

Mozilla gibt Sprachdatensammlung frei

[Originalartikel](#)

[Backup](#)

<html> <p class=„printversionback-to-article printversion-hide“><a href=„<https://www.heise.de/newsticker/meldung/Mozilla-gibt-Sprachdatensammlung-frei-4323042.html>“>zurück zum Artikel</p><figure class=„printversionlogo“><img src=„<https://1.f. ix.de/icons/svg/logos/svg/heiseonline.svg>“ alt=„heise online“ width=„180“ height=„40“></figure><figure class=„aufmacherbild“></figure><p>Die bisher durch das Common-Voice-Projekt zusammengetragenen Sprach-Samples lassen sich ab sofort frei verwenden. Es ist die größte Sammlung dieser Art.</p> <p>Mozilla hat seine Sprachdatensammlung <a href=„<https://voice.mozilla.org/>“ rel=„external noopener“ target=„_blank“>Common Voice [1] öffentlich freigegeben. Mit 1361 Stunden – das entspricht knapp zwei Monaten – transkribierten Audiodaten ist es nach Angaben von Mozilla die größte frei zugängliche Sammlung der Welt. Ebenso wichtig wie die Größte ist Mozilla die Vielseitigkeit der Samples: 42.000 Sprecher wirkten daran mit und sprachen kurze Texte in 18 verschiedenen Sprachen ein.</p> <p>Die Veröffentlichung steht unter CC0-Lizenz – der freizügigsten Variante von Creative Commons („No rights reserved“) – der Öffentlichkeit zur Verfügung. Hauptziel der Sammlung ist es, hochwertige und frei verfügbare Sprachdatensätze zum Training für Spracherkennungssysteme zu schaffen – ein Gebiet, das bisher Cloud-Anwendungen großer Konzerne mit riesigen Sprachdatensammlungen dominieren. Mit <a href=„<https://www.heise.de/meldung/Mozilla-Common-Voice-Sprachsteuerung-fuer-alle-und-ohne-Rueckgriff-auf-die-Cloud-3904454.html>“>DeepSpeech [2] entwickelt Mozilla eine eigene Open-Source-Spracherkennung, die bereits in Produkten wie <a href=„<https://mycroft.ai/blog/deepspeech-update/>“ rel=„external noopener“ target=„_blank“>Mycroft [3] oder <a href=„<https://getleon.ai/>“ rel=„external noopener“ target=„_blank“>Leon [4] eingesetzt oder getestet wird.</p> <p>Das Projekt startete <a href=„<https://www.heise.de/meldung/Mozilla-sammelt-Stimmaufzeichnungen-fuer-offene-Spracherkennungs-Software-3780795.html>“>Mitte 2017 [5] mit einer englischsprachigen Textsammlung; ein Jahr später öffnete sich Common Voice für andere Sprachen. Für Englisch hat Mozilla 685 Stunden von fast 36.000 Sprechern aufgezeichnet; Deutsch folgt auf Platz zwei mit 254 Stunden, an denen knapp 4000 Freiwillige mitwirkten.</p> <div class=„inread“> <h3 class=„subheading“ id=„nav_kabylisch0“>Kabylisch, Tatarisch, Walisisch</h3> <p>Wärend sich kommerzielle Anbieter auf die Sprachen der wichtigsten Märkte konzentrieren, finden sich bei Common Voice auch viele, die sonst kaum im Internet repräsentiert sind, etwa Kabylisch (eine algerische Berbersprache), Tatarisch oder Walisisch.

Hier treiben oft wenige Enthusiasten das Projekt voran. Neuerdings kooperiert Mozilla mit der [**Deutschen Gesellschaft für Internationale Zusammenarbeit \[6\]**](https://www.giz.de/en/html/), um zum Beispiel Sprecher in dem afrikanischen Land Ruanda zu erreichen. Einige der großen Weltsprachen hinken dagegen noch hinterher, etwa Spanisch, Arabisch oder Russisch. Seit der Release finalisiert wurde, wuchs die Zahl der Sprachen in der Aufnahmephase auf 22 an; fast 200 Stunden Aufzeichnungen kamen hinzu. Bei 70 weiteren Sprachen liegt die Vorbereitungsphase, in der die Freiwilligen setzen und die Website bersetzen. Auch wenn Deutsch in Common Voice gut vertreten ist, sucht das Projekt weiterhin Sprecher; das erklärte Ziel ist, für jede Sprache 1200 Stunden Material zu sammeln. Die Mitwirkung erfordert keine besonderen Kenntnisse und dauert nur wenige Minuten. Siehe dazu auch c't 18/2018: [**Spracherkennung**](https://www.heise.de/select/ct/2018/18/1535420071631846) für alle: Mozillas Projekte Common Voice und DeepSpeech [7]

Artikel: <http://www.heise.de/-4323042>

Artikel: <https://voice.mozilla.org/>

Artikel: <https://www.heise.de/meldung/Mozilla-Common-Voice-Sprachsteuerung-fuer-alle-und-ohne-Rueckgriff-auf-die-Cloud-3904454.html>

Artikel: <https://getleon.ai/>

Artikel: <https://www.heise.de/meldung/Mozilla-sammelt-Stimmaufzeichnungen-fuer-offene-Spracherkennungs-Software-3780795.html>

Artikel: <https://www.giz.de/en/html/>

Artikel: <https://www.heise.de/select/ct/2018/18/1535420071631846>

Copyright © 2019 Heise Medien

From: <https://schnipsl.qgelm.de/> - **Qgelm**

Permanent link: <https://schnipsl.qgelm.de/doku.php?id=wallabag:mozilla-gibt-sprachdatensammlung-frei>

Last update: 2021/12/06 15:24

