

## AOR Phase 1 Transcription Data Summary and Analysis

There are a total of 5,468 individual image files in the AOR corpus (1,146, 21% of that total, comprising images of Livy’s *Ab urbe condita* alone). For each image bearing a visible sign of Harvey’s manuscript interventions within the printed text—be it a verbal marginal comment, a non-verbal mark or element of Harvey’s person symbolic system of annotation, or simple underlining of the printed text—the three transcribers (the AOR postdoctoral researcher and two research assistants) jointly produced a transcription file in XML format. There are 2,501 of these XML files, each associated with its image file counterpart.

Table 1.1 shows how these transcription files are distributed among the 12 volumes in the AOR Phase 1 corpus.

*Table 1.1 Number of transcription files per book in corpus*

<b>Title book</b>	<b>Number of files</b>
Buchanan, <i>Ane detectioun of the duinges of Marie Quene of Scottes</i>	50
Buchanan, <i>De Maria Scotorum regina</i>	17
Castiglione, <i>Il libro del cortegiano</i>	91
Castiglione, <i>The covrtyer of Covnt Baldessar Castilio</i>	389
Domenichi, <i>Facetie, motti, et burle &amp; Guicciardini, Detti et fatti piacevoli et gravi</i>	421
Freigius, <i>Paratitla</i>	60
Frontinus, <i>The strategems</i>	224
Livy, <i>Ab urbe condita</i>	652
Machiavelli, <i>Art of warre</i>	238
Melanchton, <i>Selectarum declamationum</i>	133
Olaus, <i>Historia de gentibus septentrionalibus</i>	121
Smith, <i>De recta &amp; emendata linguæ Anglicæ</i>	105
<b>Total</b>	<b>2,501</b>

The XML files comprise several “elements”, each element corresponding to a type of annotation as described above: marginalia (a verbally written note); underlining (words and passages underlined or highlighted in some way); and a non-verbal mark or symbol. Additional tagging protocols were developed, as well, to deal with more exceptional cases of annotations such as visual drawings; numerals (passages marked up, for example, “1, 2, 3”); and errata (where the annotator appears to have corrected or specifically altered the printed text). In addition, we created the ability to quantify the specific number of words written by Harvey as marginalia, and the number of printed words underlined by Harvey within each book and across the full corpus of books.

As detailed in Table 1.2, below, the 2,501 transcription files contain a total of 102,627 individual elements. More specifically, they contain a total of 80,009 words of manuscript annotations and 226,710 underlined words, comprising a total body of annotated materials amounting to 306,719 words. To our knowledge, this dataset constitutes the largest body of information ever gathered together in a systematic, machine-readable form from manuscript annotations in early printed books.

**Table 1.2 General overview of corpus**

Title book	Marginalia	Underline	Mark	Symbol	Drawing	Numeral	Errata	Total	Marginalia_words	Underline words	Total
Buchanan, <i>Ane detectioun...</i>	13	635	221	4	0	0	5	<b>878</b>	178	2089	<b>2267</b>
Buchanan, <i>De Maria Scotorum regina</i>	8	48	3	0	0	0	2	<b>61</b>	44	167	<b>211</b>
Castiglione, <i>Il libro del cortegiano</i>	92	116	48	0	0	0	6	<b>262</b>	718	393	<b>1111</b>
Castiglione, <i>The covrtyer of Covnt Baldessar Castilio</i>	397	11981	5118	633	2	29	20	<b>18180</b>	2605	44457	<b>47062</b>
Domenichi, <i>Facetie, motti.. &amp; Guicciardini, Detti et fatti..</i>	3094	10876	5820	693	1	13	11	<b>20508</b>	39095	33288	<b>72383</b>
Freigius, <i>Paratitla</i>	50	468	331	31	0	13	2	<b>895</b>	1728	1045	<b>2773</b>
Frontinus, <i>The strategems</i>	869	3721	2549	138	0	18	3	<b>7298</b>	10351	11864	<b>22215</b>
Livy, <i>Ab urbe condita</i>	851	25919	10317	1144	1	66	54	<b>38352</b>	21080	95779	<b>116859</b>
Machiavelli, <i>Art of warre</i>	214	7288	4036	92	0	51	76	<b>11757</b>	2630	27393	<b>30023</b>
Melanchton, <i>Selectarum declamationum</i>	45	1129	327	13	1	8	9	<b>1532</b>	429	3206	<b>3635</b>
Olaus, <i>Historia de gentibus septentrionalibus</i>	53	1552	293	18	0	20	3	<b>1939</b>	798	4935	<b>5733</b>
Smith, <i>De recta &amp; emendata linguæ Anglicæ</i>	42	687	224	6	0	3	3	<b>965</b>	353	2094	<b>2447</b>
<b>Total</b>	<b>5728</b>	<b>64420</b>	<b>29287</b>	<b>2772</b>	<b>5</b>	<b>221</b>	<b>194</b>	<b>102627</b>	<b>80009</b>	<b>226710</b>	<b>306719</b>

