

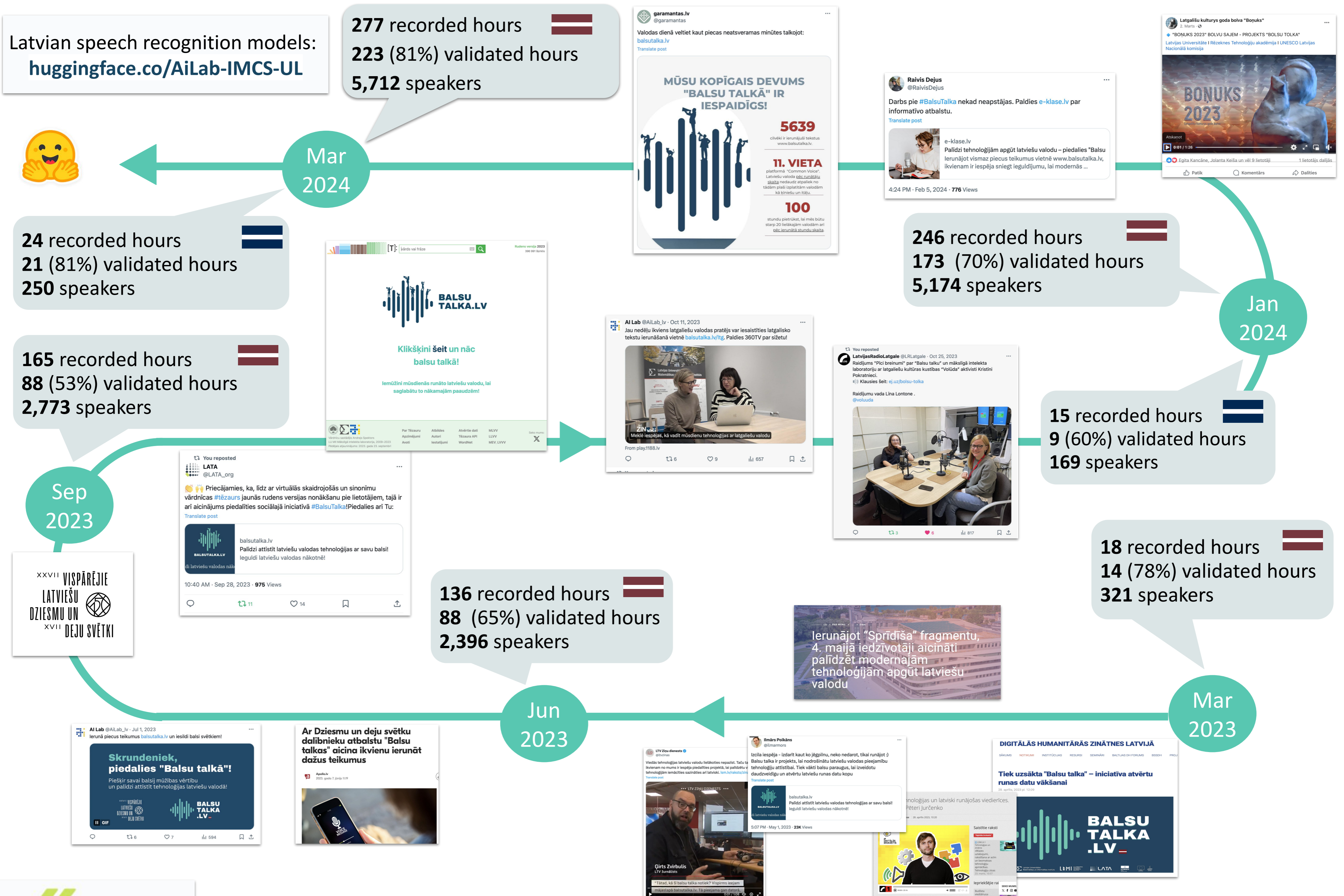
BALSU TALKA

Crowdsourcing Open Speech Corpora in Latvian and Latgalian

Baiba Saulīte¹, Ilze Auziņa¹, Kristīne Pokratniece¹,
Sanita Reinsone², Normunds Grūzītis¹, Roberts Darģis¹,
Artūrs Znotiņš¹, Ilze Ziņģe¹, Antra Kļavinska³, Raivis Dejus⁴

¹ Institute of Mathematics and Computer Science, UL
² Institute of Literature, Folklore and Art, UL
³ Rēzekne Technology Academy
⁴ Latvian Open Technologies Association

At the beginning of 2023, there were only a few open Latvian speech corpora available, the 18-hour Common Voice (CV) 13 dataset being the largest one. In the result of a national crowdsourcing initiative, organised jointly by several institutions, the size and speaker diversity of the Latvian CV 17 release have increased more than tenfold in less than a year. A follow-up initiative was also launched for Latgalian.



BalsuTalka
Balsutalka.lv Speech Corpus (Common Voice 17.0)
2024, 277 hours (1.3M tokens)
Developers: IMCS, UL, ILFA UL, LATA

1	☐	☺	vīrietis	<>	Tupēsi mājās līdz pavasara brīvdienām . </s>	▶
2	☐	☺	sieviete	<>	Ausa saulains rīts un zeme smaržoja pēc pavasara . </s>	▶
3	☐	☺	sieviete	<>	Kad sadzird pirmo pavasara mušu ielidojam istabā . </s>	▶
4	☐	☺	sieviete	<>	Sākās pavasara lauku darbi . </s>	▶
5	☐	☺	nav norādīts	<>	Tāds pavasaris ! Un , paskat , strazds ! </s>	▶
6	☐	☺	vīrietis	<>	Skatījāties Gauju , pīles un pavasari . Saldējuma mašīna Carnikavā . </s>	▶
7	☐	☺	sieviete	<>	Prasās vairāk svaiguma , vairāk pavasara garšas . </s>	▶

BolsuTalka
Bolsutalka.lv Speech Corpus (Common Voice 17.0)
2024, 24 hours (131k tokens)
Developers: RATA, IMCS, UL, ILFA UL, LATA

1	☐	☺	sieviete	<>	Ej tolkā ! </s>	▶
2	☐	☺	nav norādīts	<>	Vyssy bujbs nūkast ar tolku . </s>	▶
3	☐	☺	vīrietis	<>	Sasaitk puļņerj jau varātu i taipat , bez ituos tolks – voi ta imesis naatsarass . </s>	▶
4	☐	☺	sieviete	<>	Ka teik reikuota kaida tolka , kuozys , krystobys voi , nadūd Dīvs , bēris , vysy ir kluot ! </s>	▶
5	☐	☺	sieviete	<>	Šudīn juonūkaš vysy buļvi , deļtuo ka pareit syudu tolka . </s>	▶
6	☐	☺	sieviete	<>	Tev tik četrīs ! Tūlāik nūsasmieju , otkon kuo ta tolkā vaicuodama . </s>	▶
7	☐	☺	nav norādīts	<>	Nu tevis tolka kai nu solta pojaņnika . </s>	▶