Mosaic-IT: Cost-Free Compositional Data Synthesis for Instruction Tuning

Anonymous ACL submission

Abstract

Finetuning large language models with a variety of instruction-response pairs has enhanced their capability to understand and follow instructions. Current instruction tuning primarily relies on teacher models or human intervention to generate and refine the instructions and responses for training, which are costly, non-sustainable, and may lack diversity. In this paper, we introduce Mosaic Instruction Tuning (Mosaic-IT), a human/model-free compositional data synthesis method that can efficiently create rich and diverse augmentations from existing instruction tuning data to enhance the LLMs. Mosaic-IT randomly concatenates multiple instruction data into one and trains the model to produce the corresponding responses with predefined higher-level meta-instructions to strengthen its multi-step instruction-following and format-following skills. Our extensive evaluations demonstrate a superior performance and training efficiency of Mosaic-IT, which achieves consistent performance improvements over various benchmarks and a 80% reduction in training costs compared with original instruction tuning. Our codes and data are available at https://anonymous. 4open.science/status/mosaic-0554.

1 Introduction

003

007

014

017

027

042

The emergence of Large Language Models (LLMs) (Scao et al., 2022; OpenAI, 2023; Touvron et al., 2023a) along with their remarkable performance in down-stream tasks (Zhao et al., 2023; Xu et al., 2024a), has revolutionized the domains of Artificial Intelligence and Natural Language Processing. A key component of the recipe to unlock the exceptional ability of LLMs in understanding and following instructions is the technique of Instruction Tuning (IT) (Mishra et al., 2021; Wei et al., 2022; Chung et al., 2022), which involves the fine-tuning of LLMs on datasets comprising corresponding instruction-response pairs.



Figure 1: The brief illustration of our Mosaic-IT with different strategies. Given the original dataset, our method randomly samples and concatenates them together into more complex samples, simulating the multi-instruction-following scenarios *in a cost-free manner*.

043

044

045

047

051

054

057

060

061

062

063

065

To ensure the quality of instruction tuning data, earlier efforts (Brown et al., 2020; OpenAI, 2023; Touvron et al., 2023a; Jiang et al., 2023) carefully curate extensive, diverse, and high-quality datasets manually. Although these datasets encompass a wide range of instructions to improve instruction tuning, they require the responses to be meticulously curated by human experts (Khashabi et al., 2020; Ye et al., 2021; Wei et al., 2022; Wang et al., 2022; Du et al., 2022). Alternatively, some approaches (Wang et al., 2023b; Taori et al., 2023; Xu et al., 2023; Li et al., 2023a) leverage more capable teacher LLMs to reduce the labor-intensive process of data generation. For example, the Alpaca (Taori et al., 2023) utilizes self-instruct (Wang et al., 2023b) to automatically generate diverse instruction tuning datasets. Building on this trend and the widely acknowledged notion that more complicated instructions are more beneficial for LLMs' instruction-following ability (Xu et al., 2023; Zhao et al., 2024), numerous strategies (Zhao et al., 2024; Wu et al., 2024; Ding et al., 2023; Li et al., 2023a; Liu et al., 2023a; Li et al., 2024c,b;

077

084

091

100

101

102

104

105

106

107

109

110

111

112

113

114

115

116

117

066

Guo et al., 2024; Xu et al., 2024a) have been proposed to further diversify and complexify the instruction-response pairs, utilizing teacher models like ChatGPT-3.5 and GPT-4 (OpenAI, 2023).

Despite the enhanced performance in instructionfollowing ability offered by these existing methods, they face **Two** major issues: (1) They heavily rely on teacher models or human annotators to rewrite instruction-response pairs, which highlights the resource-intensive nature and their constraints on scalability; (2) They only increase the complexity within the scope of a single instruction, which limits the potential improvement in LLMs' instructionfollowing capabilities. Motivated by the Dense and Aligned Captions (Doveh et al., 2023) proposed for vision language (VL) models and the mosaic data augmentation proposed in Yolov4 (Bochkovskiy et al., 2020), we hypothesize that denser instructions benefit the LLM alignment, i.e. the process of instruction tuning should not be constrained by one single instruction but be extended to follow several instructions at a time, which represents a higher level of instruction-following ability that is beneficial to the training process. A similar concept during the inference phase is proposed by batch prompting (Cheng et al., 2023; Lin et al., 2024), where multiple samples are grouped in one batch allowing LLMs to generate multiple responses at one inference, while its performances are sub-optimal.

As orthogonal to the existing instruction tuning methods, we introduce Mosaic Instruction Tuning (Mosaic-IT), an innovative and model/human-free compositional approach that augments existing instruction tuning datasets, which concurrently improves the LLM performances and lowers the training expenses. As shown in Figure 1, in our method, multiple instructions and corresponding responses from the original dataset are concatenated into a single sample for fine-tuning, simulating the multi-instruction-following scenarios at no cost. Without applying any additional strategies, we term this simple process as the **Primary Mosaic** Strategy. We posit that this mosaic strategy process significantly improves the complexity and density of the original instructions, learning from which directly benefits LLMs in their instructionfollowing ability. Additionally, this method offers the advantage of directly reducing the total count of instruction-response pairs, thereby cutting down on training iterations, and accelerating the training process significantly by approximately 80% reduction. Though effective, the Primary Mosaic strategy

constrains LLMs in responding to the instructions in the original order and format, potentially limiting its further potential. Thus we further introduce three Advanced Mosaic Strategies aimed at enhancing the diversity and complexity of the mosaicked instruction-response pairs: Format, Permute, and Maskout, in which an additional metainstruction is provided as a higher-level guideline for LLMs to follow the given instructions. Illustrative examples are presented in Figure 2. Specifically, in the Format strategy, some arbitrary parsing formats will be defined in the meta-instruction thus forcing LLMs to follow these formats, which notably enhances the LLMs' capacity to follow formats. In the Permutation strategy, an arbitrary permuted order is defined thus forcing LLMs to respond in a desired order. In the Maskout strategy, some arbitrary instructions are sampled which meta-instruction forces LLMs to ignore. Moreover, the use of these Advanced strategies not only boosts the performance in several evaluation metrics but also keeps our method free of additional costs.

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

In summary, our primary contributions can be illustrated as follows:

- We propose the *first cost-free data Synthesis method*, **Mosaic-IT**, which extends existing instruction tuning from handling one single instruction at a time to following multiple instructions in diverse forms. This approach significantly enhances the potential utilization of existing high-quality datasets.
- Mosaic-IT improves the instruction-following abilities of LLMs compared to training on original data, as evidenced by consistent performance gains across a wide range of benchmarks, model families, and datasets, demonstrating strong generalization capabilities.
- Mosaic-IT substantially increases training efficiency by reducing the required number of training iterations, resulting in an approximate 80% reduction in training time, as confirmed by experimental results.

2 Methodology

2.1 Preliminaries

The instruction tuning dataset, defined as D, consists of n data samples, each represented by a triplet (*Instruction*, *Input*, *Response*). For simplicity, we define x = map(Instruction, Input) as the unified instruction, and y as the corresponding



Figure 2: Illustrative examples of Mosaic-IT. Given 3 simple data instances, our method concatenates them into data samples with diverse forms. Texts in red represent the meta-instructions that define the formats or orders for LLMs to respond. Texts in yellow are major response differences of each strategy. The **Primary Strategy** only concatenates data. The **Format Strategy** requires LLMs to respond in predefined formats. The **Permute Strategy** requires LLMs to respond in specific orders and the **Maskout Strategy** requires LLMs to ignore some instructions.

response. Therefore, D can be represented as 167 $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, denoting a set of 168 *n* instruction-response pairs. Let $p_{\theta}(\cdot)$ denote 169 170 the LLMs to be trained, with parameters θ . In the instruction tuning setting, p_{θ} is typically 171 fine-tuned by maximizing the following objective 172 on each data (x_i, y_i) , $y_{i,j}$ represents the j_{th} token 173 of response $y_i, y_{i,<j}$ represents the tokens prior to 174 175 $y_{i,j}$, and l_i represents the token length of $y_{i,j}$:

$$\max_{\theta} \sum_{j=1}^{l_i} \log p_{\theta} \left(y_{i,j} | x_i, y_{i,$$

2.2 Mosaic-IT

176

178

Motivated by the success of the existing data-179 centric instruction tuning methods, a line of 180 approaches is proposed to further enhance the instruction-response pairs utilizing extra teacher LLMs (Xu et al., 2024a). Though effective, all existing methods for instruction tuning restrict training samples to just one instruction, which severely 186 limits the potential of the existing high-quality data and the instruction-following ability of the mod-187 els to be trained. Motivated by the Dense and Aligned Captions (Doveh et al., 2023) for VL, we hypothesize that denser instructions benefit the 190

LLM alignment, thus the process of instruction tuning should not be constrained by one single instruction but be extended to follow several instructions at a time, which represents a higher level of instruction-following ability that is beneficial to the training process. Thus, we propose the cost-free data synthesis method, Mosaic Instruction Tuning (Mosaic-IT) as shown in Figure 1.

191

192

193

194

195

196

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

2.2.1 Primary Mosaic strategy

Exploring the concept of concatenating random instruction-response pairs into a unified instructionresponse pair for training remains largely unexplored. The primary challenge lies in crafting a coherent overall instruction and obtaining its corresponding response. Most existing methods utilize a strong teacher model to rewrite the instructions with prompting techniques and generate corresponding responses, introducing more cost by actually re-generating new data samples. To harness the full potential of existing data rather than directly discarding them, we introduce a simple cost-free compositional approach as shown in Figure 2, in which instructions are randomly concatenated with serial digits to form an *overall instruc*- *tion.* The concatenated overall instruction is denoted as $[x_1, ..., x_k]$, with the corresponding overall response concatenated as $[y_1, ..., y_k]$. Here, k denotes the number of original data samples integrated into each overall sample.

215

216

217

218

221

227

228

229

230

231

234

235

237

240

241

243

244

246

247

249

250

251

255

256

260

261

264

In this framework, the fundamental instructionfollowing capability is triggered by the existing instruction-response pairs, and the mosaic strategy extends this capability to a higher level in which LLMs are forced to follow multiple instructions. It represents a much more complicated scenario that benefits LLMs compared with traditional singletask instructions. Consequently, the objective function for each concatenated overall data sample can be formulated as follows:

$$\max_{\theta} \sum_{j=1} \log p_{\theta} \left([y_1, ..., y_k]_j | [x_1, ..., x_k], [y_1, ..., y_k]_{< j} \right),$$
(2)

Here, $[y_1, ..., y_k]_j$ denotes the *jth* token of the overall response, $[y_1, ..., y_k]_{<j}$ denotes the tokens prior to *jth* token, and *l* represents the length of overall response. This formulation encapsulates the essence of our approach, optimizing the model parameters θ to maximize the likelihood of generating the correct sequence of responses for the given overall instruction.

2.2.2 Advanced Mosaic Strategies

l

Though effective, this primary mosaic strategy constrains LLMs in responding to the instructions with the original order and format, potentially limiting its generalization and practical usage. In our method, the instructions and corresponding responses from the original dataset can be viewed as atomic components and our method randomly combines these elements together to form new instructions and responses. This nature allows us to further complicate this process with fancier strategies thus forcing LLMs to follow more complicated overall instructions. Hence, we propose three Advanced Mosaic Strategies to complicate and diversify the mosaicked samples as shown in Figure 2, including Format, Permute, and Maskout, with metainstructions guiding them. These strategies are still purely rule-based (Li et al., 2024a) and do not incorporate the additional human/LLM generation.

