

# Mozilla DeepSpeech: Speech-to-Text Schritt für Schritt

[Originalartikel](#)

[Backup](#)

<html> <header class=„article-header“><h1 class=„articleheading“>Mozilla DeepSpeech: Speech-to-Text Schritt f&#252;r Schritt</h1><div class=„publish-info“> Pascal Moll</div></header><figure class=„aufmacherbild“><img src=„[https://heise.cloudimg.io/width/700/q75.png-lossy-75.webp-lossy-75.foil1/\\_www-heise-de\\_/imgs/18/3/1/0/7/5/4/6/shutterstock\\_754914430-99670d63283369b8.jpeg](https://heise.cloudimg.io/width/700/q75.png-lossy-75.webp-lossy-75.foil1/_www-heise-de_/imgs/18/3/1/0/7/5/4/6/shutterstock_754914430-99670d63283369b8.jpeg)“ srcset=„[https://heise.cloudimg.io/width/700/q75.png-lossy-75.webp-lossy-75.foil1/\\_www-heise-de\\_/imgs/18/3/1/0/7/5/4/6/shutterstock\\_754914430-99670d63283369b8.jpeg](https://heise.cloudimg.io/width/700/q75.png-lossy-75.webp-lossy-75.foil1/_www-heise-de_/imgs/18/3/1/0/7/5/4/6/shutterstock_754914430-99670d63283369b8.jpeg) 700w,[https://heise.cloudimg.io/width/1050/q75.png-lossy-75.webp-lossy-75.foil1/\\_www-heise-de\\_/imgs/18/3/1/0/7/5/4/6/shutterstock\\_754914430-99670d63283369b8.jpeg](https://heise.cloudimg.io/width/1050/q75.png-lossy-75.webp-lossy-75.foil1/_www-heise-de_/imgs/18/3/1/0/7/5/4/6/shutterstock_754914430-99670d63283369b8.jpeg) 1050w,[https://heise.cloudimg.io/width/1500/q75.png-lossy-75.webp-lossy-75.foil1/\\_www-heise-de\\_/imgs/18/3/1/0/7/5/4/6/shutterstock\\_754914430-99670d63283369b8.jpeg](https://heise.cloudimg.io/width/1500/q75.png-lossy-75.webp-lossy-75.foil1/_www-heise-de_/imgs/18/3/1/0/7/5/4/6/shutterstock_754914430-99670d63283369b8.jpeg) 1500w,[https://heise.cloudimg.io/width/2300/q75.png-lossy-75.webp-lossy-75.foil1/\\_www-heise-de\\_/imgs/18/3/1/0/7/5/4/6/shutterstock\\_754914430-99670d63283369b8.jpeg](https://heise.cloudimg.io/width/2300/q75.png-lossy-75.webp-lossy-75.foil1/_www-heise-de_/imgs/18/3/1/0/7/5/4/6/shutterstock_754914430-99670d63283369b8.jpeg) 2300w“ alt=„ class=„img-responsive“ referrerpolicy=„no-referrer“ /><figcaption class=„akwa-caption“>(Bild:&#160;petrmalinak/Shutterstock.com)</figcaption></figure><p><strong>Dieses Tutorial zeigt anhand eines Praxisbeispiels, wie sich ein Sprachassistent mit DeepSpeech auf einem Raspberry Pi erstellen l&#228;sst.</strong></p><p>Intelligente Lautsprecher oder Sprachassistenten sind auf dem Vormarsch. <a href=„<https://www.heise.de/meldung/Fast-jeder-Dritte-in-Deutschland-nutzt-Sprachassistenten-4443365.html>“<strong>Mehr als ein Drittel [1]</strong></a> aller Deutschen nutzt sie. Doch was passiert eigentlich mit den dadurch erzeugten Daten? Es gibt verschiedene Vermutungen, dennoch l&#228;sst sich nicht sicher sagen, was mit den gesprochenen Informationen geschieht. Um die volle Kontrolle &#252;ber diese Daten zu behalten, entschied sich der Autor, selbst eine Offline-Anwendung mit Java zu entwickeln.</p><p>Dieses Tutorial soll den Speech-to-Text (STT)-Anteil des Entwicklungsprozesses ausschnittsweise vorstellen. Es folgen Einblicke in die Installation der DeepSpeech Engine mithilfe eines Raspberry Pi. Diese Komponenten nehmen Gesprochenes entgegen und &#252;berf&#252;hrt das Gesprochene in Text. Im n&#228;chsten Schritt erfolgt die Einf&#252;hrung des zugeh&#246;rigen Java-Programmcodes, um mit diesem Text zu interagieren und Sprachbefehle zu designen. Der Ausblick zeigt zahlreiche Erweiterungsm&#246;glichkeiten und Einsatzgebiete.</p><h3 class=„subheading“ id=„nav\_einf&#252;hrung\_in0“>Einf&#252;hrung in Mozilla DeepSpeech</h3><p>Mozillas DeepSpeech ist eine freie Speech-to-Text-Engine. Sie arbeitet mit einem durch Maschine Learning erstellten Sprachmodell, basierend auf den <a href=„<https://arxiv.org/abs/1412.5567>“ rel=„external noopener“ target=„\_blank“><strong>Forschungsergebnissen von Baidu&#8217;s Deep Speech Research Paper [2]</strong></a> und ermöglicht somit, Gesprochenes in Text umzuwandeln. Hier kommt unter anderem Googles Machine-Learning-Framework TensorFlow zum Einsatz, <a href=„<https://deepspeech.readthedocs.io/en/r0.9/>“ rel=„external noopener“ target=„\_blank“><strong>was die Implementierung vereinfachen soll [3]</strong></a>. TensorFlow Lite eignet sich besonders f&#252;r mobile Ger&#228;te, da es deutlich kleinere Sprachmodelle bei vergleichbaren Ergebnissen und Geschwindigkeit mitbringt. Mozilla arbeitet daran, verschiedene Modelle in anderen Sprachen zu generieren. Darüber nutzen sie das Projekt Common Voice. Das Ziel dieses Crowdsourcing-Projekts ist, eine freie Datenbank mit verschiedenen Sprachen und Sprechern aufzubauen.</p><p>Mitwirkende können entweder andere Spracheingaben validieren oder selbst Sprachaufnahmen beisteuern. Zur Teilnahme wird lediglich ein Mikrofon

