Challenging the Copia

Ways to a Successful Big Data Analysis of Eighteenth-Century Magazines and Treatises on Art Connoisseurship¹

Joris Corin Heyder

Introduction

Being able to compute big data in a relatively short span of time is one of the greatest advantages of DH projects. While it is almost impossible to critically read and analyze a vast corpus of c. 550 titles with an average of 300 pages, a fully digitized corpus ideally combined with linguistic indexing is relatively easy to examine in respect of a particular semantic field and structure, specific notions, word co-occurrences, and more. However, what if the corpus is not available as machine-encoded text, but only in form of text-images from scanned documents in varying, sometimes very poor qualities? As conventional optical character recognition (OCR) has a high error rate particularly in cases of the older literature used in my project, where it is for instance necessary to differentiate for example between the "long s" (I) and "f" characters, the machine encoded results are often hardly useable for further operations. Particular difficulties arise out of the diversity of the underlying material that comprises connoisseurial and philosophical French and German treatises published between 1670 and 1850 as well as magazines like the Mercure de France or the Teutsche Merkur of the same period. This is true on different levels, first, the quality and resolution of the scanned text-images, second,

¹ This article has been written within the framework of the Collaborative Research Center SFB 1288 "Practices of Comparing. Changing and Ordering the World", Bielefeld University, Germany, funded by the German Research Foundation (DFG), subproject C 01, "Comparative Viewing. Forms, Functions and Limits of Comparing Images". My heartfelt thanks to Julia Becker for proofreading my manuscript for any linguistic and stylistic errors.

the typographical wide range from German black letter print to French Antiqua, and, third, the high quantity of texts in magazines that has been written on other topics than art connoisseurship. In this paper I am seeking to propose a heuristic how to tackle both, the diverse conditions of the once digitized material as well as the potentials offered by open source corpus analysis toolkits, such as AntConc.² I will argue for a combination of a "quick and dirty" approach to mass digitization with a complementary reflecting and recontextualizing close reading.

1. Research design and problems in collecting a data corpus

When I started my project on the practices of comparing in eighteenth-century French and German art connoisseurship by first gathering a huge amount of resources from well-known and less well-known connoisseurial treatises, exhibition guides, magazines, etc., from time to time, I had to think about the etymological roots of the word "copy" that originates from the Latin copia, which meant "abundance" or "richness of material". In medieval manuscript culture the French term "copie" has been extended to its current meaning as "duplication, imitation, or reproduction." However, both sides of the coin are present in looking at the plentiful material which I had to tackle after only a short span of time. On the one hand, I had the great opportunity to find even the rarest printed material in searching engines like Gallica³ or the digital library of the Bayerische Staatsbibliothek.⁴ On the other hand, this *copia* (in the sense of "abundance") was just too vast to scour for the examples I was seeking to find, namely, explicit and implicit reflections on the role of comparing visually in art connoisseurship. The many thousands of scanned pages, or to be more precise, digital copies of distinctive genera of publications were a curse and a savior at the same time. In pre-digital times, before the tempting offer of a veritable Babylonian

² The software can be downloaded here: https://www.laurenceanthony.net/software.html [accessed: 08.05.2019].

³ https://gallica.bnf.fr/accueil/fr/content/accueil-fr?mode=desktop [accessed: 08.05.2019].

⁴ https://www.digitale-sammlungen.de/index.html?c=sammlungen_kategorien&l=de [accessed: 08.05.2019].

library,⁵ 'natural' limitations were stipulated by the possessions of the libraries themselves. Of course, Gallica represents nothing else than the immense possessions of the Bibliothèque nationale de France and 270 other French institutions like Bibliothèques municipals or Instituts de recherche,6 but everyone who ever worked in such institutions knows how difficult it sometimes can be to receive each volume one is interested in. Processes like ordering the copy, waiting for the keeper to bring it, as well as institutional frames like delimited opening hours, restricted use, and so on, limit the opportunities to see a great amount of material. Therefore, research designs are usually reduced to a certain corpus of texts and this corpus may grow and change over time, but this only happens within an adequate framework. The contrary is believed to be true with digital material that promises the user a prolific disposability of every thinkable relevant source from all over the world. This, of course, is no more than a phantasmagory. However, the potential to make every thinkable information instantaneously available led Paul Virilio (1932–2018) to the idea of the "digital ages as 'the implosion of real time".7 Silke Schwandt has taken the idea a little further and asked: "Do we live in a time of 'eternal now'"?8 What almost brings a transcendent taste to the discussion is yet worth being applied to a concrete example. A time of 'eternal now' would mean, amongst other things, that - from a digital perspective a treatise written by Roger de Piles (1635–1709) has to be considered as present as the latest breaking news. Always entwined into a multiplicity of time layers – for instance their time of origin, their proper time or their historical time -, the tremendously synchronicity of information in itself epitomizes a huge problem. The user might get overwhelmed by the potential presence of the 'eternal now' but has to choose after all, which "now" comes first. Even worse, some of the 'nows' are more present than others, some even remain invisible: One could suggest that the claimed universality of resources will never become more than a chimera.

⁵ Cf. Borges, Jorge Luis, The Library of Babel, in: Ficciones, transl. by Anthony Kerrigan et al., New York: Grove Press, 1962.

⁶ Although 'Gallica' is labeled as a service of the Bibliothèque national de France, its partners are distributed over all regions of France, cf. https://gallica.bnf.fr/html/decouvrirnos-partenaires [accessed: 23.04.2019].

⁷ Schwandt, Silke, Looking for 'Time' and 'Change': Visualizing History in the Digital Age (draft version, forthcoming), 12.

⁸ Ibid.

