Upcycling Instruction Tuning from Dense to Mixture-of-Experts via Parameter Merging

Anonymous ACL submission

Abstract

Mixture-of-Experts (MoE) shines brightly in large language models (LLMs) and demonstrates outstanding performance in plentiful natural language processing tasks. However, existing methods transforming LLMs from dense to MoE face significant data requirements and typically rely on large-scale post-training. In this paper, we propose Upcycling Instruction Tuning (UpIT), a data-efficient approach for tuning a dense pre-trained model into a MoE instruction model. Specifically, we first point out that intermediate checkpoints during instruction tuning of the dense model are naturally suitable for specialized experts, and then propose an expert expansion stage to flexibly achieve models with flexible numbers of experts, where genetic algorithm and parameter merging are introduced to ensure sufficient diversity of new extended experts. To ensure that each specialized expert in the MoE model works as expected, we select a small amount of seed data that each expert excels to preoptimize the router. Extensive experiments with various data scales and upcycling settings demonstrate the outstanding performance and data efficiency of UpIT, as well as stable improvement in expert or data scaling. Further analysis reveals the importance of ensuring expert diversity in upcycling.

1 Introduction

004

011

012

014

018

023

035

040

042

043

Large Language Models (LLMs) have demonstrated remarkable performance on various NLP tasks and are gradually becoming part of our daily lives through chatbot applications such as Chat-GPT, Copilot, etc (Ouyang et al., 2022; Touvron et al., 2023; OpenAI, 2024). As LLMs become increasingly prevalent, the high computational of traditional dense architecture with high computational costs in the inference phase poses significant obstacles to downstream deployment. How to improve the model performance without proportionally increasing computing resources become a hot topic in the field (Muennighoff et al., 2024; Xue et al., 2024). In response to this challenge, Mixture-of-Experts (MoE) receives extensive attention due to its excellent scalability, which expands model capacity with almost no extra inference overhead (Fedus et al., 2022; Zoph et al., 2022). Recently, many MoE-based LLMs have emerged in various scenarios with outstanding effectiveness and efficiency (Dai et al., 2024; Jiang et al., 2024; Zhu et al., 2024a). 044

045

046

047

051

055

058

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

081

Upcycling is garnering increasing attention for converting dense LLMs through a series of processes, including expanding experts, integrating routers, and subsequent post-training, ultimately yielding MoE-style models. As depicted in Figure 1, current solutions are broadly classified into two categories: (a) Vanilla Upcycling, which directly upcycle a dense model to a MoE model by replicating FFN layers, followed by a large-scale post-training to optimize the additional experts and corresponding routers (Komatsuzaki et al., 2023). Due to the homogeneity of experts in the initial stage, a large amount of post-training data is usually necessary, such as ~1T tokens for full parameter training or ~5M instruction data for parameter efficient fine-tuning (Dou et al., 2024; Zhu et al., 2024a). (b) Specialized Upcycling, which first trains specialized experts based on meticulously designed domain-specific datasets and then proceeds with upcycling and post-training (Sukhbaatar et al., 2024). Despite having superior performance, it still requires hundreds of billions of elaborately constructed domain data and lacks flexibility in scaling the number of domain-specific experts. To sum up, although expert specialization slightly reduces data requirements, the current approach to upcycling from dense to MoE heavily relies on a large amount of training data, how to efficiently and flexibly train a MoE instruction model based on a dense pre-trained model is still an open problem.



Figure 1: Workflow of vanilla upcycling, specialized upcycling, and the proposed upcycling instruction tuning (UpIT) solutions. UpIT achieves specialized experts with various checkpoints, increases the expert number during the expert expansion stage, and maintains discrepancy among experts through router initialization, thereby achieving efficient and flexible upcycling.



Figure 2: The performance of various checkpoints saved during an instruction tuning process, with a red star indicating the best performance on each benchmark. Checkpoints saved at different epochs excel in different benchmarks, demonstrating the potential as specialized experts.

lot experiment on dense models to observe the characteristics of models at different epochs during standard instruction tuning. Figure 2 shows checkpoints saved at different epochs exhibit interleaved performance across benchmarks in various domains. In practical terms, we categorize nine benchmarks into four domains: factual knowledge, reasoning, coding, and world knowledge, where the performance on each benchmark generally shows an upward trend followed by a downward trend, but the position of the maximum value varies. For example, the model trained at epoch 2 demonstrates superior performance on HellaSwag and Natual Question, whereas the model trained at epoch 0.25 performs best on MMLU. In other words, models with different training steps demonstrate varying expertise in handling distinct domains. This phenomenon inspires us to consider that different checkpoints during instruction tuning are inherently suitable for constructing specialized experts.

087

091

093

097

100

101

102

103

104

In light of the above findings, we propose Upcycling Instruction Tuning (UpIT), which starts from a dense pre-trained model and trains a MoE instruction model with a flexible number of experts, following the basic thought of specialized upcycling. Figure 1 illustrates the four stages of UpIT. Specifically, (1) Expert Preparation. Considering the differences among checkpoints in the pilot experiment, it is sufficient to fine-tune the dense pre-trained model and save checkpoints at fixed intervals to prepare for specialized experts, without undertaking meticulous checkpoint selection. (2) Expert Expansion. Given the fixed checkpoints, we extend a flexible number of new experts based on genetic algorithms. In each iteration, we select two experts with the greatest differences, merge their parameters and obtain a new expert. We also perform parameter scaling and dropping before merging to simulate the mutation and further promote the discrepancy of experts. (3) Router

105

107

108

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

Initialization. Traditional routers are randomly ini-125 tialized and insensitive to expert capabilities. Here, 126 we assign each expert their skilled data and intro-127 duce an auxiliary binary classification loss to pre-128 optimize the corresponding routing vector, ensuring that all experts are capable of fully exhibiting 130 their strengths in the MoE model. (4) Model Upcy-131 cling. Before post-training, we merge the param-132 eters of multiple dense models into a MoE model. 133 Unlike existing methods, the pre-optimized routing 134 vectors are merged into a routing matrix, serving 136 as the final router.

138

139

140

141

142

143

144

145

147

148

149

151

152

153

154

155

156

157

159

160

161

163

165

166

168

169

172

173

174

175

176

Overall, by leveraging the differences in existing dense checkpoints and introducing the expert expansion stage, UpIT comprehensively reduces the cost of acquiring specialized experts and improves the flexibility of expert numbers. Router initialization further maintains expert diversity, thereby encouraging more effective utilization of data characteristics during the post-training of the MoE model. From an implementation perspective, UpIT divides standard instruction tuning into two parts, where the first part is responsible for expert preparation, and the second is post-training after upcycling. Between these two stages, we find that only 1% of the training data (approximately 500 to 5,000 samples) is enough to pre-optimize the routing vectors, which means that UpIT efficiently upcycles from dense to MoE without significantly increasing the overall data requirements.

We conduct extensive experiments under LoRA and FFN-based upcycling settings, with LoRA and FFN as experts. For a fair comparison, we train all models on IDEA dataset (Wu et al., 2024a), considering data sizes ranging from 50K to 500K, and evaluate the performance of nine benchmarks. Experimental results show that UpIT is consistently better under both settings than dense instruction tuning and other upcycling baselines. Especially in situations with small amounts of training data, existing upcycling methods often can not work well, while UpIT utilizes the discrepancy of experts and achieves better results than dense baselines. Moreover, UpIT exhibits excellent scalability, with expected performance improvements when increasing training data, total expert number, or activated expert number. The router visualization and ablation study also verify the overall promoting effect of expert diversity on upcycled MoE models. In summary, the main contributions are as follows:

(1) We propose UpIT, an efficient specialized upcycling method via parameter merging, which

Algorithm 1: Workflow of UpIT

- **Input:** Dense pre-trained model Θ_d , training dataset \mathcal{D} , target number of experts n.
- 1 // Expert Preparation
- ² Fine-tune the dense model Θ_d on \mathcal{D} and obtain a series of checkpoints $\mathcal{C} = \{\Theta_d^1, \dots, \Theta_d^m\}$.
- 3 Initialize expert models $\mathcal{E} = {\mathbf{E}_1, \dots, \mathbf{E}_m}$ from checkpoints \mathcal{C} .
- 4 // Expert Expansion
- 5 Merge new expert model with Algorithm 2 and put them into the set $\mathcal{E} = {\mathbf{E}_1, \dots, \mathbf{E}_m, \dots, \mathbf{E}_n}$.
- 6 // Router Initialization
- 7 Initialize routing vectors $\mathcal{R} = {\mathbf{r}_1, \dots, \mathbf{r}_n}$ for expert models \mathcal{E} .
- 8 Construct expert-specific data $\mathcal{D}_s = \{\mathcal{D}_s^1, \dots, \mathcal{D}_s^n\}$ with Algorithm 3.
- 9 for i = 1 to n do
- ¹⁰ Fine-tune expert layer \mathbf{e}_i in \mathbf{E}_i and corresponding routing vector \mathbf{r}_i on \mathcal{D}_s^i with Equation 1.
- 11 // Model Upcycling
- 12 Initialize router ${f R}$ by concatenating all routing vectors ${\cal R}.$
- 13 Initialize the MoE model Θ_m with expert models \mathcal{E} and router **R**.
- 14 Fine-tune the MoE model Θ_m on \mathcal{D} .
- **Output:** MoE Instruction model Θ_m .

can achieve an instruction model with a flexible number of experts. To best of our knowledge, it is the first attempt to utilize intermediate dense checkpoints for model upcycling. 177

178

179

180

181

182

183

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

201

203

(2) We emphasize the importance of expert discrepancy in upcycling and incorporate the idea into the entire design of UpIT. The innovative router initialization stage ensures that all specialized experts play to their strengths in the final MoE model.

(3) Extensive experiments under LoRA and FFNbased settings show that UpIT significantly outperforms existing methods in scenarios with sufficient and insufficient data, exhibiting outstanding flexibility, scalability, and performance upper bound.

2 Methodology

In this section, we provide a detailed exposition of UpIT. Generally speaking, UpIT extends the concept of specialized upcycling. Instead of carefully curating domain-specific datasets, it further utilizes intermediate checkpoints to reduce data requirements. It also expands certain aspects related to the experts to adapt to the flexible number of experts and pre-optimizes routing vectors. This is done to ensure that each expert in the instruction MoE model can leverage its unique strengths. The basic concepts of MoE and Upcycling are described in Appendix A.2. We also conduct an in-depth discus-

- 204
- 205

208

209

210

211

212

213

214

216

217

219

221

228

229

233

241

242

243

245

246

247

248

251

254

sion on the effectiveness of UpIT in Appendix A.1.

2.1 Workflow of UpIT

Starting from the dense pre-trained model, UpIT achieves a MoE instruction model. Algorithm 1 provides the working sketch. In this section, we provide a detailed explanation of each process.

Expert Preparation. In the instruction tuning of LLMs, as shown in Figure 2, the performance of intermediate checkpoints varies across different benchmarks, and different checkpoints exhibit unique strengths in different domains, highlighting the potential to serve as specialized experts. Compared to the labour-intensive method of training diverse expert models with massive domain-specific data (Sukhbaatar et al., 2024), we believe that the natural variations among checkpoints provide a more efficient pathway to developing specialized expert models. By training dense models to generate multiple checkpoints and saving them at regular intervals during training, we can easily obtain a series of expert models proficient in different domains, resulting in a more cost-effective method for preparing specialized expert models.

Expert Expansion. Given that the fixed number of checkpoints only sometimes corresponds with the flexible requirements of expert number, acquiring additional checkpoints entails redundant training if the number of experts exceeds the saved checkpoints. Here, we propose to address these challenges by generating distinct experts from existing ones without extensive retraining (see also Algorithm 2). Specifically, we draw inspiration from genetic algorithms, where two experts with the greatest differences are selected as parents in each iteration. We simulate the mutation process by randomly assigning weights to the parents and apply DARE (Yu et al., 2024) to introduce mutations into the newly created expert further, enhancing its discrepancy and adaptability. Such an expansion process not only eliminates the need for additional retraining but also facilitates the flexible expansion of the number of experts, ultimately improving the scalability of UpIT.

Router Initialization. Since routers remain randomly initialized after upcycling, which leads to the misallocation of tokens in the early post-training stages, in UpIT, such misallocation will weaken the expert differences in the previous stage and impact the learning efficiency of MoE models. To solve this problem, we propose a data selection approach to curate expert-specific data tailored



Figure 3: Performance comparison of UpIT and vanilla upcycling methods under different size of training data. Detailed results in Section A.8

255

256

257

258

260

261

262

263

264

266

267

268

269

270

271

272

273

274

275

276

277

278

279

281

to each expert model and pre-optimize additional routing vectors to ensure the discrepancy among experts (see also Algorithm 3). Specifically, we initially embed one-dimensional routing vectors \mathcal{R} before the MoE layer in each transformer block and participate in the training process as expert-specific routers. Next, we introduce an auxiliary loss \mathcal{L}_{aux} intending to maximize the output probability of corresponding routing vectors. This ensures that the likelihood of tokens being assigned to appropriate experts increases when they pass through the router. The pre-optimizing objective of *i*-th expert model is formulated as follows,

$$\mathcal{O}_{i} = \min_{\mathbf{E}_{i}} (\alpha \mathcal{L}_{lm}(\mathbf{E}_{i}) + (1 - \alpha) \mathcal{L}_{aux}(\mathbf{E}_{i}))$$
(1)

where α is the balance coefficient, which we set to 0.5 in our experiments, and $\mathcal{L}_{lm}(\cdot)$ is the causal language model loss. The auxiliary loss $\mathcal{L}_{aux}(\cdot)$ is defined as follows,

$$\mathcal{L}_{aux}(\mathbf{E}_i) = \text{CrossEntropy}(\text{Sigmoid}(\mathbf{h}_{\mathbf{r}_i}), \mathbf{I})$$
 (2)

where $\mathbf{h}_{\mathbf{r}_i}$ is the output of routing vector and \mathbf{I} is the unit matrix. We use Sigmoid function to scale the output to (0, 1) and minimize its difference from \mathbf{I} , which is equivalent to maximizing the output probability of the routing vector on the data that current expert model excels at.

Model Upcycling. Finally, we upcycle the dense model Θ_d to MoE model Θ_m by merging all the

	HumanE.	GSM8K	HellaS.	BBH	MMLU	NQ	Tri.QA	ARC-c	ARC-e	Avg.
LoRA-based Models										
Llama 2 7B	14.63	13.95	26.58	34.73	39.84	10.06	62.06	37.29	50.26	32.16
LoRA	22.56	45.72	65.36	37.14	49.33	14.99	50.15	61.36	81.13	47.53
Self-MoE(8E, A2)	28.05	46.70	64.27	38.67	49.63	21.11	48.67	64.41	82.19	49.30
PESC(8E,A2)	28.05	46.55	63.14	37.59	46.12	16.68	49.58	61.36	74.60	47.07
$LoRAMoE_{PT}(8E,A2)$	<u>34.15</u>	47.61	60.89	37.40	46.61	17.62	46.33	60.68	74.60	47.32
$LoRAMoE_{SFT}(8E,A2)$	28.66	49.81	67.62	<u>38.88</u>	50.54	20.55	<u>50.16</u>	62.37	81.31	<u>49.99</u>
UpIT(8E,A2)	35.37	<u>49.51</u>	<u>66.00</u>	40.27	<u>50.31</u>	24.52	55.27	65.08	83.60	52.21
Self-MoE(16E, A2)	30.20	47.61	65.36	37.14	49.33	<u>24.52</u>	51.11	62.71	82.19	50.02
PESC(16E,A2)	31.10	47.62	63.14	37.59	49.08	20.83	49.58	63.05	77.62	48.85
$LoRAMoE_{PT}(16E,A2)$	40.24	46.55	65.89	36.39	48.53	19.36	46.19	61.69	76.01	48.98
$LoRAMoE_{SFT}(16E,A2)$	30.12	49.62	66.77	40.21	50.96	20.83	52.63	<u>63.41</u>	80.67	<u>50.58</u>
UpIT(16E,A2)	40.62	<u>48.37</u>	<u>66.62</u>	<u>39.43</u>	<u>50.70</u>	25.62	56.61	67.46	84.66	53.34
FFN-based Models										
Sheared Llama 2.7B	5.49	1.74	25.09	26.62	26.98	6.43	38.89	22.37	24.69	19.81
SFT	26.22	29.19	38.01	26.46	33.93	8.42	18.61	42.37	58.55	31.31
Self-MoE(4E, A2)	6.71	8.87	32.11	27.65	28.81	18.45	42.27	33.22	47.44	27.28
Upcycle _{PT} (4E, A2)	31.71	35.10	43.40	30.23	37.93	13.74	34.72	45.08	58.73	36.74
Upcycle _{SFT} (4E,A2)	23.17	33.97	50.27	29.50	<u>38.90</u>	<u>15.18</u>	34.20	48.14	<u>65.08</u>	<u>37.60</u>
UpIT(4E,A2)	<u>31.34</u>	33.81	48.97	<u>29.53</u>	40.84	14.71	<u>36.99</u>	47.80	65.96	38.88
Self-MoE(8E, A2)	10.62	22.73	34.69	28.95	30.10	15.68	37.68	40.00	50.37	30.09
$Upcycle_{PT}(8E,A2)$	<u>26.22</u>	<u>34.04</u>	51.57	28.95	<u>39.84</u>	13.57	33.86	53.22	66.49	<u>38.64</u>
$Upcycle_{SFT}(8E,A2)$	22.56	33.66	46.26	<u>29.25</u>	39.19	14.76	35.18	49.15	67.72	37.53
UpIT(8E,A2)	32.19	35.64	<u>49.15</u>	30.23	40.38	14.57	37.93	49.10	68.43	39.74

Table 1: Performance comparison under Lora-based and FFN-based upcycling settings, where (xE,Ay) indicates that y out of x experts are activated, Lora-based UpIT(16E,A2) and FFN-based UpIT(8E,A2) are expanded from Lora-based UpIT(8E,A2) and FFN-based UpIT(4E,A2), respectively. Bold text and underlined text denote the best and second-best results in each group.

expert models \mathcal{E} and routing vectors \mathcal{R} . Specifically, for the initialization of experts, we utilize pre-optimized expert models from \mathcal{E} . In terms of router initialization, we concatenate all routing vectors from \mathcal{R} to form a complete router $\mathbf{R} \in \mathbb{R}^{d_h,n}$, where d_h is the dimension of hidden states, This way, the obtained MoE block could allocate different tokens to experts skilled in processing them. Finally, we continue to utilize \mathcal{D} for post-training to achieve the final MoE model.

