

Overview of the NTCIR-12 QA Lab-2 task

Hideyuki Shibuki (YNU), Kotaro Sakamoto (YNU, NII), Madoka Ishioroshi (NII), Akira Fujita (NII),
Yoshinobu Kano (Shizuoka Univ.), Teruko Mitamura, (CMU) Tatsunori Mori (YNU), Noriko Kando (NII, SOKENDAI)



Goal

Investigation of the **real-world complex QA** technologies
using Japanese **university entrance exams** and their English
translation on the subject of "World history"

Highlights

1. Multiple-choice and **free-description** type questions
2. Understanding of the surrounding context
3. Evaluation of free-description, especially, **essay questions**
4. Construction of a hierarchy of question formats (see the paper)
5. Competition with high-school students from all over Japan

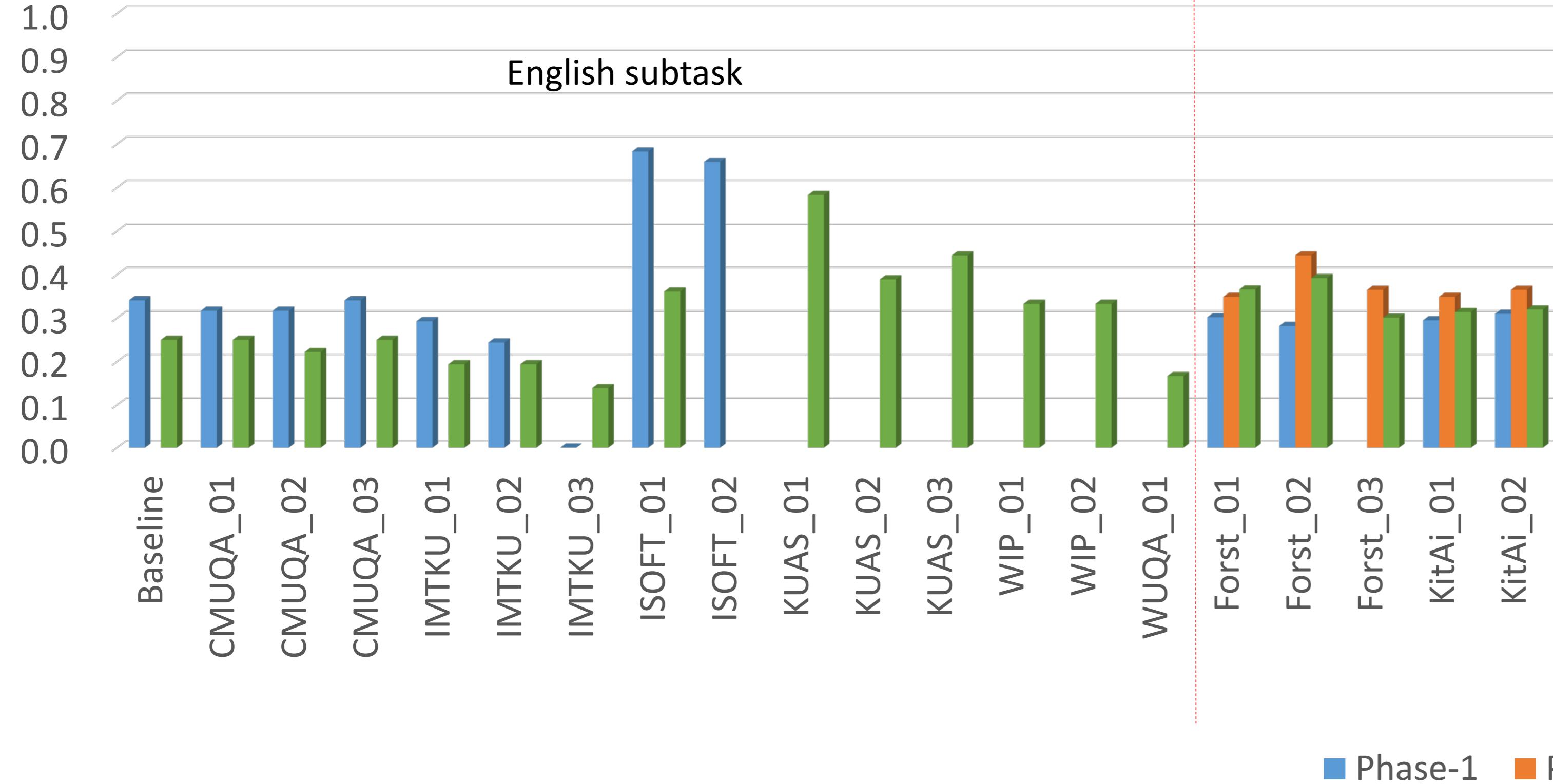
Data Set

Exam name	Type	Training	Phase-1 (EN & JA)	Phase-2 (JA only)	Phase-3 (EN & JA)
Center Test	Multiple choice	97,01,03,05, 07,09	99	N/A	11
Secondary Exam	Free description	00,05,07,09	01,03,06,10	N/A	02,04,08,11
Yozemi (JA only)	Multiple choice	13b,13c	12,13a	N/A	13d,14a
Benesse (JA only)	Multiple choice	14Jun	14Nov	15Jun	14Sep
Sundai (JA only)	Free description	14Aug,14Nov	13Nov	15Aug	13Aug

Participants are free to participate any particular phase and either of exams.

