

CommonMorph

A Crowd-Sourcing Platform for Morphological Resources

Aso Mahmudi¹ Sina Ahmadi² Kemal Kurniawan¹ Rico Sennrich² Eduard Hovy¹ Ekaterina Vylomova¹

¹The University of Melbourne ²University of Zurich



Universität
Zürich UZH

LREC
2026
Palma

The Gap

How can morphological data collection scale through structured collaboration between linguists and speakers?

1,500+

languages at risk of extinction by 2100

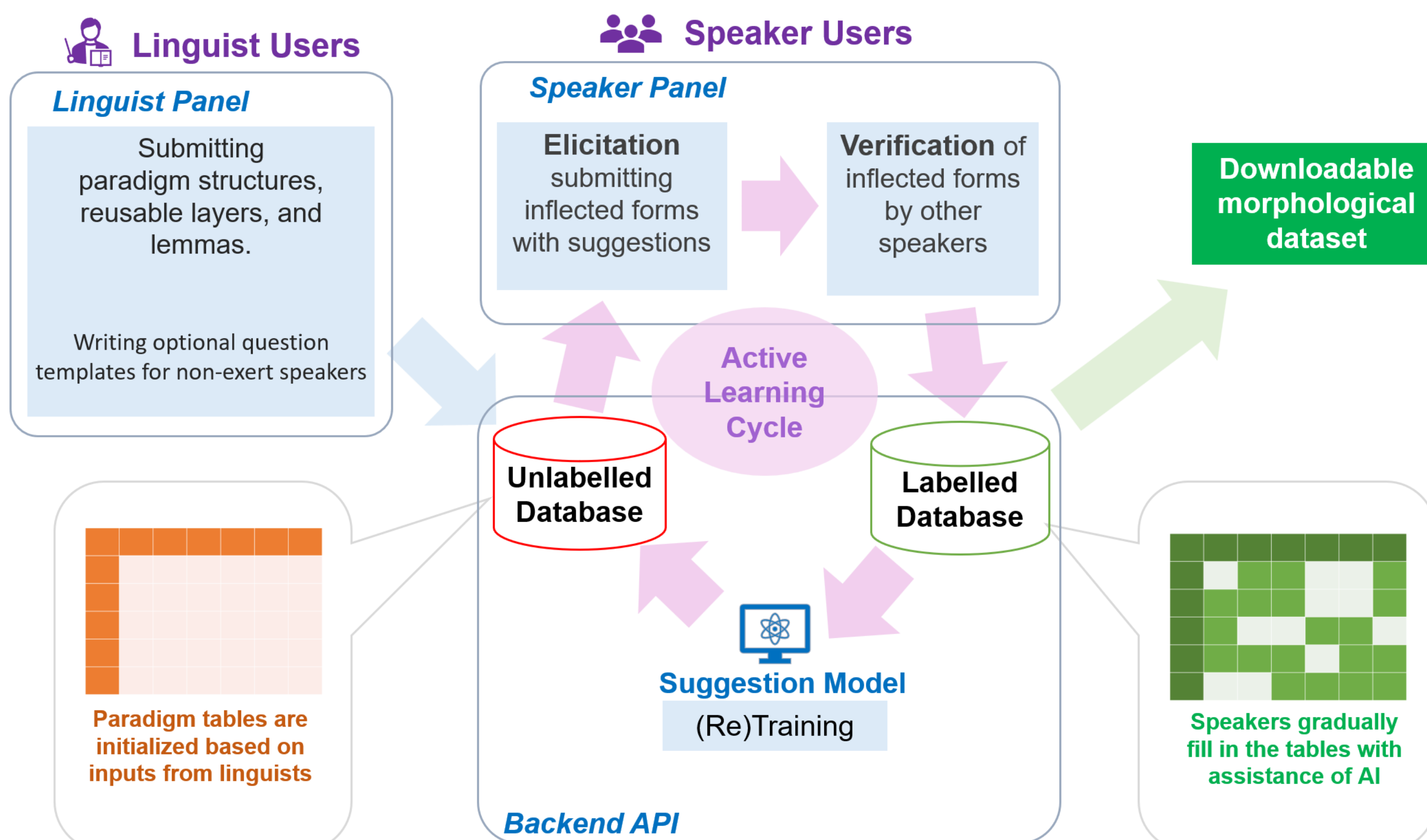
169

in UniMorph today (< 2.5% of the world's)

40 hours

per language to build a morphological grammar

What is CommonMorph?



Linguist Interface

Linguists define paradigms, lexicon, reusable layers, and optional elicitation prompts. Materials can be **imported and adapted** from related varieties.

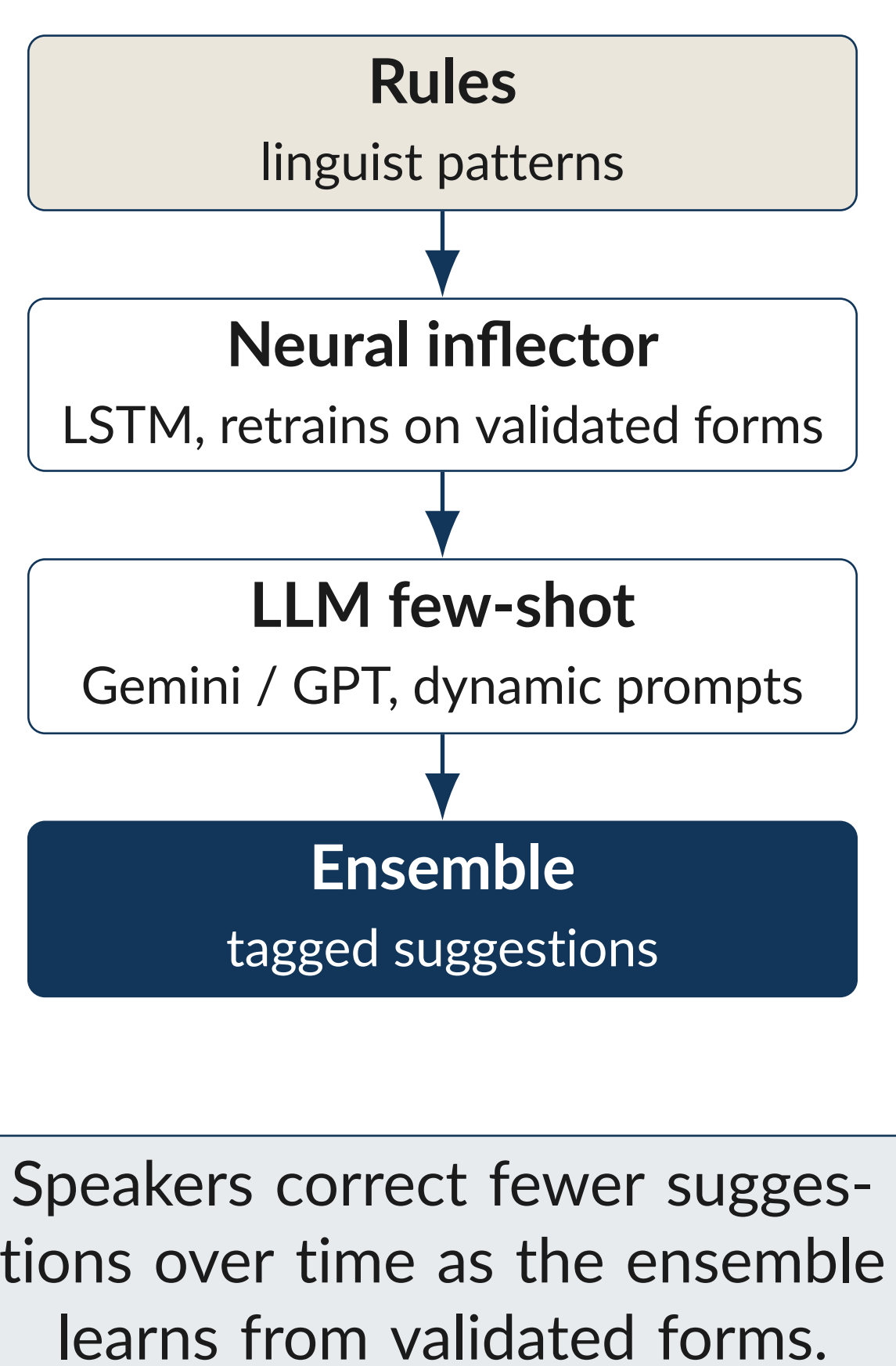
The Linguist Interface includes four main sections: **Paradigm structures** (Structure Title/Alias, Morphosyntactic Feature, Add feature, Formula, Reusable Layer), **Reusable layers** (Reusable Morpheme Title/Alias, Morpheme, Morphosyntactic Feature, Add feature, Priority), **Lexicon** (Inflection Class, Gloss (Meaning), Lemma (Dictionary Entry), Stem1, Stem2, Description, Priority), and **Question design** (Insert a new question form, Please write a template question, Need inspiration?).