2.2 Training Details

286

287

290

291

295

300

303

To comprehensively evaluate the effectiveness of UpIT, we utilize two types of upcycling settings:

(1) **FFN-based Upcycling:** Initially, we fully fine-tune all parameters of the dense pre-trained model to accumulate several expert models. In the expert expansion stage, we apply the genetic algorithm to construct expert modules (i.e. FFN layers), average the parameters of backbone modules (i.e. all layers except FFN) in candidate expert models, and result in new diverse expert models. We select expert-specific data to pre-optimize the FFN layers and routing vectors during the router initialisation stage. Finally, in the model upcycling stage, we average the backbone parameters of all expert models and concatenate the routing vectors, integrating FFN layers to produce the final MoE models. 304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

(2) **LoRA-based Upcycling:** The key difference from FFN-based upcycling is that in parameterefficient fine-tuning, parameters of backbone modules remain unchanged. Instead, we augment FFN layers with LoRA matrices, and operate on the values of LoRA matrices during expert expansion.

Following previous work (Fedus et al., 2022), during post-training, we also use load balancing loss, $\mathcal{L}_{\text{load}} = n \cdot \sum_{i=1}^{n} f_i \cdot P_i$, where *n* is the expert number, f_i is the fraction of tokens dispatched to expert E_i , P_i is the average fraction of the router probability allocated for expert E_i .

3 Experiments

3.1 Experimental Setup

Baselines. To assess the effectiveness of UpIT, we compare its performance against several baselines. For LoRA-based settings, we consider the



Figure 4: Performance comparison of UpIT and vanilla upcycling methods under different total and activated experts. Detailed results in Section A.9.

following baselines. (1) LoRA (Hu et al., 2021), (2) Self-MoE (Kang et al., 2024), (3) PESC (Wu et al., 2024a), (4) LoRAMoE_{PT} (Dou et al., 2024), and (5) LoRAMoE_{SFT}. For FFN-based settings, we compare UpIT with (1) SFT, (2) Self-MoE (Kang et al., 2024), (3) Upcycle_{PT} (Komatsuzaki et al., 2023), and (4) Upcycle_{SFT}. More detailed descriptions of the baselines can be found in Appendix A.5.

Dataset. Following (Wu et al., 2024a), we simultaneously train UpIT and compared baselines on a diverse set of datasets, encompassing Magicoder (Wei et al., 2023) for coding, Meta-MathQA (Yu et al., 2023) for mathematical and SlimORCA (Lian et al., 2023) for general abilities from various subjects. We do not perform quality filtering or other operations on the three datasets. We randomly sample data in a 1:1:3 ratio to create the final training dataset with 500K samples.

More training details and additional experimental results can be found in Appendix A.4 and A.7.

3.2 Main Results

328

329

331

332

334

341

351

354

Table 1 presents an analysis of benchmark results for LoRA and FFN-based models across various domains, revealing the following key insights.

(1) The UpIT demonstrates remarkable performance across various benchmarks, highlighting its effectiveness compared to existing upcycle solutions. Specifically, when compared with LoRAMoE_{SFT}(8E,A2), the LoRA- based UpIT(8E,A2) achieves an average performance improvement of 2.22%. When the number of experts is expanded to 16, UpIT(16E,A2) sustains a competitive edge over LoRAMoE_{SFT}(16E,A2), exhibiting a lead of 2.76%. Similar trends are observed in FFN-based scenarios, where UpIT(4E,A2) and UpIT(8E,A2) outperform Upcycle_{SFT}(4E,A2) and Upcycle_{SFT}(8E,A2) by 1.28% and 2.21%. This comprehensive analysis further corroborates the applicability of UpIT across diverse MoE architectures, consistently yielding optimal performance. 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

383

387

388

389

391

392

393

394

395

397

399

400

401

402

403

404

405

(2) In comparisons with PESC(8E,A2) and PESC(16E, A2), which take adapter structure as experts, $LoRAMoE_{PT}(8E, A2)$ and LoRAMoE_{PT}(16E,A2) display respective advantages of 0.25% and 0.13%, thereby underscoring a slightly superiority of LoRA-based MoE models over adapter-based counterparts. More than that, in FFN-based upcycling, two Self-MoE models experience a collapse in performance, a phenomenon not observed in LoRA-based settings. We posit that this is due to the excessive number of expert parameters introduced in FFN-based upcycling, and the small data in instruction tuning is insufficient to differentiate the experts sufficiently, which hinders the ability of only training routers to fit the diverse data effectively.

3.3 Scaling the Training Dataset

To assess the data-efficient nature of UpIT, we validate UpIT and vanilla upcycling approaches by randomly sampling 50K, 100K, and 200K samples from the full 500K dataset, enabling experiments across four data sizes. As illustrated in Figure 3, we have several intriguing findings.

(1) In the context of LoRA-based scenarios, UpIT(8E,A2) demonstrates comparable performance to LoRAMoE(8E,A2) trained on 500K samples, with only 50K training samples. When scaling up to 16 experts, UpIT(16E,A2) outperforms Lo-RAMoE(16E,A2) trained on the full 500K dataset again with 100K training samples. These findings extend to FFN-based settings, underscoring the data-efficient essence of UpIT and the ability to diminish the data demands of upcycling notably.

(2) Both existing LoRA and FFN-based models face performance growth saturation issues with traditional SFT (LoRA) and vanilla upcycling strategies. Specifically, while a noticeable performance increase occurs as the dataset scales from 50K to



Figure 5: Proportion of tokens dispatched to each expert on different benchmarks, where experts in UpIT exhibit stronger diversity than LoRAMoE.

	Avg.		Avg.		Avg.
UpIT(Front.Half)	51.37	UpIT(w/o EE)	53.31	UpIT(w/o Init)	49.96
UpIT(Uniform)	<u>51.71</u>	UpIT(Random)	52.41	UpIT(Random)	49.30
UpIT(Back.Half)	52.21	UpIT(Genetic)	53.34	UpIT(Skilled)	52.21

Table 2: Further analysis containing different checkpoint selection approaches (left), different expert expansion methods (middle), and different router initialization strategies (right). Detailed results are shown in Section A.11.1, A.11.2, and A.11.3

200K, the performance growth stabilizes as it continues to expand from 200K to 500K, with the average performance exhibiting a log-like curve. In contrast, UpIT shows nearly linear growth trends in FFN-based models and even exhibits accelerated performance gains as the dataset size increases in LoRA-based models. This strongly indicates that MoE models trained using UpIT could efficiently capture the principles of token dispatch more and possess a higher performance upper bound.

3.4 Scaling the Number of Experts

406

407

408

409

410

411

412

413

414

415

416

To examine the impact of scaling the number of 417 total experts and activated experts on MoE mod-418 els, we investigate the effects of UpIT and vanilla 419 420 upcycling methods with different expert numbers. The first conclusion drawn from Figure 4 is that 421 UpIT demonstrates superior performance across 422 all configurations. Furthermore, as the number of 423 activated experts increases, the growth trend of av-494 erage performance gradually slows down, which is 425 attributed to the fact that the evaluation benchmark 426 is domain-specific, and simply increasing the num-427 ber of activated parameters does not consistently 428 yield substantial improvements. We also find that 429 under the same activated parameters, as the number 430 of experts increases, vanilla upcycling even expe-431 riences several performance drops, whereas UpIT 432 433 consistently shows improvements in performance. Due to the inefficiency of data utilization in vanilla 434 upcycling, increasing the number of experts during 435 training leads to a reduction in data allocation for 436 each expert, and the router fails to dispatch tokens 437

to experts appropriately, results in unpredictable model performance.

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

3.5 Router Analysis

To assess the efficiency and interoperability of UpIT, understanding its token dispatch mechanism is essential. We comprehensively analyze the distribution patterns of designated experts across four representative benchmarks: HumanEval, GSM8K, NQ, and MMLU. The results of this examination are illustrated in Figure 5, focusing specifically on the 15th layer of LoRAMoE and UpIT with 8 total experts and 2 activated experts. Significantly, Expert 4 exhibits significantly higher activation within the HumanEval benchmark compared to the other datasets, while Expert 3 demonstrates a substantial activation rate in MMLU compared to other experts. The analysis reveals that, aside from exhibiting slight routing preferences in GSM8K, LoRAMoE dispatches tokens evenly among the experts across the other three benchmarks. In contrast, UpIT accurately allocates tokens from different domains to specific experts, highlighting the significant differentiation among routers and experts, resulting in a more data-efficient model upcycling routine.

