

Automatic Manuscript Comparison Prototype

Andrew Edmondson

Introduction

Inspired by DNS sequences, I decided to explore the idea of encoding NT manuscripts in a similar way: ABCABDCD... etc.

The concept is as follows: Each verse will count as a unit, including the title (verse 0) of a work. Each unique version of that verse then gets assigned a letter, a, b, c etc. Then any text can be encoded and compared quickly with any other.

Texts

Transcriptions of manuscripts are downloaded from <http://www.igntp.org/> in XML form. These are then parsed and each verse is indexed and assigned a letter. 07 (Codex Basiliensis) was chosen as a base text, somewhat arbitrarily, as it is Byzantine and therefore likely to be more similar to the majority of other manuscripts than, say, 01 (Codex Sinaiticus).

The base text of the first verses of John 1 is then represented as:

MS 07 aaaaaaaaaaaaaa...

Comparing 04 (Codex Ephraemi rescriptus) with 07 we see:

MS 07 aaaaaaaaaaaaaa...

MS 04 dabaaaaaaa...

The first two verses of 04 are missing, and verse 3 is just a fragment. It can then be seen that v4 is the same in 07 and 04. In the web-representation of the output, hovering over a letter shows the text of the verse in that manuscript. So, for example, v5, is as follows:

07: και το φως εν τη σκοτια φαινει και η σκοτια αυτο ου κατελαβεν

04: και το φως εν τη σκοτεια φαινει και η σκοτεια αυτο ου κατελαβεν

Correctors as states of the text

The XML transcriptions not only show the final state of the text, but where possible also show the original and details of corrections made by different hands. 01, for example, can be represented like this:

MS 01 (firsthand)	a b b b a b b b b b b c b c b
MS 01 (S1)	c d d c c c c
MS 01 (S2)	e
MS 01 (ca)	a c a d g c c d c c c c c d c c c a c d c
MS 01 (cb2)	f a c c
MS 01 (corrector)	c

This shows that v3 has been corrected (by "ca") to be the same as 07 (although this, of course, says nothing about a direct relationship between 01 and 07). The corrector changed "ουδε εν" to "ουθεν".

MS 07 (firsthand)	aaa
MS 07 (2)	a
MS 07 (corrector)	b b
MS 01 (firsthand)	abbbabbbbbbbcbcbbbbbbbbbbbcbcbcbcbcbcbcbcbcbcbcbcb
MS 01 (S1)	c d d c c c c
MS 01 (S2)	e
MS 01 (ca)	a c a d g cc d cc c c cd cc ca cd c
MS 01 (cb2)	f a c c
MS 01 (corrector)	c
MS 02 (firsthand)	aacaadaacaaecbaaeccacacccddadaadcdcadecddcdacdeaccd
MS 03 (firsthand)	aaacaecababcbfdhbdddcddceddedcedddbefdedededdefcdde
MS 03 (2)	d e b
MS 03 (secunda_manu)	a d
MS 04 (firsthand)	dabaaaaaadfibefedfeabeffefedfeeacfgef
MS 04 (2)	j g egc ac g gfegf gif
MS 05 (firsthand)	aaebafdc dabdgkc
MS 05 (A)	c
MS 05 (B)	eh l
MS 05 (K)	f
MS 09 (firsthand)	b da ed cafihbaaafahcaeaahaaahaaacadjaaaefeafeafaf
MS 09 (corrector)	e a
MS 011 (firsthand)	aaaaaaaaaaaadibafhgaicacagigaafgafaddkggafgfafeagef
MS 011 (corrector)	g g f
MS 013 (firsthand)	aaaacafaad gdjbaaiafafaaaaahaaaaaacadaaaafghaefgaaaf
MS 017 (firsthand)	aaaaaagaaeadfmaeahgacaaaadhghaiaaaahlhaghiaffhefff
MS 017 (corrector)	d i
MS 019 (firsthand)	caaadaaaeaaajkbbejidfcdfdhjadhgjgeeafmjihiaaggifhgg
MS 019 (corrector)	ji

MS 021 (firsthand)	d a a a a a g a a a c a d f n a a k a h a g a a f a d a i a h k a e a a d j a a j e j a h f j a f h a
MS 022 (firsthand)	j h e g c i k i c a i l h e a a d n k
MS 024 (firsthand)	j a g m a a c e d o l a
MS 028 (firsthand)	a a a f a a g a a a a a e f m a e g a a k c a a a l a d a a d a e a a d k i a f j d a i f k a i a f
MS 028 (corrector)	a
MS 029 (firsthand)	a g a m j k a g n

It is quite straightforward to spot Byzantine manuscripts by the amount of “a” entries (Although this actually means that the verse is identical to the text in 07. A different base text might yield different results).

Several questions or points of interest are raised by considering the data above. For example, what was 028's corrector's reason for changing v39 to the “a” text. Why is the “d” text so popular in v37? Or the “c” text in v9? Etc.

Next steps

The next steps in this prototype may include:

1. Further automatic regularisation – e.g. nomina sacra. The lack of such regularisation is currently corrupting the results somewhat, so while I can develop a method of analysis some of the conclusions may be spurious at this stage.
2. Finding a way to represent the degree of difference for each text form of a verse – is it radically different or has just a single letter changed?
3. Finding a way to represent the type of difference – addition, omission, etc.
4. Integrate with the complete set of IGNTP manuscripts – e.g. thus including papyri
5. ?