



Autophon user guide Norwegian Bokmål

Engine: Montreal Forced Aligner 1.0

Model: NoFA 1.0

1 Introducing Autophon and Forced Alignment

Autophon is a free, user-friendly tool for phoneticians that performs forced alignment (FA) – the automated process of converting speech recordings and their transcriptions into phonetically time-stamped annotations.

Autophon leverages widely used alignment engines developed by the phonetics community, including:

- FAVE¹
- faseAlign²
- Montreal Forced Aligner, version 1.0³
- Montreal Forced Aligner, version 2.0³

The tool produces time-aligned phonetic annotations compatible with Praat⁴, based on two user inputs: (1) a speech recording and (2) its orthographic transcript.

This user guide is specifically for **Norwegian Bokmål**, using the **Montreal Forced Aligner 1.0** engine with the **NoFA 1.0** model. Autophon may support additional engine-model combinations for this language; therefore, ensure you are using the best option for your needs.

While many forced aligners exist, they often require command-line usage and are tied to outdated or incompatible operating systems. **Autophon offers a platform-independent, intuitive alternative for phoneticians worldwide.**

2 Using the app

2.1 Aligning files without registering To align smaller files, go to the main page and click Add files at the bottom. A box titled *Transcription Mode: change transcription mode* will appear. Click the heading to choose one of four *Transcription Modes* (see below), then select your files.

2.2 Registering and logging in To align larger files or access the full suite, click Sign up to create a free account. This helps us monitor usage for funders and guard against bots. After signing up, check your email for a verification link. If it doesn't arrive, check your spambox and wait 15 minutes before contacting tech support.

2.3 Cost Autophon is free of charge.

2.4 Aligning files in a registered account Once registered and verified, go to the Aligner tab and click Add files. A box titled *Transcription Mode: change transcription mode* will appear. Click the heading to choose one of four *Transcription Modes*, then select your files.

2.5 Transcription modes Autophon supports four *Transcription Modes*, named for the fields they're commonly used in: *Experimental Linguistics A*, *Experimental Linguistics B*, *Computational Linguistics*, and *Variationist Linguistics*. Each mode can be selected via the corresponding box in Figure 1, which illustrates expected file structures and links to instructional videos.

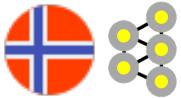
Video instructions for each transcription mode can be viewed. In addition, sample templates for each mode are available for [download here](#).

¹FAVE was built by Rosenfelder, Fruehwald, Brickhouse, Evanini, Seyfarth, Gorman, Prichard, and Yuan (2022). It relies on the Hidden Markov Toolkit (Young, Woodland, and Byrne 1993).

²faseAlign was built by Wilbanks (2022). Like FAVE, it relies on the Hidden Markov Toolkit (Young, Woodland, and Byrne 1993).

³The Montreal Forced Aligner was developed by McAuliffe, Socolof, Mihuc, Wagner, and Sonderegger (2017). It uses the Kaldi toolkit (Povey, Ghoshal, Boulianne, Burget, Gembek, Goel, Hannemann, Motlicek, Qian, Schwarz, et al. 2011).

⁴Praat is a speech analysis tool developed by Boersma and Weenink (2017).



Experimental Ling A (click to see video guide)	Experimental Ling B (click to see video guide)	Computational Ling (click to see video guide)	Variationist Ling (click to see video guide)
<pre>yourzip.zip - yourtrans.xlsx/tsv/txt - file0001.wav - file0002.wav - file0003.wav ... - file9999.wav</pre> <p><i>Transcriptions in a master file absent of time stamps - as separate rows with separate audio* files for each transcription.</i></p>	<pre>yourzip.zip - yourtrans.xlsx/tsv/txt - file01.wav - file02.wav - file03.wav ... - file99.wav</pre> <p><i>Transcriptions in a master file with start and end time stamps with more than one row per audio* file.</i></p>	<pre>yourzip.zip - file0001.lab - file0001.wav - file0002.lab - file0002.wav - file0003.lab - file0003.wav ... - file9999.lab - file9999.wav</pre> <p><i>Transcriptions as separate same-name lab and audio* files, absent of time stamps.</i></p>	<pre>yourzip.zip - file01.TextGrid - file01.wav - file02.eaf - file02.wav - file03.tsv - file03.wav - file04.xlsx - file04.wav ... - file99.txt - file99.wav</pre> <p><i>Longer transcription files in TextGrid, eaf, tsv, txt, or xlsx format with same-name audio* files.</i></p>