3.6 Explore the Upper Bound of UpIT

In the main experiment, we aligned total training amounts fairly to compare UpIT with baseline methods. Here, we extend the training epochs to better understand the performance upper bound of UpIT. For the baselines, we continue training for an additional 4 epochs, but they do not show sig-



Figure 6: Performance comparison of UpIT and vanilla upcycling methods with different training epochs where x represent training epochs. For UpIT, we consider different allocations of the training epochs between expert preparation and model upcycling. Detailed results are shown in Section A.10

nificant performance gains. Instead, we expand the training for UpIT in two ways: UpIT(2,6) in-470 cludes 2 epochs for expert preparation followed by 6 epochs for post-training, while UpIT(4,4)472 comprises 4 epochs for expert preparation followed 474 by 4 epochs for post-training. Figure 6 shows that UpIT demonstrates continuous performance im-475 provement with more training epochs, indicating 476 greater potential than the baselines. Besides, it is interesting that UpIT(4,4) experiences longer iter-478 479 ation epochs during expert preparation, with greater divergences between expert models, leading to a more rapid upward trend. In the post-training stage, after only 4 epochs, it achieves comparable results 482 to UpIT(2,6) and is expected to reach higher performance in the upper bound.

3.7 Further Analysis

469

471

473

477

480

481

483

484

485

486

487

488

489

490

491

492

493

494

Different Checkpoint Selection Strategies during Expert Preparation. In this section, we would like to answer the question of how to select checkpoints if the number of saved checkpoints exceeds the required number of experts. As shown in Table 2 (left), the latter-half selection performs better. Detailed results in Table 13 indicate that the primary performance differences stem from improvements in mathematical reasoning and coding

abilities as training progresses, and selecting later checkpoints might enhance these capabilities.

495

496

497

498

499

500

501

502

503

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

535

536

537

538

539

540

541

542

543

544

Different Parameter Merging Strategies during Expert Expansion. Next, we investigate various expert merging strategies. Table 2 (mid) reveals that the genetic algorithm-based expert expansion achieves performance comparable to the method of constructing experts with more checkpoints, highlighting the effectiveness of our merging strategy in generating diverse experts. Additionally, merging two randomly selected expert models in each round results in a performance decline, again validating the importance of maintaining expert diversity.

Different Data Selection Strategies during Router Initialization. We compare the expertspecific data selection approach with two alternatives: without the router initialization stage and randomly selecting data to pre-optimize routing vectors. Table 2 (right) illustrates a significant performance decline that occurs without router initialization, underscoring the importance of this stage and the effectiveness of utilizing checkpoints as experts. Furthermore, performance diminishes when the training data is randomly selected, due to the loss of inherent diversity among checkpoints when using similar data in pre-optimizing different expert models. Overall, the phenomenon in this paper is similar to the findings of discussing expert diversity in existing work (Lo et al., 2024)

4 Conclusion

In this paper, we present a novel, flexible, scalable, and data-efficient approach, Upcycling Instruction Tuning (UpIT), for transforming dense pre-trained models into MoE instruction models. By leveraging intermediate checkpoints as specialized experts and implementing an expert expansion stage with genetic algorithms and parameter merging, UpIT successfully enhances expert diversity while allowing for a flexible number of experts. The strategic selection of seed data ensures that each expert and router performs optimally within the MoE framework. Our extensive experiments demonstrate that UpIT not only achieves superior performance across various benchmarks but also maintains stability in expert and data scaling. Further analysis emphasizes the critical importance of expert diversity in the upcycling process. Overall, UpIT offers a promising pathway for enhancing the efficiency and effectiveness of MoE models, paving the way for future advancements in the field.

Limitations

545

553

554

557

560

561

562

565

568

572

573

574

575

576

577

578

579

584

585

586

587

588

590

592

594

595

597

546 Our proposed UpIT method innovatively uses in-547 termediate checkpoints to prepare experts. Addi-548 tionally, during the router initialization stage, it 549 strengthens the discrepancy among experts and pre-550 optimizes the routers, endowing UpIT with pow-

erful data-efficient nature. Meanwhile, an expert expanding method based on the genetic algorithm is employed to effectively enhance scalability and alleviate the situation of insufficient checkpoints. Although UpIT has already addressed the deficiencies of existing Upcycling approaches, there are still the following limitations and aspects that can be further explored in the future.

(1) In order to verify the universality of the UpIT method, we have chosen the basic LoRA and FFN-based MoE architectures. The performance of UpIT in many more advanced MoE architectures is worthy of in-depth study.

(2) UpIT adopts an expert expansion method based on the genetic algorithm, and we have also observed the potential of expert expansion. More advanced expert expansion methods are worthy of future exploration.

(3) Due to the limitations of computing resources, the overall scale of our MoE model does not exceed 10 billion parameters. Therefore, it is also worth attempting to apply the UpIT framework in larger-scale MoE models.

References

- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating large language models trained on code. Preprint, arXiv:2107.03374.
 - Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind

Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv:1803.05457v1*.

598

599

601

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- OpenCompass Contributors. 2023. Opencompass: A universal evaluation platform for foundation models. https://github.com/open-compass/ opencompass.
- Damai Dai, Chengqi Deng, Chenggang Zhao, R. X. Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding Zeng, Xingkai Yu, Y. Wu, Zhenda Xie, Y. K. Li, Panpan Huang, Fuli Luo, Chong Ruan, Zhifang Sui, and Wenfeng Liang. 2024. Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models. *Preprint*, arXiv:2401.06066.
- Shihan Dou, Enyu Zhou, Yan Liu, Songyang Gao, Jun Zhao, Wei Shen, Yuhao Zhou, Zhiheng Xi, Xiao Wang, Xiaoran Fan, Shiliang Pu, Jiang Zhu, Rui Zheng, Tao Gui, Qi Zhang, and Xuanjing Huang. 2024. Loramoe: Alleviate world knowledge forgetting in large language models via moe-style plugin. *Preprint*, arXiv:2312.09979.
- William Fedus, Barret Zoph, and Noam Shazeer. 2022. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Preprint*, arXiv:2101.03961.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. *Preprint*, arXiv:2009.03300.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *Preprint*, arXiv:2106.09685.
- Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2023. Editing models with task arithmetic. *Preprint*, arXiv:2212.04089.
- Robert A. Jacobs, Michael I. Jordan, Steven J. Nowlan, and Geoffrey E. Hinton. 1991. Adaptive mixtures of local experts. *Neural Computation*, 3(1):79–87.
- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang,

764

707

657

655

- 671
- 673
- 677
- 678 679

691

703

706

Timothée Lacroix, and William El Sayed. 2024. Mixtral of experts. Preprint, arXiv:2401.04088.

- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. triviaqa: A Large Scale Distantly Supervised Challenge Dataset for Reading Comprehension. arXiv e-prints, arXiv:1705.03551.
- Junmo Kang, Leonid Karlinsky, Hongyin Luo, Zhen Wang, Jacob Hansen, James Glass, David Cox, Rameswar Panda, Rogerio Feris, and Alan Ritter. 2024. Self-moe: Towards compositional large language models with self-specialized experts. Preprint, arXiv:2406.12034.
- Aran Komatsuzaki, Joan Puigcerver, James Lee-Thorp, Carlos Riquelme Ruiz, Basil Mustafa, Joshua Ainslie, Yi Tay, Mostafa Dehghani, and Neil Houlsby. 2023. Sparse upcycling: Training mixture-of-experts from dense checkpoints. Preprint, arXiv:2212.05055.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N. Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: a benchmark for question answering research. Transactions of the Association of Computational Linguistics.
- Chen Li, Weiqi Wang, Jingcheng Hu, Yixuan Wei, Nanning Zheng, Han Hu, Zheng Zhang, and Houwen Peng. 2024a. Common 7b language models already possess strong math capabilities. Preprint, arXiv:2403.04706.
- Dengchun Li, Yingzi Ma, Naizheng Wang, Zhengmao Ye, Zhiyuan Cheng, Yinghao Tang, Yan Zhang, Lei Duan, Jie Zuo, Cal Yang, and Mingjie Tang. 2024b. Mixlora: Enhancing large language models finetuning with lora-based mixture of experts. *Preprint*, arXiv:2404.15159.
- Wing Lian, Guan Wang, Bleys Goodson, Eugene Pentland, Austin Cook, Chanvichet Vong, and "Teknium". 2023. Slimorca: An open dataset of gpt-4 augmented flan reasoning traces, with verification.
- Ka Man Lo, Zeyu Huang, Zihan Qiu, Zili Wang, and Jie Fu. 2024. A closer look into mixture-ofexperts in large language models. arXiv preprint arXiv:2406.18219.
- Tongxu Luo, Jiahe Lei, Fangyu Lei, Weihao Liu, Shizhu He, Jun Zhao, and Kang Liu. 2024. Moelora: Contrastive learning guided mixture of experts on parameter-efficient fine-tuning for large language models. Preprint, arXiv:2402.12851.
- Michael Matena and Colin Raffel. 2022. Merging models with fisher-weighted averaging. Preprint, arXiv:2111.09832.

