formula used for text_bits is not the same as the one in the paper [1], but it is equivalent (implementing the formula as writte in the paper would hve resulted in an ugly software).

`charset_size` is calculated by subtracting the number of 'rare' characters which are not considered by the grammar from the total number of characters found in the text.

## Chunkified grammar and asemic writer

The main class `ChunksGrammar`, in PackGrammarAndWrite.cs, contains the slots of the chunkified grammar and the transitions table:

```csharp
public class ChunksGrammar
{
    public List<Dictionary<string, EValueOcc>> slots_table = new List<Dictionary<string, EValueOcc>>();

    public List<Dictionary<string, SlotTransition>> transitions_table = new List<Dictionary<string, SlotTransition>>();
}
```

Given this is out of the scope of the paper [1] , no further explanations are given here.

## References

[1] Mauro Lanzini "The evaluation of slot grammars and their application in the study of the Voynich" Manuscript, 2025. https://www.academia.edu/