Format In the Format strategy, some arbitrary formats are defined in the meta-instruction to force LLMs to follow these formats in the response. The formats mainly contain two categories: 1) *Serial Digit Format* and (2) *Response Parsing Format*. The serial digits establish the initial instruction order that guides LLMs to follow sequentially. We manually define 10 types of serial digit format, which will be randomly sampled during each mosaic process. For response parsing, we simulate the scenario where the users try to extract specific information from the responses. We define 27 types of parsing brackets and 17 types of parsing text pairs, which will be randomly sampled and assembled during each mosaic process. Examples can be found in Appendix G, which can be easily extended for customized training settings. We denote responses with specific formats as $y'_i = wrap(y_i, s_{format})$, and l as the token length of the overall response. An additional metainstruction s_{format} specifying the required format will be included in the overall instruction. Thus, the objective function for each mosaic data point:

265

266

267

269

270

271

272

273

274

275

276

277

278

279

285

286

287

289

290

291

292

293

294

295

296

297

298

299

300

301

303

304

305

306

307

308

309

310

311

$$\max_{\theta} \sum_{j=1}^{l} \log p_{\theta} \Big([y'_1, ..., y'_k]_j \Big| [x_1, \tag{3}$$

$$x_2, ..., x_k, s_{format}], [y'_1, ..., y'_k]_{< j}$$

Permute and Maskout Building upon the Format strategy, we further introduce two strategies for our Mosaic-IT, Permutation and Maskout.

In the **Permute** strategy, an arbitrary permuted order is defined in the meta-instructions, forcing LLMs to follow. Moreover, several high-level rules are defined to ensure the complexity and diversity of meta-instructions, e.g., forcing LLMs to respond to each instruction in the randomly generated permutation list, forcing LLMs to respond in the alphabetical order of each instruction, forcing LLMs to respond according to the length of instructions, etc. The detailed rule types and descriptions are depicted in Appendix G. These various meta-instructions not only provide higher-level guidelines for LLMs to follow multiple instructions but also inherently enhance the instruction perception ability of LLMs. In our settings, LLMs are required to generate responses selectively conditioned on some critical parts of the overall instruction, forcing them to first understand the formats and other requirements, indicating a more comprehensive understanding of the context given. The meta-instruction is denoted as $s_{permute}$ and is included in the overall instruction. The permuted response list is denoted as $[y'_{1'}, ..., y'_{k'}] =$ $Permute([y'_1, ..., y'_k], s_{permute})$. Thus the objective function can be formulated as below:

312
$$\max_{\theta} \sum_{j=1}^{l} \log p_{\theta} \left([y'_{1'}, ..., y'_{k'}]_j \middle| [x_1, (4) \right)$$

314

315

316

317

320

321

322

324

325

326

327

329

331

333

335

336

337

341

344

347

354

 $x_2, ..., x_k, s_{format}, s_{permute}], [y'_{1'}, ..., y'_{k'}]_{< j}$

In the **Maskout** strategy, some arbitrary instructions are selected in the meta-instructions forcing LLMs to ignore them. Several high-level rules are also defined similarly to the permute strategy, including forcing LLMs to ignore the instructions with given random digits, forcing LLMs to ignore the longest one/several instructions, forcing LLMs to ignore odd-numbered instructions, etc. The details are provided in Appendix G. Similarly, the meta-instruction is denoted as $s_{maskout}$ and the response list is denoted as $[y'_1, ..., y'_m] =$ $Maskout([y'_1, ..., y'_k], s_{maskout})$, where m is the count of responses after masking out. Thus the objective function can be formulated as below:

$$\max_{\theta} \sum_{j=1}^{l} \log p_{\theta} \left([y'_{1}, ..., y'_{m}]_{j} \middle| [x_{1}, (5) \\ x_{2}, ..., x_{k}, s_{format}, s_{maskout}], [y'_{1}, ..., y'_{m}]_{< j} \right)$$

It's important to note that our mosaic strategies entail *no supervision cost*, and the predefined rules are flexible and have the potential for further extension. We utilize the version with three Advanced strategies as our default Mosaic-IT.

How to decide the Number of Instructions k: Number of Instructions denotes the number of original data samples that are integrated into an overall sample. In addition to the detailed mosaic strategies being used, this count also dramatically affects the effect of Mosaic-IT. Our experiments reveal that larger and more diverse numbers of instructions will benefit LLM training. By default, we set the maximum number of instructions as $k_{max} = 10$, and randomly sample an integer that is smaller or equal to k_{max} under a uniform distribution. If the number causes the data sample to be longer than the max length, it will be automatically reduced to the max number which remains the sample length within the limits.

3 Experimental Results

3.1 Main Results

In this section, we present the evaluation results comparing our methods with the baseline methods on *6* baseline models (Mistral-7B (Jiang et al., 2023), Llama2-7B (Touvron et al., 2023b), Llama2-13B, Llama-3-8B (Dubey et al., 2024), Phi-3 (Abdin et al., 2024), Gemma2-2B (Team et al., 2024)) and *4* instruction tuning datasets (Alpaca-GPT4 (Peng et al., 2023), Alpaca (Taori et al., 2023), WizardLM-70k (Xu et al., 2023)), Magpie (Xu et al., 2024b), on *5* commonly used evaluation metrics and additional Human Evaluation. Detailed experimental setup and descriptions of evaluation metrics can be found in Appendix B. 355

356

357

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

Table 1 shows the results on 2 general evaluation settings (Pair-Wise Comparison and Open LLM leaderboard). Pair-wise Winning Score indicates the result directly comparing our models with the corresponding baseline models, which is calculated as (Num(Win)-Num(Lose))/Num(All) +1. These values that are greater than 1.0 represent better responses generated by our models. The performances on the Huggingface Open LLM Leader**board** are also presented, and we bold the greater average values for each comparison. The consistent outperforming results on different base models and datasets represent the effectiveness and robustness of our methods. Results on more advanced baseline models and datasets can be found in Table 2. Our method shows consistent improvements compared with baseline models. The performance gains in pair-wise comparison indicate our method helps LLMs generate more detailed and accurate responses, and the performance gains in the open leaderboard indicate our method helps alleviate harm to the intrinsic capabilities of base LLMs.

To better understand how our method improves the instruction-following abilities of LLMs, we further compare the performance on the other 3 benchmarks for fine-grained analysis as shown in Table 2. On the Alpaca Eval 2 benchmark, our method has a consistent improvement with or without the Length Control (LC), indicating that the improvement of response qualities does not directly originate from the length of responses. On the MT-Bench, the 1-round scores of our method are consistently higher, while the 2-round scores slightly fluctuate, indicating that our method mainly improves the response quality for single-round conversations since our meta instructions only focus on single-round scenarios in this version. On the IF Eval benchmark, our method consistently improves the performances on different settings, both Prompt-level and Instruction-level. Compared with the previous benchmarks, IF Eval mainly focuses on the constraint-following ability of LLMs. The

Model	Detect	Mathad	Pair-wise ↑	Hu	ggingfa	ce Open LLN	1 Leaderb	Leaderboard ↑		
Widdel	Dataset	Method	Winning Score	Average	ARC	HellaSwag	MMLU	TruthfulQA		
		Baseline	1.000	59.70	55.03	78.87	56.01	48.88		
	Alpaca-Gr 14	Mosaic-IT	1.349	63.65	59.04	81.85	60.09	53.62		
Mistral-7B	Alpaga	Baseline	1.000	55.15	51.96	74.61	52.85	41.20		
	Аграса	Mosaic-IT	1.390	58.86	56.23	79.57	57.06	42.58		
	Wirord 70k	Baseline	1.000	57.86	51.88	77.93	53.76	47.89		
	Wizaru-70K	Mosaic-IT	1.161	61.11	57.85	82.13	57.42	47.08		
		Baseline	1.000	58.71	54.69	80.05	47.89	52.21		
	Alpaca-GI 14	Mosaic-IT	1.073	58.84	54.18	80.54	47.92	52.70		
Llama2-7B	Alpaga	Baseline	1.000	55.25	54.35	78.65	47.02	40.98		
	Аграса	Mosaic-IT	1.096	55.32	53.75	78.65	46.88	41.98		
	Wirord 70k	Baseline	1.000	57.09	54.18	79.25	46.93	48.02		
	Wizaru-70K	Mosaic-IT	1.197	57.41	54.69	79.69	48.11	47.13		
		Baseline	1.000	61.47	58.70	83.12	54.13	49.92		
	Alpaca-Gr 14	Mosaic-IT	1.110	63.26	58.87	83.54	55.75	54.87		
Llama2-13B	Almana	Baseline	1.000	57.63	57.25	81.23	54.13	37.91		
	Аграса	Mosaic-IT	1.046	58.80	56.57	81.79	54.28	52.55		
	Wigond 70k	Baseline	1.000	61.24	57.04	83.39	55.76	48.78		
	vvizai d-70k	Mosaic-IT	1.078	61.50	58.70	83.69	56.44	47.18		

Table 1: The performance comparison on the Pair-wise Comparison Winning Score and the Open LLM Leaderboard, on 3 different base models and 3 different instruction tuning datasets.

Model	Dataset	Method	Pair-wise ↑ Score	Open LLM ↑ Average	Alpaca Ev Rate (LC)	7 al 2 ↑ Rate	MT-Be 1-round	ench↑ 2-round	IF E P(L)	val↑ I(L)
Llama-3-8B	Magpie	Baseline Mosaic-IT	1.000 1.133	56.15 60.13	9.22 12.23	13.74 16.05	8.10 8.36	7.08 7.49	35.67 40.67	47.72 52.76
Phi-3	Magpie	Baseline Mosaic-IT	1.000 1.014	62.90 63.54	13.82 14.04	17.68 17.67	7.78 7.89	6.42 6.16	44.36 50.83	55.52 62.35
Gemma2-2B	Magpie	Baseline Mosaic-IT	1.000 1.032	46.37 48.36	5.35 5.66	7.77 8.54	4.57 5.16	3.23 3.96	21.81 22.18	32.49 34.77
Mistral 7B	Alpaca-GPT4	Baseline Mosaic-IT	1.000 1.349	59.70 63.65	3.98 5.00	7.28 7.81	6.44 7.11	5.26 4.69	35.86 38.08	45.92 50.23
	Wizard-70k	Baseline Mosaic-IT	1.000 1.161	57.86 61.11	4.13 4.44	6.46 7.56	6.21 6.95	4.70 4.32	41.96 45.47	53.00 56.47

Table 2: The performance comparison across multiple model families and datasets on five evaluation metrics. Rate(LC) in Alpaca Eval represents length-controlled win rates. In IF Eval, P(L) and I(L) represent Prompt-level and Instruction-level accuracy in the Loose setting.

consistent improvement in this benchmark represents that our method not only improves the response qualities of the LLMs but also improves their controllability regarding formats.

Further **Human Evaluations** are conducted on Mistral-7B with Alpaca-GPT4 and WizardLM dataset. For the comparison on (1) Alpaca-GPT4: the model using Mosaic-IT wins on 68 out of 100 instruction, ties on 3, and losses on 29 instructions; on (2) WizardLM: the model using Mosaic-IT wins on 63 out of 100 instruction, ties on 6, and losses on 31 instructions. This human evaluation also further verifies the effectiveness of our Mosaic-IT.

To conclude, our Mosaic-IT shows consistent improvement in instruction-following and constrainfollowing ability and response quality with a reduction of approximately 80% of the training time cost. Given that our method is a cost-free data synthesis technique that does not rely on any additional human/LLM generation, the observed improvements are remarkable.