Thus, working with digital texts needs clear frameworks, too. Often, they are set more or less hazardously by the resources themselves: How easy are they to find? Is the access to the digitized resources restricted? Is it possible to download the files, and if yes, is the quality of the scans sufficient for the needs of OCR or not? Of course, these limitations should not and cannot be decisive to handle the flood of information. How, then, is it possible to cope with the *copia*? An answer to this allegedly simple question can potentially be found in Petrarch's request to pledge yourself to *sufficientia* ('moderation'). In chapter 43 entitled "De librorum copia" of his work *De remediis utriusque fortunae*9 the allegorical figures *Gaudium* ('Joy') and *Ratio* ('Reason') dispute the challenges of the *copia* of books a long time before the invention of the letter press. ¹⁰ It was not the excess of accessible information but rather the unsuspecting use that motivated Petrarch to formulate his media critical remarks. He warns against the growing quantity of information and recommends being moderate in its use.

2. Digitizing the corpus

"Tene mensuram" ('be moderate') perhaps should also have been the motto in the beginning of my project, but it was not. Instead, I optimistically started the research – as has already been mentioned – by bringing together hundreds of PDF files with more or less relevant content in expectation of a quick OCR process. My naïve optimism was guided by the experiences with commercial OCR software that I used for almost all of my scanned secondary literature. As the average of the results appeared to be valuable, I presumed that an improved OCR engine developed by our IT-team would produce comparable results with regard to the texts I was working on. The idea was to establish a quick and dirty approach that allowed me to find relevant

⁹ A fully digitized version of c. 1490 probably printed by the editor Heinrich Knoblochtzer can be found here: http://mdz-nbn-resolving.de/urn:nbn:de:bvb:12-bsb11303461-1 [accessed: 21.04.2019].

¹⁰ Siegel, Steffen, Tabula: Figuren der Ordnung um 1600, Berlin: Akademie-Verlag, 2009, 31.

¹¹ Motto of German Emperor Maximilian I (1493-1519).

¹² Kohle, Hubertus, Digitale Bildwissenschaft, Glückstadt: Verlag Werner Hülsbusch. Fachverlag für Medientechnik und -wirtschaft, 2013, 37–38. For a full digitized version, cf. http://archiv.ub.uni-heidelberg.de/artdok/volltexte/2013/2185 [accessed: 08.04.2019].

words, co-occurrences, as well as characteristic semantic markers for comparisons, such as "plus ancienne/plus jeune", "copié d'après", "pareillement/different", "moins avantageuse que", etc. I intuitively estimated that it would be sufficient to reach an average text recognition of c. 80 percent to be able to run first tests with applications like *AntConc*. Only much later I learned that in machine learning based projects that use handwritten text recognition like Transkribus¹³ a character error rate (CER) of less than five percent on average could be achieved.

Several months went by until the IT-team could first establish an optimized version of the free OCR software Tesseract.¹⁴ It implemented specific requirements, for instance, the recognition of old font types like Gothic/ Fraktur fonts as well as layout detection. 15 However, a major problem consisted in the scans itself: As Sven Schlarb of the Österreichische Nationalbibliothek (Austrian National Library) has shown in a talk at the final IMPACT conference (Improving Access to Text)¹⁶ in 2010, an optimal scan is crucial to obtain optimum results. Of course, besides the font question, typical challenges of historical material appeared in most of the PDFs like warped book pages, curved text lines, different print intensities, distortions and contaminants, handwritten annotations, complex layouts, and of course time-specific orthography. 17 From a digital analytical perspective the PDFs had to be parsed. One item or one 'now' – the PDF – had not only to be dissolved into images or n-'nows' but also by help of a functional extension parser into a multiplicity of elements, like the text block, page numbers, characteristic layout features, and so forth. These items, then, had to be enhanced by different steps like border detection, geometric correction ('unwarping'), as well as binarization, i.e., the conversion of a picture into black and white. Although those processes are automatable to a certain degree, they maybe take the longest time. It is worth illustrating that with the different outcomes gained within the digitization process: A characteristic quality acces-

¹³ Cf. https://transkribus.eu/Transkribus/ [accessed: 13.06.2019].

¹⁴ https://en.wikipedia.org/wiki/Tesseract_(software) [accessed: 08.04.2019].

¹⁵ For a precise description of these operations, cf. the paper by Patrick Jentsch and Stephan Porada in this volume.

¹⁶ https://impactocr.wordpress.com/2010/05/07/an-overview-of-technical-solutions-in-impact/ [accessed: 08.04.2019].

¹⁷ https://impactocr.wordpress.com/2010/05/07/an-overview-of-technical-solutions-in-impact/ [accessed: 08.04.2019].

sible at Gallica, present in an example taken from the magazine *Mercure de France* from August 1729 (ill. 1a) promises prosperous results.

Ill. 1a & 1b

A O U S T. 1729. 1786 0 1 QUESTION proposée dans le Mer-1 L'" cure de May 1729. page 986. La perfection est-elle plus difficile a acquerir dans la Peinture que dans la Sculpture? RE'PONSE. L'Experience décide; elle nous apprend que la moitié des Etudes d'un Peintre, suffit pour devenir Sculpteur. L'étenduë de la Peinture est immense, la Sculpture a ses bornes. Tout ce qui est visible est du ressort de la Peinture : combien de parties dans la Nature se sont-elles soustraites de la puissance de la Sculpture ? L'air, les variations infinies que chaque saison y produit, sa sérenité, son opacité, ses éloignemens, ses Phénomenes, les Nuës, les differens degrez des lumieres du Soleil & de la Lune, les Arbres, les Campagnes, les Eaux, &c. toutes parties considerables dont la Sculpture est exempte de faire les études ; que de temps, que de peines épargnées : austi voyons - nous dans la Sculpture trente Eleves arriver en peu de temps à un éminent degré, lorsque dans la Peinture un seul d'un grand nom-

The scan is clearly legible, the page is not warped, the page layout appears to be not too difficult to "read". A first test with a commercial OCR performing PDF-reader shows, however, the limited capacity of such applications ever more drastically to us. In fact, by using the OCR tool with the options

"French" and "dpi dissolution 600" (ill. 1b) and by transferring the PDF to .txt-format afterwards zero result is achieved.

