



Multilingual benchmarks should reflect the regional and cultural knowledge of where languages are spoken or used, not be translated from western-centric resources

**The largest multilingual exam collection to date with ~200K question-answer pairs focusing on regional & cultural knowledge**

**Limitations of the current translated benchmarks**

- They primarily reflect US / Western-centric knowledge
- They may contain errors introduced during translation
- They often exhibit translationese artifacts

**A balanced subset is available on Hugging Face**

**CohereForAI/include-base-44**

Tasks: Multiple Choice Languages: 44 Languages

Size: 22,637 ArXiv: arxiv: 2411.19799

**44 LANGUAGES**

- Azerbaijani, Bulgarian, Greek, Armenian, Basque
- Croatian, Hungarian, Nepali, Macedonian, Tagalog
- Serbian, Albanian, Lithuanian, Malayalam, Georgian
- Bengali, Estonian, Turkish, Belarusian, Telugu, Dutch
- Hebrew, Hindi, Malay, Urdu, Kazakh, Arabic
- Korean, Ukrainian, Tamil, Uzbek, Russian, Spanish
- Chinese, Italian, French, German, Finnish, Persian
- Indonesian, Vietnamese, Portuguese, Japanese, Polish

**Motivation: The same questions can have different answers depending on where they are asked**

**REGIONAL KNOWLEDGE** (Law & Regulations)

Ποιο είναι το επιτρεπό όριο αλκοόλ ανά λίτρο αίματος στην οδήγηση? (Какой уровень алкоголя в крови допустим при вождении?)

What is the Blood Alcohol Limit (BAC%) for driving? (What is the Blood Alcohol Limit (BAC%) for driving?)

0 %, 0.03 %, 0.05 %, 0.08 %

57 TASKS: Arts & Humanities, Social Sciences, STEM, Business & Commerce, Health education, Professional Licenses, Occupational Licenses +

### HOW DO LLMs PERFORM ON INCLUDE ?

**LEADERBOARD**

Model	5-shot	COT
GPT-4o	76.2	78.5
Llama3.1-70B Inst	70.6	55.6
Aya-expans-32B	60.0	57.8
Qwen2.5-14B	61.0	51.6
Aya-expans-8B	47.8	
Llama-3.1-8B	51.6	
Llama-3.1-8B Inst	54.4	
Gemma-7B	54.7	
Gemma-7B Inst	39.2	
Qwen2.5-7B	54.5	
Qwen2.5-7B Inst	53.9	

Performance on GPT-4o with different generation windows on high vs low resource languages

Multilingual LLMs do not follow instructions the same way in all languages. GPT-4o shows 3.1% gain when increasing the generation length window

Open models transfer ~20% more to unseen languages with shared scripts than to those with different ones

Increased model size significantly enhances multilingual capabilities (with the same pre-training data)

Instruction tuning may hurt the multilingual ability likely due to English-heavy post-training

Models struggle more on tasks requiring regional knowledge than those assessing universal knowledge

English prompts offer modest improvements of 1-2% compared to In-language prompts