- Niklas Muennighoff, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Jacob Morrison, Sewon Min, Weijia Shi, Pete Walsh, Oyvind Tafjord, Nathan Lambert, et al. 2024. Olmoe: Open mixture-of-experts language models. arXiv preprint arXiv:2409.02060.
- OpenAI. 2024. Gpt-4 technical report. Preprint, arXiv:2303.08774.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. Preprint, arXiv:2203.02155.
- Sainbayar Sukhbaatar, Olga Golovneva, Vasu Sharma, Hu Xu, Xi Victoria Lin, Baptiste Rozière, Jacob Kahn, Daniel Li, Wen tau Yih, Jason Weston, and Xian Li. 2024. Branch-train-mix: Mixing expert llms into a mixture-of-experts llm. Preprint, arXiv:2403.07816.
- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, , and Jason Wei. 2022. Challenging big-bench tasks and whether chain-of-thought can solve them. arXiv preprint arXiv:2210.09261.
- 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.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and finetuned chat models. Preprint, arXiv:2307.09288.
- Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. 2023. Magicoder: Source code is all you need. arXiv preprint arXiv:2312.02120.

765

766

768

- 778 779 780 781 782 783 783
- 7 7 7
- 787 788 789 790
- 791 792 793
- 793 794
- 795 796
- 797
- 800 801
- 802 803 804

8

- 807 808
- 8
- 810 811

812 813 814

- 815
- 816 817
- 817 818

- Mitchell Wortsman, Gabriel Ilharco, Samir Yitzhak Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S. Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, and Ludwig Schmidt. 2022. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. *Preprint*, arXiv:2203.05482.
 - Haoyuan Wu, Haisheng Zheng, Zhuolun He, and Bei Yu. 2024a. Parameter-efficient sparsity crafting from dense to mixture-of-experts for instruction tuning on general tasks. *Preprint*, arXiv:2401.02731.
- Xun Wu, Shaohan Huang, and Furu Wei. 2024b. Mixture of lora experts. *Preprint*, arXiv:2404.13628.
 - Fuzhao Xue, Zian Zheng, Yao Fu, Jinjie Ni, Zangwei Zheng, Wangchunshu Zhou, and Yang You. 2024. Openmoe: An early effort on open mixture-of-experts language models. *Preprint*, arXiv:2402.01739.
 - Prateek Yadav, Derek Tam, Leshem Choshen, Colin Raffel, and Mohit Bansal. 2023. Ties-merging: Resolving interference when merging models. *Preprint*, arXiv:2306.01708.
 - Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. 2024. Language models are super mario: Absorbing abilities from homologous models as a free lunch. *Preprint*, arXiv:2311.03099.
 - Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2023. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*.
 - 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.
 - Jinghan Zhang, Shiqi Chen, Junteng Liu, and Junxian He. 2023. Composing parameter-efficient modules with arithmetic operations. *Preprint*, arXiv:2306.14870.
 - Shanshan Zhong, Shanghua Gao, Zhongzhan Huang, Wushao Wen, Marinka Zitnik, and Pan Zhou. 2024.
 Moextend: Tuning new experts for modality and task extension. *Preprint*, arXiv:2408.03511.
- Tong Zhu, Xiaoye Qu, Daize Dong, Jiacheng Ruan, Jingqi Tong, Conghui He, and Yu Cheng. 2024a. Llama-moe: Building mixture-of-experts from llama with continual pre-training. *Preprint*, arXiv:2406.16554.
- Xingkui Zhu, Yiran Guan, Dingkang Liang, Yuchao Chen, Yuliang Liu, and Xiang Bai. 2024b. Moe jetpack: From dense checkpoints to adaptive mixture of experts for vision tasks. *Preprint*, arXiv:2406.04801.

Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du, Yanping Huang, Jeff Dean, Noam Shazeer, and William Fedus. 2022. St-moe: Designing stable and transferable sparse expert models. *Preprint*, arXiv:2202.08906. 819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

A Appendix

A.1 Discussion on the Effectiveness of UpIT

In this section, we mainly want to discuss why the UpIT framework is effective.

In this paper, we aim to emphasize the importance of expert discrepancy in upcycling and incorporate the idea into the entire design of UpIT. Figure 1 highlights that there is a certain discrepancy among intermediate checkpoints, and different checkpoints are proficient in different domains. Both our PPL-based data selection method and the router initialization stage make use of this characteristic among intermediate checkpoints, strengthening the discrepancy among experts and enabling the routers to correctly allocate tokens. Therefore, the motivation of the UpIT method originates from the discrepancy among checkpoints and effectively reinforces this discrepancy. Eventually, the MoE model after upcycling can obtain differentiated experts in advance and routers that can correctly allocate tokens. As a result, it can effectively reduce the data requirements and improve the performance.

A.2 Basic Concepts of MoE

In this section, we briefly review some basic concepts of MoE models.

Mixture-of-Experts (MoE). MoE significantly scale up the total parameter number and increases the knowledge capacity of language models, by selectively activating some of the parameters during inference, it does not proportionally increase the computational workload. In transformers-based models, the feed-forward neural network (FFN) layer in each transformer block is typically replaced with a MoE layer, which comprises N identical and independent experts $\{E_i\}_{i=1}^N$, along with a router $R(\cdot)$ for assigning experts, where each expert generally corresponds to an FFN layer, and in the scenario of parameter-efficient fine-tuning (PEFT), the expert may also be LoRA matrices. Formally, for the hidden states h of the attention layer, the output o of the MoE layer is represented as $\mathbf{o} = \sum_{i=1}^{N} R(\mathbf{h}) E_i(\mathbf{h})$. Here, $E_i(\mathbf{h})$ is the output of the *i*-th expert, $R(\mathbf{h})$ denotes the score for all experts, where experts with the highest scores are usually selected to calculate the final output.

Upcycling. Upcycling seeks to avoid training MoE models from scratch by transforming an existing dense model into MoE, followed by a posttraining stage to integrate all the parameters into an organic whole. It starts with dense models, forming experts by expanding the FFN layer or creating new LoRA branches, and then adding routers to control the dispatch of input tokens. In this process, the remaining embedding layers, attention blocks, normalization modules, and output layers are directly transferred from the initial dense model to the ultimate MoE model, and the router is randomly initialized and optimized in the post-training stage.

A.3 Related Work

870

871

874

875

878

879

882

890

891

894

900

901 902

903

904

905

906

907

908

909

910

911

912

913

914 915

916

917

918

919

Mixture of Experts. Mixture of Experts (MoE) (Jacobs et al., 1991) modifies the FFN layers or inserts additional branches to construct experts and activates them sparsely, thereby significantly enlarging the model capacity without noticeably increasing computational costs. The exploration of MoE has increasingly captured attention in recent years. Vanilla upcycling (Komatsuzaki et al., 2023) copies FFN layers, followed by post-training, have achieved a more convenient MoE training strategy. LoRAMoE (Dou et al., 2024), MoELoRA (Luo et al., 2024), MixLoRA (Li et al., 2024b) and MoLE (Wu et al., 2024b) develop an MoE model by incorporating several LoRA branches as experts, utilizing sparse activation or linear weighting for model construction. PESC (Wu et al., 2024a) introduces adapter-based structures after the FFN layers, exploring a parameter-efficient MoE model that diverges from the LoRA paradigm. MoExtend (Zhong et al., 2024) adapt to new tasks by expanding the MoE layer during the training process, mitigating catastrophic forgetting. MoE Jetpack (Zhu et al., 2024b) introduces checkpoint recycling, which leverages checkpoints to enhance the flexibility and diversity of expert initialization. In contrast, Branch-Train-MiX (Sukhbaatar et al., 2024) and Self-MoE (Kang et al., 2024) explore a Specialized Upcycling method by introducing specialized experts. Despite superior performance, they still require considerable domain data to acquire specialized experts. In this paper, we integrate the advantages of the work above and utilize intermediate checkpoints for expert preparation, innovatively propose an expert expansion strategy and an stage of pre-optimizing routing vectors to enhance flexibility, scalability and data efficiency. Model Merging. Model merging has emerged

as a prominent research direction in recent years, focusing on consolidating multiple task-specific models into a unified model with diverse capabilities (Wortsman et al., 2022; Ilharco et al., 2023). Model merging usually considers the combination of model parameters without accessing the original training data. Average Merging (Wortsman et al., 2022) is one typical model merging approach, which utilizes averaged parameters to construct the merged model. Task Arithmetic (Zhang et al., 2023) employs a pre-defined scaling term to distinguish the importance of various models. Fisher Merging (Matena and Raffel, 2022) performs automatic weighted merging of parameters, where the Fisher information matrix calculates the weights. TIES-Merging (Yadav et al., 2023) tackles the task conflicts in (Zhang et al., 2023) by trimming low-magnitude parameters, resolving sign disagreements, and disjointly merging parameters with consistent signs. DARE (Yu et al., 2024) first sparsifies delta parameters of several SFT homologous models and then merges them into a single model. We innovatively integrate the model merging concept into the MoE model, leveraging genetic algorithms and DARE to expand and evolve new experts. This approach enhances the scalability of our framework.