Figure 1: The Transcription Mode selection menu for Autophon.

daDK_small

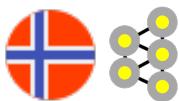
- ↳ X0297
 - > 1
 - ↳ 2
 - X0297-dk15-09082000-1715_u0295139-1.wav
 - X0297-dk15-09082000-1715_u0295140-1.lab
 - X0297-dk15-09082000-1715_u0295140-1.wav
 - X0297-dk15-09082000-1715_u0295141-1.lab
 - X0297-dk15-09082000-1715_u0295141-1.wav
 - X0297-dk15-09082000-1715_u0295142-1.lab
 - X0297-dk15-09082000-1715_u0295142-1.wav
 - X0297-dk15-09082000-1715_u0295143-1.lab
 - X0297-dk15-09082000-1715_u0295143-1.wav
 - X0297-dk15-09082000-1715_u0295144-1.lab
 - X0297-dk15-09082000-1715_u0295144-1.wav
 - X0297-dk15-09082000-1715_u0295145-1.lab
 - X0297-dk15-09082000-1715_u0295145-1.wav
 - X0297-dk15-09082000-1715_u0295146-1.lab
 - X0297-dk15-09082000-1715_u0295146-1.wav
 - X0297-dk15-09082000-1715_u0295147-1.lab
 - X0297-dk15-09082000-1715_u0295147-1.wav

daDK_small

- ↳ X0297
 - > 1
 - ↳ 2
 - X0297-dk15-09082000-1715_u0295140-1.TextGrid
 - X0297-dk15-09082000-1715_u0295141-1.TextGrid
 - X0297-dk15-09082000-1715_u0295142-1.TextGrid
 - X0297-dk15-09082000-1715_u0295143-1.TextGrid
 - X0297-dk15-09082000-1715_u0295144-1.TextGrid
 - X0297-dk15-09082000-1715_u0295145-1.TextGrid
 - X0297-dk15-09082000-1715_u0295146-1.TextGrid
 - X0297-dk15-09082000-1715_u0295147-1.TextGrid

- ↳ X0298
 - > 1
 - X0298-dk17-09082000-1822_u0296002-1.TextGrid
 - X0298-dk17-09082000-1822_u0296003-1.TextGrid
 - X0298-dk17-09082000-1822_u0296004-1.TextGrid
 - X0298-dk17-09082000-1822_u0296005-1.TextGrid
 - X0298-dk17-09082000-1822_u0296006-1.TextGrid
 - > 2
 - > 3
 - X0298-dk17-09082000-1822_u0296230-1.TextGrid
 - X0298-dk17-09082000-1822_u0296231-1.TextGrid
 - X0298-dk17-09082000-1822_u0296232-1.TextGrid
 - X0298-dk17-09082000-1822_u0296233-1.TextGrid

Figure 2: Autophon outputs finished TextGrids using the same subfolder structure as the uploaded files.



Experimental linguistics A: Upload a two-column spreadsheet (Excel **xlsx**, or tab-delimited **txt/tsv**) with audio filenames in column 1 and transcriptions in column 2. No time stamps allowed. This format suits short clips and resembles CommonVoice⁵. Use zip or individual file upload.

Experimental linguistics B: Same structure as A, but with four columns: audio file name, start time, end time, and transcription. Designed for longer recordings requiring segmentation. Time stamps must be in real-number seconds (e.g., **1.23** or **1,23**); no colons or hour-minute markers are permitted (e.g., you may not use something like **00:00:01.23**).

Computational Linguistics: Upload matching audio and **lab** files (containing only the corresponding transcription). Files may be zipped with nested folders—Autophon preserves the hierarchy (Figure 2). No time stamps permitted.

Variationist Linguistics:⁶ Upload paired transcription and audio files (individually or zipped). Transcriptions may be in Praat **TextGrid**, ELAN **eaf**, or tabular format (**xlsx**, **txt**, **tsv**). Use either three or four columns:

- Four-column: speaker, start time, end time, transcription
- Three-column: start time, end time, transcription

Time stamps must be real-number seconds (comma or period decimal separators); formats with colons (e.g., **00:00:01.23**) are not supported.

2.6 File formats and codecs

If you encounter errors during upload, it's often due to unsupported file formats or codecs. The simplest fix is to re-save your files in a common format using tools like Praat or ELAN.

Transcription file formats: Autophon accepts transcription files in most standard encodings, including **UTF-8** and **UTF-16** (**Windows CRLF**). If you encounter issues, try re-saving the file or email a sample to tech support.

Audio file formats: Autophon supports a wide range of audio formats, including: **WAV**, **FLAC**, **MP3**, and more. Stereo files are not currently accepted. Therefore, convert all audio to mono first.

2.7 Transcription preparation

Regardless of the transcription mode, each entry should contain between one and 20 words. Boundary demarcations must include at least 0.01 seconds of silence before and after the speech. Figure 3 shows a five-word phrase with a 0.03-second pre-boundary and a 0.25-second post-boundary. This sort of variability is expected and handled well by Autophon.⁷

2.8 Select a language

After uploading files into the aligner, Autophon will suggest a language and language model. You may override this suggestion using the dropdown menu.

2.9 Task list

The task list displays all uploads along with file name, upload date, language, tier count, file size, word count, and an inventory of missing words. You can delete the task and start over, add words to your custom pronunciations box (described below), or proceed by clicking **Align**.

2.10 Missing words

To understand this feature, it helps to know that forced alignment maps phonemic pronunciations – defined in language-specific dictionaries – onto the speech stream using statistical models. These dictionaries contain a finite set of words. The missing words feature lists items not found in Autophon's dictionary and provides suggested pronunciations. Autophon will use these suggestions by default, but you can reject them and enter your own. The next section explains how.

2.11 Your custom pronunciations

If you disagree with either (a) Autophon's pronunciation suggestions for missing words or (b) the default dictionary entries, you can override them here. Enter your own phonemic transcriptions in this box, which will take precedence over both.

Pronunciations must be entered using the alphanumeric string specific to the language model at hand – in this case, the **NoFAbet**. Section 4 provides a key that maps the NoFAbet to its IPA⁸ equivalents.

⁵<https://commonvoice.mozilla.org>

⁶This field originally drove the development of forced alignment in the early 2000s.

⁷If your transcriptions are segmented with exact start and end times, performance may degrade and boundary shifts may occur. If you're working with such data, contact tech support—we are interested in designing a fifth transcription mode for these cases.

⁸International Phonetic Alphabet

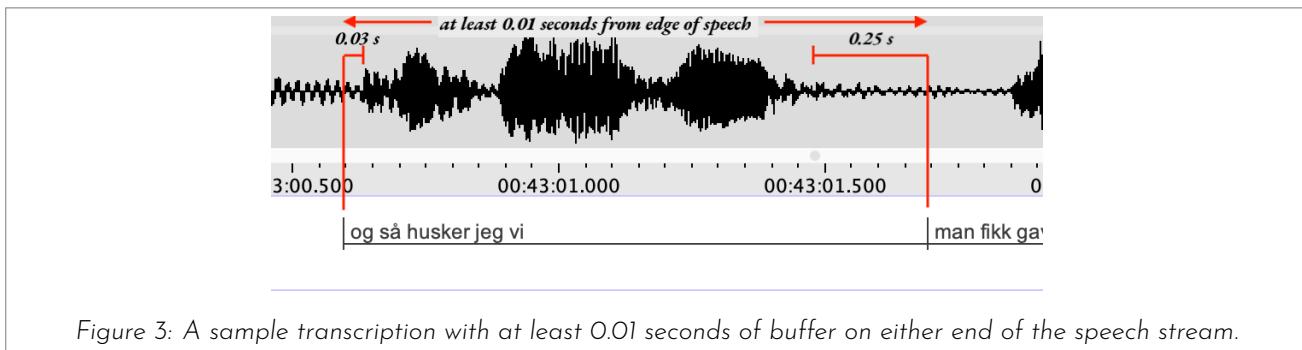


Figure 3: A sample transcription with at least 0.01 seconds of buffer on either end of the speech stream.