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

3.2 Ablation Studies

Detailed ablations are in Appendix C.

Ablation on Mosaic Strategies investigates how different mosaic strategies affect LLM performances. We find out that further implementing Advanced Strategies, (Format, Permute, Maskout), improves LLM performance as they largely diversify and complicate the instructions only implementing Primary Strategy. Ablation on the Max Number of Instructions investigates the trend between the number of instructions that are concatenated to-

gether and the LLMs' performances. We find out 439 that in the uniform distribution setting, more of 440 the instructions are concatenated together, as this 441 process makes the concatenated instruction com-442 plex. Ablation on the Distribution of Number 443 of Instructions investigates how different distribu-444 tions affect the LLMs' performance. It is revealed 445 that except for the maximum number of instruc-446 tions concatenated together, the distribution is also 447 important. Ablation on Semantic Grouping inves-448 tigates how semantic grouping, i.e., group instruc-449 tions with similar semantic meaning rather than 450 pure random grouping, affects performance. We 451 show that semantic grouping has its own benefits 452 while random selection is much more convenient. 453

4 Further Discussion

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

4.1 Improving Efficiency

One of the benefits of our method is the efficiency of the training process. Given an existing dataset, our mosaic processes largely decrease the number of total overall instructions and the total number of gradient descents, leading to a reduction in the training process. The detailed comparison is shown in Table 3, which is based on the Mistral-7B model on the Alpaca-GPT4 dataset. The time is calculated based on four NVIDIA A100 Graphic Cards. As shown, our method greatly decreases the training time to approximately 16% to 25% while achieving better performances, especially when there are mosaic samples with larger permutation counts. Please note that the time reduction is *ratio-based*, when larger datasets or models are trained, the absolute time reduction gap between baseline methods and our methods will be much more obvious¹.

4.2 Alleviating Memorizing

In the original instruction tuning process, each data sample will be trained several times for LLMs without changes to the instructions and responses. This training process poses risks to the potential memorizing effects on training samples, which can be partially indicated by the "*stair-like*" training loss curves as shown in Figure 3. In the figure, all the training settings are kept the same between the baseline models and Mosaic-IT models, including the Learning Rate, Warm-up Ratio, Learning Rate Schedule (Cosine), Batch Size, etc. For the baseline methods, the training loss hardly decreases within each epoch of training 486 but drops dramatically when the LLMs meet the 487 same training samples again, which indicates a 488 potential memorizing effect of training samples 489 and potential overfitting. However, when utilizing 490 our method, the random mosaics of original 491 instructions with diverse and complex meta-492 instructions largely diversify the overall training 493 instructions. Although each original data sample 494 will still be seen by LLMs several times during 495 training, the overall context varies dramatically 496 as each original sample is only an atomic element 497 of the overall mosaic sample, indicating that there 498 will be no identical overall instructions during the 499 whole training process. Thus this augmentation 500 largely alleviates the potential memorizing and 501 overfitting problems as shown in the figure, where 502 the training loss decreases smoothly, representing 503 the gradual learning process. 504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

4.3 Why It Works

Detailed analysis are in Appendix A.

4.3.1 Quantitive Analysis

Compared with the original instruction tuning, our Mosaic-IT largely increases the complexity and difficulty of the original instructions. InsTag (Lu et al., 2023) proposes a ChatGPT-based method (Number of InsTag), while Cherry (Li et al., 2024e) proposes a perplexity-based Instruction-Following Difficulty (IFD) score. To quantitatively evaluate the difficulty and complexity of instruction-tuning data, we compute these two metrics on the Alpaca and WizardLM70k datasets to verify the effectiveness of our method in improving the difficulty/complexity as shown in Table 5, the large increase of our method on Number of InsTag (2.62 to 10.93) and IFD scores (0.60 to 0.76) verifies our assumption that Mosaic-IT benefits LLMs by diversify and complicate the original instructions. This is a widely accepted manner for data synthesis (Zhao et al., 2024; Wu et al., 2024; Ding et al., 2023; Li et al., 2023a; Liu et al., 2023a; Li et al., 2024c; Guo et al., 2024), while our method accomplishes it in a cost-free manner.

4.3.2 Qualitative Analysis

Generally, following multiple instructions at the same time is difficult for most of the existing LLMs, especially when more constraints are required like specific formats or orders. As shown in Table 4, even the powerful GPT4 model does not do well

¹For example, fine-tuning Llama-3-8B on 1 million data for 2 epochs requires approximately 140 hours, while our method can reduce the time to approximately 20 hours.

Settings	Baseline	Fix	Exponential	Pareto	Log-normal	Logistic	Uni-2	Uni-4	Uni-6	Uni-8	Uni-10	Uni-12
Time (min)	827	121	129 15.6%	133	133	143	716	426	305 36.0%	245 20.6%	202	173
	100.0%	14.0%	13.0%	10.1%	1.421	17.5%	0.0%	31.3%	1 202	29.0%	24.4%	20.9%
Winning Score	1.000	0.982	0.995	1.417	1.431	1.417	0.989	1.142	1.303	1.294	1.349	1.376

Table 3: The training time comparison of different settings, and the pair-wise winning scores are also provided for better illustration. "Uni-2" represents uniform distribution with max count as 2. Mosaic-IT reduces the training time to 16% - 25% while achieving better performance.



Figure 3: The training loss curve comparisons between the original instruction tuning process and our Mosaic-IT with w datasets on (a) Mistral-7B, (b) Llama2-7B, and (c) Llama2-13B. The "stair-like" loss curves for the original training process indicate potential memorizing effects, while our loss curves are smoother. All the training settings are kept the same between the baseline models and Mosaic-IT models, including the Learning Rate, Warm-up Ratio, Learning Rate Schedule (Cosine), Batch Size, etc.

in this condition. Thus we regard the capability to follow multiple instructions at the same time as a higher level of instruction-following capabilities.

Previous methods train LLMs to produce a response for a single instruction or query. Instead, our method produces compositional data synthesis to train LLMs to generate multiple responses for multiple instructions in diverse forms (e.g., order, mask, format) specified by different metainstructions. It also enforces LLMs to partition the input context correctly and manage the interference and dependencies among multiple instructions. These are critical to developing and improving the compositional reasoning capabilities of LLMs, which have not been covered by mainstream instruction-tuning frameworks.

4.3.3 Qualitative Examples

535

536

540

541

542

544

545

546

548

549

550

551

552

553

554

555

556

558

559

562

To better understand the differences between the responses generated by the baseline model and the Mosaic-IT model, a pair of qualitative examples are presented in Figure 4 (baseline model) and Figure 5 (Mosaic-IT model). The example instruction is selected from the WizardLM test sets and the models are Mistral-7B trained on Alpaca-GPT4 with or without our method.

Some of the real-world instructions can be complex and hard to answer all at once but require LLMs to "decompose" the original overall instruction into pieces to respond. Due to the lack of this kind of difficult data in the training set, the capability of most LLMs is largely restricted as shown in Figure 4. The baseline LLM directly ignores the "sub-query" of explaining "*protocols and standards such as TCP/IP, HTTP, FTP, DNS, DHCP, and ARP*". On the contrary, as a cost-free data synthesis method, Mosaic-IT simulates this kind of instruction by using original easy instructions, thus equipping LLMs with the capability to respond to these difficult instructions without ignoring some parts of it, as shown in Figure 5.

563

564

565

566

567

568

569

570

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

589

5 Conclusion

We introduce Mosaic Instruction Tuning (Mosaic-IT), a novel, human/model-free method to enhance instruction tuning for LLMs. By concatenating multiple instruction-response samples and using higher-level meta-instructions, Mosaic-IT improves multi-step and format-following capabilities. Our evaluations show superior performance and an 80% reduction in training costs compared to the original methods. Mosaic-IT's simplicity and efficiency make it a scalable solution for improving LLMs without extensive human intervention or resource-intensive teacher models. Our results highlight the potential of innovative data synthesis techniques in advancing LLM capabilities.

592

597

598

612

613

614

615

617

618

622

625

627

632

634

635

637

639

641

642

643

645

Limitations

The potential limitations of our work: (1) Currently, three Advanced Mosaic Strategies with corresponding high-level rules are proposed and utilized in our method, however, we believe more strategies and predefined rules can be further introduced. (2) The optimal distribution of the number of instructions for the mosaic process still needs further justification in future studies. (3) It is unknown whether the inclusion of extra models or careful curation/selection of instructions for concatenation will further improve the performance of Mosaic-IT largely.

References

- Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin Jin, Nikos Karampatziakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. 2024. Phi-3 technical report: A highly capable language model locally on your phone. Preprint, arXiv:2404.14219.
 - Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. 2020. Yolov4: Optimal speed

and accuracy of object detection. *Preprint*, arXiv:2004.10934.

646 647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In Advances in Neural Information Processing Systems, volume 33, pages 1877–1901. Curran Associates, Inc.
- Alexander Bukharin and Tuo Zhao. 2023. Data diversity matters for robust instruction tuning. *Preprint*, arXiv:2311.14736.
- Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay Srinivasan, Tianyi Zhou, Heng Huang, and Hongxia Jin. 2023. Alpagasus: Training a better alpaca with fewer data. *Preprint*, arXiv:2307.08701.
- Zhoujun Cheng, Jungo Kasai, and Tao Yu. 2023. Batch prompting: Efficient inference with large language model APIs. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 792–810, Singapore. Association for Computational Linguistics.
- Cheng-Han Chiang and Hung-yi Lee. 2023. Can large language models be an alternative to human evaluations? In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 15607–15631, Toronto, Canada. Association for Computational Linguistics.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An opensource chatbot impressing gpt-4 with 90%* chatgpt quality.
- Hyung Won Chung, Le Hou, S. Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Wei Yu, Vincent Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed Huai hsin Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling instruction-finetuned language models. *ArXiv*, abs/2210.11416.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question

- 703 704 705 706 710 711 712 713 714 715 716 717 719 720 721 723 724 725 727 728 729
- 731 732 734 735 736 737 740 741 742 743 744 745 746 747 748 749 750 751 752 753 754 755 756

answering? try arc, the ai2 reasoning challenge. Preprint, arXiv:1803.05457.

- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. Flashattention: Fast and memory-efficient exact attention with io-awareness. Preprint, arXiv:2205.14135.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. Preprint, arXiv:2305.14314.
- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. Enhancing chat language models by scaling high-quality instructional conversations. arXiv preprint arXiv:2305.14233.
- Sivan Doveh, Assaf Arbelle, Sivan Harary, Roei Herzig, Donghyun Kim, Paola Cascante-Bonilla, Amit Alfassy, Rameswar Panda, Raja Giryes, Rogerio Feris, Shimon Ullman, and Leonid Karlinsky. 2023. Dense and aligned captions (DAC) promote compositional reasoning in VL models. In Thirty-seventh Conference on Neural Information Processing Systems.
- Qianlong Du, Chengqing Zong, and Jiajun Zhang. 2023. Mods: Model-oriented data selection for instruction tuning. Preprint, arXiv:2311.15653.
- Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. GLM: General language model pretraining with autoregressive blank infilling. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 320–335, Dublin, Ireland. Association for Computational Linguistics.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi,

Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, 761 Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph 762 Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, 763 Kalyan Vasuden Alwala, Kartikeya Upasani, Kate 764 Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, 765 Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz 768 Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, 769 Lukas Blecher, Lukas Landzaat, Luke de Oliveira, 770 Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, 771 Manohar Paluri, Marcin Kardas, Mathew Oldham, 772 Mathieu Rita, Maya Pavlova, Melanie Kambadur, 773 Mike Lewis, Min Si, Mitesh Kumar Singh, Mona 774 Hassan, Naman Goyal, Narjes Torabi, Nikolay Bash-775 lykov, Nikolay Bogoychev, Niladri Chatterji, Olivier 776 Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, 779 Punit Singh Koura, Puxin Xu, Qing He, Qingxiao 780 Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon 781 Calderer, Ricardo Silveira Cabral, Robert Stojnic, 782 Roberta Raileanu, Rohit Girdhar, Rohit Patel, Ro-783 main Sauvestre, Ronnie Polidoro, Roshan Sumbaly, 784 Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, 786 Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, 788 Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer 790 Whitman, Sten Sootla, Stephane Collot, Suchin Gu-791 rurangan, Sydney Borodinsky, Tamar Herman, Tara 792 Fowler, Tarek Sheasha, Thomas Georgiou, Thomas 793 Scialom, Tobias Speckbacher, Todor Mihaylov, Tong 794 Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor 795 Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent 796 Gonguet, Virginie Do, Vish Vogeti, Vladan Petro-797 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whit-798 ney Meers, Xavier Martinet, Xiaodong Wang, Xiao-799 qing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei 800 Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine 801 Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue 802 Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng 803 Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam 805 Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva 806 Goldstand, Ajay Menon, Ajay Sharma, Alex Boesen-807 berg, Alex Vaughan, Alexei Baevski, Allie Feinstein, 808 Amanda Kallet, Amit Sangani, Anam Yunus, An-809 drei Lupu, Andres Alvarado, Andrew Caples, An-810 drew Gu, Andrew Ho, Andrew Poulton, Andrew 811 Ryan, Ankit Ramchandani, Annie Franco, Apara-812 jita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, 813 Ashwin Bharambe, Assaf Eisenman, Azadeh Yaz-814 dan, Beau James, Ben Maurer, Benjamin Leonhardi, 815 Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi 816 Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Han-817 cock, Bram Wasti, Brandon Spence, Brani Stojkovic, 818 Brian Gamido, Britt Montalvo, Carl Parker, Carly 819 Burton, Catalina Mejia, Changhan Wang, Changkyu 820 Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, 821 Chris Cai, Chris Tindal, Christoph Feichtenhofer, Da-822 mon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, 823 Danny Wyatt, David Adkins, David Xu, Davide Tes-824

tuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria

825

826

836

857

861

867

870

871

872

873

874

875

876

877

878

879

887

Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. 2024. The Ilama 3 herd of models. *Preprint*, arXiv:2407.21783.

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

- Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Alpacafarm: A simulation framework for methods that learn from human feedback. *Preprint*, arXiv:2305.14387.
- Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, Jason Phang, Laria Reynolds, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2021. A framework for few-shot language model evaluation.
- Weidong Guo, Jiuding Yang, Kaitong Yang, Xiangyang Li, Zhuwei Rao, Yu Xu, and Di Niu. 2024. Instruction fusion: Advancing prompt evolution through hybridization. *Preprint*, arXiv:2312.15692.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. In *International Conference on Learning Representations*.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. *Preprint*, arXiv:2310.06825.
- Daniel Khashabi, Sewon Min, Tushar Khot, Ashish Sabharwal, Oyvind Tafjord, Peter Clark, and Hannaneh Hajishirzi. 2020. UNIFIEDQA: Crossing format boundaries with a single QA system. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1896–1907, Online. Association for Computational Linguistics.
- Diederik P. Kingma and Jimmy Ba. 2017. Adam: A method for stochastic optimization. *Preprint*, arXiv:1412.6980.
- Miyoung Ko, Jinhyuk Lee, Hyunjae Kim, Gangwoo Kim, and Jaewoo Kang. 2020. Look at the first sentence: Position bias in question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1109–1121, Online. Association for Computational Linguistics.

- 946 947
- 94
- .
- 950 951
- 951
- 95
- 955
- 956
- 957 958
- 960 961
- 9
- 964
- 9
- 967
- 969
- 970 971 972
- 973
- 974 975
- 976
- 977 978
- 980 981
- 982 983

98

987

- 9
- 991
- 9
- 993 994

9

996 997 998

999

1000

Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring how models mimic human

ence on Learning Representations.

Ming Li, Han Chen, Chenguang Wang, Dang Nguyen,

Ming Li, Jiuhai Chen, Lichang Chen, and Tianyi Zhou.

2024b. Can LLMs speak for diverse people? tuning

LLMs via debate to generate controllable controversial statements. In *Findings of the Association for*

Computational Linguistics ACL 2024, pages 16160-

16176, Bangkok, Thailand and virtual meeting. As-

Ming Li, Lichang Chen, Jiuhai Chen, Shwai He, Jiuxi-

ang Gu, and Tianyi Zhou. 2024c. Selective reflection-

tuning: Student-selected data recycling for LLM

instruction-tuning. In Findings of the Association for

Computational Linguistics ACL 2024, pages 16189-

16211, Bangkok, Thailand and virtual meeting. As-

Ming Li, Lichang Chen, Jiuhai Chen, Shwai He, and

Tianyi Zhou. 2023a. Reflection-tuning: Recycling

data for better instruction-tuning. In NeurIPS 2023

Workshop on Instruction Tuning and Instruction Fol-

Ming Li, Yong Zhang, Shwai He, Zhitao Li, Hongyu

Zhao, Jianzong Wang, Ning Cheng, and Tianyi Zhou. 2024d. Superfiltering: Weak-to-strong data filtering

for fast instruction-tuning. In Proceedings of the

62nd Annual Meeting of the Association for Compu-

tational Linguistics (Volume 1: Long Papers), pages

14255-14273, Bangkok, Thailand. Association for

Ming Li, Yong Zhang, Zhitao Li, Jiuhai Chen, Lichang

Chen, Ning Cheng, Jianzong Wang, Tianyi Zhou, and

Jing Xiao. 2024e. From quantity to quality: Boosting

LLM performance with self-guided data selection

for instruction tuning. In Proceedings of the 2024

Conference of the North American Chapter of the

Association for Computational Linguistics: Human

Language Technologies (Volume 1: Long Papers),

pages 7595–7628, Mexico City, Mexico. Association

Xian Li, Ping Yu, Chunting Zhou, Timo Schick, Luke

Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori,

Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and

Tatsunori B. Hashimoto. 2023c. Alpacaeval: An

automatic evaluator of instruction-following models.

https://github.com/tatsu-lab/alpaca_eval.

Jianzhe Lin, Maurice Diesendruck, Liang Du, and

Robin Abraham. 2024. Batchprompt: Accomplish

more with less. In The Twelfth International Confer-

translation. arXiv preprint arXiv:2308.06259.

Zettlemover, Omer Levy, Jason Weston, and Mike

Lewis. 2023b. Self-alignment with instruction back-

Computational Linguistics.

for Computational Linguistics.

sociation for Computational Linguistics.

sociation for Computational Linguistics.

lowing.

arXiv preprint arXiv:2406.15938.

Diangi Li, and Tianyi Zhou. 2024a. Ruler: Improv-

ing llm controllability by rule-based data recycling.

falsehoods. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics* (*Volume 1: Long Papers*), pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.

1002

1003

1005

1006

1007

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

- Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. 2023a. What makes good data for alignment? a comprehensive study of automatic data selection in instruction tuning. *arXiv preprint arXiv:2312.15685*.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023b. G-eval: Nlg evaluation using gpt-4 with better human alignment. *Preprint*, arXiv:2303.16634.
- Keming Lu, Hongyi Yuan, Zheng Yuan, Runji Lin, Junyang Lin, Chuanqi Tan, Chang Zhou, and Jingren Zhou. 2023. #instag: Instruction tagging for analyzing supervised fine-tuning of large language models. *Preprint*, arXiv:2308.07074.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2021. Cross-task generalization via natural language crowdsourcing instructions. *arXiv preprint arXiv:2104.08773*.
- OpenAI. 2023. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. *Preprint*, arXiv:2304.03277.
- Teven Le Scao, Angela Fan, Christopher Akiki, 1029 Elizabeth-Jane Pavlick, Suzana Ili'c, Daniel Hesslow, 1030 Roman Castagn'e, Alexandra Sasha Luccioni, Franc-1031 cois Yvon, Matthias Gallé, Jonathan Tow, Alexan-1032 der M. Rush, Stella Rose Biderman, Albert Web-1033 son, Pawan Sasanka Ammanamanchi, Thomas Wang, 1034 Benoît Sagot, Niklas Muennighoff, Albert Villanova 1035 del Moral, Olatunji Ruwase, Rachel Bawden, Stas 1036 Bekman, Angelina McMillan-Major, Iz Beltagy, Huu 1037 Nguyen, Lucile Saulnier, Samson Tan, Pedro Ortiz 1038 Suarez, Victor Sanh, Hugo Laurenccon, Yacine Jer-1039 nite, Julien Launay, Margaret Mitchell, Colin Raf-1040 fel, Aaron Gokaslan, Adi Simhi, Aitor Soroa Etx-1041 abe, Alham Fikri Aji, Amit Alfassy, Anna Rogers, 1042 Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, 1043 Chris C. Emezue, Christopher Klamm, Colin Leong, 1044 Daniel Alexander van Strien, David Ifeoluwa Ade-1045 lani, Dragomir R. Radev, Eduardo Gonz'alez Pon-1046 ferrada, Efrat Levkovizh, Ethan Kim, Eyal Bar 1047 Natan, Francesco De Toni, Gérard Dupont, Germán 1048 Kruszewski, Giada Pistilli, Hady ElSahar, Hamza 1049 Benyamina, Hieu Trung Tran, Ian Yu, Idris Abdul-1050 mumin, Isaac Johnson, Itziar Gonzalez-Dios, Javier 1051 de la Rosa, Jenny Chim, Jesse Dodge, Jian Zhu, 1052 Jonathan Chang, Jorg Frohberg, Josephine L. To-1053 bing, Joydeep Bhattacharjee, Khalid Almubarak, 1054 Kimbo Chen, Kyle Lo, Leandro von Werra, Leon 1055 Weber, Long Phan, Loubna Ben Allal, Ludovic Tan-1056 guy, Manan Dey, Manuel Romero Muñoz, Maraim 1057 Masoud, Mar'ia Grandury, Mario vSavsko, Max 1058 Huang, Maximin Coavoux, and Mayank Singh. 2022. 1059

Bloom: A 176b-parameter open-access multilingual language model. *ArXiv*, abs/2211.05100.

1060

1061

1062

1063

1064

1066

1067

1069

1074

1075

1076

1077

1078

1079

1080

1081

1082

1084

1085

1086

1087

1088

1089

1090

1091

1092

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114 1115

1116

1117

1118

1119

- Andrea Sottana, Bin Liang, Kai Zou, and Zheng Yuan. 2023. Evaluation metrics in the era of GPT-4: Reliably evaluating large language models on sequence to sequence tasks. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 8776–8788, Singapore. Association for Computational Linguistics.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https:// github.com/tatsu-lab/stanford_alpaca.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Mikuła, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. 2024. Gemma: Open models based on gemini research and technology. Preprint, arXiv:2403.08295.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. Llama: Open and efficient foundation language models. *Preprint*, arXiv:2302.13971.
 - Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay

Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti 1120 Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton 1121 Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, 1122 Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, 1123 Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-1124 thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan 1125 Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, 1126 Isabel Kloumann, Artem Korenev, Punit Singh Koura, 1127 Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-1128 ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-1129 tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-1130 bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-1131 stein, Rashi Rungta, Kalyan Saladi, Alan Schelten, 1132 Ruan Silva, Eric Michael Smith, Ranjan Subrama-1133 nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-1134 lor, Adina Williams, Jian Xiang Kuan, Puxin Xu, 1135 Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, 1136 Melanie Kambadur, Sharan Narang, Aurelien Ro-1137 driguez, Robert Stojnic, Sergey Edunov, and Thomas 1138 Scialom. 2023b. Llama 2: Open foundation and 1139 fine-tuned chat models. Preprint, arXiv:2307.09288. 1140

Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. 2023a. Large language models are not fair evaluators. *Preprint*, arXiv:2305.17926.

