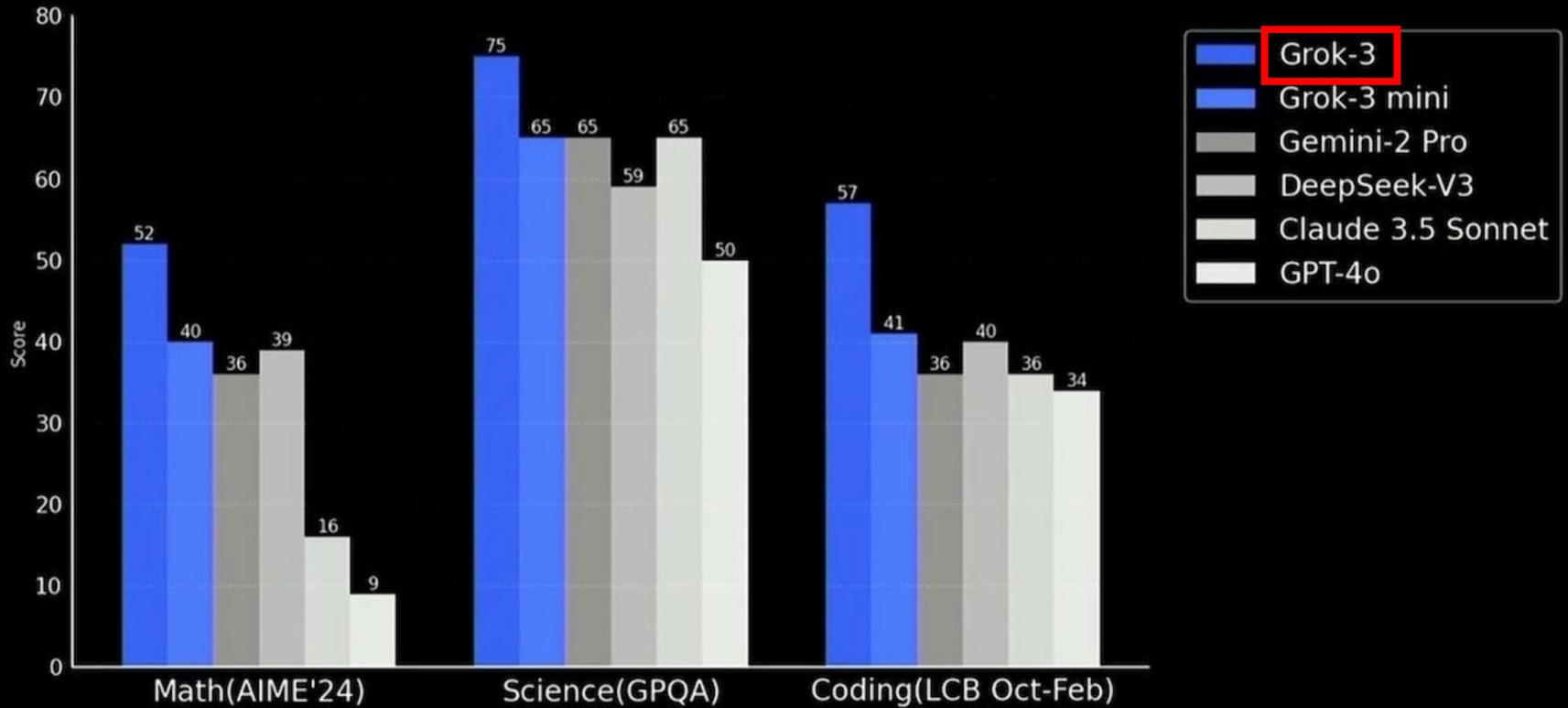
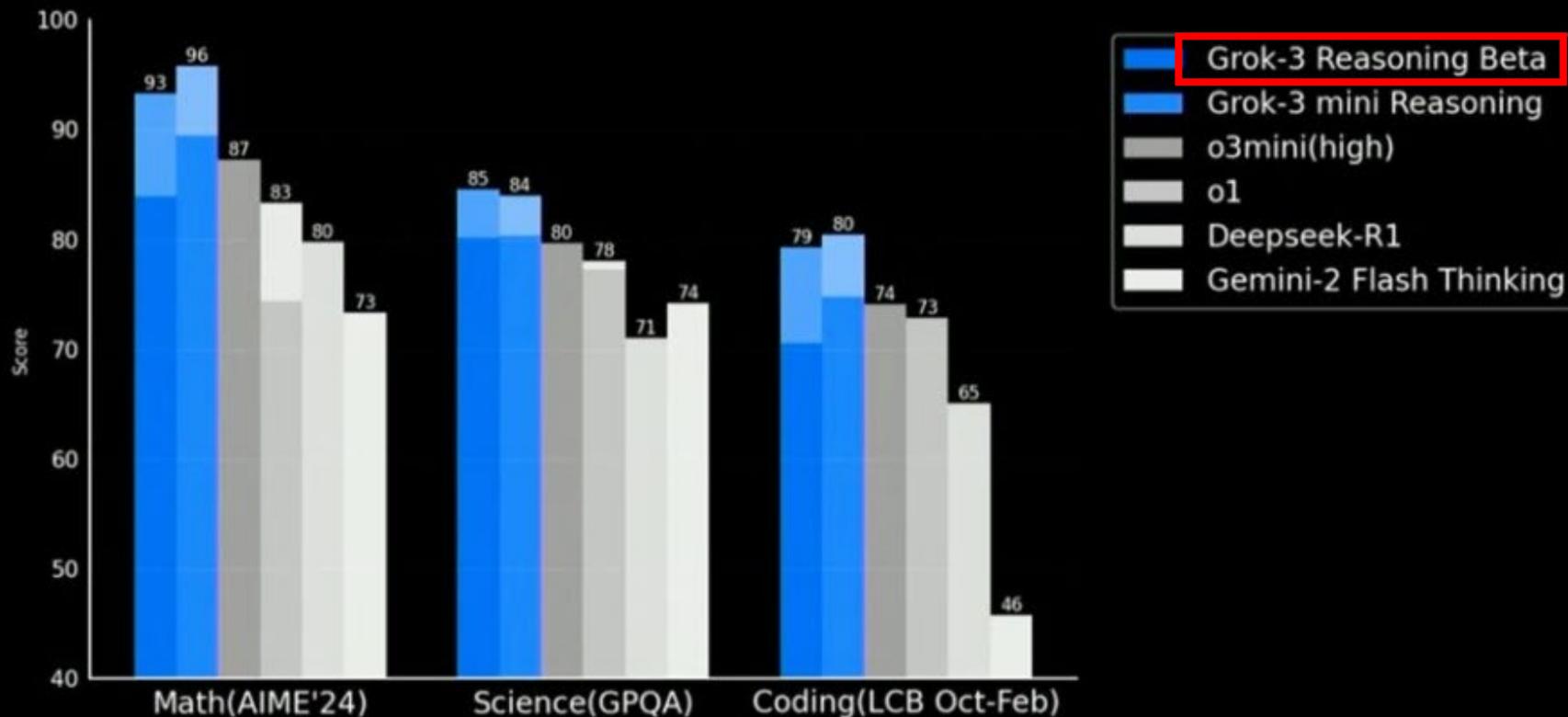