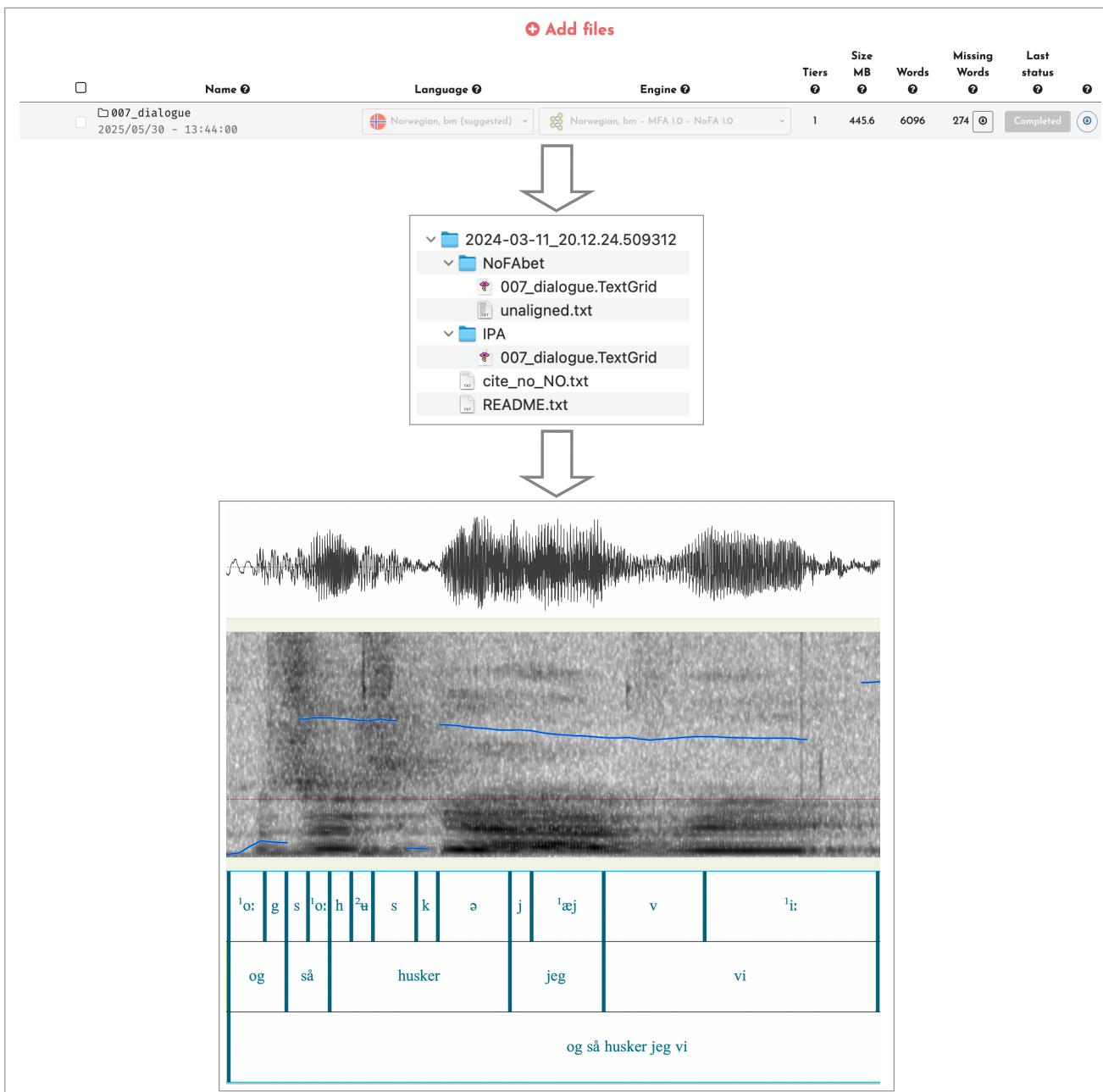
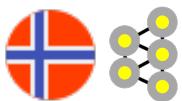


Figure 4: The alignment process, including task list, folder structure, and Praat TextGrid.



You can enter pronunciations directly in the dictionary box or upload them from a **txt** file. The maximum input length is 50 000 characters.

Entries must follow the format:

- word[space]phoneme[space]phoneme OR word[tab]phoneme[space]phoneme

Each phoneme must be separated by a space, and the lookup may not include two or more words – Autophon will interpret the second word as a phone and produce an error. You may submit multiple pronunciations for the same word by repeating the word on separate lines with different phoneme strings. Autophon will evaluate the best match for each speech instance. Refer to Figure 5 and the examples below.