Speaker Interfaces

The interface adapts to the speaker's expertise. **Experts** verify whole tables; **non-experts** answer one plain-language question at a time. Suggestions are **tagged by source**: rule, neural net, or LLM.

The Speaker Interfaces are divided into four types: **Elicit (non-expert)** (Please answer the following question in your own language), **Verify (non-expert)** (How would you say 'love/amó' when you talk about yourself?), **Elicit (expert)** (My grammar knowledge: Basic (high school grammar)), and **Verify (expert)** (My grammar knowledge: Expert (familiar with linguistics terms)).

Active Learning

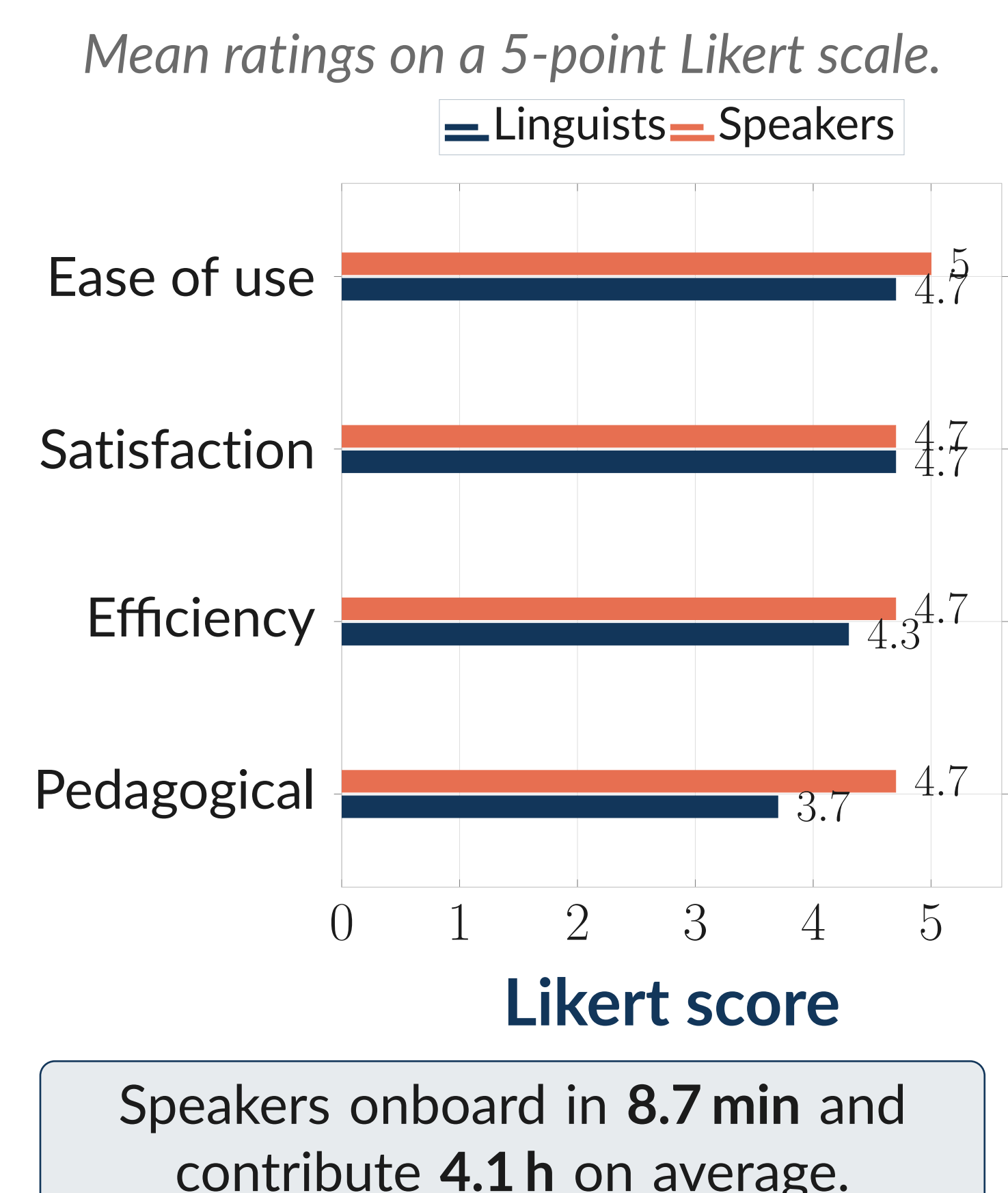


Suggestion Accuracy across Case Studies

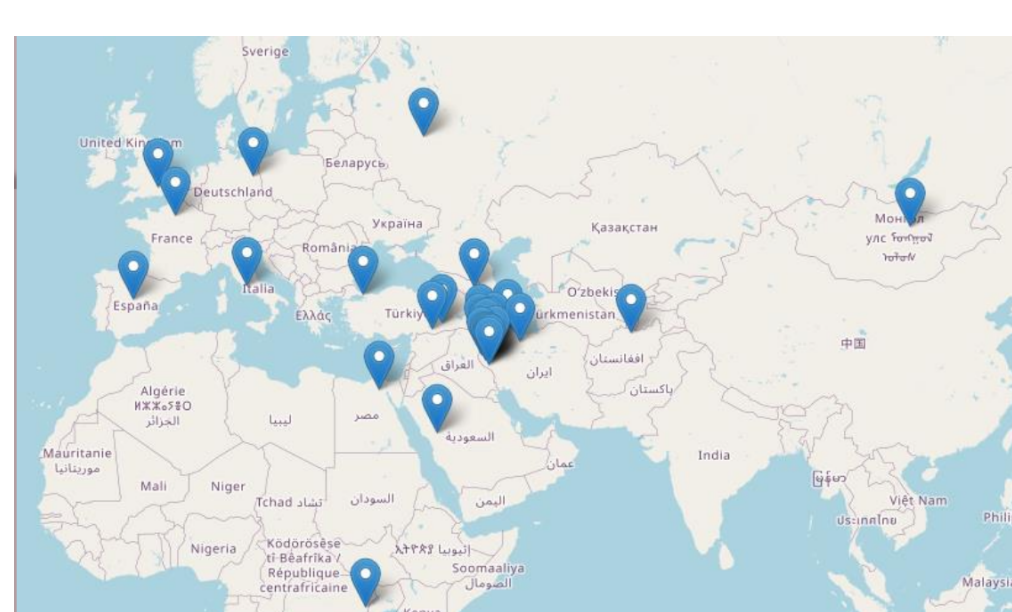
Character Error Rate (%) on 3,000 held-out forms – lower is better.

Model	Hawrami C.	Kurdish	Farsi	Turkish	Spanish	Arabic
Linguist rules						
Patterns only	12.33	1.39	0.86	31.74	0.82	6.62
+ morphophon.	5.78	0.74	0.85	2.97	0.82	2.24
Neural inflector						
100 samples	48.41	54.23	37.25	23.31	20.98	41.93