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023b. Self-instruct: Aligning language models with self-generated instructions. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 13484–13508, Toronto, Canada. Association for Computational Linguistics.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krima Doshi, Kuntal Kumar Pal, Maitreya Patel, Mehrad Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang Karia, Savan Doshi, Shailaja Keyur Sampat, Siddhartha Mishra, Sujan Reddy A, Sumanta Patro, Tanay Dixit, and Xudong Shen. 2022. Super-NaturalInstructions: Generalization via declarative instructions on 1600+ NLP tasks. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 5085-5109, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. 2022. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*.
- Minghao Wu, Abdul Waheed, Chiyu Zhang, Muham-
mad Abdul-Mageed, and Alham Fikri Aji. 2024.1176Lamini-Im: A diverse herd of distilled models from
large-scale instructions. *Preprint*, arXiv:2304.14402.1179

Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. Wizardlm: Empowering large language models to follow complex instructions. *Preprint*, arXiv:2304.12244.

1180

1181

1182

1183 1184

1185

1186

1187

1188

1189

1190

1191

1192

1193 1194

1195 1196

1197

1198

1199 1200

1201

1203

1204

1205

1206 1207

1208

1210

1211

1212 1213

1214

1215

1216

1217 1218

1219

1220

1221

1222

1224

1225

1226 1227

1228

1229

1230

1231 1232

1233

1234

1235

1236

- Xiaohan Xu, Ming Li, Chongyang Tao, Tao Shen, Reynold Cheng, Jinyang Li, Can Xu, Dacheng Tao, and Tianyi Zhou. 2024a. A survey on knowledge distillation of large language models. *ArXiv*, abs/2402.13116.
- Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yuntian Deng, Radha Poovendran, Yejin Choi, and Bill Yuchen Lin. 2024b. Magpie: Alignment data synthesis from scratch by prompting aligned llms with nothing. *Preprint*, arXiv:2406.08464.
- Qinyuan Ye, Bill Yuchen Lin, and Xiang Ren. 2021. CrossFit: A few-shot learning challenge for crosstask generalization in NLP. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 7163–7189, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. HellaSwag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, Florence, Italy. Association for Computational Linguistics.
- Zhiyuan Zeng, Jiatong Yu, Tianyu Gao, Yu Meng, Tanya Goyal, and Danqi Chen. 2024. Evaluating large language models at evaluating instruction following. In *The Twelfth International Conference on Learning Representations*.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. A survey of large language models. *Preprint*, arXiv:2303.18223.
- Yingxiu Zhao, Bowen Yu, Binyuan Hui, Haiyang Yu, Fei Huang, Yongbin Li, and Nevin L. Zhang. 2024. A preliminary study of the intrinsic relationship between complexity and alignment. *Preprint*, arXiv:2308.05696.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zhuohan Li, Zi Lin, Eric P. Xing, Joseph E. Gonzalez, Ion Stoica, and Hao Zhang. 2024a. Lmsyschat-1m: A large-scale real-world llm conversation dataset. *Preprint*, arXiv:2309.11998.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Preprint*, arXiv:2306.05685.

- Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan 1238 Ye, Zheyan Luo, Zhangchi Feng, and Yongqiang Ma. 1239 2024b. Llamafactory: Unified efficient fine-tuning 1240 of 100+ language models. In Proceedings of the 1241 62nd Annual Meeting of the Association for Compu-1242 tational Linguistics (Volume 3: System Demonstra-1243 tions), Bangkok, Thailand. Association for Computa-1244 tional Linguistics. 1245
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao1246Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu,1247Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis,1248Luke Zettlemoyer, and Omer Levy. 2023. Lima: Less1249is more for alignment. Preprint, arXiv:2305.11206.1250

1252 1253

1254

1255

1256

1257

1258

1259

1260

1263

1264

1265

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1288

1289

1291

1292

1293

1294

1295

1296

1298

A Why It Works

A.1 Preliminary Experiments: Performance Degradation

The motivation of our Mosaic-IT is also rooted in the observation that when handling multiple instructions simultaneously, a performance degradation will incurred for even strong LLMs like GPT-4-turbo. While LLMs generally perform well when responding to single instructions, their capability to follow multiple instructions at once tends to decline noticeably. BatchPrompt has shown the uncertainty when LLMs are requested to answer multiple formatted questions at one time. Moreover, in some cases, e.g for general open-domain instructions, LLMs might directly ignore some of the instructions, especially when the LLMs are required to respond to the instructions in a random pre-defined order, which is exactly simulating our *Permute* strategy.

To quantitatively analyze this phenomenon, experiments using GPT-3.5-turbo and GPT-4-turbo are conducted on the WizardLM test set. Specifically, we compare the models' performance when responding to multiple instructions concurrently versus responding to a single instruction at each time, by utilizing LLM-based Pair-Wise comparison, as shown in Table 4. All the win rates are lower than 1.0, demonstrating a clear and significant reduction in response quality when these models are required to respond to multiple instructions at one time. Moreover, the possibility of missing instructions (Miss Rate) increases further when they are required to respond to the instructions in a predefined random order rather than a sequential order. These results clearly demonstrate the difficulties of following several instructions at a time and why it can be regarded as a higher level of instructionfollowing capability.

A.2 Why It Works

Mosaic-IT trains LLMs to follow metainstructions for compositional reasoning.

Previous methods train LLMs to produce a response for a single instruction or query. Instead, our method produces a compositional data synthesis method to train LLMs to generate multiple responses for multiple instructions in diverse forms (e.g., order, mask, format) specified by different meta-instructions. It also enforces LLMs to partition the input context correctly and manage the interference and dependencies among multiple instructions. These are critical to developing and improving the compositional reasoning capabilities of LLMs, which have not been covered by mainstream instruction-tuning frameworks.

Mosa	aic-IT	creates	m	ore chall	lenging ar	nd com-		
plex	instru	ictions	to	further	improve	LLMs'		
instruction-following capabilities.								

Mosaic-IT's composition of multiple instructions and the diverse meta-instructions create more challenging and complex instruction-tuning data for LLMs. Moreover, since we do not rely on data synthesis using LLMs but solely apply some rules to existing data, the correctness and quality of the augmented data are guaranteed. As shown in Table 4, even powerful LLMs like GPT4 can not follow concatenated instructions. Besides, it has been widely accepted that such challenging and complex instructions improve LLMs' instructionfollowing capability (Zhao et al., 2024; Wu et al., 2024; Ding et al., 2023; Li et al., 2023a; Liu et al., 2023a; Li et al., 2024c,b; Guo et al., 2024). Mosaic-IT follows this intuition by making the instruction more challenging and complex in order to improve LLMs. Different from previous methods relying on humans or stronger teacher LLMs to create the challenging samples, Mosaic-IT does not require any humans/models to create the augmentations.

To quantitatively evaluate the difficulty and complexity of instruction-tuning data, InsTag (Lu et al., 2023) proposes a ChatGPT-based method (Number of InsTag), while Cherry (Li et al., 2024e) proposes a perplexity-based Instruction-Following Difficulty (IFD) score. We compute these two metrics on the Alpaca and WizardLM70k datasets to verify the effectiveness of our method in improving the difficulty/complexity:

Number of InsTag: The number of InsTag is used to measure the complexity of the instructions. A larger value of the Number of InsTag indicates the intentions of the instruction are complex and benefit the LLM instruction tuning process. For the experiments below, we prompt GPT40 with the exact prompt provided in the paper to generate the Instags. As shown in Table 5, Mosaic-IT largely increases the average number of InsTag, indicating a large increase in instruction intention complexity, further leading to better performance.

IFD score: IFD score is a perplexity-based metric used to evaluate the instruction-following diffi-

1300 1301

1299

1302 1303

1304 1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

Dain Wice (Multi ve Single)	3 Instr	uctions	5 Instr	ructions	7 Instructions		
ran-wise (white vs. Single)	Win Rate ↑	Miss Rate \downarrow	Win Rate ↑	Miss Rate ↓	Win Rate \uparrow	Miss Rate ↓	
GPT-3.5-turbo (Sequential)	0.357	0.014	0.336	0.055	0.303	0.064	
GPT-3.5-turbo (Random)	0.315	0.124	0.330	0.156	0.198	0.312	
GPT-4-turbo (Sequential)	0.176	0.000	0.137	0.000	0.140	0.000	
GPT-4-turbo (Random)	0.139	0.000	0.153	0.014	0.101	0.005	

Table 4: Pair-wise win rate of performances when responding to multiple instructions concurrently versus responding to a single instruction each time, and miss rate when responding to multiple instructions concurrently. "3 Instructions" represents the setting where 3 random instructions are concatenated together for inference. "Sequential" and "Random" represents the setting where the models are asked to respond to each instruction sequentially, or in a random pre-defined order.

culty of a given instruction-response pair. A higher IFD score indicates that it is hard for the current model to build a connection between the instruction and the corresponding response, so it can be used to select training data beneficial for LLM instruction tuning. For the experiments below, we utilized the IFD score computed on GPT2. As shown in Table 5, Mosaic-IT increases IFD scores, indicating an increase in the instruction-following difficulty, which leads to an improvement in performance.

1348

1349

1350

1352

1353

1354

1355

1356

1357

1358 1359

1360

1361

1363

1364

1365

1366

1368

1369

1372

1373

1374

1375

1376

Method	Iı	ns Tag	IFD		
	Alpaca	Wizard-70k	Alpaca	Wizard-70k	
Original Mosaic-IT	2.62 9.75	4.20 10.93	0.60 0.76	0.67 0.79	

Table 5: The comparison between the original dataset and Mosaic-IT enhanced dataset with respect to the **Number of InsTag** and the **IFD score**.

A.3 Qualitative Example

To better understand the differences between the responses generated by the baseline model and the Mosaic-IT model, a pair of qualitative examples are presented in Figure 4 (baseline model) and Figure 5 (Mosaic-IT model). The example instruction is selected from the WizardLM test sets and the models are Mistral-7B trained on Alpaca-GPT4 with or without our method.

Some of the real-world instructions can be complex and hard to answer all at once but require LLMs to "decompose" the original overall instruction into pieces to respond. Due to the lack of this kind of difficult data in the training set, the capability of most LLMs is largely restricted as shown in Figure 4. The baseline LLM directly ignores the "sub-query" of explaining "*protocols and standards such as TCP/IP, HTTP, FTP, DNS, DHCP, and ARP*". On the contrary, as a cost-free data synthesis method, Mosaic-IT simulates this kind of instruction by using original easy instructions, thus equipping LLMs with the capability to respond to these difficult instructions without ignoring some parts of it, as shown in Figure 5. 1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

The consistent better performances on Alpaca Eval 2 with length controls show that the performance gains of Mosaic-IT do not purely originate from longer responses. On the contrary, Mosaic-IT actually makes the lengths of training data more diverse compared with the original data, e.g. some of the samples are the concatenation of only 1 original sample and some of them are the concatenation of 10 original samples. This ensures diverse generation lengths according to the given instructions rather than dominated by the original length distribution. Identify some of the main components of a computer network such as hosts, routers, switches, hubs, links, and interfaces. Explain how they communicate and coordinate with each other using protocols and standards such as TCP/IP, HTTP, FTP, DNS, DHCP, and ARP. Describe how data is transmitted and received over a network using concepts such as packets, frames, headers, addresses, ports, sockets, and checksums.