Results

Correct rates of multiple-choice questions



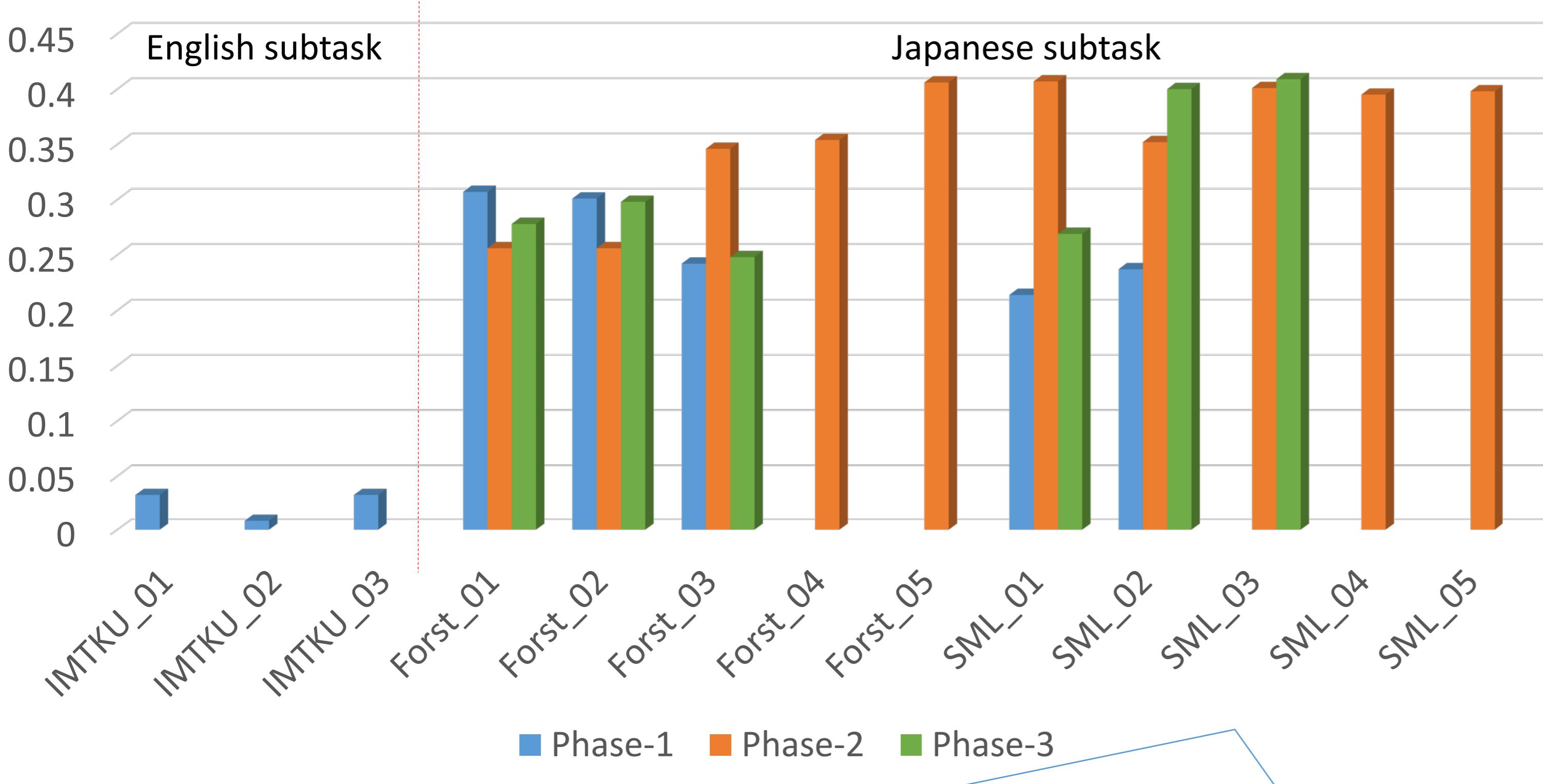
English subtask

Japanese subtask

Combination runs

$$\text{Correct Rate} = \frac{\text{number of correct answers}}{\text{total of inputted questions}}$$