ben&#246;tigt. F&#252;r die englische Sprache ist dieses Projekt schon weit entwickelt, so ist die Fehlerquote sehr gering und selbst f&#252;r nicht Muttersprachler mit Dialekt vielseitig anwendbar. Das deutsche Sprachmodell funktioniert ebenfalls gut, hat aber noch Schwierigkeiten bei W&#246;rtern, die &#228;hnlich klingen. Beispielsweise wird „Licht an“ oft als „Lichter“ erkannt.

</p><h3 class=„subheading“ id=„nav\_ben&#246;tigte1“>Ben&#246;tigte Komponenten und Voraussetzungen</h3><p>F&#252;r das STT-Projekt sind diese Komponenten notwendig:</p><ul class=„rtelist rtelist-unordered“><li>Raspberry Pi 4 (ein PI3 B+ ist ausreichend, kann aber zu Verz&#246;gerungen bei der Ergebnisr&#252;ckgabe f&#252;hren)</li><li>Java 8 mit Maven</li><li>Python 3.5 auf dem Raspi</li><li>Mikrofon (je nach Entfernung empfiehlt sich ein Ringmikrofon)</li></ul><h3 class=„subheading“ id=„nav\_deepspeech\_als2“>DeepSpeech als Server aufsetzen</h3><p>Zun&#228;chst sollte das Standardbetriebssystem „Raspberry PI OS“ installiert sein. Anschlie&#223;nd empfiehlt es sich, die gesamte Speicherplatte als Speicher &#252;ber das Konfigurationsmen&#252; bereitzustellen. Somit entfallen m&#246;gliche Speicherprobleme zu einem späteren Zeitpunkt, beispielsweise bei der Verwendung eines gr&#246;&#223;eren Sprachmodells.</p><figure class=„a-inline-image a-u-inline“><div><figcaption class=„a-caption“>Aufruf der Config mit sudo raspi-config (Abb. 1)</figcaption></div></figure><p>&#220;ber Advanced Options (Punkt 7) und A1 Expand Filesystem l&#228;sst sich die gesamte Gr&#246;&#223;e der Speicherplatte nutzen. Dar&#252;ber hinaus ist das Aktivieren des Secure File Transfer Protocol (SSH) oder Virtual Network Computing (VNC) zu empfehlen, um nicht immer wieder externe Peripherie f&#252;r die Bedienung anschlie&#223;en zu m&#252;ssen.</p><p>Nach erfolgreichem Abschluss der Grundinstallation gilt es, DeepSpeech zu installieren, was &#252;ber den Konsolenbefehl

```
pip3 install deepspeech
```