The effects of a first binarization (ill. 2a) demonstrate two things: first, an optimized image enormously increases the success of the 'reading' process, and, second, despite such achievements the effects appear to remain difficult to optimize (ill. 2b).

Ill. 2a & 2b

```
A O U S T. 1729. 1786
1 51 • ge 9·.-.6 •
QUESTION proposée dans le Mer-
      cure de May 1729. page 986.
La perfection est-elle plus difficile a ac-
                                                                       per ion tif-e ! plu1 difficile a t'C... u \simir d tils 1 P intur qu dans l Scitlptur ?
  querir dans la Peinture que dans la
                  Sculpture ?
                                                                       F. 'P F.
               RE'PONSE.
L'Experience décide; elle nous apprend
                                                                       · Ile no~ t1S a ~ re11
                                                                       e 1
' P' 1e 1ce ~ décJ
    que la moitié des Etudes d'un Pein-
tre, suffit pour devenir Sculpteur. L'éten-
duë de la Peinture est immense, la Sculp-
ture a ses bornes. Tout ce qui est visible est
                                                                       s Ettl es · un 1ein1
1 jri '
du ressort de la Peinture : combien de
parties dans la Nature se sont-elles sous-
traites de la puissance de la Sculpture ?
L'air, les variations infinies que chaque
                                                                       ( u pr L r. 1 r -n~~

o r J *

u · · 1 ; j 11 tL r i

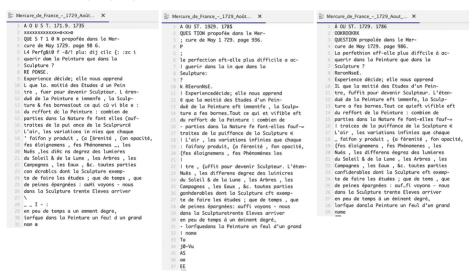
[e, la culpr

-. s o îl s · _o t c e ui fi illble
saison y produit, sa sérenité, son opacité,
ses éloignemens, ses Phénomenes, les
Nuës, les differens degrez des lumieres
du Soleil & de la Lune, les Arbres, les
Campagnes, les Eaux, &c. toutes parties
                                                                       o t c la <u>inrur</u> · cot _i n e
a r · _ 11 s la r t-1 -s {ou{:
ic s e la i n d la Scu : ture ~
considerables dont la Sculpture est exemp-
te de faire les études ; que de temps , que
de peines épargnées : aussi voyons - nous
                                                                       J'air le a <u>iations</u> 1t1fini s qu\sim cl1 .qu n i , ( (ér \cdot ' , on pactt ', '1 i n 1 s " 1 'not 1 s les
dans la Sculpture trente Eleves arriver
en peu de temps à un éminent degré.
lorfque dans la Peinture un feul d'un grand
                                                                       , 1 es 1 \cdot n s \cdot e re z s \underline{lutn} i re il e \underline{\mbox{un}} , 1es ; . r \cdot , les
                                                                       c 1 ..
```

With the OCR function of the commercial PDF reader at least (and of all things) the word "difficile" was recognizable for the machine. Keeping those unsuccessful data in mind, the results of the optimized Tesseract OCR

(ill. 3a-c) go far beyond any expectations. But this is only true for anyone knowing the enormous difficulties visible in illustrations 1b and 2b. When I studied the outcomes of our first OCR tests I was shocked. Even if many words (ill. 3a), as for instance "Aoust" or "Peinture" in line one and six, were readable, others like "visible" in line thirteen became illegible because of the disassembling of the readable and the unreadable syllables and/or signs into two or more parts: "vi ble".

Ill. 3a & 3b & 3c



What was even worse was the recognition of cursive letters; the disintegration of the word "perfection" – that forms part of line five – into the heap of signs: "Perfg&i0" ensued a concrete problem for the spotting of words. This is true because one would only gain a decisive outcome by considering the word stem "Perf" that is far too vague for a target-oriented research. Better results could be achieved by a second, revised Tesseract-version (ill. 3b).

While the italics caused fewer problems, as the correctly recognized example "perfection" proves, also the word "vifible" is now represented as one integral word. Here, the long "s" is 'read' as "f", which is a common but manageable problem in the underlying text genera of eighteenth-century French prints. Remarkably, the integrality of the text is distracted by a bizarre line adjustment. In reality, the twelfth line "e que la moitié des Etudes d'un Pein"

should have continued with line twenty two: "| tre, {uffit pour devenir Sculpteur. L'éten-", but instead is followed by line thirteen "duë de la Peinture eft immentfe, la Sculp=", which actually should have succeeded line twenty two. For the spotting of words, such an error is hardly decisive, but it is for word co-occurrences. We could gain the best possible result with a third, improved version of Tesseract (ill. 3c).