Table 1.2 therefore provides a fairly objective overview of the reader's interaction with the printed texts in the Phase 1 Harvey corpus. It also reveals several interesting reading patterns undertaken by Harvey across these books. For example, both of Buchanan's books may be regarded as lightly annotated, relative to the rest of the corpus, and yet in Buchanan's *Ane Detectioun*, Harvey underlined far more words of the printed text than were underlined in the second Buchanan book, *De Maria Scotorum*. Although further analysis needs to address the relative length of a book, and the average number of words captured per annotation tag, applied statistical analysis of such major variations within and across books in the corpus (in addition to further exploration of the relationship of common themes that appear both in underlined printed texts and allied marginal annotations) may provide new ways to evaluate

more fundamental elements of historical reading practices and strategies than have been possible in a purely analog environment. These analyses may offer a basis for further and perhaps unexpected lines of enquiry and investigation into historical reading practices.

The humanities content team has, however, also determined that such statistical analysis should not be treated as the endpoint of scholarly investigation of the AOR resource, but rather may be used as a tool to support the identification of specific research questions that might not otherwise occur to the researcher in the analog environment of the library reading room. That research may lead to the identification of demonstrable, even statistically significant relationships only made clear through the analysis of the data extracted from the books in the XML transcription and encoding process, and lead to additional qualitative as well as quantitative scholarly inquiries.

It is also important to note, in the context of the data summarized below in Table 1.2, that the large majority of printed texts and marginal manuscript annotations in the Phase 1 Harvey corpus were produced in languages other than the English, most notably Latin, and in several modern European vernacular languages (Italian, French, etc.). This repeatedly demonstrates Harvey's general practice of inscribing annotations into books in the same language in which they were originally printed. The XML schema includes a <language> tag to identify which language is being used where. There are 6,243 language tags, of which 1,283 (20.5%) were used to identify passages in English and 4,388 language tags (70.2%) to identify those in Latin. The number of language tags actually exceeds the number of marginalia tags (5,728), since some marginal notes were written in more than one language, and thus one marginalia tag can consist of several language tags. These data may reveal information of particular use to corpus linguistics scholars and to students of the histories of translation and polyglot composition.

This process of transcription was not a straightforward mechanical process, either: in many cases external textual and historical research was required to confirm a given transcription based on contextual evidence contained outside the books in the corpus, or, in some cases, across books within the Harvey corpus, wherever similarities appeared. This was particularly notable in the case of identifying the specific historical personages he cited, identifying specific texts referenced by him in brief short titles, and particular geographic places and locations cited in reference to the printed texts. The Phase 1 corpus transcription files contain 3,950 people tags, 709 book tags, and 221 location tags, which collectively refer (due to repetitions) to 970 individuals, 341 book titles, and 88 locations.