Stress and/or accent **must** accompany every vowel and diphthong with a number. Consult Figure 1 for the specific digits used in this model; refer to Figure 5 to see how these are operationalized.

The figure shows the Autophon interface. On the left, there is a decorative icon of an open book with a sound wave graphic. A button labeled "Click to open" is next to it. An arrow points from this icon to the right side of the interface. The right side is titled "Your Custom Pronunciations" and includes a "Swedish - MFA 1.0 - SweFA 2.0" dropdown. Below this is a text input field with the placeholder "Type directly into the field below or upload a text file here". Underneath the input field, there is a table showing four entries:

	word	stress	phones
1	ackompanjera	AH0	K OAH0 M P AH0 N J EE1 R AH0
2	ackompanjerade	AH0	K OAH0 M P AH0 N J EE1 R AH0 D EH0
3	ackompanjerades	AH0	K OAH0 M P AH0 N J EE1 R AH0 D EH0 S
4	ackompanjerads	AH0	K OAH0 M P AH0 N J EE1 R AH0 D S

Below the table, there are two sections: "Correct:" and "Incorrect:" with examples of valid and invalid phoneme strings.

Correct:

dababy	D	AA ₀	B	EE ₁	B	II ₀
dababy	D	B	EJ ₁	B	II ₀	
da_baby	D	AA ₀	B	EE ₁	B	II ₀

Incorrect:

dababy	D	AA	B	EE	B	II	(vowel-stress numbering is missing)
dababy	D	AA ₀	B	EE ₁	B	II ₀	(phones missing a space between them)
da_baby	D	AA ₀	B	EE ₁	B	II ₀	(two look-ups on a single line)

Figure 5: Interface with dictionary entry (left) and phoneme string input (right).

2.12 Aligning files To begin alignment, click *Align* to the far right of the upload list. Alignment typically takes a few minutes, depending on server load.

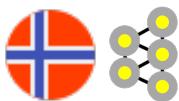
2.13 Downloading the annotations When alignment is complete, you can download the annotations as Praat TextGrids by clicking the downward arrow beside the task list. See Figure 4 for an illustration.

3 How to cite

Any dissemination or publication that makes use of this Autophon package for **Norwegian Bokmål**, which uses **NoFA 1.0** within **The Montreal Forced Aligner 1.0** for its engine, should cite the relevant references listed below. Proper citation is essential: not only to acknowledge the “daisy chain” of technical and academic work underpinning Autophon, but also to reinforce the incentives for sharing tools with the broader community.

While space constraints may tempt you to remove references to software, we strongly encourage prioritizing these citations. If trimming is necessary, consider reducing peripheral citations in the literature review instead.

- Boersma, P., & Weenink, D. (2017). Praat: Doing phonetics by computer [softw.], ver.6.0.36. www.praat.org
- McAuliffe, M., Socolof, M., Mihuc, S., Wagner, M., & Sonderegger, M. (2017). Montreal Forced Aligner: Trainable text-speech alignment using Kaldi. *Proceedings of Interspeech*, 498–502.
- Young, N. J. (2020). NoFA 1.0 - norsk modell for forced alignment, ver. 1.0. <https://www.nb.no/sprakbanken/ressurskatalog/oai-nb-no-sbr-59/>
- Young, N. J., & Anikwe, K. H. (2025). Autophon - Automatic phonetic annotation and online forced aligner. www.autophon.org



4 Phoneme key

Autophon will output two versions of the same TextGrid for every file you align: (1) a TextGrid in the NoFAbet specific to NoFA 1.0 for The Montreal Forced Aligner 1.0 and (2) a TextGrid in the International Phonetic Alphabet (IPA). The phoneme key is located in Table 1.