While numbers can now be 'read' with more precision by the machine, and the lines almost exactly imitate the actual content of the page, other aspects – like the recognition of italics – came out worse again. All in all, one can nevertheless say that such a result permits the application of software solutions like AntConc or Voyant¹⁸ in the intended scope. On the level of time alone, this positive outcome came at an expensive cost, as the process of enhancing, reading and evaluating meanwhile took approximately six minutes per page, which brought me to the following calculation: at that time, given a sum of 470 PDFs with an average of 300 pages, the working process for one PDF alone would take 1800 minutes or 30 hours. Of course, an average processing power of approximately 846.000 minutes or 14.100 hours was simply unfeasible in light of all the other projects maintained by our IT-team. Therefore, we had to decide, whether or not the project would have to be continued. What we were looking for was a return to a manageable number of 'nows'.

3. Strategies to reduce and to expand the corpus

The guiding question was how the extensiveness of the text corpus could be reduced as much as possible. One part of the solution was as trivial as it was analogue. It turned out that for our purposes only a smaller part of the magazine's text corpus was truly useful. Texts on art typically constitute only a very small percentage of journals like the *Mercure de France*¹⁹ or the

¹⁸ https://voyant-tools.org [accessed: 24.04.2019].

¹⁹ An overview on the digitized versions of the Mercure the France between 1724 and 1778 can be found here: http://gazetier-universel.gazettes18e.fr/periodique/mercure-de-france-1-1724-1778 [accessed: 27.05.2019]. The issues printed between 1778 and 1791 can be found here: http://gazetier-universel.gazettes18e.fr/periodique/mercure-de-france-2-1778-1791 [accessed: 27.05.2019].

Teutsche Merkur.²⁰ We, therefore, decided to let a student assistant²¹ preselect all art relevant passages in the magazines by usually starting from the table of contents. Of course, not every magazine still had its table of contents, which is why skimming the entire magazine was in any case obligatory. Nevertheless, we preferred to run the risk of overlooking texts to gather a relevant text quantity within a reliable period of investigation rather than being stuck with only a few digitized magazines that would have been the output of the limited resources we had for the whole project. Finally, we could bring together relevant material from the Mercure the France, printed between 1724 and 1756,²² and from the Teutsche Merkur, published between 1773 and 1789. Admittedly, these periods of time do not illustrate art connoisseurship over a whole century, but they can, nonetheless, exemplify the prolific stages of connoisseurial practices in eighteenth-century France and Germany. Needless to say, the choice of the time spans also depends on capacity factors such as available digital resources and the workforce offered by the SFB.

Which other methods were available to reduce the corpus? It unquestionably appeared to be rather counter-intuitive to preselect passages in treatises on art and aesthetics. Although greater parts of, for example, Kant's *Critique of Judgement*²³ may have no relevance for the importance of practices of comparing in eighteenth-century connoisseurship, others could be all the more inspiring. A preselection here means to skip maybe the most interesting references. Then again, also in the field of significant art treatises it is possible to value more and less relevant material; to start the digitization with the most discussed ones proved to be a convenient approach.

A useful strategy not to reduce but to expand the corpus, however, is the concentration on eighteenth-century art treatises and/or magazines that have already been digitized referring to TEI-guidelines (Text Encod-

²⁰ An overview on the digitized versions of the *Teutsche Merkur* between 1773 and 1789 can be found here: http://ds.ub.uni-bielefeld.de/viewer/toc/1951387/1/LOG_0000/ [accessed: 17.04.2019].

²¹ I would like to thank our student assistant Felix Berge for his help.

²² The Mercure the France usually comprises six issues per year. In consequence, thirty years make a sum of almost two hundred issues with an average of four hundred pages per issue

²³ Kant, Immanuel, Critique of Judgement, ed. by Nicholas Walker, transl. by James Creed Meredith, Oxford [u. a.]: Oxford University Press, 1952.

ing Initiative). ²⁴ They allow the user to analyze works in depth, for instance, in terms of a semantical structure, the recognition of a specific typology of comparisons, ²⁵ and so forth. Unfortunately, only a small percentage of seventeenth- and eighteenth-century treatises and magazines with a focus on art and aesthetics available as scanned images or in PDF version is also accessible in fully digitized and linguistically marked .txt- or .xml-formats. For the research subject in question, databases like the *Observatoire de la vie littéraire*²⁶ or the *Deutsches Textarchiv*²⁷ offer a number of fully digitized texts that are central for eighteenth-century connoisseurship, as for instance Jean-Baptiste Dubos's *Réflexions critiques sur la poésie et la peinture*²⁸ or Johann Joachim Winckelmann's *Geschichte der Kunst des Alterthums*. ²⁹ It may heuristically make sense to use such versions to exemplarily shed light on the reflections of practices of comparing itself, while the material digitized by means of dirty OCR could rather be used for quantitative questioning or to track down passages that could be of interest for the project. ³⁰

Of course, the more digitized, evaluated and corrected resources are at disposal, the more comprehensive the outcome will after all probability be. In case of the *Mercure the France*, for instance, the *Observatoire de la vie littéraire* offers a fully digitized and evaluated version of the forerunner magazine named *Mercure galant*. This magazine was published between 1672 and 1710, and today 465 issues are available for a deeper analysis. Given the hypothesis that reflections on practices of comparing grew in the course of the eighteenth-century, one could expect an increase of certain key notions, as for instance "comparaison/comparer" or "jugement/juger", or the co-occurrence of such notions. Therefore, it seems utterly useful to expand the corpus of the *Mercure de France* by the corpus of the *Mercure galant* to visualize the keyword's use over the years. Such a juxtaposition, however, brings

²⁴ https://tei-c.org/guidelines/[accessed: 25.04.2019].

²⁵ A good example for such an approach is Olga Sabelfeld's contribution in this volume.

²⁶ http://obvil.sorbonne-universite.site [accessed: 23.05.2019].

²⁷ http://www.deutschestextarchiv.de [accessed: 23.05.2019].

²⁸ https://obvil.sorbonne-universite.fr/corpus/critique/dubos_critiques [accessed: 23.05.2019].