ROUGE-1 scores of essay questions



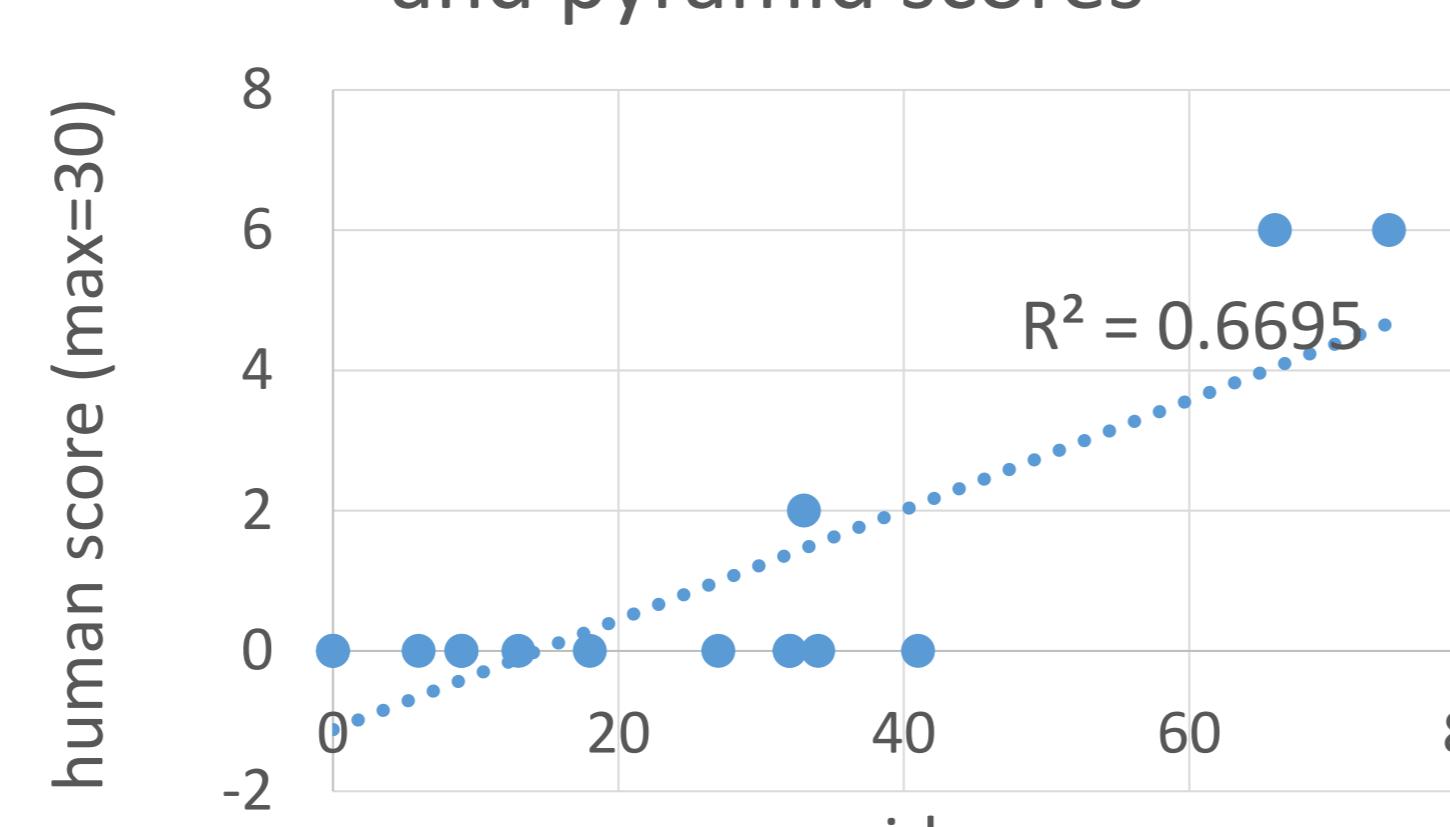
How much does ROUGE-1 score affect the evaluation of essay questions?

A *human expert evaluated outputs with top priority*.