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

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

A.4 Implementation Details.

We utilize Llama 2 7B and Sheared Llama 2.7B to train LoRA-based and FFN-based models. We adopt a constant learning rate of 2e-4 and 2e-5 for LoRA-based and FFN-based settings, respectively. All models are trained for 4 epochs in total. For UpIT, we first train the dense model for 2 epochs and prepare specialized expert models using intermediate checkpoints saved between epochs 1 and 2. In router initialization, we randomly select 1% of the training data and train for 4 epochs to pre-optimize routing vectors. In model upcycling, UpIT does not introduce additional training by training the upcycling MoE models for 2 epochs. Therefore, the total training duration also amounts to 4 epochs, including 2 epochs for expert preparation and 2 for post-training. All experiments are conducted using a global batch size of 128 and a context length of 2048, running on 8 NVIDIA A800 GPUs. For evaluation, we use various benchmarks to validate comprehensively our method. Please refer to Appendix A.6 for detailed evaluation benchmarks and metrics.

Algorithm 2: Genetic Algorithm-Based Ex-								
pert Expansion								
Input: Existing <i>m</i> expert models								
$\mathcal{E} = \{\mathbf{E}_1, \cdots, \mathbf{E}_m\}$, target number of								
experts n.								
1 for $i = 1$ to $n - m$ do								
2 for $j = 1$ to $len(\mathcal{E})$ do								
3 for $k = j + 1$ to $len(\mathcal{E})$ do								
4 Find two expert models \mathbf{E}_{j^*} , \mathbf{E}_{k^*} with								
the smallest similarity.								
5 Setting weights α and β randomly, s.t.,								
$\alpha + \beta = 1.$								
6 Merge new expert model via								
$\mathbf{E}_{m+i} = \text{DARE}(\alpha \mathbf{E}_{j^*}, \beta \mathbf{E}_{k^*})$ and put \mathbf{E}_{m+i}								
Output: Expanded expert models								
$\hat{\mathcal{E}} = \{\mathbf{E}_1, \cdots, \mathbf{E}_m, \cdots, \mathbf{E}_n\}.$								

Algorithm 3: Expert-Specific Data Selection for Router Initialization

Input: Training dataset D, Expert models
E = {E₁, ··· , E_n}, and data capacity for
each expert C.

1 Initialize n empty expert-specific data buckets
{D_s¹,..., D_sⁿ}.
2 Construct seed dataset D_s by randomly selecting 1%

- of the data in \mathcal{D} .
- 3 foreach data d_i in \mathcal{D}_s do

10

970

972

974

975

977

978

979

981

985

- 4 Calculate the perplexity list \$\mathcal{P}_i = [p_i^1, \dots, p_i^n]\$ of each expert on \$d_i\$.
 5 Sort \$\mathcal{P}_i\$ in order (with small perplexity at the
- 5 Sort \mathcal{P}_i in order (with small perplexity at the beginning), denoted as \mathcal{P}'_i .

6 **foreach** p_i^i in \mathcal{P}'_i **do** 7 Get the expert index j. 8 **if** $len(\mathcal{D}_e^j) < C$ **then** 9 Append data d_i to bucket

Append data d_i to bucket \mathcal{D}_e^j . Break.

Output: Expert-specific data buckets $\{\mathcal{D}_s^1, \ldots, \mathcal{D}_s^n\}$.

A.5 Detailed Description of Baselines

For LoRA-based settings, we compare several baselines with our proposed UpIT. (1) LoRA (Hu et al., 2021): It adds low-rank matrix branches for parameter-efficient fine-tuning. (2) Self-MoE (Kang et al., 2024): It employs specialized experts to build MoE model, and only train the routers during the post-training stage. We solely reuse the training strategy of Self-MoE, and only train the routers after upcycling with intermediate checkpoints. (3) PESC (Wu et al., 2024a): It adds several adapter structures after the FFN block as experts. (4) LoRAMoE_{PT} (Dou et al., 2024): It employs the same structure as UpIT which insert several LoRA branches as experts, and (5) LoRAMoE_{SFT}: It copies the final-step checkpoint to form MoE

blocks.

Similarly, for FFN-based settings, we compare UpIT with (1) SFT, It is the standard fine-tuning solution. (2) Self-MoE (Kang et al., 2024), It is similar with the LoRA-based method. (3) Upcycle_{PT} (Komatsuzaki et al., 2023): It is the vanilla upcycling approach for transforming a dense pre-trained model to the MoE model. (4) Upcycle_{SFT}. It copies the final-step checkpoint to form MoE blocks. 986

987

988

989

990

991

992

993

994

995

996

A.6 Evaluation Metrics.

To validate the effectiveness of our method, we em-997 ploy comprehensive evaluation benchmarks, which 998 contain various abilities. (1) Factual Knowledge: 999 To assess the LLMs' factual knowledge, we employ the Massive Multitask Language Understand-1001 ing (MMLU) (Hendrycks et al., 2021), ARC-e and 1002 ARC-c (Clark et al., 2018) datasets. MMLU com-1003 prises questions across 57 subjects from elemen-1004 tary to professional difficulty levels. ARC-e and 1005 ARC-c contain questions for science exams from 1006 grade 3 to grade 9. We report the 0-shot accu-1007 racy based on answer perplexity for MMLU and 1008 ARC. (2) **Reasoning**: We utilize the test split of 1009 the Grade School Math (GSM8K) (Cobbe et al., 2021), HellaSwag (HellaS.) (Zellers et al., 2019) 1011 and Big-Bench-Hard (BBH) (Suzgun et al., 2022) 1012 benchmarks to evaluate reasoning abilities. We re-1013 port the 8-shot accuracy for GSM8K and the 3-shot 1014 accuracy for HellaSwag. (3) Coding: To probe the 1015 LLMs' ability to generate functionally correct pro-1016 grams from docstrings, we utilize HumanEval (Hu-1017 manE.) (Chen et al., 2021) and report the pass@1 1018 performance. (4) World Knowledge: We adopt 1019 Natural Question (NQ) (Kwiatkowski et al., 2019) 1020 and TriviaQA (Joshi et al., 2017) to evaluate the 1021 commonsense question-answering ability. All of the above evaluations are performed using open-1023 compass (Contributors, 2023) framework, and to 1024 expedite evaluation, we enable batch-padding with 1025 a batch size of 32. 1026

A.7 Further Experiments

A.7.1 Ablation Study of Existing Experts

As shown in Table 3, we explore the impact of the1029number of different existing experts (m in Algo-1030rithm 2) on the performance. We find that as the1031number of existing experts increases, the performance of the model continues to improve.1032

	HumanE.	GSM8K	HellaS.	BBH	MMLU	NQ	TriviaQA	ARC-c	ARC-e	Avg.
UpIT(4E,A2)	34.21	47.20	63.92	38.79	49.15	24.07	46.38	63.72	81.93	49.93
UpIT(8E,A2)	35.37	49.51	66.00	40.27	50.31	24.52	55.27	65.08	83.60	52.21
UpIT(12E,A2)	37.51	50.33	66.18	39.90	50.48	25.09	55.61	65.77	83.91	52.75
UpIT(16E,A2)	40.62	48.82	65.58	40.60	51.59	25.24	57.00	67.12	83.25	53.31

Table 3: Results of ablation study on existing experts.

	HumanE.	GSM8K	HellaS.	BBH	MMLU	NQ	TriviaQA	ARC-c	ARC-e	Avg.
UpIT(8E,A2)	35.37	49.51	66.00	40.27	50.31	24.52	55.27	65.08	83.60	52.21
UpIT(10E,A2)	36.32	48.76	66.08	39.82	50.43	24.77	54.85	66.10	83.91	52.34
UpIT(12E,A2)	37.82	49.03	65.67	40.03	50.75	25.13	55.93	66.38	84.10	52.76
UpIT(14E,A2)	39.40	48.41	66.39	39.38	50.91	25.48	56.38	67.04	84.10	53.05
UpIT(16E,A2)	40.62	48.37	66.62	39.43	50.70	25.62	56.61	67.46	84.66	53.34

Table 4: Results of ablation study on target experts.

1061

1062

1063

1064

1065

1067

1034

A.7.2 Ablation Study of Target Experts

As shown in Table 4, we explore the impact of the number of different target experts (*n* in Algorithm 2) on the performance. We find that as the number of expanded experts increases, the performance of the model also continues to grow. In particular, when the number of experts is 12 and 16, that is, the number of checkpoints is 8 and the number of expanded experts is 4 and 8 respectively, compared with directly using 12 and 16 checkpoints for expert initialization, the performance of the former is slightly better than that of the latter. This further proves the effectiveness of our expert expansion method.