# Benchmarks



x1

# Reasoning + Test-Time Compute



推論モードでのテスト結果

<https://shift-ai.co.jp/blog/16730/>

Total #models: 205. Total #votes: 2,668,091. Last updated: 2025-02-16.

Code to recreate leaderboard tables and plots in this [notebook](#). You can contribute your vote at [lmarena.ai](#)!

Category: Overall

Apply filter

Style Control

Show Deprecated

**Overall Questions**

#models: 205 (100%) #votes: 2,668,091 (100%)

Rank* (UB)	Rank (StyleCtrl)	Model	Arena Score	95% CI	Votes	Organization	License
1	1	<a href="#">chocolate (Early Grok-3)</a>	1402	+7/-6	7829	xAI	Proprietary
2	4	<a href="#">Gemini-2.0-Flash-Thinking-Exp-01-21</a>	1385	+5/-5	13336	Google	Proprietary
2	2	<a href="#">Gemini-2.0-Pro-Exp-02-05</a>	1379	+5/-6	11197	Google	Proprietary
2	1	<a href="#">ChatGPT-4o-latest (2025-01-29)</a>	1377	+5/-6	10529	OpenAI	Proprietary

# More Statistics for Chatbot Arena (Overall)

Figure 1: Confidence Intervals on Model Strength (via Bootstrapping)