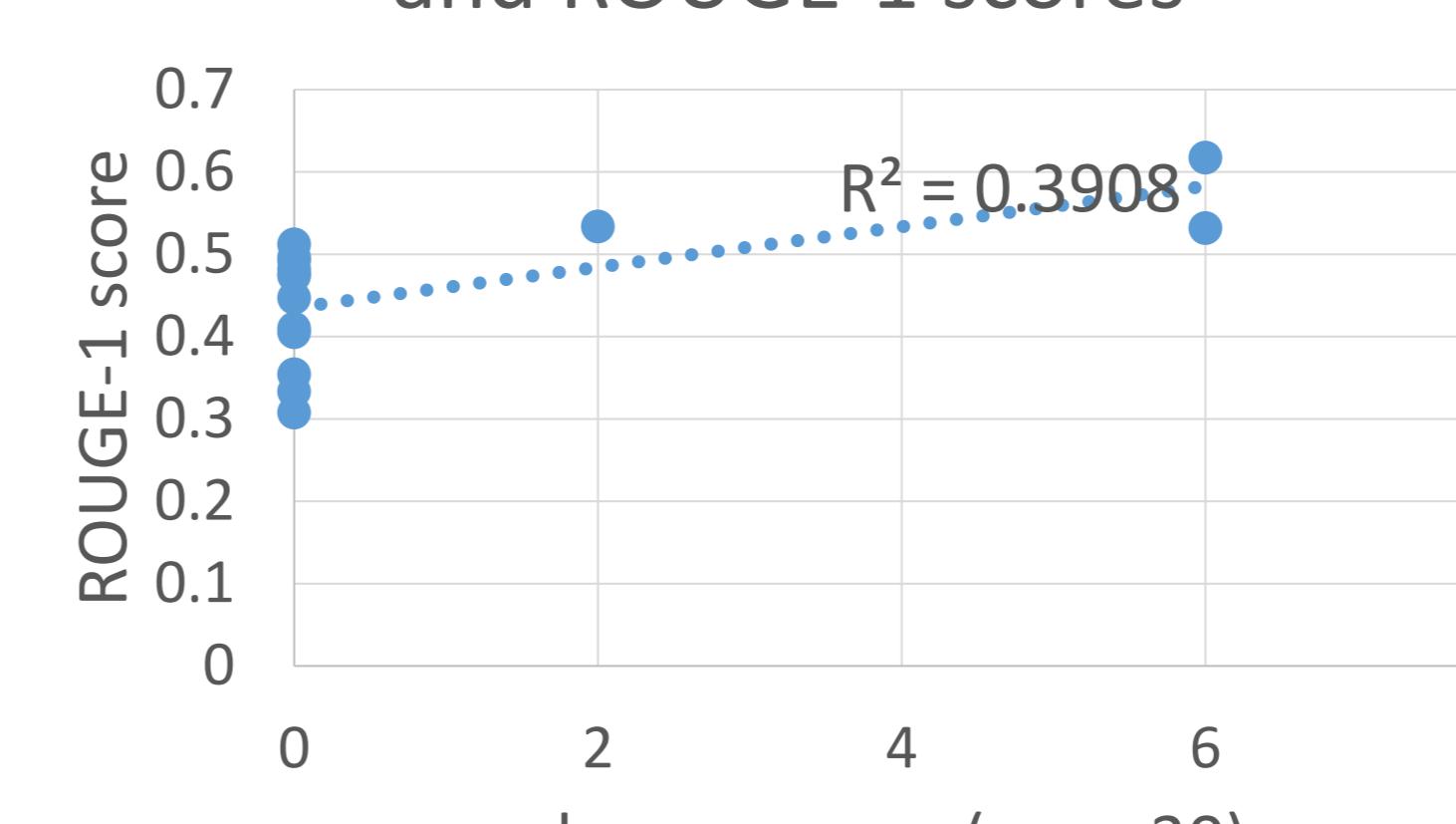
Three human experts made nuggets for pyramid method.

A weak/moderate positive correlation was found, but ROUGE-1 is not enough to get logical consistency.

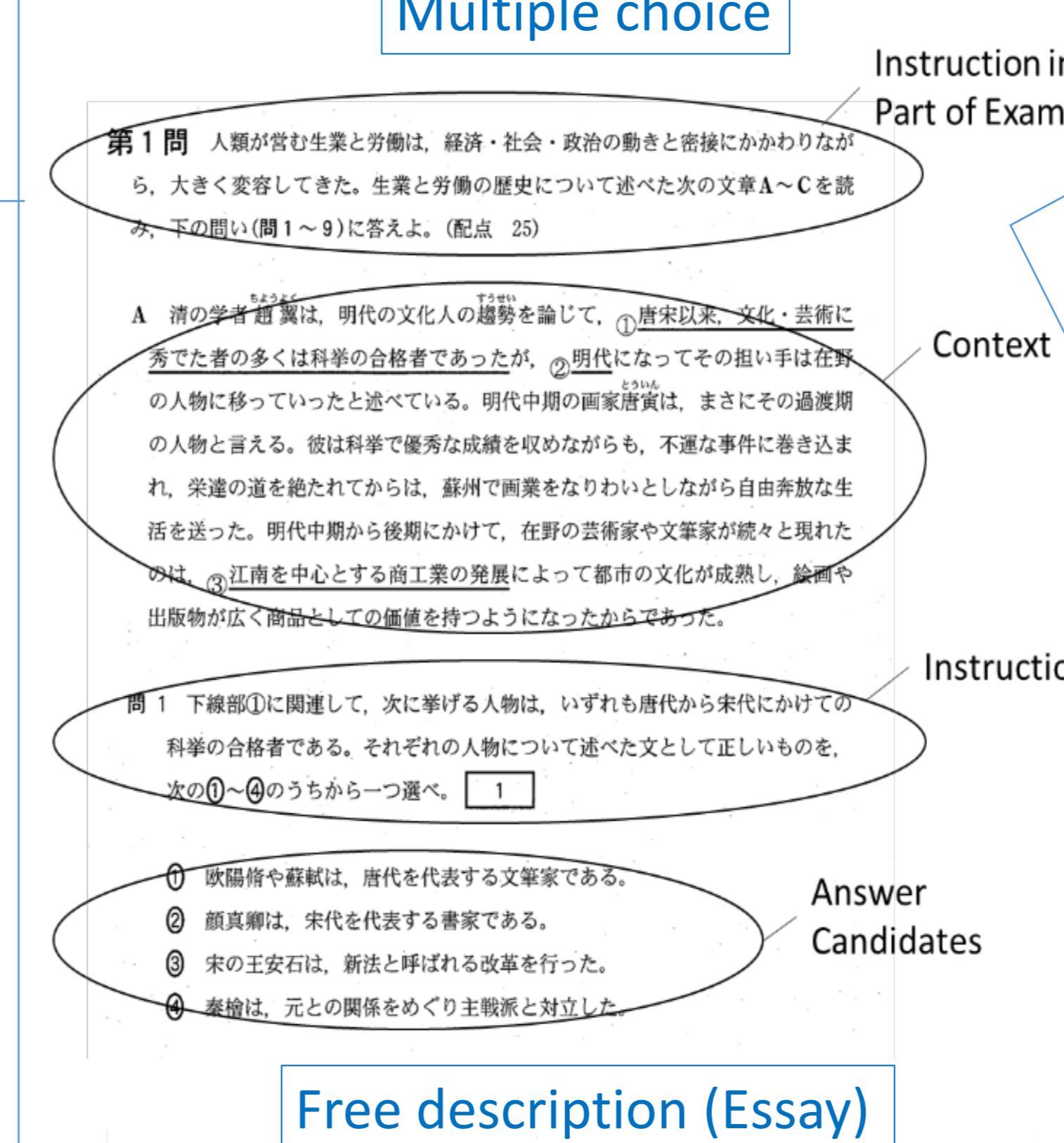
Correlation between human and pyramid scores



