

# OCR4all: Modern tool for old texts

April 24 2019

---



Ce qui fut iadis viciuey  
 Criminel ignominieuy  
 Aboli laisse reprouue  
 Est par les aultres approuue  
 Maintenant en ce present aage  
 Par les aultres nouuel vsaige  
 Nouveau rit coustume nouuelle  
 Est gardee la chose est telle  
 Mais ie ne puis pas bien penser  
 En mon cueur ne contrepenser  
 Lequel est le plus fol des deux  
 Usant du nouueau rit et vieuy  
 Du cil qui prent les grans coudeyz  
 Comme font vng tas de lourdeyz  
 Du cil qui porte manches larges  
 Comme font maintenant les paiges  
 Si non de dire quil me semble  
 Que mectre se peuent ensemble

Mat. xvij.

Comme deux folz car cest tout vng  
 Si lung est bien noir lautre est brun  
 Se lung est fol lautre lest plus  
 Si lung boiteuy lautre perclus  
**Jadis** estoit grande louange  
 Qui maintenant seroit estrange  
 Aux anciens peres porter  
 Grande barbe et devez noter  
 Que a le temple de socrates  
 Tous les philosophes apres  
 Et auant quil fust mort portoient  
 La barbe grant lesquelz estoient  
 Rempliz de grande sapience  
 Mais eulz decedez leur science  
 Sop voyant par nous contempnee  
 Sen est lassus au ciel volee  
 Nous laissant ca bas tous inbertes  
 Dont sont innartables les pertes  
**Libidinite** corumpue  
 A commence faire repue  
 Et tenir son cours par le monde  
 Toutes vertuz et chose monde  
 Dont les haultz cieulz sont decorez  
 Ny sont plus mais sont demorez  
 Tous maulz tous vices et pechez  
 Dont les humains sont entachez  
**Tout** le monde se contrefait  
 Et veullent ce que dieu a fait  
 Par presumption contrefaire  
 En cuidant mieulz que dieu les faire  
 Qui est vng peche par trop grant  
 Honteuy sont et honte les prant  
 De porter grant barbe au vsaige  
 De peur de monstret leur dieulz aage  
 Mais leurs corps et viz si bien gardent  
 Si bien les acoutrent et fardent  
 Que iamais ne deuiennent dieulz  
 Se semble et aussi leurs cheueuy  
 Les vngz comme sicambriens  
 Et comme les ethiopiens  
 Les portent tous crasspes et tors  
 Faisant a nature grans tors

ff. Regi. x.

Socrates

Sicambri.  
Ethiopes.

can be reliably converted into computer-readable text with OCR4all. Credit: Dresden State and University Library, CC BY-SA 4.0

Historians and other humanities' scholars often have to deal with difficult research objects: centuries-old printed works that are difficult to decipher and often in an unsatisfactory state of conservation. Many of these documents have now been digitized—usually photographed or scanned—and are available online worldwide. For research purposes, this is already a step forward.

However, there is still a challenge to overcome: bringing the digitized old fonts into a modern form with text recognition software that is readable for non-specialists as well as for computers. Scientists at the Center for Philology and Digitality at Julius-Maximilians-Universität Würzburg (JMU) in Bavaria, Germany, have made a significant contribution to further development in this field.

With OCR4all, the JMU research team is making a new tool available to the scientific community. It converts digitized historical prints with an [error rate](#) of less than one percent into computer-readable texts. And it offers a [graphical user interface](#) that requires no IT expertise. With previous tools of this kind, user-friendliness was not always a given, as the users mostly had to work with programming commands.

## **Developed in cooperation with the humanities**

The new OCR4all tool was developed under the direction of Christian Reul together with his computer science colleagues Professor Frank Puppe (Chair of Artificial Intelligence and Applied computer science) and Christoph Wick as well as Uwe Springmann (Digital Humanities expert) and numerous students and assistants.

OCR4all originates from the JMU Kallimachos project, which is funded by the German Federal Ministry of Education and Research. This cooperation between the humanities and computer science will be continued and institutionalized in the newly founded JMU Center for Philology and Digitality.

In developing OCR4all, computer scientists have collaborated with the humanities at JMU—including German and Romance studies and literature studies in the project "Narragonien digital." The aim was to digitize the "Narrenschiff," a moral satire by Sebastian Brant, a bestseller of the 15th century that was translated into many languages. Furthermore, OCR4all has been frequently used in the JMU's Kolleg "Medieval and Early Modern Times."

OCR4all is freely available to the public on the GitHub platform (with instructions and examples): <https://github.com/OCR4all>

## **Each print shop had its own font**

Christian Reul explains the challenges involved in the development of OCR4all: Automatic text recognition (OCR = Optical Character Recognition) has been working very well for modern fonts for some time now. However, this has not yet been the case for historical fonts.

"One of the biggest problems was typography," says Reul. One of the reasons for this is that the first printers of the 15th century did not use uniform fonts. "Their printing stamps were all carved by themselves, each printing house practically had its own letters."

## **Error rates below one percent**

Whether "e" or "c," whether "v" or "r"—it is often not easy to

distinguish in old prints, but software can learn to recognize such subtleties. To do so, it has to be trained on sample material. In his work, Reul has developed methods to make training more efficient. In a case study with six historical prints from the years 1476 to 1572, the average error rate in automatic text recognition was reduced from 3.9 to 1.7 percent.

Not only was the methodology improved, JMU computer scientist Christoph Wick has also decisively further refined the technical component by developing the Calamari OCR tool, which is also freely available and has since been fully integrated into OCR4all, promising even better results. Now, even for the oldest printed works, error rates of less than one percent can be achieved in general.

## **Lexical projects**

Reul has also convinced external partners of the quality of Würzburg's OCR research. In cooperation with the "Zentrum für digitale Lexikographie der deutschen Sprache" (Berlin), Daniel Sanders' "Wörterbuch der deutschen Sprache" (Dictionary of the German Language) has been digitally indexed, and a scientific publication on this work is currently being prepared. The various lines of this text often contain different fonts, representing different semantic information. Here, the existing approach to character recognition was extended in such a way that not only the text but also the typography and thus the complex content structure of the lexicon may be reproduced very precisely.

The computer scientist from Würzburg will soon complete his doctoral thesis, but he is also willing to continue working with OCR in the future: "The computer science behind OCR is extremely exciting," he says. A possible project in the near future: the creators of the "Idiotikon," a dictionary of the Swiss-German language, have indicated their interest in

collaboration since they might well need the Würzburg's specialist knowledge.

**More information:** [github.com/OCR4all](https://github.com/OCR4all)  
[github.com/Calamari-OCR](https://github.com/Calamari-OCR)

[jlcl.org/content/2-allissues/1 ... 18/jlcl\\_2018-1\\_1.pdf](https://jlcl.org/content/2-allissues/1...18/jlcl_2018-1_1.pdf)

[jlcl.org/content/2-allissues/1 ... 18/jlcl\\_2018-1\\_4.pdf](https://jlcl.org/content/2-allissues/1...18/jlcl_2018-1_4.pdf)

Provided by University of Würzburg

Citation: OCR4all: Modern tool for old texts (2019, April 24) retrieved 27 April 2024 from <https://phys.org/news/2019-04-ocr4all-modern-tool-texts.html>

This document is subject to copyright. Apart from any fair dealing for the purpose of private study or research, no part may be reproduced without the written permission. The content is provided for information purposes only.