A.7.3 Exploration of Catastrophic Forgetting

We conduct experiments by incorporating the pretrained model as an expert to explore the forgetting issue in training MoE models. Specifically, we conduct experiments under two configurations, one without an expert expansion stage and the other with an expert expansion stage. As shown in Table 5, Compared to Llama 2-7B, the fine-tuned models exhibit a performance decline on the TriviaQA benchmark. By incorporating the pre-trained model into the UpIT pipeline (note that only intermediate checkpoints are used in our paper), there is a considerable recovery in performance on TriviaQA. However, the performance of other benchmarks are declined, especially on HumanEval and GSM8K, as the pre-trained model is not satisfactory on these tasks. It raises the issue of the trade-off between downstream performance and pre-trained knowledge, and we believe is a topic worthy of in-depth exploration in future endeavors.

A.7.4 Checkpoints Discrepancy on Other Dataset

In order to explore that this discrepancy existing among the intermediate checkpoints is not only present in our training data, we conduct experiments using the Alpaca (Taori et al., 2023) dataset, as shown in Table 6. We find that during the training process with the Alpaca dataset, the checkpoints at different steps show the same conclusion, that is, different checkpoints are proficient in different domains. 1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1094

1095

1096

A.7.5 Compared to Specialized Upcycling

We compare the performance of the UpIT method and the Specialized Upcycling method. Specifically, we divide our data into four parts, including 100K from MetaMathQA, 100K from Magicoder, and two parts of 150K from SlimORCA, to train four domain-specific experts separately, followed by upcycling and post-training. As shown in Table 7, it can be observed that UpIT does not have a significant performance gap compared to specialized upcycling, but specialized upcycling relies on meticulously constructed domain-specific data and a more complex training process, which limits its application. In contrast, UpIT can directly learn from mixed data and reinforce discrepancies among experts.

A.8 Detail Results of Scaling the Training Dataset

Table 8 and 9 presents the detailed results of LoRA-1097based and FFN-based models under 50K, 100K,1098200K and 500K training samples which correspond1099to Figure 3. The detailed results further substan-1100

	HumanE.	GSM8K	HellaS.	BBH	MMLU	NQ	TriviaQA	ARC-c	ARC-e	Avg.
Llama 2-7B	14.63	13.95	26.58	34.73	39.84	10.06	62.06	37.29	50.26	32.16
UpIT(8E,A2)	35.37	49.51	66.00	40.27	50.31	24.52	55.27	65.08	83.60	52.21
UpIT(8E,A2) w/ pre-trained	34.15	49.33	65.36	40.78	50.96	23.75	58.73	64.71	82.26	52.23
UpIT(16E,A2)	40.62	48.37	66.62	39.43	50.70	25.62	56.61	67.46	84.66	53.34
UpIT(16E,A2) w/ pre-trained	40.13	47.67	66.76	38.20	51.32	25.47	58.11	66.49	83.82	53.11

Table 5: Results of incorporating pre-trained models as experts.

Steps	HumanE.	GSM8K	HellaS.	BBH	MMLU	NQ	TriviaQA	ARC-c	ARC-e
300	0.00	7.28	34.01	29.68	33.60	21.80	51.60	36.61	58.02
600	12.80	5.61	25.30	31.38	35.79	21.41	51.54	42.71	66.14
900	3.66	6.29	28.15	35.53	39.83	19.14	45.26	50.51	70.02
1200	14.02	7.20	40.19	32.39	33.97	20.06	50.13	24.41	35.98
1500	12.80	9.63	31.29	35.45	38.67	18.25	47.07	39.32	56.97
1800	13.41	8.87	35.85	34.51	39.90	17.40	47.40	35.93	48.85
2100	16.46	8.34	36.92	34.71	38.51	17.31	47.60	37.97	49.21
2400	14.63	8.26	35.87	34.85	39.16	17.04	47.54	37.29	52.03

Table 6: Results of different checkpoints on Alapca dataset.

tiate our findings, demonstrating that UpIT exhibits higher data utilization efficiency. It achieves
stronger performance with less data, showcasing
the data-efficient nature of UpIT.

1105 A.9 Detail Results of Scaling the Experts

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

Table 10 shows the detailed results of different numbers of experts and different activated parameters which correspond to Figure 4. The detailed results indicate that UpIT consistently maintains the desired growth trend during the scaling of experts and activated parameters, whereas the baselines exhibit an unstable performance growth trend during scaling, making it difficult to reliably predict performance expectations.

A.10 Detailed Results of Upper Bound

Table 11 and 12 show the detailed results of perfor-1116 mance upper bound during continuous training for 1117 LoRA-based and FFN-based models, respectively. 1118 The detailed results demonstrate that UpIT main-1119 tains a consistent trend of gradual performance 1120 improvement during ongoing training, while SFT 1121 (LoRA) and Upcycle_{PT} (LoRAMoE) exhibit per-1122 formance instability throughout this process. Fur-1123 1124 thermore, compared to UpIT(2), UpIT(4) benefits from a broader selection range of expert models 1125 during expert preparation, resulting in greater di-1126 versity among the expert models and thus a faster 1127 rate of performance growth. 1128

A.11 Detail Results of Further Analysis

A.11.1 Detail Results of Different Strategies during Expert Preparation

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

Table 13 shows the detailed results of different checkpoint selection approaches during expert preparation. We find that utilizing the latter half of the checkpoints as expert models yields stronger performance, particularly in mathematical reasoning and code generation capabilities. We believe this is due to the continuous improvement of mathematical reasoning and code generation as the volume of data or training increases. This observation aligns with the conclusions in (Li et al., 2024a), which indicate that selecting the latter half of the checkpoints can enhance mathematical reasoning and code generation abilities.

A.11.2 Detail Results of Different Strategies during Expert Expansion

Table 14 shows the detailed results of different expert expanding strategies during expert expansion. We are pleasantly surprised to find that our genetic algorithm-based method demonstrates highly competitive performance compared to using all expert models, indicating that the merged expert models indeed possess sufficient diversity. In contrast, the approach of randomly selecting expert models to construct new experts results in a decline in performance, which can be attributed to the insufficient diversity among the randomly chosen expert models.

	HumanE.	GSM8K	HellaS.	BBH	MMLU	NQ	TriviaQA	ARC-c	ARC-e	Avg.
Specialized(4E,A2)	34.82	46.93	63.81	39.06	48.37	23.94	47.02	62.88	83.05	49.99
UpIT(4E,A2)	34.21	47.20	63.92	38.79	49.15	24.07	46.38	63.72	81.93	49.93

	HumanE.	GSM8K	HellaS.	BBH	MMLU	NQ	TriviaQA	ARC-c	ARC-e	Avg.
Results of 50K training se	amples									
LoRA	18.29	37.00	51.84	39.71	48.61	15.46	56.27	55.93	76.90	44.45
$LoRAMoE_{PT}(8E,A2)$	17.07	34.87	52.21	37.33	47.24	22.66	53.64	54.24	75.13	43.82
$LoRAMoE_{PT}(16E,A2)$	21.34	35.48	46.01	38.39	<u>47.97</u>	23.38	52.69	56.95	74.60	44.09
UpIT(8E,A2)	18.90	33.74	58.19	39.95	47.26	27.04	60.42	<u>58.64</u>	80.07	<u>47.13</u>
UpIT(16E,A2)	18.63	<u>35.67</u>	60.11	<u>39.87</u>	47.35	<u>26.59</u>	61.37	59.96	80.76	47.81
Results of 100K training	samples									
LoRA	21.95	40.94	60.50	37.32	<u>50.97</u>	14.88	55.41	62.37	77.95	46.92
$LoRAMoE_{PT}(8E,A2)$	18.90	36.54	57.06	37.84	48.93	22.99	54.05	60.00	78.31	46.07
$LoRAMoE_{PT}(16E,A2)$	<u>20.12</u>	39.73	58.31	36.50	47.49	23.10	53.62	62.03	77.60	46.50
UpIT(8E,A2)	18.90	39.58	60.01	<u>38.17</u>	51.18	<u>25.79</u>	<u>58.76</u>	64.41	80.07	<u>48.54</u>
UpIT(16E,A2)	<u>20.12</u>	<u>40.41</u>	61.50	39.73	50.43	26.59	59.35	<u>63.73</u>	80.78	49.18
Results of 200K training	samples									
LoRA	20.73	46.32	58.86	38.64	48.79	16.07	53.93	61.69	78.84	47.10
$LoRAMoE_{PT}(8E,A2)$	26.22	42.23	58.36	37.63	47.07	20.97	51.79	60.75	77.45	46.94
$LoRAMoE_{PT}(16E,A2)$	30.49	40.03	60.59	36.67	48.42	21.52	50.80	62.03	76.37	47.44
UpIT(8E,A2)	20.12	43.97	65.16	<u>40.46</u>	51.18	<u>25.54</u>	<u>58.54</u>	66.10	79.89	<u>50.11</u>
UpIT(16E,A2)	22.82	<u>44.96</u>	<u>64.89</u>	40.64	<u>51.02</u>	25.69	60.26	<u>65.97</u>	83.25	51.06
Results of 500K training	samples									
LoRA	22.56	45.72	65.36	37.14	49.33	14.99	50.15	61.36	81.13	47.53
$LoRAMoE_{PT}(8E,A2)$	34.15	47.61	60.89	37.40	46.61	17.62	46.33	60.68	74.60	47.32
$LoRAMoE_{PT}(16E,A2)$	40.24	46.55	65.89	36.39	48.53	19.36	46.19	61.69	76.01	48.98
UpIT(8E,A2)	35.37	49.51	<u>66.00</u>	40.27	<u>50.31</u>	<u>24.52</u>	55.27	<u>65.08</u>	83.60	52.21
UpIT(16E,A2)	40.62	<u>48.37</u>	66.62	<u>39.43</u>	50.70	25.62	56.61	67.46	84.66	53.34

Table 7: Comparison between specialized upcycling and UpIT.