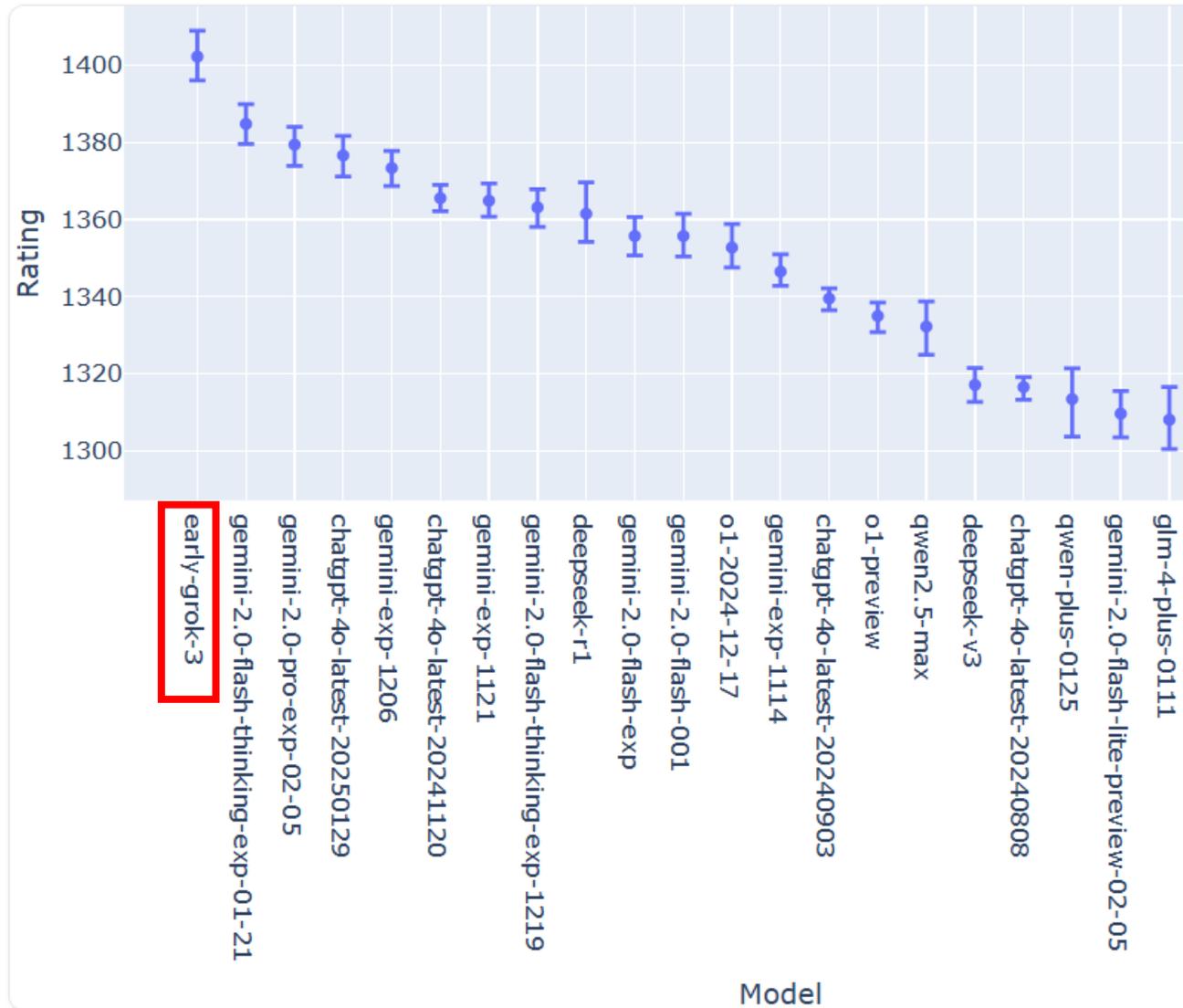
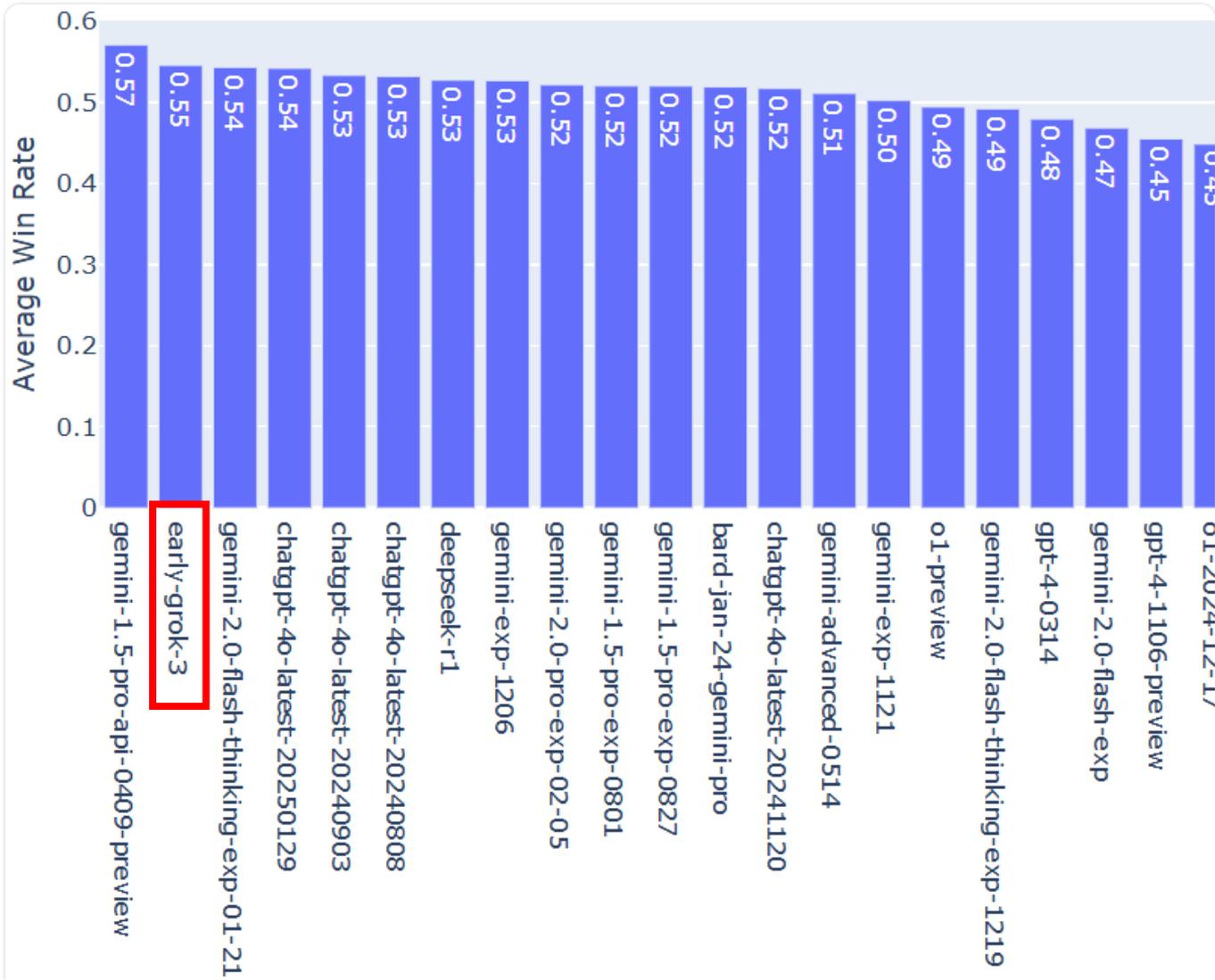


Figure 2: Average Win Rate Against All Other Models (Assuming Uniform Sampling and No Ties)



# Chatbot Arena Overview

(Task)

Sort by Rank

Sort by Arena Score

Model	Overall	w/ Style Control	Hard Prompts	Prompts w/ Style Control	Coding	Math	Creative Writing	Instruction Following	Longer Query	Multi- Turn
early-grok-3	1	1	1	1	1	1	1	1	1	1
gemini-2.0- flash- thinking-exp- 01-21	2	4	1	1	2	1	1	1	1	1
gemini-2.0- pro-exp-02-05	2	2	1	1	1	2	1	1	1	1
chatgpt-4o- latest- 20250129	2	1	4	3	2	10	1	2	1	1
deepseek-r1	5	2	2	1	2	1	4	2	2	1
gemini-2.0- flash-001	5	8	4	7	2	1	4	6	4	4

<https://lmarena.ai/?leaderboard>

# Chatbot Arena Overview

(Task)

Sort by Rank

Sort by Arena Score

Model ▲	Overall ▲	Overall w/ Style Control ▲	Hard Prompts ▲	Hard Prompts w/ Style Control ▲	Coding	Math ▲	Creative Writing	Instruction Following	Longer Query	Multi-Turn
early-grok-3	1402	1333	1397	1327	1399	1352	1421	1395	1424	1414
gemini-2.0-flash-thinking-exp-01-21	1384	1311	1385	1316	1368	1349	1397	1385	1395	1398
gemini-2.0-pro-exp-02-05	1379	1320	1381	1318	1371	1330	1391	1381	1398	1404
chatgpt-4o-latest-20250129	1376	1339	1353	1310	1359	1303	1419	1371	1395	1415
deepseek-r1	1361	1320	1368	1339	1361	1351	1357	1365	1361	1396
gemini-2.0-flash-001	1355	1284	1354	1285	1352	1331	1363	1353	1361	1374

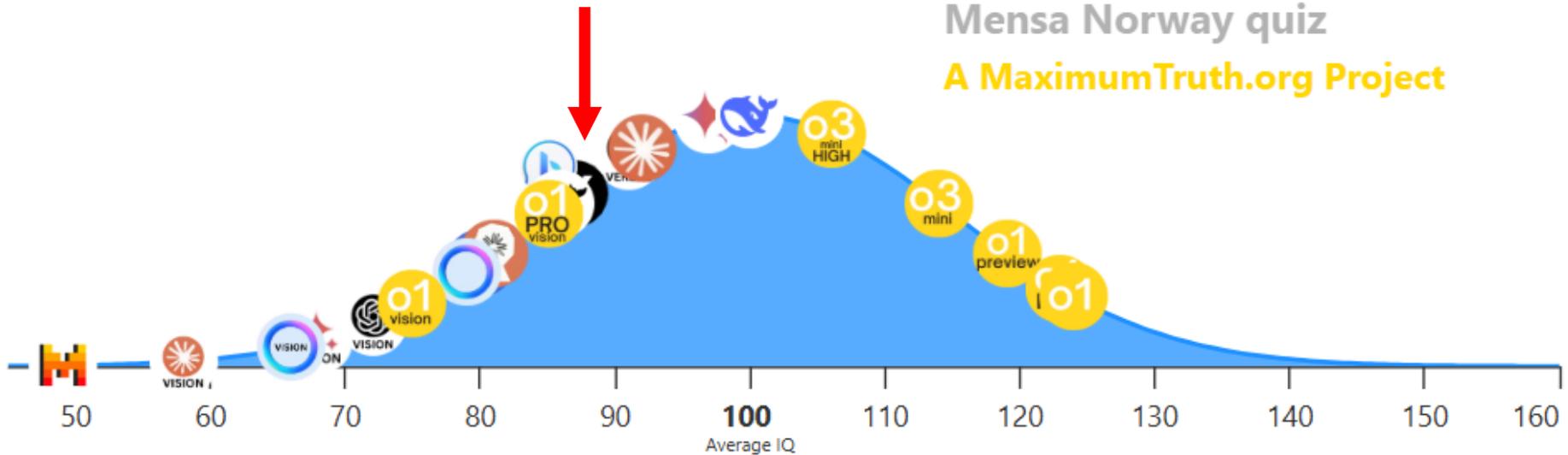
<https://lmarena.ai/?leaderboard>

# IQ Test Results

**Reset**   

Mensa Norway IQ Scores (Average of last 7 tests)

**TrackingAI.org**  
Mensa Norway quiz  
A MaximumTruth.org Project



- |  |   |  |
|--|---|--|
|  Gemini Advanced (Vision) |  Gemini Advanced       |  GPT4 Omni (Vision)   |
|  GPT4 Omni               |  ChatGPT-4            |  Grok-3 ←            |
|  Mistral                |  Claude-3.5 Sonnet   |  Llama-3.3          |
|  Claude-3 Opus (Vision) |  Bing Copilot        |  Llama-3.2 (Vision) |
|  Claude-3 Opus          |  OpenAI o1 preview   |  DeepSeek V3        |
|  OpenAI o3 mini         |  OpenAI o3 mini high |  DeepSeek R1        |

# Rank by Test Source

Reset

Show Offline Test

Show Mensa Norway



Rank is sorted by Mensa Norway IQ Scores

Mensa Norway