NoFAbet	IPA example	NoFAbet	IPA example	NoFAbet	IPA example	NoFAbet	IPA example
Vowels		AX	ə behage	KJ	ç kino	Syllabic consonants	
AA	a: bad	L	l land	M	m man	LX	l djevelsk
AH	a hatt	M	n nord	N	ŋ eng	MX	m landsomfattande
AE	æ vær	AEW	æw sau	NG	p pil	NX	n avtalen
AEH	æ vært	AJ	aj kai	P	r rose	RLX	l varsel
EE	e: lek	OEJ	œj køye	R	d rekord	RNX	n turen
EH	ɛ penn	OJ	ɔj konvoy	RD	l perle	RX	r søker
II	i: vin	OU	oʊ show	RL	ɳ barn	Paralinguistic features	
IH	ɪ sitt			RN	t stort	EXH	<exhale>
OA	o: rå	Consonants		S	s sil	INH	<inhale>
OAH	ɔ gått	B	b bil	SJ	ʂ sju	NHES	<nasal>
OE	ø lok	D	d dag	T	t tid	VHES	<vowel>
OEH	œ høst	DH	ð this*	TH	θ thin*	LG	<laughter>
OO	u: bod	DJ	dʒ George*	TSJ	tʃ church*	Lexical stress	
OH	u fort	F	f fin	V	v vase	○0	○ advarer
UU	ʉ hus	G	g gul	W	w Washington	○1	''○ advarer
UH	ʉ russ	H	h hes	X	x ach**	○2	''○ advarsel
YY	y ny	J	j ja	Z	z zigzag*	○3	○ advaring
YH	y nytt	K	k kost				
UX	ꝝ girl*						

* English word

** German word

Table 1: NoFAbet, IPA, and lexical examples. The denotation for lexical stress means that any NoFAbet vowel or diphthong must always be followed by the numbers 1 (toneme one), 2 (toneme two), 3 (secondary accent), or 0 (unstressed).

Every NoFAbet vowel and diphthong is followed by a numerical code that denotes suprasegmental information. ○0 refers to lexically unstressed vowels, ○1 - vowels with **tonem 1**, ○2 - vowels with **tonem 2** vowels, and ○3 - vowels with secondary stress. We encourage you to inform us of errors and provide suggestions for changes.

5 Acoustic model and pronunciation dictionary

This specific Autophon package for **Norwegian Bokmål** uses NoFA 1.0 within The Montreal Forced Aligner 1.0, which was trained on spontaneous and read-aloud Norwegian Bokmål from the RUNDKAST⁹ and NB Tale corpora¹⁰. The pronunciation dictionary was adapted from The NST Pronunciation Lexicon for Norwegian Bokmål¹¹.

6 Performance metrics

NoFA 1.0 is quite accurate, measured here by comparing alignments of 1000 phonemes, each from one speaker from three different dialect regions, respectively: Southeastern (Oslo), Western (Bergen), and Central (Trøndelag). The alignments are compared in Table 2 with manual segmentation.

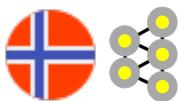
7 Data security and GDPR compliance

Files uploaded to Autophon are encrypted and transmitted to a secure server hosted by Digital Ocean within the European Union (Frankfurt and Amsterdam). Transcriptions and audio files are automatically deleted immediately

⁹Amdal, Strand, Almberg, and Svendsen (2008)

¹⁰<https://www.nb.no/sprakbanken/en/resource-catalogue/oai-nb-no-sbr-31/>

¹¹<https://www.nb.no/sprakbanken/en/resource-catalogue/oai-nb-no-sbr-23/>



Dialect	Speaker	n boundaries	median onset difference (ms)	pct 20 ms or less	pct 10 ms or less
Southeastern (Oslo)	Anniken Hauglie	1000	10	79%	49%
Western (Bergen)	Audun Lysbakken	1000	10	77%	48%
Central (Trøndelag)	Trine Skei Grande	1000	12	72%	40%

Key

n boundaries number of boundaries tested against the manual gold standard (g.s.)
median onset difference (ms) median difference between aligner boundaries and manual g.s. boundaries
pct 20 ms or less percentage of aligner boundaries that fall within 20 milliseconds of manual g.s. boundaries
pct 10 ms or less percentage of aligner boundaries that fall within 10 milliseconds of manual g.s. boundaries

Table 2: Accuracy metrics for NoFA 1.0

after alignment. This approach enhances privacy while also reducing storage costs. By contrast, finished TextGrids remain available in your account until you choose to delete them. Once deleted, they are permanently removed from our servers.

If you upload files but do not initiate alignment by clicking *Align*, the files will be automatically purged at 3:00 AM GMT¹².

Autophon adheres to the principles of the European Union’s General Data Protection Regulation (GDPR). We collect only four pieces of user information: name, title, affiliation, and email address. Once a file is aligned, the corresponding audio is permanently deleted. Deleting a file from your task list also permanently removes the transcription and filename metadata. You may delete your account at any time, which will erase all associated personal data. However, we **do** retain anonymized alignment metadata – such as a randomly assigned alphanumeric user ID and summary usage statistics – to demonstrate the platform’s utility to funders.

8 Features and limitations

What Autophon is: Autophon is a web-based application designed to simplify forced alignment workflows and expand access for users with minimal technical background. It is particularly useful for research on under-resourced languages and non-standard varieties, and emphasizes ease of use, format flexibility, and language model diversity. The backend relies on existing forced alignment technologies developed over the past decades, wrapped in a modern frontend that facilitates fast, OS-independent processing.

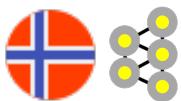
Key features include:

1. Fully web-based and platform-independent (OS-agnostic).
2. No programming or installation required.
3. Accepts a wide range of transcription and audio formats.
4. Capable of processing low-resource and non-standard language varieties.
5. Supports user-defined pronunciation dictionaries and multiple transcription modes.

What Autophon is not: Important caveats to bear in mind:

1. Alignment quality depends on transcription accuracy, recording quality, and language characteristics.
2. Performance may vary across languages, dialects, and speaking styles.
3. Benchmarking accuracy is ongoing and not available for all models.

¹²Users working near this cutoff time—e.g., at 2:55 AM GMT—should be aware that their files may disappear if alignment is not initiated in time.



4. Core updates to underlying alignment engines may not be immediately reflected, due to the complexity of the Autophon backend.

9 Budget and funding

Autophon costs approximately 25 000 SEK (2 300 EUR) per year to run. Founded by Dr. Nate Young (who is the sole copyright holder), the project has since received support from the University of Helsinki, Linnaeus University, The Swedish Academy, the Department of Linguistics and Scandinavian Studies at the University of Oslo, and the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No. 892963. Additional funding for language model development has come from The National Library of Norway¹³.

We continue to seek funding and welcome collaboration. If you are experienced in grant writing or interested in supporting the project, please reach out via the support page.

Acknowledgements

Numerous individuals have contributed to Autophon. We especially thank Michael McGarrah for strategic guidance and Kaosi Anikwe for extensive backend and frontend development. Ismail Raji Damilola helped implement a bootstrapping function to expand phoneme inventories. Additional contributions in the early stages came from Nabil Al Nazi, Zamanat Abbas Naqvi, and Santiago Recoba.

We also wish to acknowledge the people who helped make the NoFA model possible. Per Erik Solberg and The National Library of Norway supported the development of NoFA 1.0.

References

- Amdal, I., Strand, O. M., Almberg, J., & Svendsen, T. (2008). RUNDKAST: An Annotated Norwegian Broadcast News Speech Corpus. *Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC'08)*. http://www.lrec-conf.org/proceedings/lrec2008/pdf/486_paper.pdf
- Boersma, P., & Weenink, D. (2017). Praat: Doing phonetics by computer [softw.], ver.6.0.36. www.praat.org
- McAuliffe, M., Socolof, M., Mihuc, S., Wagner, M., & Sonderegger, M. (2017). Montreal Forced Aligner: Trainable text-speech alignment using Kaldi. *Proceedings of Interspeech*, 498–502.
- Povey, D., Ghoshal, A., Boulianne, G., Burget, L., Glembek, O., Goel, N., Hannemann, M., Motlicek, P., Qian, Y., Schwarz, P., et al. (2011). The Kaldi speech recognition toolkit (tech. rep.). IEEE Signal Processing Society. Piscataway.
- Rosenfelder, I., Fruehwald, J., Brickhouse, C., Evanini, K., Seyfarth, S., Gorman, K., Prichard, H., & Yuan, J. (2022). FAVE (Forced Alignment and Vowel Extraction) Program Suite v2.0.0 [Zenodo].
- Wilbanks, E. (2022). faveAlign (Version 1.14). <https://github.com/EricWilbanks/faveAlign>
- Young, S. J., Woodland, P. C., & Byrne, W. J. (1993). HTK: Hidden Markov Model Toolkit V1. 5. Cambridge Univ. Eng. Dept. Speech Group; Entropic Research Lab. Inc.

¹³<https://www.nb.no/sprakbanken/ressurskatalog/oai-nb-no-sbr-59/>