As Table 1.2 also suggests the dataset created in the course of AOR Phase 1 may rightly be regarded as a "big data" set in the sense that the complexity of the data could never be adequately encompassed by an individual researcher working within the constraints of traditional "analog" scholarship. At the same, the average size (in terms of bytes) of the transcriptions is actually is small, as demonstrated below in Table 1.3, relative to data sets generated, for example, in the sciences. This raises an important point in the broader conversation that brings humanistic content in direct conjunction with quantifiable data; namely, that humanists must speak responsibly about the quantitative size of the data that humanities projects generate. The qualitative arguments and analyses that may be generated within humanistic scholarly discourse should not seek to draw greater authority or preponderance simply from the fact that they sit atop a relatively massive pile of quantifiable data (piles of data that would seem, by contrast, remarkably modest in the natural sciences or in the world of finance). It is incumbent upon humanists to consider foremost the qualitative

aspects of the dataset, since they, and not the actual size of the dataset itself, helps the broader scholarly community to address particular research questions. Both forms of data are exceedingly useful, but simply must be treated proportionally.

**Table 1.3 Average size of XML file per book**

<b>Title book</b>	<b>Number of files</b>	<b>Average length of files (in bytes)</b>
Buchanan, <i>Ane detectioun of the duinges of Marie Quene of Scottes</i>	50	2293
Buchanan, <i>De Maria Scotorum regina</i>	17	1047
Castiglione, <i>Il libro del cortegiano</i>	91	1181
Castiglione, <i>The covrtyer of Covnt Baldessar Castilio</i>	389	5222
Domenichi, <i>Facetie, motti, et burle &amp; Guicciardini, Detti et fatti piacevoli et gravi</i>	421	8283
Freigius, <i>Paratitla</i>	60	2818
Frontinus, <i>The strategems</i>	224	5493
Livy, <i>Ab urbe condita</i>	652	6972
Machiavelli, <i>Art of warre</i>	238	5414
Melanchton, <i>Selectarum declamationum</i>	133	1710
Olaus, <i>Historia de gentibus septentrionalibus</i>	121	2360
Smith, <i>De recta &amp; emendata linguæ Anglicæ</i>	105	1537
<b>Total</b>	<b>2501</b>	<b>5465</b>

Table 1.3 demonstrates, further, that the AOR XML dataset, no matter its objective size, may generate a range of meaningful insights. For example, it corroborates a pattern that is also visible in Table 1.2; namely, that some books are much more lightly annotated than others, in terms of the relative number of tags generated and words captured. In table 1.3, this is also demonstrated by the average size of the transcriptions within a given book. A quantitative analysis, in this latter context, may prove helpful, but it cannot formulate the basis of an argument that relies on the seeming “size” of the dataset to carry off or prove by default a qualitative point of analysis on that score. It is the quality of the data, rather, and the ways in which they are structured to support user interface with the data, that constitute the major asset of AOR as a digital asset.

During the first year of AOR Phase 1, transcribers have not only generated comprehensive transcriptions and translations of all manuscript annotations in the corpus of 12 books, but have also checked one another’s transcriptions independently in order to ensure the accuracy and quality of the data provided in the encoding process. Because of the relatively large size and complexity of this data within a humanities context, that final, secondary review process for checking transcriptions will require further effort within the second year of Phase 1. Currently there are 969 transcriptions (38.7% of the total number of transcriptions) which still need to be checked in this way, as demonstrated in Table 1.4, below.

**Table 1.4 Number of files that remain to be checked in AOR Phase 1, year 2**

<b>Title book</b>	<b>Number of files which need to be checked</b>
Castiglione, <i>The covrtyer of Covnt Baldessar Castilio</i>	389
Domenichi, <i>Facetie, motti, et burle &amp; Guicciardini, Detti et fatti piacevoli et gravi</i>	119
Frontinus, <i>The strategems</i>	224
Livy, <i>Ab urbe condita</i>	237
<b>Total</b>	<b>969</b>

The items in Table 1.4 reflect those books in the Harvey corpus for which the transcriptions have been completed, but still require independent cross-checking. There are an additional c. 75 remaining questions regarding humanistic content within completed transcriptions of annotations that are still unresolved due to ambiguous or complex translations, the presence of as yet indecipherable, or otherwise unknown non-verbal marks or symbols, as well as particular references to books or historical figures that are still difficult to identify at present. Our goal in AOR Year 2 is to resolve all these remaining tasks of the humanities team, and complete this portion of the annotations work by March 31, 2015.