²⁹ http://www.deutschestextarchiv.de/book/show/winckelmann_kunstgeschichteo1_ 1764 [accessed: 23.05.2019].

³⁰ See for instance the projects by Christine Peters and Malte Lorenzen in this volume.

³¹ https://obvil.sorbonne-universite.fr/corpus/mercure-galant/ [accessed: 23.05.2019].

its own challenges, because data repositories like the *Observatoire de la vie littéraire* are equipped with their own search function but cannot be extrapolated to download .txt-files (or other data formats). Although it is possible to obtain the data by means of the source code, such an indirect download process, once again, takes quite some time. Another problem lies in the comparability of the data itself because on the one hand, in addition to art critiques the downloaded data in the *Mercure galant* also comprises articles on poetry, theatre, history, politics and disputes discussed in the *Académie*. The repository of texts from the *Mercure the France* used in our project, on the other hand, includes only those parts on art critique. While it was necessary to reduce the corpus for processing dirty OCR, it now proves itself a disadvantage, because the much bigger corpus of the *Mercure galant* would have to be reduced accordingly. It is questionable, though, whether or not such a time exposure is worth the result.

This is why a tentatively carried out first comparison between the hits for the root word "compar" in the *Mercure galant* and the *Mercure the France* has only proceeded by its total numbers. Unfortunately, little can be concluded from the "concordance plot" for the regular expression search with the root word "compar". From the 465 issues of the *Mercure galant*, evidently, some issues have much more hits than others but a steady increase can not be established (ill. 4a, b).

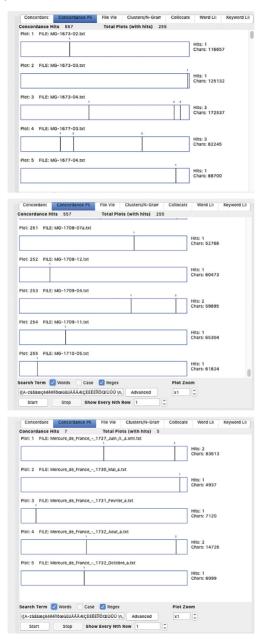
Evidently, the number of hits in the first years around 1673 is as high as around 1710. The same is true for the search in the extracted contributions on art critique in the *Mercure de France* (ill. 4c).

Here, only five out of thirty-six issues have a hit for the root word in question.³³ The situation in the *Mercure galant* is more or less similar given a hit rate of 255 out of 465 issues and considering that only a smaller part of the content is devoted to art critiques as such. Only one out of nine hits evident in the five plots of ill. 4a represents the root word "compar" in a text on art. All other hits are related to different topics such as theatre or music. At first

³² The regular expression used for the search in AntConc was provided by Stephan Porada. It is: "([A-zàâäæçèéêëîīôœùûüÀÂÄÆÇÈÉÊËÎÎÔŒÙÛÜ \n,;;'-\d\(\)><-]*|compar[a-zàâäæçèéêëîīôœùûüÀÂÄÆÇÈÉËËÎÎÔŒÙÛÜ\n]*)compar[a-zàâäæçèéêëîîôœùûüÀÂÄÆÇÈÉËËÎÎÔŒÙÛÜ]* [A-zàâäæçèéêëîîôœùûüÀÂÄÆÇÈÉËÏÎÔŒÙÛÜ \n,;;'-\d(\)><-]*". The root word "compar" can be exchanged by any other root word.

³³ When I tested the extracted text passages in AntConc, I could only refer to the first thirty-six out of altogether 110 articles on art published between June 1724 and October 1754.

Ill. 4a & 4b & 4c



sight, the result of the comparison appears to be disappointing but it is not given the astonishing outcome of a stability of hits in the course of time. This can either mean that the starting hypothesis on the increase of reflections on practices of comparing in the course of the eighteenth-century is probably wrong or that it cannot be shown by means of the use of words that indicate such a reflection. Another error source could perhaps be linked to the length of the examined periods, which might have been too short for substantial statements.

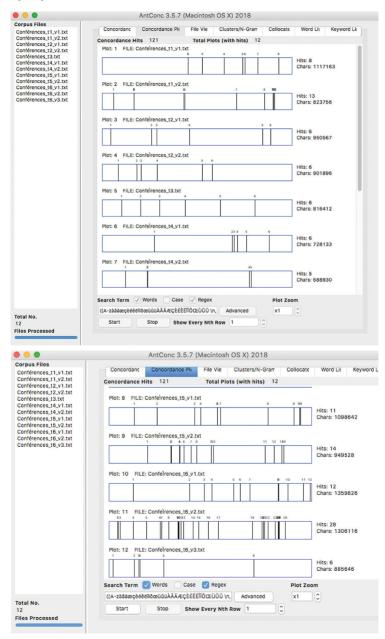
In other cases, it is possible to quantify the increasing use of notions, as can be seen in the following example. In the course of a contribution to the conference "Media of Exactitude" held in Basel in 2018, 34 I firstly asked for the occurrence of the word "exactitude" (and its derivates) and, secondly, the co-occurrence³⁵ of the words "comparer/comparaison" (and its derivates) and "exactitude" (and its derivates) in the Conférences de l'Académie Royale de Peinture et de Sculpture published between 1648 and 1792. My decision for choosing the lectures of the Conférence was due to the fact that they reflect in a way interests, fashions and focal points of well-known French connoisseurs for more than a century. Moreover, they are accessible as digitized, edited, and optical-character-recognized data.36 Although the result is far from being representative, in the chosen span of time the word "exactitude" (ill. 5a, b) was increasingly used, as can be proved by comparing the hits in the early years (ill. 5a) with those in the later years (ill. 5b). The same is true for the co-occurrence of the word "comparison" (ill. 5c, d), which is optionally followed by the word "exactitude". What is remarkable is the fact that in both cases only by the middle of the century an increase becomes visible.