A computer network is composed of several components that work together to facilitate communication and data transfer between devices. These components include:

1. Hosts: A host is a device that is connected to the network, such as a computer, smartphone, or printer. Each host has a unique IP address, which is used to identify it on the network.

2. Routers: A router is a device that connects multiple networks together and directs traffic between them. It uses protocols such as IP to route data packets to their destination.

3. Switches: A switch is a device that connects multiple devices together on a single network. It uses MAC addresses to direct data packets to their destination.

4. Hubs: A hub is a device that connects multiple devices together on a single network. It broadcasts data packets to all connected devices, which allows them to communicate with each other.

5. Links: A link is the physical connection between two devices on a network. It can be a cable, a wireless connection, or a combination of both.

6. Interfaces: An interface is the part of a device that connects it to the network. It can be a network card, a USB port, or a wireless card.

All of these components communicate and coordinate with each other using protocols and standards such as TCP/IP, HTTP, FTP, DNS, DHCP, and ARP. These protocols define how data is transmitted and received over a network, how devices are addressed, and how traffic is routed.

Data is transmitted over a network in packets or frames. Each packet or frame contains a header, which includes information such as the destination IP address, the source IP address, and the type of data being transmitted. The data is then transmitted over the network, and the receiving device uses the information in the header to identify the data and direct it to the appropriate application or program.

In summary, a computer network is composed of several components that work together to facilitate communication and data transfer. These components communicate and coordinate with each other using protocols and standards, and data is transmitted over the network in packets or frames, using headers to identify the destination, source, and type of data.

Figure 4: An example from the WizardLM test set with the response generated by the Baseline model.

A computer network is made up of several components that work together to facilitate communication and data transfer between devices. Some of the main components include:

1. Hosts: These are the devices that are connected to the network and use its resources, such as computers, printers, and servers.

2. Routers: Routers are devices that direct traffic between networks, and are used to connect multiple networks together.

3. Switches: Switches are used to direct traffic within a network, and are used to connect multiple devices together on a single network.

4. Hubs: Hubs are devices that connect multiple devices together on a single network, and are used to distribute data to multiple devices at once.

5. Links: Links are the physical connections between devices on a network, such as Ethernet cables or wireless connections.

6. Interfaces: Interfaces are the points of connection between the network and the devices on it, such as network cards or wireless adapters.

For these components to communicate and coordinate with each other, they use a variety of protocols and standards. Protocols are sets of rules that dictate how data is transmitted and received, while standards are agreed-upon methods or formats that allow for interoperability between different systems. Some common protocols and standards used in computer networking include:

- TCP/IP: This is a suite of protocols that is used to communicate between devices on a network. It includes protocols for addressing, routing, and transporting data.

- HTTP: This is a protocol used for accessing and transferring data over the web. It is used to communicate between web servers and clients.

- FTP: This is a protocol used for transferring files over a network. It is used to communicate between FTP servers and clients.

- DNS: This is a protocol used for translating domain names to IP addresses. It is used to coordinate communication between devices on a network.

- DHCP: This is a protocol used for assigning IP addresses to devices on a network. It is used to coordinate communication between devices on a network.

- ARP: This is a protocol used for finding the MAC address of a device on a network. It is used to coordinate communication between devices on a network.

When data is transmitted over a network, it is broken down into small packets, each of which contains information about where it needs to go and what it contains. These packets are sent over the network, hopping from one device to the next until they reach their destination. At each hop, the device examines the packet's header to see where it needs to go next, and then forwards it along the appropriate link.

Once the packet reaches its destination, the receiving device examines the header to see if it is the intended recipient. If it is, the device will extract the data from the packet and use it, otherwise, the packet will be discarded.

To ensure that the data in the packets is not corrupted during transmission, checksums are used. A checksum is a mathematical calculation performed on the data in the packet, which produces a value that can be used to verify that the data has not been altered during transmission. If the checksum of the data in a packet is incorrect when it arrives at its destination, it is discarded, and the sender is notified of the error.

In summary, data is transmitted over a network by being broken down into small packets, each of which contains information about where it needs to go and what it contains. These packets are sent hopping from one device to the next, using a variety of protocols and standards to coordinate their transmission and reception. Checksums are used to ensure that the data in the packets has not been corrupted during transmission, and any errors are detected and corrected.

Figure 5: An example from the WizardLM test set with the response generated by the Mosaic-IT model.

1395

1396

1397

1398

1399

1400

1401

1402

1403

1404

1406

1407

1408

1409

1410

1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

1421

1422

1423

1424

1425

1426

1427

1428

1429

1430

1431

1432

1433

1434

1435

1436

1437 1438

1439

1440

1441

1442

B Experimental Setup

B.1 Implementation Details

For the three base pre-trained models, Llama2-7B, Llama2-13B (Touvron et al., 2023b), and Mistral-7B (Jiang et al., 2023), we utilize the prompt and code base from Vicuna (Chiang et al., 2023) and flash attention (Dao et al., 2022). The overall training arguments are aligned with the common training configuration. The Adam optimizer (Kingma and Ba, 2017) is utilized with the batch size of 128 and with the max token length of 2048. When training the baseline models Llama2-7B and Llama2-13B, the maximum learning rate is set to 2×10^{-5} with the warmup rate as 0.03 for 3 epochs. When training the baseline models on Mistral-7B, the maximum learning rate is set to 1×10^{-5} with the warmup rate as 0.1 for 3 epochs. For the three models, Llama-3-8B (Dubey et al., 2024), Phi-3 (Abdin et al., 2024), and Gemma2-2B (Team et al., 2024), we utilize the code base from LLaMA-Factory (Zheng et al., 2024b). The max token length is set with 4096 following the modern settings and we train the model for 2 epochs. Other parameters are kept the same as the above.

When training with Mosaic-IT, we run the mosaic process n times for each experiment to simulate n epochs of training, n represents the number of epochs trained on baseline models, to ensure the alignment of overall data sample counts. Then these augmented data are mixed together and used for training 1 epoch while all other configurations are kept the same as baselines.

B.2 Training Dataset

The Alpaca dataset (Taori et al., 2023) comprises 52,000 instruction-following samples and is constructed utilizing the self-instruct paradigm (Wang et al., 2023b). This dataset was produced by employing OpenAI's text-davinci-003 model. Characterized as a classical dataset with moderate quality attributes, the Alpaca dataset serves as an initial platform to validate our methodology. To further substantiate our approach using a dataset of inherently high quality, we also applied our method to the Alpaca-GPT4 dataset (Peng et al., 2023), which features responses generated by GPT4. The WizardLM dataset (Xu et al., 2023) is also utilized in our method, which contains 70,000 samples created by the evolution algorithm proposed by them. With ChatGPT-3.5 utilized, the data quality on WizardLM is largely guaranteed. The Vicuna 1M

dataset (Zheng et al., 2024a) is a large-scale dataset1443containing one million real-world conversations1444with 25 state-of-the-art LLMs, due to the computa-1445tion budget, 300k instances are randomly sampled1446for our experiments. Magpie dataset (Xu et al.,14472024b) is a most recent SOTA synthetic dataset1448with 300k samples.1449

1450

1451

1452

1453

1454

1455

1456

1457

1458

1459

1460

1461

1462

1463

1464

1465

1466

1467

1468

1469

1470

1471

1472

1473

1474

1475

1476

1477

1478

1479

1480

1481

1482

1483

1484

1485

1486

1487

1488

1489

1490

1491

1492

1493

B.3 Evaluation Metrics

Pair-wise Comparison by using powerful LLMs like GPT-4 is recently widely accepted and becoming a common practice (Touvron et al., 2023b; Chiang et al., 2023; Dettmers et al., 2023; Liu et al., 2023b; Chiang and Lee, 2023). The evaluation of responses from LLMs, especially in open-domain contexts where definitive ground truth is hard to establish, continues to be an intricate and evolving research domain. Recent studies, however, have indicated a notable alignment between GPT-4's performance evaluations and human assessments (Zheng et al., 2023; Li et al., 2023c; Sottana et al., 2023), thereby establishing a credible foundation for this evaluative methodology. We adopted test instruction sets from WizardLM (Xu et al., 2023), comprising 218 diverse, human-curated instructions for pair-wise comparison. We directly follow the evaluation framework proposed by (Chen et al., 2023; Li et al., 2024e), which evaluates responses on a scale spanning from 1 to 10 across multiple dimensions. To further address positional bias, as discussed by (Ko et al., 2020; Wang et al., 2023a), the comparison is conducted in two distinct sequences, LLM1's response first and then LLM2's response first, ensuring a fair assessment of model performance. Evaluation outcomes are categorized into 'win-tie-loss' for each instruction. The detailed evaluation prompt is provided in Appendix H.

Open LLM Leaderboard, employing the evaluation framework from Eval Harness (Gao et al., 2021), offers a detailed and systematic approach to assessing the capabilities of generative language models through a set of diverse evaluation tasks. This methodology zeroes in on four pivotal benchmarks: ARC (Clark et al., 2018), HellaSwag (Zellers et al., 2019), MMLU (Hendrycks et al., 2021), and TruthfulQA (Lin et al., 2022). These benchmarks collectively provide a comprehensive evaluation of the models' reasoning abilities, their grasp of common-sense knowledge, and their accuracy in presenting factual information. Consequently, the leaderboard presents valuable insights. Alpaca-Eval Leaderboard, leveraging the AlpacaFarm evaluation dataset, presents a dependable and efficient automated evaluation tool for LLMs (Li et al., 2023c; Dubois et al., 2023). This tool benchmarks the responses generated by LLMs against those from Davinci003, focusing on the models' ability to adhere to generic user instructions.

1494

1495

1496

1497

1498

1499

1500

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

1511

1512

1513

1514

1516

1517

1518

1519

1520

1521

1522

1523

1524

MT-Bench (Multi-turn Benchmark) (Zheng et al., 2023) is a benchmark tool designed for automated evaluating LLMs in multi-turn dialogue settings. It focuses on analyzing conversation flow and the model's ability to follow instructions with 80 high-quality, multi-turn questions.

IFEval (Instruction-Following Eval) (Zeng et al., 2024) is a straightforward and easy-to-produce evaluation benchmark focusing on a set of "verifiable instructions". It contains 25 types of verifiable instructions and 541 prompts, with each prompt containing one or multiple verifiable instructions.

Human Evaluation is further implemented to substantiate the superiority of our approach based on the WizardLM test set. The test set contains 100 samples randomly sampled from the original WizardLM test set. Three human evaluators were tasked with comparing the outputs generated by the models under consideration, using the same criteria as in the previous pairwise evaluation. Each evaluator was presented with three response options: Win, Tie, and Loss. The final outcomes were determined by a majority vote.

1527

1528

1530

1531

1532

1533

1535

1536

1537

1538

1540

1541

1542

1543

1544

1545

1546

1547 1548

1550

1552 1553

1554

1555

1556

1558

1559

1560

1561

1562

1563

1565

1566

1567

1568

1569

1570

1571

1572

1573

C Detailed Ablation Studies

In this section, extensive ablation experiments are conducted on Mistral-7B using with the Alpaca-GPT4 dataset to verify our method. We utilize Pair-wise comparison for evaluation.