geschieht. Nach wenigen Minuten Wartezeit sollte die Erfolgsmeldung erscheinen:</p><figure class=„a-inline-image a-u-inline“><div><figcaption class=„a-caption“>Erfolgsmeldung: DeepSpeech erfolgreich installiert (Abb. 2)</figcaption></div></figure><p>Mit

```
pip3 install deepspeech-server
```

wird nun der Server installiert. Dieser Vorgang kann einige Minuten in Anspruch nehmen.</p><figure class=„a-inline-image a-u-inline“><div><figcaption class=„a-caption“>Ben&#246;tigte Sprachmodellkomponenten (Abb. 3)</figcaption></div></figure><p>Anschlie&#223;nd ist DeepSpeech installiert, jedoch ist noch kein Sprachmodell hinterlegt. Pre-Trained-Modelle lassen sich <a href=„<https://github.com.mozilla/DeepSpeech/releases>“ rel=„external noopener“ target=„\_blank“><strong>via GitHub herunterladen [4]</strong></a>. F&#252;r das deutsche Modell <a href=„[https://drive.google.com/drive/folders/1o-N-VH\\_0P89fcRKWEUIVDm-\\_z18Kbkb](https://drive.google.com/drive/folders/1o-N-VH_0P89fcRKWEUIVDm-_z18Kbkb)“ rel=„external noopener“ target=„\_blank“><strong>empfiehlt sich jedoch diese Adresse [5]</strong></a>. Dieses Modell beinhaltet 1582 Sprach-Stunden. Ben&#246;tigt werden das Modell TensorFlow lite und der Scorer.</p><figure class=„a-inline-image a-u-inline“><div><figcaption class=„a-caption“>Scorer und Tflite Datei in /home/pi/.local/bin abgelegt (Abb. 4)</figcaption></div></figure><p>Beide Dateien gilt es herunterzuladen und in der Konfiguration zu hinterlegen. Am einfachsten ist dies mit dem VNC Viewer m&#246;glich.</p><p>Im letzten Schritt ist es notwendig, die Konfiguration zu hinterlegen. Hierzu dient eine neue Datei namens <em>config.json</em>.</p><pre class=„rtetx-listing listing“><code>Config.json{ „deepspeech“: { „model“ : „output\_graph\_de.tflite“, „scorer“ : „kenlm\_de.scorer“, „beam\_width“: 500, „lm\_alpha“: 0.931289039105002, „lm\_beta“: 1.1834137581510284 }, „server“: { „http“: { „host“: „0.0.0.0“ , „port“: 8080, „request\_max\_size“: 1048576 } }, „log“: { „level“: [ { „logger“: „deepspeech\_server“, „level“: „DEBUG“} ] } }</code></pre><p>Zum Abschluss der Konfiguration ist nun ein Neustart des

Raspi erforderlich.</p><h3 class=„subheading“ id=„nav\_starten\_des3“>Starten des DeepSpeech-Servers</h3><p>Um den Server zu starten, ist im Verzeichnis <em>cd /home/pi/.local/bin</em> der Steuerbefehl

```
deepspeech-server --config config.json
```

einzugeben.</p><figure class=„a-inline-image a-u-inline“><div><figcaption class=„a-caption“>DeepSpeech-Server erfolgreich gestartet (Abb. 5)</figcaption></div></figure><p>Sollte der Fehler <em>Original error was: libf77blas.so.3: cannot open shared object file: No such file or directory</em> auftauchen, hilft der Befehl:

```
sudo apt-get install libatlas-base-dev
```

