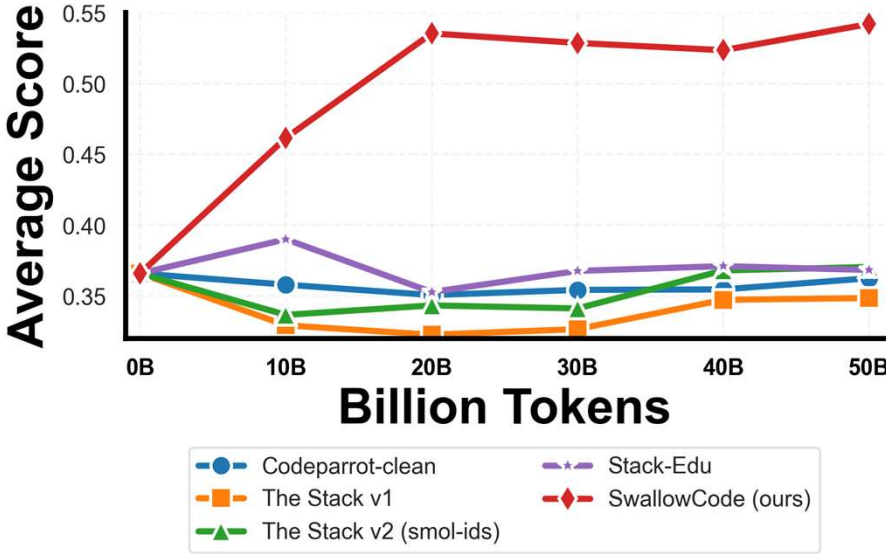


ABCI 3.0開発加速利用（2024年度）成果概要（公開用）

| 課題名：<br>継続事前学習に基づく汎用性および有用性の高い日本語大規模言語モデルの構築   | 実施時期：2025年1月から3月<br>所属機関名：産業技術総合研究所<br>代表者氏名：高村大也  |                |                         |              |                         |           |                    |    |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |
|--|--|----------------|-------------------------|--------------|-------------------------|-----------|--------------------|----|------|------|------|------|------|-----|------|------|------|------|------|-----|------|------|------|------|------|-----|------|------|------|------|------|-----|------|------|------|------|------|-----|------|------|------|------|------|
| 成果概要：数学とコードに強いLlama-3.1-Swallow-8B-v0.5構築のためのコードデータセット作成を行った。オープンなコードデータセットであるThe Stack v2にはSyntax Errorやコーディング規約を守らないコードが多数存在している。そこでLlama-3.3-70B-Instructを利用し、高品質なコードデータセットに作り直すデータ合成を行った。これにより、高品質なコードデータセットを構築することができた。 |  |                |                         |              |                         |           |                    |    |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |
| 成果のポイント：<br>NeurIPS 2025 Dataset trackに投稿予定の論文でも報告するように、作成したデータセットで学習を行うことで、他の主要なオープンコードデータセットで学習したLLMを著しく上回る性能を発揮するLLMを構築することができた。<br>このデータセットは、Llama-3.3 Licenseのもとでhuggingfaceにて公開する予定である。                                |  <table><caption>Average Score vs Billion Tokens</caption><tr><th>Billion Tokens</th><th>Codeparrot-clean</th><th>The Stack v1</th><th>The Stack v2 (smol-ids)</th><th>Stack-Edu</th><th>SwallowCode (ours)</th></tr><tr><td>0B</td><td>0.37</td><td>0.37</td><td>0.37</td><td>0.37</td><td>0.37</td></tr><tr><td>10B</td><td>0.36</td><td>0.33</td><td>0.34</td><td>0.39</td><td>0.46</td></tr><tr><td>20B</td><td>0.35</td><td>0.32</td><td>0.34</td><td>0.35</td><td>0.53</td></tr><tr><td>30B</td><td>0.35</td><td>0.32</td><td>0.34</td><td>0.37</td><td>0.52</td></tr><tr><td>40B</td><td>0.36</td><td>0.35</td><td>0.37</td><td>0.37</td><td>0.52</td></tr><tr><td>50B</td><td>0.36</td><td>0.35</td><td>0.37</td><td>0.37</td><td>0.54</td></tr></table> | Billion Tokens | Codeparrot-clean        | The Stack v1 | The Stack v2 (smol-ids) | Stack-Edu | SwallowCode (ours) | 0B | 0.37 | 0.37 | 0.37 | 0.37 | 0.37 | 10B | 0.36 | 0.33 | 0.34 | 0.39 | 0.46 | 20B | 0.35 | 0.32 | 0.34 | 0.35 | 0.53 | 30B | 0.35 | 0.32 | 0.34 | 0.37 | 0.52 | 40B | 0.36 | 0.35 | 0.37 | 0.37 | 0.52 | 50B | 0.36 | 0.35 | 0.37 | 0.37 | 0.54 |
| Billion Tokens   | Codeparrot-clean   | The Stack v1   | The Stack v2 (smol-ids) | Stack-Edu    | SwallowCode (ours)      |           |                    |    |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |
| 0B   | 0.37   | 0.37           | 0.37                    | 0.37         | 0.37                    |           |                    |    |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |
| 10B  | 0.36   | 0.33           | 0.34                    | 0.39         | 0.46                    |           |                    |    |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |
| 20B  | 0.35   | 0.32           | 0.34                    | 0.35         | 0.53                    |           |                    |    |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |
| 30B  | 0.35   | 0.32           | 0.34                    | 0.37         | 0.52                    |           |                    |    |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |
| 40B  | 0.36   | 0.35           | 0.37                    | 0.37         | 0.52                    |           |                    |    |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |
| 50B  | 0.36   | 0.35           | 0.37                    | 0.37         | 0.54                    |           |                    |    |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |
| 成果についてより詳細な情報を提供しているWebページ、発表論文などの情報：  |  |                |                         |              |                         |           |                    |    |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |     |      |      |      |      |      |