³⁴ For the conference, cf. https://www.genauigkeit.ch/talk/conference/[accessed:12.05.2019].

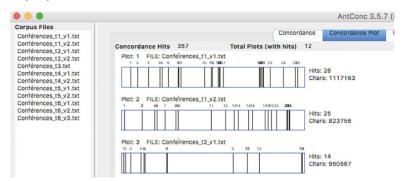
³⁵ The regular expression used for the search in AntConc was provided by Stephan Porada. It is: "([A-zàâāæçèéêëîîôœùûüÀÂÄÆÇÈÉÊËÎÎÔŒÙÛÜ \n,;;-\d\(\)><-]*|compar[a-zàâāæçèéêëîîôœùûüÀÂÄÆÇÈÉÊËÎÎÔŒÙÛÜ\n,;;-\d\(\)><-]*|compar[a-zàâāæçèéêëîîôœùûüÀÂÄÆÇÈÉÊËÎÎÔŒÙÛÜ\n,;;-\d\(\)><-]*)?(juge[a-zàâāæçèéêëîîôœùûüÀÂÄÆÇÈÉÊËÎÎÔŒÙÛÜ\n]*|[A-zàâāæçèéêëîîôœùûüÀÂÄÆÇÈÉÊËÎÎÔŒÙÛÜ\n],;;-\d\(\)><]*)? [A-zàâāæçèéêëîîôœùûüÀÂÄÆÇÈÉÊËÎÎÔŒÙÛÜ\n],;;-\d\(\)><]*". The root words "compar" and "juge" can be exchanged by any other root words.

³⁶ The entire body of lectures given at the Parisian Académie royale de peinture et de sculpture were made available by an editing research project at the German Center for Art History (DFK Paris). For the fully digitized corpus, cf. https://dfk-paris.org/en/research-project/editing-and-publishing-conférences-de-l'académie-royale-de-peinture-et-de [accessed: 23.05.2019].

Ill. 5a & 5b



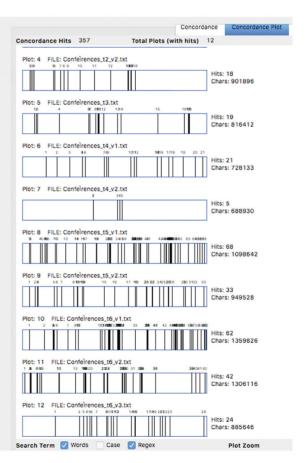
Ill. 5c & 5d



Corpus Files

Conférences_t1_v2.txt
Conférences_t2_v2.txt
Conférences_t2_v1.txt
Conférences_t2_v2.txt
Conférences_t3.txt
Conférences_t3.txt
Conférences_t4_v2.txt
Conférences_t4_v2.txt
Conférences_t5_v2.txt
Conférences_t5_v2.txt
Conférences_t5_v2.txt
Conférences_t6_v2.txt
Conférences_t6_v2.txt
Conférences_t6_v2.txt

Conférences_t6_v3.txt



4. How to handle first results?

First examinations of the isolated and OCR-processed passages in the Teutsche Merkur from 1773 to 1789 have shown that, for instance, the word "Vergleich/Vergleichung/vergleichen" can hardly be found. This holds particularly true because of the Fraktur, whose majuscule letters "B" and "V", for example, look very similar. Therefore, using string characters like "rgleich" instead of root words like "Vergleich" or "Bergleich" undoubtedly yields more results. While doing so, I stumbled upon a text, which had only one hit for "Bergleichung", but proved being a revealing example for practices of comparison in the close reading. The text "Ueber Christus und die zwölf Apostel, nach Raphael von Mark=Anton gestochen, und von Herrn Prof. Langer in Düsseldorf kopiert" was published in the fourth issue of the *Teutsche Merkur* in 1789.37 As an apparent advertisement, the text of an anonymous author stresses the advantage of a series of engravings after Marcantonio Raimondi's Apostles that shall help the beholder to refresh the vision of Raphael's ingenious inventions.³⁸ The print series was copied by the artist Johann Peter von Langer (1756-1824) shortly before 1789. While the greater part of the review describes the prints and inventions themselves, in a subsequent passage the author discusses the value of the copies. They would inspire a fresh appraisal of the prints:

"These sheets arguably give us insight into the notion of the value of the originals in regard to invention, posture, drapery, character of hairs and faces. We can safely claim that no amateur of arts should fail to purchase these copies by Langer, even if he as an exception already possessed the originals [i. e., the prints by Marcantonio Raimondi]. In that case, the copies would still give some food for thought like a good translation." 39

³⁷ Anonymus, Ueber Christus und die zwölf Apostel, nach Raphael von Mark=Anton gestochen, und von Herrn Prof. Langer in Düsseldorf kopiert, in: Teutscher Merkur 4 (1789), 269–277.

³⁸ The connection between Raphael and Raimondi was demonstrated in-depth by: *Bloemacher, Anne,* Raffael und Raimondi, Berlin [et al.]: Deutscher Kunstverlag, 2016.