C.1 Ablation on Mosaic Strategies

Ablation on Mosaic Strategies is presented in Table 6a. "Primary" represents the Primary Mosaic Strategy. The winning score of this setting is greater than 1.0, indicating a better performance compared with the baseline method. This comparison directly verifies the effectiveness of the idea of introducing multiple instructions during training, which complicates the instructions at no cost and improves the instruction-following ability of LLMs. "Format" represents the Format Strategy. Although the winning score is only slightly greater than the naive version, this version makes it possible for LLMs to follow the customized user-defined formats, indicating great potential for the controllability of LLMs. Moreover, the format version can be easily used with other types of meta instructions, showing great extensibility. "Permute" represents the Permute Strategy that builds on the Format Strategy with a probability of 1/2, similar to "Maskout". "Permute/Maskout" represents our default setting, where the Permute or Maskout Strategies are utilized together with the Format Strategie with a probability of 1/3. All these 3 settings show higher performance than the format version, indicating the effectiveness of Advanced Mosaic Strategies which define more complicated meta instructions.

C.2 Ablation on the Max Number of Instructions

Ablation on the Max Number of Instructions is presented in Table 6b, including the pair-wise comparison values. As shown in the table, when the max number is set as 2, i.e. at most 2 instructions/responses are concatenated together, the performance is almost the same as the baseline, indicating the ineffectiveness. However, when the max number grows, the corresponding winning scores also grow consistently. This trend shows that the more instructions concatenated together, the better the instruction-following ability. We hypothesize that, with the growth of the number of instructions, the overall instruction becomes much harder to follow, especially for the permute

	Winning Score	Win	Tie	Lose				
Primary	1.261	110	55	53				
Format	1.284	109	62	47				
Permute	1.334	118	55	45				
Maskout	1.376	121	58	39				
Permute/Maskout	1.349	123	48	47				
(a) Ablation on Mosaic-IT strategies.								

	Winning Score	Win	Tie	Lose
Max Count = 2	0.989	70	75	73
Max Count = 4	1.142	92	65	61
Max Count = 6	1.303	111	62	45
Max Count = 8	1.294	112	58	48
Max Count = 10	1.349	123	48	47
Max Count = 12	1.376	124	52	42

(b) Ablation on the Max Number of Instructions.

Table 6: Ablation on (a) Mosaic-IT strategies and (b)Max Number of Instructions.

and maskout strategies, which benefits LLMs' instruction-following capability.

	Winning Score	Win	Tie	Lose	$Mix \le 5$
Fix	0.982	90	34	94	2.39%
Exponential	0.995	94	29	95	2.58%
Pareto	1.417	129	51	38	8.94%
Log-normal	1.431	136	40	42	6.83%
Logistic	1.417	123	49	46	15.84%
Uniform	1.349	123	48	47	51.45%

Table 7: Ablation on the **Distribution of Number of Instructions**. The distribution formula and data counts for different settings are shown in Appendix E. "Mix ≤ 5 " represents the percentage of samples with the number of instructions less or equal to 5.



Figure 6: Ablation on the **Distribution of Number of Instructions**, the visualization of distributions.

C.3 Ablation on the Distribution of Number of Instructions

Ablation on the Distribution of Number of1578Instructions is presented, including the pair-wise1579comparison values in Table 7 and detailed number1580distribution comparisons in Figure 6, which aims1581at identifying how this count distribution affects1582the performance of our method.The detailed

1576

distribution formula and data counts are provided 1584 in the Appendix E. "Fix" represents the setting 1585 where all the overall instructions are concatenated 1586 with a fixed number of instructions, which we set 1587 as 10 unless the overall instructions exceed the 1588 max length limit. "Exponential" represents the 1589 setting where the number of instructions is sampled 1590 following the exponential distribution. Under 1591 these two settings, less than 3% of the overall 1592 instructions are concatenated by less or equal to 1593 5 original instructions. The lack of few-instruction 1594 concatenated samples negatively affects the LLMs' 1595 ability to follow the single instruction, which is 1596 employed by most of the existing evaluation meth-1597 ods, leading to worse performances. "Pareto", 1598 "Log-normal", and "Logistic" represents the corresponding distribution that are utilized for 1600 sampling. Different from the above two settings, 1601 approximately 10% of the overall instructions 1602 are composed of fewer original instructions, thus 1603 ensuring the LLMs are trained with samples with 1604 sufficiently diverse lengths, resulting in optimal performances. "Uniform" is our default setting, 1606 representing using the uniform distribution where 1607 1608 different numbers are sampled evenly. In this situation, the LLMs are trained with samples with the most diverse lengths, thus avoiding the LLMs 1610 overfit to simple lengthy responses. 1611

Ablation on Semantic Grouping

To demonstrate that including other concatenation

methods could further strengthen our argument,

we further conduct experiments using a seman-

tic grouping approach based on the Alpaca-GPT4

Specifically, we utilize a sentence transformer

model *all-mpnet-base-v2*² to obtain the semantic

embedding for each sample in the dataset, and then

we applied the K-means algorithm to group these data samples into multiple clusters. To ensure

enough samples per cluster, we set K=52 as the

dataset contains 52k samples in total. Given the

clusters, each concatenated sample is composed of multiple samples randomly drawn from the same cluster. We keep using the same training hyperpa-

rameters as in our main experiments. As shown

in Table 8 we report the performance on two eval-

uation metrics: pair-wise comparison and Alpaca

dataset to fine-tune Mistral.

1612

C.4

Eval.

10

1613 1614

1615 1616

1618 1619 1620

1617

1623 1624 1625

1622

1627 1628 1629

1630 1631

The semantic concatenation can still outperform 1632 the non-mosaic baseline by a large margin, indicat-1633 ing the effectiveness and potential of our Mosaic-IT 1634 augmentations and tasks. The semantic concate-1635 nation method has a slightly lower performance 1636 than the pure-random concatenation method, on 1637 pair-wise comparison and Alpaca Eval 2 scores. 1638 However, it achieves a much higher Alpaca Eval 2 (LC) score. This result suggests that the response 1640 quality of the model trained with semantic concate-1641 nation is on par with pure-random but the response 1642 length is shorter and more condensed. We found 1643 that semantic grouping leads to clusters with highly 1644 different average lengths of samples: The largest 1645 average length is 316.7 tokens while the smallest 1646 is 31.4 tokens. This discrepancy makes the lengths 1647 of Mosaic-IT concatenated samples more diverse, 1648 resulting in a better trade-off between the quality 1649 and length of the responses. 1650

²https://huggingface.co/sentence-transformers/ all-monet-base-v2

Method	Alpaca Eval 2 (LC)	Alpaca Eval 2	Pair-wise Compare (with non-mosaic)	Pair-wise Compare (with pure-random)
Pure-random Concatenation	5.00	7.81	1.349	1.000
Concatenation with Semantic Groups	7.80	6.51	1.275	0.936

Table 8: Comparison with Semantic grouping

D Multi-Instruction Evaluation

To verify our trained models' capability to follow the multiple instructions and meta-instructions in one inference, we create a test set of compositional instructions from WizardLM test sets using Mosaic-IT. For simplicity, we name this new test setting as *Mosaic Task*, which evaluates LLMs' capability to follow multiple instructions with additional diverse constraints (meta-instructions). One example of *Mosaic Task* is shown as follows.

Example of Mosaic Task

System Prompt

You are a helpful and precise assistant for providing the answer.

User Prompt

Respond to each of the following instructions in reverse of the original order.

[Ins1] [Ins2] [Ins3]

1651

1652

1653

1654

1655

1656

1657

1658

1659

1660

1661 1662

1663

1664

1665

1666

1667

1668

1669

1670

1671

1672

1673

1674

1676

1677

1678

1680

1681

1682

1684

Figure 7: The prompt we used to request GPT4-Turbo to evaluate the responses.

We use the success rate(%) to evaluate the performance of models on the Mosaic task. A response is successful if it follows the meta-instruction and no instruction is ignored (unless the meta-instruction masks it). In the table below, we report the success rate (%) of LLMs following three meta-instruction strategies, i.e., Format, Permute, and Maskout, on compositional augmentations of different numbers of instructions (i.e., 3, 5, 7 instructions). We report the success rates of GPT40, two base models, and their Mosaic-IT finetuned versions, as shown in Table 9.

The results expose the weaknesses of existing LLMs on Mosaic-IT tasks and show that training on Mosaic-IT augmentations can significantly improve performance. Specifically, Existing LLMs, even GPT40, can not perfectly follow multiple instructions with diverse constraints, not to mention other open-source models like Llama3 finetuned on datasets such as Magpie. These results further demonstrate the difficulty and complexity of Mosaic-IT tasks for existing LLMs, indicating the novelty of our method. The compositional reasoning capability required by Mosaic-IT tasks cannot be covered by the capabilities of base LLMs and1685existing instruction-tuning datasets. For example,1686the success rates of Mistral + Alpaca-GPT4 (base-1687line) and Llama3 + Magepie (baseline) are similar,1688although Llama3 + Magepie has relatively better1689general instruction-following capabilities among1690them.1691

1693

1694

1695

1696

Our method can bridge the significant gap and enhance LLMs' capability to follow multiple instructions with diverse constraints. Moreover, our data augmentation is cost-free and does not take any effort from humans or models.

	3	Instructio	ons	5	Instructio	ons	7	Instructio	ons
Model	Format	Permute	Maskout	Format	Permute	Maskout	Format	Permute	Maskout
GPT4o	59.17	55.05	41.46	56.88	51.38	26.13	29.82	37.16	24.27
Mistral + Alpaca-GPT4 (baseline) Mistral + Alpaca-GPT4 (mosaic)	20.18 98.32	3.67 66.51	3.25 69.11	10.09 95.87	2.75 60.55	5.41 67.57	7.34 97.25	0.92 64.68	0.97 66.02
Llama3 + Magepie (baseline) Llama3 + Magepie (mosaic)	16.06 97.71	8.26 79.82	7.32 84.55	9.63 94.95	1.38 72.94	5.41 77.48	5.50 76.61	2.75 61.01	3.88 85.44

Table 9: Performance comparison across multiple instructions settings.

E Detailed Distribution for Ablation on Mixture Distribution

E.1 Distribution description

The detailed distribution descriptions and formulas are provided below.

Exponential Distribution³: The exponential distribution is a continuous probability distribution used to model the time or space between events in a Poisson process. The probability density function (PDF) of the exponential distribution is:

$$f(x;\lambda) = \lambda e^{-\lambda x}$$
 for $x \ge 0$,

where $\lambda = 1$ by default in our setting. We will resample with this distribution if the sampled value

 x_{sample} is greater than k_{max} .

Log-normal Distribution⁴: The log-normal distribution is a continuous probability distribution of a random variable whose logarithm is normally distributed. It is often used to model variables that are positively skewed, such as income, stock prices, and other financial data. The probability density function (PDF) for a log-normal distribution is given by:

$$f(x;\mu,\sigma) = \frac{1}{x\sigma\sqrt{2\pi}} \exp\left(-\frac{(\ln x - \mu)^2}{2\sigma^2}\right) \quad \text{for} \quad x > 0$$

where $\mu = 0$ and $\sigma = 0$ by default in our setting. We will resample with this distribution if the sampled value x_{sample} is greater than k_{max} .