.</p><p><strong>DeepSpeech-Aufruf</strong></p><p>Um DeepSpeech per Curl aufzurufen, wird die IP des Raspberry-Servers benötigt sowie der Curl-Befehl:

```
curl -X POST --data-binary @testfile.wav http://IP:8080/stt
```

</p><h3 class=„subheading“ id=„nav\_deepspeech4“>DeepSpeech verlangt eine Sounddatei</h3><p>Aktuell existieren keine nativen Java Libraries, um DeepSpeech zu nutzen. Es gibt lediglich eine Android-Bibliothek, die für Desktop-Anwendungen jedoch wenig hilfreich ist. Daher ist es erforderlich, DeepSpeech als Server einzusetzen. Somit lässt sich eine Verbindung über HTTP aufbauen und ein eigenes Java-Programm entwickeln. Der Curl-Befehl lässt es schon vermuten: DeepSpeech verlangt eine Sounddatei im WAV-Format, die entgegengenommen, interpretiert und dann als Text zurückgeliefert wird. Das nachfolgende Programm geht genauso vor.</p><h3 class=„subheading“ id=„nav\_blick\_in\_die5“>Blick in die Bibliothek</h3><p>Zunächst sind zwei Bibliotheken notwendig. Beide befinden sich im Maven Central Repository und lassen sich einfach ins Project Object Model (POM) einbinden:</p><pre class=„rtetx-listing listing“><code>&lt;dependencies&gt; &lt;dependency&gt; &lt;groupId&gt;org.apache.httpcomponents&lt;/groupId&gt; &lt;artifactId&gt;httpclient&lt;/artifactId&gt; &lt;version&gt;4.5.13&lt;/version&gt; &lt;/dependency&gt; &lt;dependency&gt; &lt;groupId&gt;commons-io&lt;/groupId&gt; &lt;artifactId&gt;commons-io&lt;/artifactId&gt; &lt;version&gt;2.8.0&lt;/version&gt; &lt;/dependency&gt; &lt;/dependencies&gt;</code></pre><h3 class=„subheading“ id=„nav\_ans\_eingemachte6“>Ans Eingemachte mit Java</h3><p>Der nachfolgende Code ist vereinfacht dargestellt: Exception Handling wurde bewusst weitestgehend vernachlässigt. Der Client besteht aus den drei Methoden

```
transcribe()
```

,

```
createAudioOutputStream(...)
```

und

```
sendDataToServer(...)
```

.</p><p><strong>übertragen des Gesprochenen</strong></p><p>Die Methode

## transcribe()

&#246;ffnet den Mikrofoneingang und legt gleichzeitig einen Buffer für das Gesprochene an.

```
<p><pre class="rtetx-listing listing"><code>AudioFormat format = new AudioFormat(16000, 16, 1, true, false);</code></pre><p>Die Abtastrate des Eingangssignals wird hiermit auf 16000, die Samplesize auf 16 und die Kanalanzahl auf Mono festgelegt. Der <code>true</code>-Wert gibt an, dass das Eingangssignal signed ist, es kann also sowohl positive als auch negative Werte beinhalten. Das ist notwendig, da Sprachsignale eben solche Werte beinhalten. Der <code>false</code>-Wert steht für <code>littleEndian</code> und gibt die Speicherreihenfolge an; <code>true</code> steht für <code>bigEndian</code>.</p><p>Diese Einstellungen sind erforderlich, da die hinterlegten Sprachmodelle von DeepSpeech ebenfalls mit diesen Werten abgelegt wurden. Ein Wechsel von mono auf stereo beeinflusst die Erkennung negativ.</p><p>Schließlich wird das Mikrofon mit dem vorgegebenen Audioformat geöffnet und ist bereit Signale zu empfangen.</p><pre class="rtetx-listing listing">
```

```
DataLine.Info info = new DataLine.Info(TargetDataLine.class, format); if ( ! AudioSystem.isLineSupported(info) ) { throw new Exception("Mikrofon ist nicht verfügbar!"); } TargetDataLine line = (TargetDataLine) AudioSystem.getLine(info); line.open(format, line.getBufferSize());
```