³⁹ Translation by the author. The original quote is: "Diese Blätter gewähren also uns streitig einen Begriff von dem Werth der Originale in Absicht auf Erfindung, Stellung, Wurf der Falten, Charakter der Haare und der Gesichter, und wir dürfen wohl sagen, daß kein Liebhaber der Künste versäumen sollte, sich diese Langerischen Copien anzuschaffen,

At a time when the concept of originality evolved into the most important scheme for artistic achievements, 40 such praise for efficacy of the copy is rather astonishing. Moreover, the author underlines the fact that it can be enlightening to compare the original prints by Raimondi with the copies by Langer to expose still more the creators' artistic understanding and their i.e., Raphael's and Raimondi's - light and fortunate nature. 41 The call for a visual comparison ("Vergleichung") is followed by a meticulous analysis of the differences between Langer's copies and the original prints (ill. 6a, 6b); it brings out the tendency of connoisseurship to not only get the most complete possible overview but particularly to fragment the objects of research into small comparable entities. Langer's prints were used as a foil to foster the ideal-typical execution and planning of the original prints right into the folds and hatchings ("In den Originalen ist keine Falte, von der wir uns nicht Rechenschaft zu geben getrauen"). This and many other aspects make the short review a promising example for applied connoisseurship, in terms of the contemporary expectations of comparisons, of medial considerations as well as of the explicit guidance for a comparative approach.

The chance is rather modest that I would have found this review in the seemingly endless issues of the *Mercure de France* and the *Teutsche Merkur*. I could, however, track it down very quickly thanks to dirty OCR and the hit list in the AntConc concordance plot. Other than quantitative assertions, a qualitative perspective necessarily requests a thorough recontextualization of the extracted file. In the discussed case, the result could not be more surprising, because the anonymously published text is well known in the domain of research on Johann Wolfgang Goethe's art critical writings.⁴² The review was

selbst in dem seltenen Falle wenn er die Originale besäße; denn auch alsdann würden ihm diese Copien, wie eine gute Uebersetzung, noch manchen Stoff zum Nachdenken geben." cf. *Anonymus*, Ueber Christus und die zwölf Apostel, nach Raphael von Mark=Anton gestochen, und von Herrn Prof. Langer in Düsseldorf kopiert, 275.

- 40 Cf. for instance: *Mortier, Roland*, L'originalité: une nouvelle catégorie esthétique au siècle des lumières. Histoire des idées et critique littéraire, Genf: Droz, 1982.
- 41 The original quote is: "bey dem größten Kunstverstand, ein so leichtes und glückliches Naturell ihrer Urheber, daß sie uns wieder unschätzbar vorkommen." cf. Anonymus, Ueber Christus und die zwölf Apostel, nach Raphael von Mark=Anton gestochen, und von Herrn Prof. Langer in Düsseldorf kopiert, 275.
- 42 Cf. Osterkamp, Ernst, Bedeutende Falten. Goethes Winckelmann-Rezeption am Beispiel seiner Beschreibung von Marcantonio Raimondis Apostelzyklus, in: Thomas W. Gaehtgens (ed.), Johann Joachim Winckelmann, 1717–1768, (Studien zum achtzehnten

Ill. 6a & 6b





written by Goethe subsequent to his Italian Journey, and, apparently, the publisher of the *Teutsche Merkur*, Christian Martin Wieland, could temporarily provide him with the original prints by Marcantonio Raimondi. ⁴³ It is stunning to see that the text was already discussed with regard to practices of comparative vision in Johannes Grave's seminal work on Goethe as a collector of prints. One could therefore say that such a positive result proves the efficacy of the approach to record relevant hits for a vocabulary reflecting comparisons as comprehensively as possible in different kinds of sources. However, it will hardly be possible to recontextualize every hit – or in other words every 'now' – towards a multiplicity of 'nows'. But it is, of course, possible to selectively densify certain "nows". Such a process would also comprise a reference to the visual resource itself (ill. 6a, 6b). Today, Langer's prints are surprisingly difficult to find within

Jahrhundert, 7), Hamburg: Meiner Verlag, 1986, 265–288; Osterkamp, Ernst, Im Buchstabenbilde: Studien zum Verfahren Goethescher Bildbeschreibungen (Germanistische Abhandlungen, 70), Stuttgart: Metzler, 1991, 54–71; Grave, Johannes, Der 'ideale Kunstkörper': Johann Wolfgang Goethe als Sammler von Druckgraphiken und Zeichnungen (Ästhetik um 1800, 4), Göttingen: Vandenhoeck+Ruprecht, 2006, 240–243.

⁴³ J. Grave, Der 'ideale Kunstkörper', 240.

the image repositories provided by libraries, museums, and so forth. This circumstance evokes the lack of any illustration in publication formats like the *Teutsche Merkur* and it underlines the fact that Goethe's review actually might have prompted readers to make the purchase. It also shows how much it differs from what Peter Bell⁴⁴ has baptized 'digital connoisseurship': while the potentials of machine learning image recognition rest on a basis of hundreds of thousands of images online, in Goethe's time for the majority of people prints were still a rare commodity that had to be assembled in cumbersome collections. It could be said that at that time every comparison had its own value.

Conclusion

The here proposed heuristic focuses not only on one specific method but seeks to combine different approaches, namely big data analysis and close readings, in reaction to both the diverse condition of the digitized material as well as the potentials offered by already fully digitized and OCR-processed open source material. One goal was to bring together as much material as possible to be able to trace changes in a long-term perspective with regard to the number of hits of vocabulary reflecting comparisons. At the same time, close reading enables recontextualization of the gathered bits and pieces and to switch from a quantitative to a qualitative argument. The application of a big data approach has been proven as a reliable 'good nose' for texts that could be crucial for a project on practices of comparison in connoisseurship. It turned out that dealing with the 'eternal nows' not only helped us to strategically preselect the material but also to encourage entirely new questions. In the long term and also in light of a cost/benefit ratio it seems to have been worth the effort, because the established digital library allows even many more ways to question the once assembled material. It could be a next step to publish the library of OCR-processed texts on art and aesthetic and to share them with the scientific community. Such plans, however, have to be implemented in accordance with legal obligations, financial aftercare requirements, and so forth, so that one is reminded anew of the motto: "Tene mensuram".