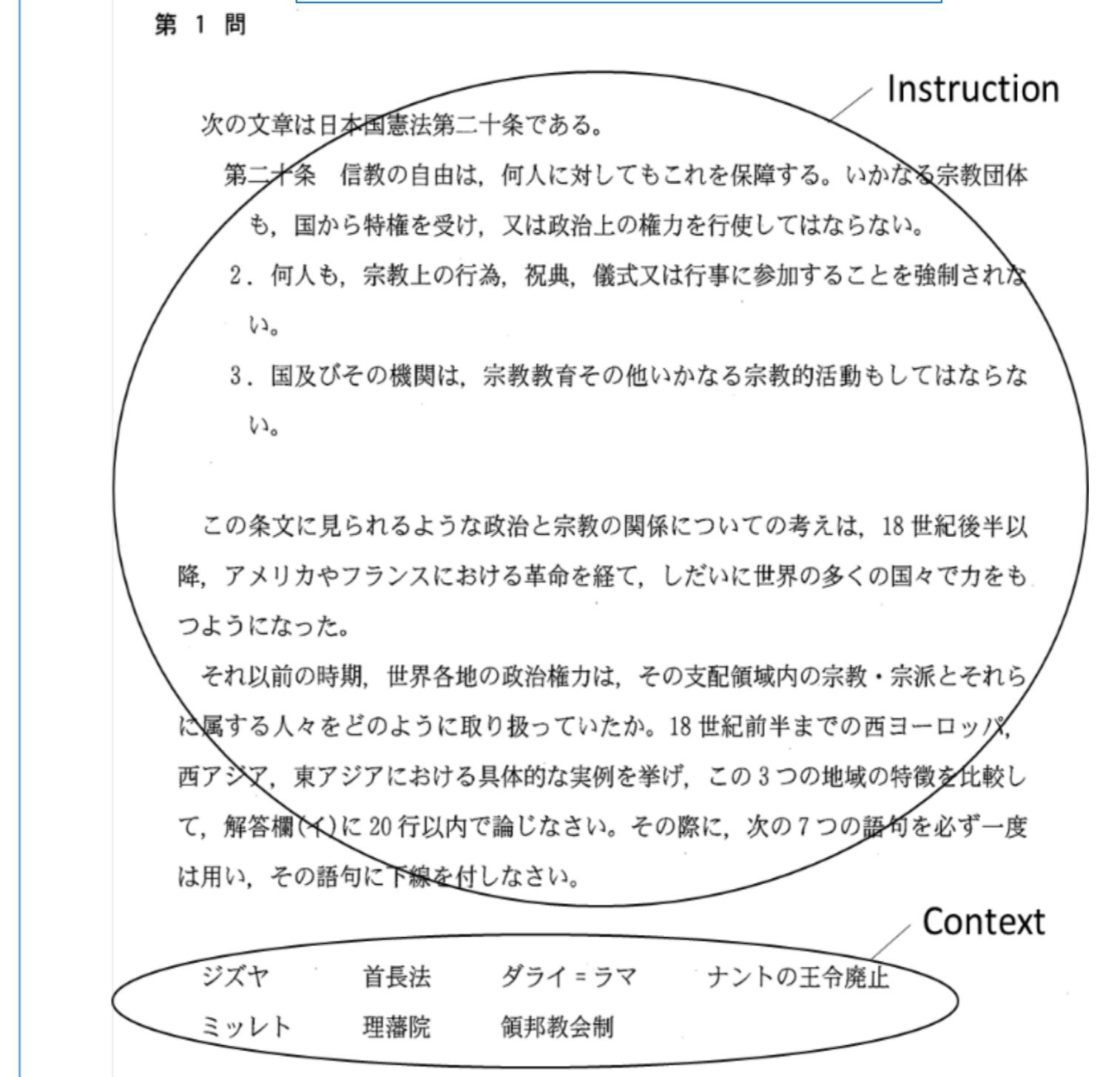
Correlation between human and ROUGE-1 scores