Zum Speichern des Eingangssignales wird noch ein Buffer zur Verfügung gestellt:

```
<p><code>ByteArrayOutputStream out = new ByteArrayOutputStream(); byte[] data = new byte[line.getBufferSize() / 4];</code></pre><p>Nachdem alle Vorbereitungen abgeschlossen sind, lässt sich das Signal aufzeichnen und abspeichern:</p><pre class="rtetx-listing listing">
```

```
System.out.println("StartLine"); line.start(); long startTime = System.currentTimeMillis(); while ( ( System.currentTimeMillis() - startTime ) < 2000 ) { int numBytesRead = line.read(data, 0, data.length); out.write(data, 0, numBytesRead); } System.out.println("ende");
```

Der Eingang hält zwei Sekunden zu und speichert das Gesprochene in dem zuvor festgelegten Buffer. In einem nächsten Schritt wird der

## ByteBuffer

an die Methode

```
createAudioOutputStream(...)
```

weitergegeben. Das Ergebnis davon geht an

```
sendDataToServer(...)
```

```
.</p><pre class="rtetx-listing listing"><code>ByteArrayOutputStream byos = createAudioOutputStream(line, out); return sendDataToServer(byos);</code></pre><p><strong>Ausgabe in die Datei</strong></p><p>Die Methode <code>createAudioOutputStream</code> wandelt den Inhalt des
```

<code>ByteBuffer</code> in eine WAV-Datei um. Dazu berechnet sie die benötigte WAV-Größe, liest den <code>ByteBuffer</code> ein und schreibt schließlich eine WAV als <code>ByteArrayOutputStream</code> raus:</p><pre class=„rttx-listing listing“>

```
private ByteArrayOutputStream createAudioOutputStream(TargetDataLine line,
ByteArrayOutputStream out) throws IOException { int frameSizeInBytes =
line.getFormat().getFrameSize(); byte audioBytes[] = out.toByteArray();
AudioInputStream audioInputStream = new AudioInputStream(new
ByteArrayInputStream(out.toByteArray()),
line.getFormat(), audioBytes.length /
frameSizeInBytes); ByteArrayOutputStream byos = new
ByteArrayOutputStream(); AudioSystem.write(audioInputStream,
AudioFileFormat.Type.WAVE, byos); audioInputStream.close();
return byos;}
```

</pre><p><strong>Kommunikation mit DeepSpeech</strong></p><p>Nachdem nun die WAV-Datei als

## ByteArrayOutputStream

vorliegt, lässt sich eine Verbindung zum DeepSpeech Server aufbauen und die Datei dorthin senden:</p><pre class=„rttx-listing listing“><code>private String
sendDataToServer(ByteArrayOutputStream byos) throws IOException, ClientProtocolException {
HttpClient httpclient = HttpClient.createDefault(); HttpPost httppost = new
HttpPost(SERVERADRESS); ByteArrayEntity bae = new ByteArrayEntity(byos.toByteArray());
httppost.setEntity(bae); HttpResponse response = httpclient.execute(httppost); HttpEntity entity =
response.getEntity(); if (entity != null) { try (InputStream instream = entity.getContent()) { return
IOUtils.toString(instream, „UTF-8“); } } return „Kein Ergebnis“;}</code></pre><p>Die gezeigte
Methode baut zunächst eine Serververbindung auf und setzt anschließend einen
<code>Post</code>-Befehl ab. Die Serveradresse sieht beispielsweise so aus:
<em><http://localhost:8080/stt></em>. </p><p>Das zuvor erstellte <code>ByteArray</code> wird
nun in eine Entity verpackt und an den Server gesendet. Zum Abschluss erfolgt das Auslesen der
zurückgegebenen Response.</p><figure class=„a-inline-image a-u-inline“><div><figcaption
class=„a-caption“>DeepSpeech-Ergebnis von „Hallo Welt“ (Abb.
6)</figcaption></div></figure><figure class=„a-inline-image a-u-inline“><div><figcaption
class=„a-caption“>Eclipse Konsolenausgabe von „Hallo Welt“ (Abb.
7)</figcaption></div></figure><h3 class=„subheading“ id=„nav\_fazit\_und7“>Fazit und
Ausblick</h3><p>Der hier vorgestellte Prozess ist nur ein Ausschnitt von dem, was mit DeepSpeech
möglich ist. Der eigene Sprachassistent des Autors ist mittlerweile in der vierten Iteration und
ermöglicht es ihm, das Licht zu steuern, Termine aus dem Kalender abzufragen und YouTube
zu nutzen.</p><p>Problematisch sind englische Begriffe und Worte, die ähnlich klingen. Auch
wurde hier bewusst das Reagieren auf unterschiedliche Lautstärkepegel vernachlässigt.
Mit dem Audiodispatcher und Silencelistener der Tarsos-Bibliothek lässt sich diese Problematik
lösen, da der Sprachassistent erst dann reagiert, wenn ein bestimmter Lautstärkepegel
überschritten wird und auch erst dann aufhört zu hören, wenn dieser Pegel
wieder unterschritten wird. Damit ist eine fixe Zeitangabe für das Zuhören nicht mehr
notwendig.</p><p>Unberücksichtigt blieben eigene Sprachmodelle. Es besteht auch die
Möglichkeit, eigene Wörter oder eigene Modelle zu entwickeln und mit TensorFlow zu
trainieren. Dabei lassen sich Entferungen zum Mikrofon oder Umgebungsgeräusche
kompensieren. Die Möglichkeiten sind grenzenlos und stehen den gegenwärtigen
Sprachassistenten auf dem aktuellen Markt in nichts nach. Mit diesem Vorgehen lassen sich die