⁴⁴ Bell, Peter/Ommer, Björn, Digital Connoisseur? How Computer Vision Supports Art History, in: Stefan Albl/Alina Aggujaro (eds.), Connoisseurship nel XXI secolo. Approcci, Limiti, Prospettive, Rome: Artemide, 2016, 187–197.

Bibliography

- Anonymus, Ueber Christus und die zwölf Apostel, nach Raphael von Mark=Anton gestochen, und von Herrn Prof. Langer in Düsseldorf kopiert, in: Teutscher Merkur 4 (1789), 269–277.
- Bell, Peter/Ommer, Björn, Digital Connoisseur? How Computer Vision Supports Art History, in: Stefan Albl/Alina Aggujaro (eds.), Connoisseurship nel XXI secolo. Approcci, Limiti, Prospettive, Rome: Artemide, 2016, 187–197.
- Bloemacher, Anne, Raffael und Raimondi, Berlin [et al.]: Deutscher Kunstverlag, 2016.
- Borges, Jorge Luis, The Library of Babel, in: Ficciones, transl. by Anthony Kerrigan et al., New York: Grove Press, 1962.
- Grave, Johannes, Der 'ideale Kunstkörper': Johann Wolfgang Goethe als Sammler von Druckgraphiken und Zeichnungen (Ästhetik um 1800, 4), Göttingen: Vandenhoeck+Ruprecht, 2006.
- Kant, Immanuel, Critique of Judgement, ed. by Nicholas Walker, transl. by James Creed Meredith, Oxford [u. a.]: Oxford University Press, 1952.
- Kohle, Hubertus, Digitale Bildwissenschaft, Glückstadt: Verlag Werner Hülsbusch. Fachverlag für Medientechnik und -wirtschaft, 2013.
- *Mortier, Roland, L'* originalité: une nouvelle catégorie esthétique au siècle des lumières. Histoire des idées et critique littéraire, Genf: Droz, 1982.
- Osterkamp, Ernst, Im Buchstabenbilde: Studien zum Verfahren Goethescher Bildbeschreibungen (Germanistische Abhandlungen, 70), Stuttgart: Metzler, 1991.
- Osterkamp, Ernst, Bedeutende Falten. Goethes Winckelmann-Rezeption am Beispiel seiner Beschreibung von Marcantonio Raimondis Apostelzyklus, in: Thomas W. Gaehtgens (ed.), Johann Joachim Winckelmann, 1717–1768, (Studien zum achtzehnten Jahrhundert, 7), Hamburg: Meiner Verlag, 1986, 265–288.
- Schwandt, Silke, Looking for 'Time' and 'Change': Visualizing History in the Digital Age (draft version, forthcoming).
- Siegel, Steffen, Tabula: Figuren der Ordnung um 1600, Berlin: Akademie-Verlag, 2009.

Illustrations

- Ill. 1a: Page from the Mercure de France, Aoust 1729, 1785 © http://gazetier-universel.gazettes18e.fr/periodique/mercure-de-france-1-1724-1778 [last access: 12.6.2019].
- Ill. 1b: .txt-file gained with a commercial OCR tool by 'reading' ill. 1a © Screenshot by the author
- Ill. 2a: Binarized page from the Mercure de France, Aoust 1729, 1785 © Screenshot by the author
- Ill. 2b: .txt-file gained with a commercial OCR tool by 'reading' ill. 2a © Screenshot by the author
- Ill. 3a: .txt-file gained with a first version of an improved Tesseract OCR tool © Screenshot by the author
- Ill. 3b: .txt-file gained with a second version of an improved Tesseract OCR tool © Screenshot by the author
- Ill. 3c: .txt-file gained with a third version of an improved Tesseract OCR tool © Screenshot by the author
- Ill. 4a: AntConc concordance plot of hits for the occurrence of the root word "compar" in the *Mercure galant* between 1673–1677 © Screenshot by the author
- Ill. 4b: AntConc concordance plot of hits for the occurrence of the root word "compar" in the *Mercure galant* between 1708–1710 © Screenshot by the author
- Ill. 4c: AntConc concordance plot of hits for the occurrence of the root word "compar" in the *Mercure de France* between 1727–1732 © Screenshot by the author
- Ill. 5a: AntConc concordance plot of hits for the occurrence of the root word "exactitude" in the *Conférence* between 1648–1746 © Screenshot by the author
- Ill. 5b: AntConc concordance plot of hits for the occurrence of the root word "exactitude" in the *Conférence* between 1747–1792 © Screenshot by the author
- Ill. 5c: AntConc concordance plot of hits for the co-occurrence of the root words "compare" and "exactitude" in the *Conférence* between 1648–1746 © Screenshot by the author

- Ill. 5d: AntConc concordance plot of hits for the co-occurrence of the root words "compare" and "exactitude" in the *Conférence* between 1747–1792 © Screenshot by the author
- Ill. 6a: St Paul, Johannes Peter von Langer (after a print by Marcantonio Raimondi), etching, 216×135 mm (sheet), c. 1789, Wolfenbüttel, Herzog August Bibliothek, Graph.A1:1470 © http://diglib.hab.de?grafik=graph-a1-1470 [last access: 12.6.2019].
- Ill. 6b: St Paul, Marcantonio Raimondi (after an invention by Raphael), etching, 215 × 141 mm (sheet), c. 1520, Wolfenbüttel, Herzog August Bibliothek, MRaimondi AB 3.26 © http://kk.haum-bs.de/?id=raim-m-ab3-0026 [last access: 12.6.2019].