Logistic Distribution⁵: The logistic distribution is a continuous probability distribution used in various fields, including logistic regression, modeling growth, and in some cases as an alternative to the normal distribution due to its heavier tails. The probability density function (PDF) for the logistic distribution is given by:

$$f(x;\mu,s) = \frac{e^{-(x-\mu)/s}}{s\left(1 + e^{-(x-\mu)/s}\right)^2}$$

where $\mu = 0$ and s = 2 by default in our setting. We will resample with this distribution if the sampled value x_{sample} is greater than k_{max} .

Pareto Distribution⁶: The Pareto II or Lomax distribution is a shifted Pareto distribution. It can be considered a simplified version of the Generalized Pareto distribution, with the scale set to one and the location set to zero. The probability density function (PDF) for the Pareto distribution is:

$$f(x; \alpha) = \frac{\alpha m^{\alpha}}{x^{\alpha+1}} \quad \text{for} \quad x \ge m,$$

³https://numpy.org/doc/stable/reference/ random/generated/numpy.random.exponential.html ⁴https://numpy.org/doc/stable/reference/ random/generated/numpy.random.lognormal.html ⁵https://numpy.org/doc/stable/reference/ where m = 1 and $\alpha = 1$ by default in our set-
ting. We will resample with this distribution if the
sampled value $x_{sample} - 1$ is greater than k_{max} .174117421743

After getting x_{sample} , a floor function will be uti-1744lized to get the corresponding integer and the final1745concatenation count $k = k_{max} - floor(x_{sample})$.1746

E.2 Distribution visualization

The detailed data counts for different distributions1748are provided in Figure 8.1749

random/generated/numpy.random.logistic.html
 ⁶https://numpy.org/doc/stable/reference/
random/generated/numpy.random.pareto.html



Figure 8: Bar plots of detailed data counts for different distributions in the Ablation on the Numbers of Instructions: (a) Fix Max Number, (b) Exponential Distribution, (c) Log-normal Distribution, (d) Logistic Distribution, (e) Pareto Distribution, (f) Uniform Distribution.

1752

1753

1754

1755

1756

1757

1758

1759

1760

1761

1763

1764

1765

1767

1768

1769

1770

1771

1772

1773

1774

1775

1776

1777

1779

1780

1781

1782

1783

1784

1785

1786

1787

1788

1789

1790

1791

1792

1793

1796

1797

1798

1799

1800

F Related Work

Earlier research in instruction tuning primarily centered on constructing expansive, high-quality datasets through intensive curation by human experts, a process both time-consuming and laborintensive (Khashabi et al., 2020; Ye et al., 2021; Wei et al., 2022; Wang et al., 2022; Du et al., 2022). Motivated by the success of Alpaca (Taori et al., 2023), recent studies have explored automated approaches for developing instructiontuning datasets.

Instruction Data Improvement: WizardLM (Xu et al., 2023) first proposes an Evol Algorithm to complicate the existing data and reach supreme performance. LaMini-LM (Wu et al., 2024) innovatively generates "Topic-Guided" instructions utilizing Wiki data. Tree-Instruct (Zhao et al., 2024) preliminarily explores the relationship between instruction complexity and Alignment and proposes adding nodes to complicate the instruction. UltraChat (Ding et al., 2023) establishes broad thematic scopes, systematically generating numerous instructions within each. Reflection-Tuning (Li et al., 2023a) sequentially refines both instructions and responses by focusing on specific evaluative criteria. DEITA (Liu et al., 2023a) utilizes Chat-GPT to diversify and then select the data. Selective Reflection-Tuning (Li et al., 2024c) proposes a teacher-student collaborative pipeline to improve and select the data. Instruction Fusion (Guo et al., 2024) proposes to utilize ChatGPT4 to merge two distinct instructions for further complexity enhancement. These advancements showcase a shift towards automating the generation and refinement of datasets, reducing reliance on human labor.

Instruction Data Selection: It is widely accepted that "quality is all you need" (Touvron et al., 2023b; Zhou et al., 2023) for instruction tuning. LIMA (Zhou et al., 2023) demonstrates that merely 1,000 human-carefully-curated, highquality training instances can substantially enhance the instruction-following performance. InsTag (Lu et al., 2023) employs the proprietary model, Chat-GPT, to tag instruction data and select data with complex tags. Alpagasus (Chen et al., 2023) utilizes proprietary LLMs chatGPT and Claude2 to directly assess the quality of instruction tuning data. Cherry LLM (Li et al., 2024e) proposes the Instruction-Following Difficulty (IFD) scores to assess the difficulty of the instructions, which is a self-guided method in which no extra LLMs are

utilized. Motivated by Humpback (Li et al., 2023b), 1801 Selective Reflection-Tuning (Li et al., 2024c) ex-1802 tends the IFD score to a reverse version, focus-1803 ing on the feasibility of responses. (Du et al., 1804 2023) and (Bukharin and Zhao, 2023) utilize re-1805 ward models as the base scores for measuring data 1806 quality. DEITA (Liu et al., 2023a) experiments on 1807 several different data selection metrics and builds a dataset with high quality. Superfiltering (Li et al., 1809 2024d) reveals the consistency between weak and 1810 strong language models in perceiving instruction 1811 difficulty, making the filtering process much more 1812 efficient. All these works are devoted to distin-1813 guishing and selecting good data samples from bad 1814 ones for instruction tuning. 1815

1816 G Predefined Rules

1817Examples of predefined formats can be found in1818Table 10 and detailed predefined rule descriptions1819can be found in Table 11.

Serial Digit	Parsing Bracket	Parsing Text	Assembled Examples
i	(text)	BEGIN, END	1. (BEGIN) <i>response</i> (END)
(<i>i</i>)	[text]	START, END	(1). [START] <i>response</i> [END]
[<i>i</i>]	$\langle text \rangle$	RESPONSE, END	[1]. (RESPONSE) <i>response</i> (END)
$\langle i \rangle$	$\ll text \gg$	RESPONSE, END OF RESPONSE	$\langle 1 \rangle$. «RESPONSE» <i>response</i> «END OF RESPONSE»
$\ll i \gg$	text	OPEN, CLOSE	≪1≫. OPEN <i>response</i> CLOSE
###i	[text]	OPEN RESPONSE, CLOSE	###1. [IOPEN RESPONSEI]response[ICLOSEI]
##i	$\langle text\rangle$	INITIATE, TERMINATE	##1. (IINITIATEI) <i>response</i> (ITERMINATEI)
##i##	#text#	START POINT, END POINT	##1##. #START POINT#response#END POINT#
<i>i</i>	*text*	RES_START, RES_END	<pre>11. *RES_START*response*RES_END*</pre>
<i>i</i>	@text@	RES, /RES	1 . @RES@response@/RES@

Table 10: Examples of predefined formats, including the Serial Digit formats and Response Parsing formats. "*i*" represents the real number serial number, "*text*" represents the replaceable parsing text, and "*response*" represents the real response of the concatenated overall instructions/responses. The response parsing formats are composed of the parsing bracket and text. In each mosaic process, random formats will be sampled simulating the real-world user-defined formats. The last column represents the assembled examples using the formats in the same row.

Strategy	Rule Name	Rule Description
Permute	FIX	Respond in the order of a provided list.
Permute	REVERSE	Respond in reverse of the original order.
Permute	ALPHA	Respond in the alphabetical order of the first letter of tasks.
Permute	REVERSE_ALPHA	Respond in the reverse alphabetical order of the first letter of tasks.
Permute	LENGTH_WORD	Respond according to the length (words) of tasks, respond to short ones first.
Permute	REVERSE_LENGTH_WORD	Respond according to the length (words) of tasks, respond to long ones first.
Permute	LENGTH_CHAR	Respond according to the length (characters) of tasks, respond to short ones first.
Permute	REVERSE_CHAR_WORD	Respond according to the length (characters) of tasks, respond to long ones first.
Permute	ODD_EVEN	First respond to the odd-numbered tasks, then the even-numbered ones.
Permute	EVEN_ODD	First respond to the even-numbered tasks, then the odd-numbered ones.
Maskout	FIX	Ignore the tasks provided in the list.
Maskout	WORD_LONG	Ignore the longest one/several task(s) according to the word count.
Maskout	WORD_SHORT	Ignore the shortest one/several task(s) according to the word count.
Maskout	ODD	Ignore the odd-numbered tasks.
Maskout	EVEN	Ignore the even-numbered tasks.

Table 11: Predefined rules for the Permute and Maskout strategy. A random rule will be sampled for each mosaic process, which largely complicates and diversifies the mosaicked instructions.

H Prompt for Evaluation

1820 1821 1822

The detailed pair-wise comparison prompt for the pair-wise comparison is in Figure 9.

Prompt for Performance Evaluation

System Prompt

You are a helpful and precise assistant for checking the quality of the answer.

User Prompt

[Question] *Question* [The Start of Assistant 2's Answer] *Answer 2* [The End of Assistant 2's Answer] *Answer 2* [The End of Assistant 2's Answer]

We would like to request your feedback on the performance of two AI assistants in response to the user question displayed above.

Please rate the helpfulness, relevance, accuracy, level of details of their responses. Each assistant receives an overall score on a scale of 1 to 10, where a higher score indicates better overall performance. Please first output a single line containing only two values indicating the scores for Assistant 1 and 2, respectively. The two scores are separated by a space. In the subsequent line, please provide a comprehensive explanation of your evaluation, avoiding any potential bias and ensuring that the order in which the responses were presented does not affect your judgment.

Figure 9: The prompt we used to request GPT4-Turbo to evaluate the responses.

I Detailed Performance Scores on Llama3, Phi3 and Gemma2

1825The detailed performance scores on the Open LLM1826Leaderboard and IFEval, for Llama-3-8B, Phi-3,1827and Gemma2-2B.

Model	Dataset	Method	Open LLM Leaderboard ^					IF Eval ↑			
			Average	ARC	HellaSwag	MMLU	TruthfulQA	Prompt (S)	Inst (S)	Prompt (L)	Inst (L)
Llama-3-8B	Vicuna	Baseline Mosaic-IT	52.51 55.62	44.54 47.78	70.66 73.77	49.68 56.11	45.18 44.83	19.04 29.76	30.70 43.17	21.26 31.42	33.45 45.56
	Magpie	Baseline Mosaic-IT	56.15 60.13	50.09 53.58	71.29 76.62	54.40 60.82	48.84 49.52	29.39 38.08	40.76 49.64	35.67 40.67	47.72 52.76
Phi-3	Vicuna	Baseline Mosaic-IT	62.06 62.30	58.96 58.45	76.48 77.66	64.89 65.24	47.89 47.87	28.47 30.13	40.29 39.57	30.50 32.35	43.17 41.85
	Magpie	Baseline Mosaic-IT	62.90 63.54	59.30 60.23	75.07 76.30	65.89 66.14	51.35 51.50	39.56 42.33	50.84 53.60	44.36 50.83	55.25 62.35
Gemma2-2B	Vicuna	Baseline Mosaic-IT	48.90 51.31	43.43 46.33	64.20 69.32	41.50 44.29	46.46 45.31	20.51 21.44	32.61 33.57	23.66 24.03	35.61 36.93
	Magpie	Baseline Mosaic-IT	46.37 48.36	39.59 39.33	60.71 64.10	35.46 39.87	49.75 50.16	19.78 19.78	29.74 31.65	21.81 22.18	32.49 34.77

Table 12: The performance comparison on more model families and datasets on all five automatic evaluation metrics. In IF Eval, P and I represent Prompt-level and Instruction-level accuracy.