Sprachkompetenzen von DeepSpeech auch für Java nutzen und nicht mehr nur für Python.

Der vorgestellte Code [findet sich auf GitHub \[6\]](https://github.com/PMO-IT). Auf dem [eigenen Blog](http://pmo-it.de/blog/) zeigt der Autor [7] unter anderem die Möglichkeiten Philips Hue Leuchten mit Java zu steuern. Für den Sommer plant er seinen Sprachassistenten auf GitHub zur Verfügung zu stellen.

Pascal Moll ist freiberuflicher Berater und Java-Entwickler. Seine Schwerpunkte liegen im Bereich des Test-Managements und Testautomatisierung von Web - und Desktopapplikationen insbesondere SAP. Neben seiner Beratertätigkeit arbeitet er auch als freiberuflicher Trainer für Java, Cucumber und Selenium Schulungen. Mehr Informationen finden sich auf [https://pmo-it.de \[8\]](https://pmo-it.de).

( )

---

URL dieses Artikels:

<https://www.heise.de/-6048698>

</small></p><p><strong>Links in diesem Artikel:</strong><br /><small>

<strong>[1]</strong> https://www.heise.de/meldung/Fast-jeder-Dritte-in-Deutschland-nutzt-Sprachassistenten-4443365.html

</small><br /><small>

<strong>[2]</strong> https://arxiv.org/abs/1412.5567

</small><br /><small>

<strong>[3]</strong> https://deepspeech.readthedocs.io/en/r0.9/

</small><br /><small>

<strong>[4]</strong> https://github.com.mozilla/DeepSpeech/releases

</small><br /><small>

<strong>[5]</strong> https://drive.google.com/drive/folders/1o0-N-VH\_0P89fcRKWEU1VDm-\_z18Kbkb

</small><br /><small>

<strong>[6]</strong> https://github.com/PMO-IT

</small><br /><small>

<strong>[7]</strong> http://pmo-it.de/blog/

</small><br /><small>

```
<strong>[8]</strong>&#160;https://pmo-it.de
</small><br /><small>
<strong>[9]</strong>&#160;mailto:mdo@ix.de
</small><br /></p><p class=„printversion__copyright“><em>Copyright &#169; 2021 Heise  
Medien</em></p> </html>
```

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

Permanent link:  
[https://schnipsl.qgelm.de/doku.php?id=wallabag:wb2mozilla-deepspeech\\_-speech-to-text-schritt-fr-schritt](https://schnipsl.qgelm.de/doku.php?id=wallabag:wb2mozilla-deepspeech_-speech-to-text-schritt-fr-schritt)

Last update: **2025/06/27 11:17**