Question Data

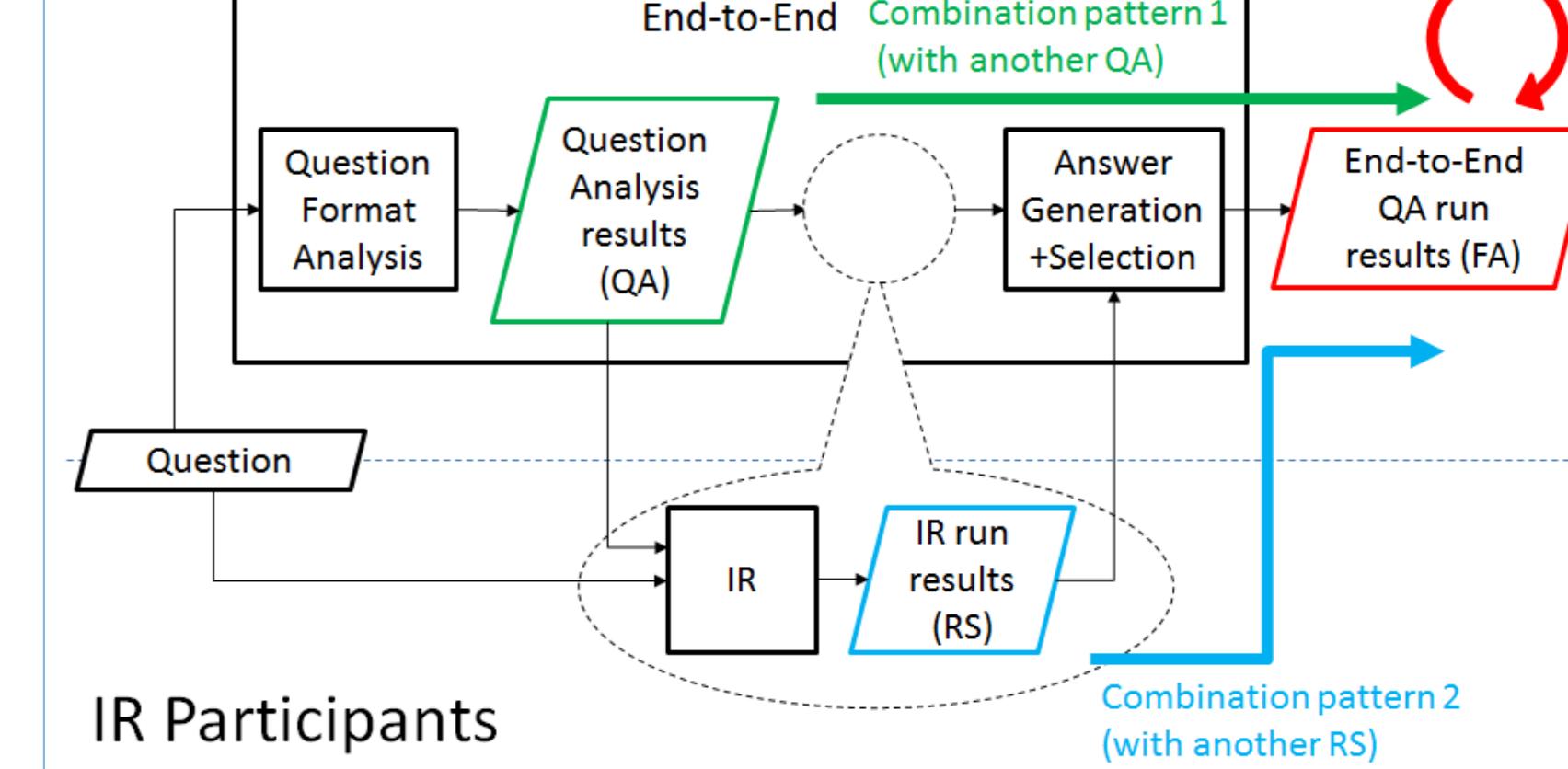


Free description (Essay)