500 samples	15.13	29.12	23.63	11.83	7.18	8.45
1,000 samples	7.88	27.37	54.02	3.33	5.91	6.40
2,000 samples	6.46	14.36	24.22	4.68	1.14	2.72
LLM few-shot						
gemini-2.5 (1)	17.00	14.86	3.71	2.36	0.50	7.25
gemini-2.5 (2)	14.45	13.43	1.86	1.50	0.50	4.59
gemini-2.5 (3)	13.03	6.86	2.32	2.15	0.50	5.07
gpt-5-mini (1)	19.26	15.14	4.64	0.43	0.00	10.14
gpt-5-mini (2)	11.61	14.57	4.64	0.21	0.25	7.97
gpt-5-mini (3)	12.18	7.43	3.25	0.21	0.00	4.83

User survey



Your language next?



Five language families, four morphological types.

Bring your own variety.

Takeaways

- **CommonMorph** is a participatory platform for documenting morphology – linguists scaffold, speakers fill in.
- **Adaptive interfaces** let both experts and non-experts contribute meaningfully.
- **Active learning + LLMs** make every annotation cycle faster.
- **Cross-variety reuse** bootstraps related languages from existing scaffolds.
- **New data released**: first comprehensive verbal morphology of Hawrami.
- **Open-source and UniMorph-compatible**.

Try it / Fork it



Aso-UniMelb/CommonMorph
common-morph.com