

Russian Al Struggles in Its Own Language

By Ben Dubow

January 28, 2025



Digital Transportation Forum in Moscow's Technology Valley. **Arthur Novosiltsev / Moskva News Agency**

In a sleek presentation hall at Sber's Moscow headquarters in November, executives showcased their latest achievement: <u>GigaChat MAX</u>, a new AI model they claimed would cement Russia's position as a leading AI power.

The model's improved 80% score on the industry-standard <u>MMLU math benchmark</u> represents a major improvement on Russia's previous efforts. But beneath the corporate optimism lay an uncomfortable reality: American models had achieved the same feat over a year ago.

This widening capability gap represents more than just benchmark numbers. While Russian developers celebrate incremental improvements in basic language processing, American and Chinese models are racing ahead with advanced capabilities like chain-of-thought reasoning

and multimodal understanding. The distance between Russian models and the global leaders is growing. Worryingly for Moscow, Russian models still lag behind in their own language, whereas American models continue to dominate.

The new model's benchmarks showcase how far behind Russia remains. GigaChat MAX bases its generations on 20 billion associations among words (called parameters), roughly equivalent to ChatGPT-3.5 released nearly two years ago. GigaChat MAX can remember 131k tokens (words and semantically significant punctuation or parts of words) in a conversation, compared to 2 million for Gemini 1.5, Google's previous generation.

Sber preferred not to draw these comparisons, however. Instead, in their <u>announcement</u> blog, Sber put MAX up against open-source models like Meta's LLaMa and Google's Gemma. The developers bragged of MAX's comparatively cheap costs. However, they said it would be able to surpass Google's Gemma 2-9B after further training, a veiled admission of its lackluster Russian language capabilities.

Math, by comparison, proved a bright spot for MAX. That said, the 80% score it achieved is the same as that of <u>DeepSeek</u>, the Chinese lab whose r1 model <u>upended</u> American AI-sector stocks after its free app made waves. DeepSeek's code forms the basis for MAX.

Related article: Embarrassingly for the Kremlin, Russian AI Isn't Good Enough for Its Own Disinformation

However, Russian developers using open-source models that originate from China and are available to anyone across the world is hardly the sort of collaboration likely to give Russia much advantage. Chinese labs seem broadly uninterested in the Russian market, even compared to those of hegemons the BRICS are supposedly countering. Alibaba's most popular open-source model, Qwen, performs about as well in English as Meta's LLaMa, according to ChatbotArena. The model is nearly unusable in Russian, according to the technical blog announcing GigaChat MAX. The most China seems willing to offer Russia are shiny medals for its researchers.

Russia's underperformance is mostly a product of its own making: international isolation brought on by the annexation of Crimea just as the AI boom began, then the flight of most tech talent and punishing sanctions with the full-scale invasion hit the same year as ChatGPT-3's release. The consultancy Tortoise Media ranks Russia 31st on their scale of AI vibrancy. Estonia, with about the same population as Kazan, is only one spot behind.

Nonetheless, Putin still <u>declared</u> at the conference that AI is "the most important resource for achieving the national development goals of the country, to ensure the strengthening of its defense capability, the qualitative development of the economy and social sectors, public administration, and the growth of innovation." He is most likely correct about that. But there is little sign that any models besides American ones can deliver those results, even in Russian.

The views expressed in opinion pieces do not necessarily reflect the position of The Moscow Times.

Original url:

