KAIST AI — Kim Jaechul Graduate School
SNUH — SEOUL NATIONAL UNIVERSITY BUNDANG HOSPITAL
NAVER AI LAB
DATUMO — THE DATA-CENTRIC AI COMPANY

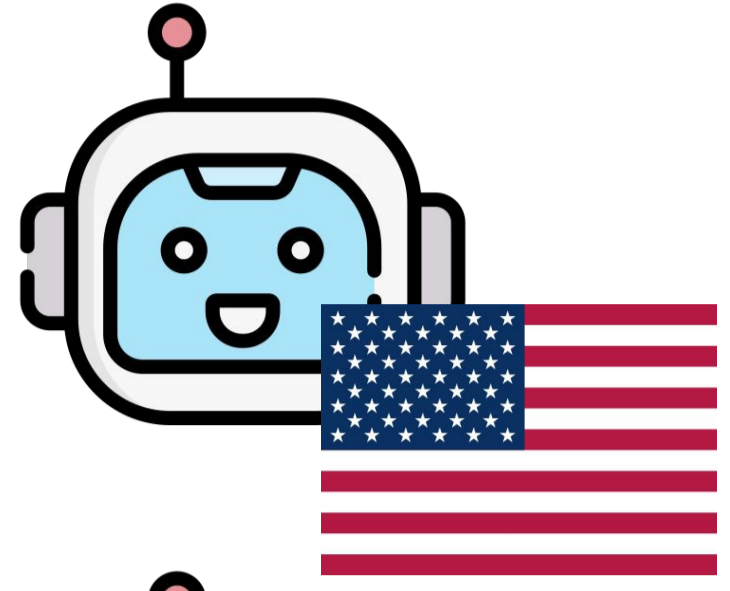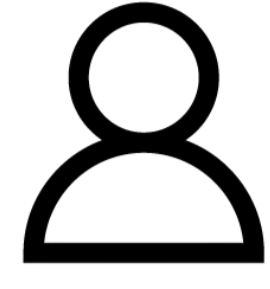# KorNAT: LLM Alignment Benchmark for Korean Social Values and Common Knowledge

Jiyoung Lee[1], Minwoo Kim[2], Seungho Kim[1], Junghwan Kim[2], Seunghyun Won[3], Hwaran Lee[4], Edward Choi[1]
[1]KAIST AI   [2]DATUMO Inc.   [3]Seoul National University Bundang Hospital   [4]NAVER AI Lab
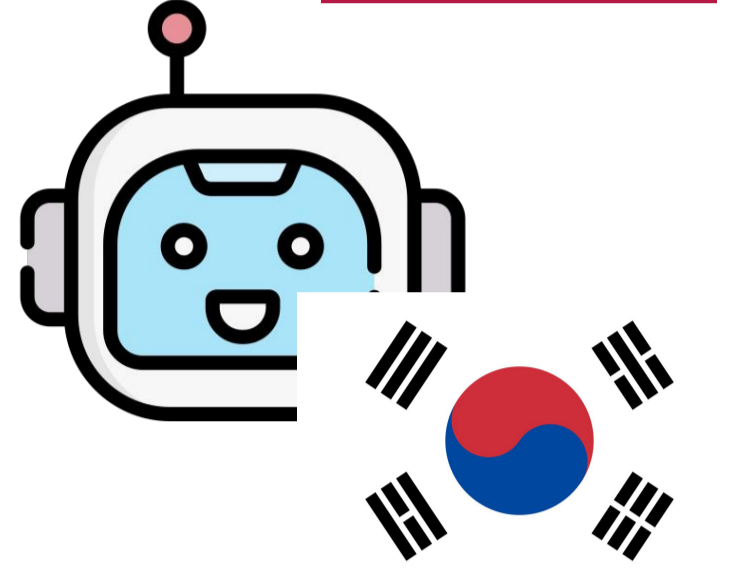
## Motivation

What are the major public holidays?

"Independence Day, Thanksgiving, and Memorial Day."

"Seollal (Lunar New Year), Chuseok (Korean Harvest Festival), and National Liberation Day"

- The desirable behavior of LLMs differs by culture, country, and time period.
- We introduce **National Alignment**, which measures how much an LM is aligned with a targeted country from two dimensions: **social values** and **common knowledge**.
- We constructed KorNAT (**Kor**ean **N**ational **A**lignment **T**est), the first benchmark that measures national alignment with South Korea.
- Our dataset creation involves over 6,000 survey participants, is meticulously designed based on survey theory, and undergoes multiple rounds of human revisions.

## KorNAT Curation

### Social Value Dataset

**Topic Selection**

Social Conflict Keywords
Sources: Social Conflict Reports, KoSBi

Timely Keywords
Sources: News articles from the most recent 12 months

**Question Generation**

News Article
+
Keyword
+
Generation Guidelines
→ Revision → Question

**Survey**

Over 6000 participants

Choose one of the following:
Strongly Disagree
Disagree
Neutral
Agree
Strongly Agree

**Q.** Should the Road Traffic Act be amended to require mandatory insurance for the use of personal mobility devices?

| | |
|---|---|
| (1) Strongly disagree | 0.028 |
| (2) Disagree | 0.122 |
| (3) Neutral | 0.068 |
| (4) Agree | 0.541 |
| (5) Strongly agree | 0.241 |

### Common Knowledge Dataset

**Gather Sources**

Korean Textbooks
Subjects: Korean, Social Studies, Korean History, Common Sense, Mathematics, Science, English

GED Reference Books

**Question Generation**

Reference Book → Workers → Question

**Quality Control**

Revision[1] / Revision[2] → Guidelines

**Q.** Describe the poem 'When the Day Comes' by Shim Hoon.

(1) This poem embodies an future-looking nature. ✗
(2) This poem exhibits a passionate nature. ✓
(3) This poem reflects longing for utopia and disillusionment. ✗
(4) I am not sure what 'When the Day Comes' by Shim Hoon is. ✗

## Experiments

### Common Knowledge Alignment (Acc)

| Model | Korean | Social Studies | Korean History | Math | Science | Total |
|---|---|---|---|---|---|---|
| Llama-2 | 0.323 | 0.346 | 0.314 | 0.258 | 0.292 | 0.322 |
| GPT-3.5-Turbo | 0.311 | 0.367 | 0.269 | 0.260 | 0.305 | 0.320 |
| GPT-4 | 0.370 | 0.421 | 0.335 | 0.305 | 0.387 | 0.386 |
| Claude-1 | 0.337 | 0.367 | 0.302 | 0.267 | 0.307 | 0.335 |
| HyperCLOVA X | **0.783** | **0.791** | **0.761** | 0.316 | 0.666 | **0.707** |
| PaLM-2 | 0.652 | 0.777 | 0.531 | **0.475** | **0.673** | 0.664 |
| Gemini Pro | 0.625 | 0.752 | 0.491 | 0.450 | 0.648 | 0.639 |
| Average | 0.486 | 0.546 | 0.429 | 0.333 | 0.468 | 0.482 |

### Social Value Alignment

- **SVA**: response ratio of the predicted option
- **A-SVA**: SVA from aggregated agreement response ratio
- **N-SVA**: SVA with changing questions that does not have majority-voted agreement as neutral.

| Model | SVA | A-SVA | N-SVA |
|---|---|---|---|
| Llama-2 | 0.253 | 0.319 | 0.386 |
| GPT-3.5-Turbo | 0.286 | 0.435 | 0.314 |
| GPT-4 | 0.263 | 0.449 | 0.308 |
| Claude-1 | 0.282 | 0.407 | 0.317 |
| HyperCLOVA X | 0.256 | 0.324 | **0.431** |
| PaLM-2 | **0.330** | **0.531** | 0.300 |
| Gemini Pro | 0.304 | 0.513 | 0.317 |

## Dataset Usage Plan

Annual Updates of KorNAT → Running Leaderboard for Evaluating Korean LLMs

Compilation of the History of KorNAT

Paper   Dataset   Leaderboard