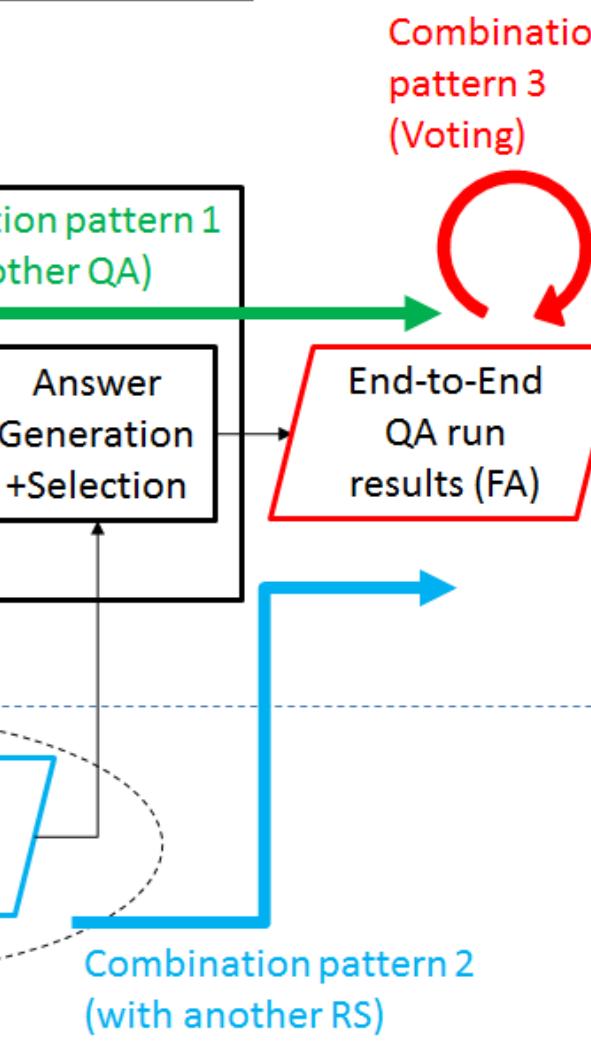


Combination Pattern

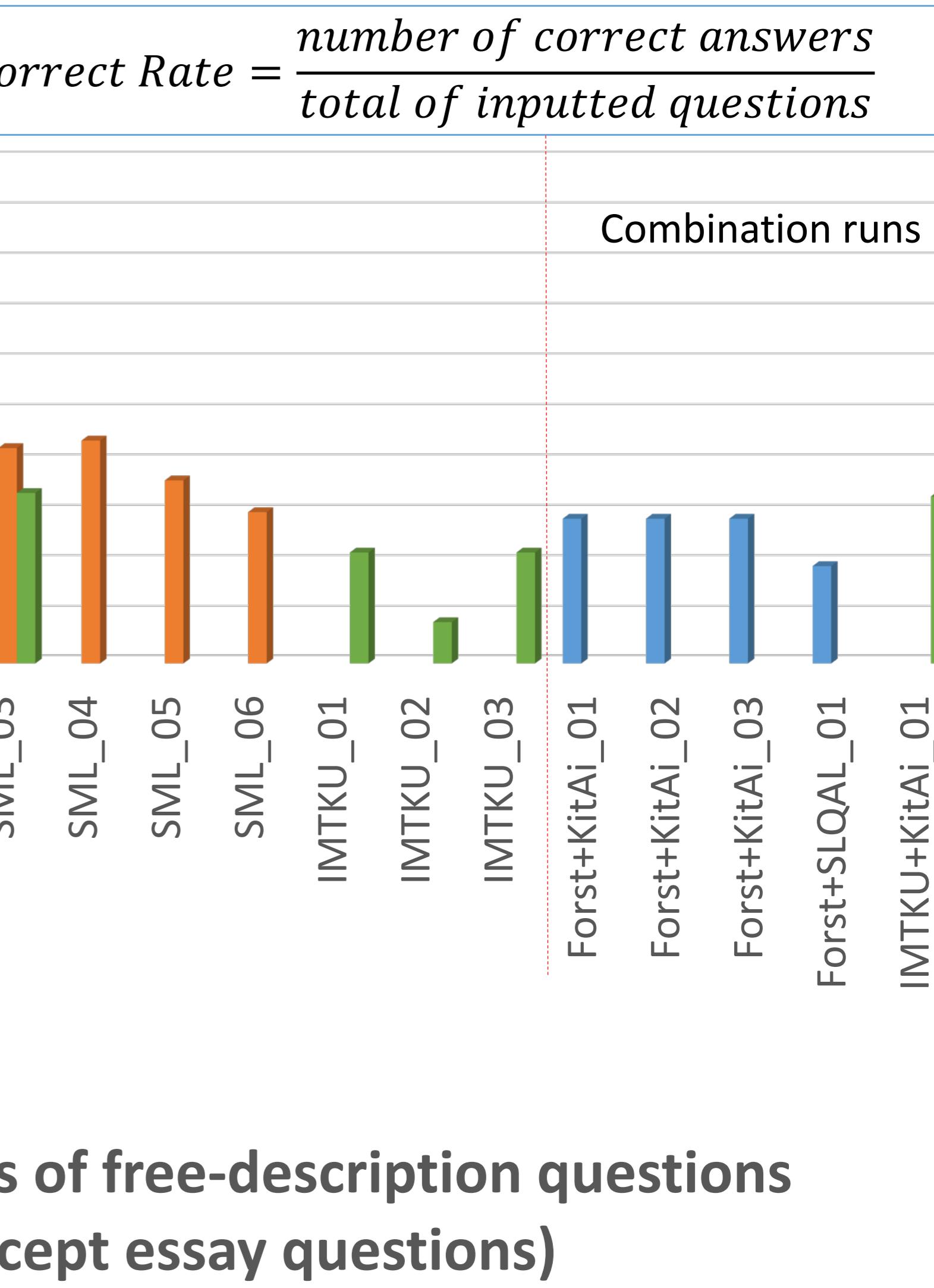
QA Participants



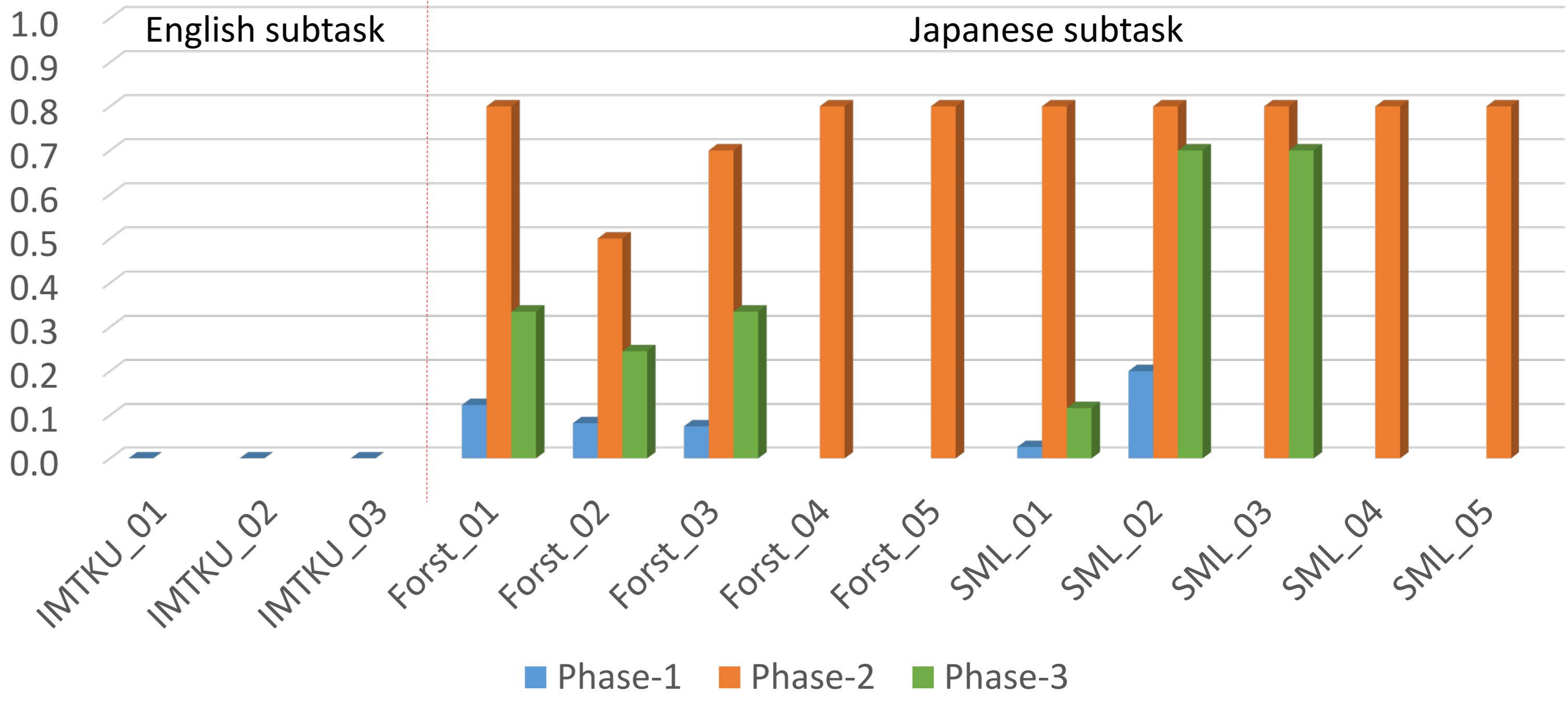
IR Participants



Combination pattern 3 (Voting)



Correct rates of free-description questions (except essay questions)



Correlation between pyramid and ROUGE-1 scores