Table 8: Detailed results of LoRA-based models under four data sizes, where (xE,Ay) indicates that y out of x experts are activated. LoRA-based UpIT(16E,A2) is expanded from LoRA-based UpIT(8E,A2). Bold text and underlined text denote the best and second-best results in each group.

A.11.3 Detail Results of Different Strategies during Router Initialization

Table 15 shows the detailed results of different expert-specific data selection approaches. We find that our proposed PPL-based data selection method achieves the best performance. Interestingly, the method of randomly selecting expert-specific data has a detrimental effect. We believe that randomly selected data leads to a reduction in the diversity of the experts, resulting in poorer performance.

1167

1168

	HumanE.	GSM8K	HellaS.	BBH	MMLU	NQ	TriviaQA	ARC-c	ARC-e	Avg.
Results of 50K traini	ng samples									
SFT	9.76	12.28	34.43	28.43	29.57	16.01	40.13	40.00	51.50	29.12
$Upcycle_{PT}(4E,A2)$	9.15	<u>15.95</u>	29.90	<u>29.20</u>	31.75	14.46	32.57	38.64	50.79	28.05
$Upcycle_{PT}(8E,A2)$	7.93	13.87	<u>38.34</u>	29.23	33.60	14.40	35.68	<u>41.36</u>	55.73	30.02
UpIT(4E,A2)	8.54	14.78	33.05	28.39	<u>35.72</u>	<u>14.54</u>	37.68	44.41	64.02	<u>31.24</u>
UpIT(8E,A2)	<u>9.62</u>	17.39	39.29	<u>29.20</u>	36.12	14.35	<u>37.93</u>	44.41	<u>63.20</u>	32.39
Results of 100K train	ing samples									
SFT	6.71	13.27	35.87	28.55	30.95	15.68	40.75	41.49	54.32	29.73
Upcycle _{PT} (4E,A2)	14.02	20.17	42.48	29.44	36.12	14.07	33.41	45.08	58.91	32.63
$Upcycle_{PT}(8E,A2)$	12.20	21.08	47.75	28.90	38.53	14.07	35.28	48.81	63.14	<u>34.42</u>
UpIT(4E,A2)	14.76	18.57	44.54	27.99	37.82	<u>15.43</u>	<u>37.93</u>	44.07	62.26	33.71
UpIT(8E,A2)	<u>14.63</u>	21.92	<u>47.04</u>	<u>29.25</u>	39.84	14.49	36.86	<u>46.78</u>	60.85	34.63
Results of 200K train	ing samples									
SFT	9.15	16.15	44.90	29.03	34.73	16.04	40.59	44.07	60.67	32.81
$Upcycle_{PT}(4E,A2)$	19.51	24.26	<u>50.84</u>	<u>29.08</u>	37.95	13.82	34.72	46.44	59.79	35.16
$Upcycle_{PT}(8E,A2)$	18.29	22.90	47.47	30.08	37.82	13.91	35.23	50.51	<u>67.72</u>	35.99
UpIT(4E,A2)	20.63	<u>25.55</u>	50.92	<u>29.08</u>	<u>38.72</u>	15.37	35.84	<u>49.49</u>	61.20	<u>36.31</u>
UpIT(8E,A2)	21.34	28.73	47.47	<u>29.08</u>	39.84	<u>15.43</u>	<u>36.19</u>	50.51	66.49	37.23
Results of 500K train	ing samples									
SFT	26.22	29.19	38.01	26.46	33.93	8.42	18.61	42.37	58.55	31.31
$Upcycle_{PT}(4E,A2)$	<u>31.71</u>	35.10	43.40	30.23	37.93	13.74	34.72	45.08	58.73	36.74
$Upcycle_{PT}(8E,A2)$	26.22	30.62	46.76	28.95	39.84	13.57	33.86	49.62	<u>66.49</u>	37.33
UpIT(4E,A2)	31.34	33.81	<u>48.97</u>	<u>29.53</u>	40.84	14.71	36.99	47.80	65.96	<u>38.88</u>
UpIT(8E,A2)	32.19	35.64	49.15	30.23	<u>40.38</u>	<u>14.57</u>	37.93	<u>49.10</u>	68.43	39.74

Table 9: Detailed results of FFN-based models under four data sizes, where(xE,Ay) indicates that y out of x experts are activated. FFN-based UpIT(8E,A2) is expanded from FFN-based UpIT(4E,A2). Bold text and underlined text denote the best and second-best results in each group.

	HumanE.	GSM8K	HellaS.	BBH	MMLU	NQ	TriviaQA	ARC-c	ARC-e	Avg.
LoRA-based upcycling m	odels									
$LoRAMoE_{PT}(8E,A1)$	29.27	41.32	48.71	34.35	38.17	16.68	41.11	55.25	72.49	41.93
$LoRAMoE_{PT}(8E,A2)$	34.15	47.61	60.89	37.40	46.61	17.62	46.33	60.68	74.60	47.32
$LoRAMoE_{PT}(8E,A4)$	39.63	48.82	63.88	37.26	49.25	<u>20.39</u>	50.10	65.42	80.42	50.57
$LoRAMoE_{PT}(8E,A6)$	41.46	50.95	66.70	37.46	50.65	22.33	50.80	64.07	80.95	51.71
$LoRAMoE_{PT}(16E,A1)$	27.44	36.09	43.58	33.05	39.58	15.57	36.79	51.96	65.61	38.85
$LoRAMoE_{PT}(16E,A2)$	40.24	46.55	65.89	36.39	48.53	19.36	46.19	61.69	76.01	48.98
$LoRAMoE_{PT}(16E,A4)$	42.68	<u>49.96</u>	58.65	<u>37.59</u>	48.64	19.39	<u>50.94</u>	<u>65.08</u>	79.19	50.24
$LoRAMoE_{PT}(16E,A6)$	42.07	51.05	<u>64.19</u>	37.77	<u>49.84</u>	22.33	51.60	62.71	80.42	<u>51.33</u>
UpIT(8E,A1)	20.73	43.59	63.51	40.01	49.13	20.86	48.10	62.03	80.78	47.64
UpIT(8E,A2)	35.37	49.51	66.00	40.27	50.31	24.52	55.27	65.08	83.60	52.21
UpIT(8E,A4)	40.12	49.05	65.82	39.98	51.15	25.96	<u>57.50</u>	67.80	84.66	53.56
UpIT(8E,A6)	43.62	49.13	65.81	39.53	50.96	25.48	55.57	67.89	85.19	53.69
UpIT(16E,A1)	21.95	40.71	61.47	40.19	49.39	24.99	57.62	64.07	81.31	49.08
UpIT(16E,A2)	40.62	48.37	66.62	39.43	50.70	<u>25.62</u>	56.61	67.46	84.66	53.34
UpIT(16E,A4)	<u>42.53</u>	<u>50.49</u>	65.85	41.08	51.13	25.29	56.80	66.78	<u>84.83</u>	<u>53.86</u>
UpIT(16E,A6)	43.62	51.25	<u>66.14</u>	<u>40.74</u>	51.33	25.48	56.95	67.12	85.19	54.20
FFN-based upcycling mo	odels									
Upcycle _{PT} (4E,A1)	13.20	26.73	42.73	29.28	38.72	13.02	32.08	49.15	62.26	34.13
$Upcycle_{PT}(4E,A2)$	<u>31.71</u>	35.10	43.40	30.23	37.93	<u>13.74</u>	<u>34.72</u>	45.08	58.73	36.74
$Upcycle_{PT}(4E,A4)$	32.93	33.89	51.95	29.62	39.16	13.91	34.16	49.83	67.55	39.22
$Upcycle_{PT}(8E,A1)$	11.68	23.28	41.68	28.11	37.93	13.49	30.73	48.92	61.39	33.02
$Upcycle_{PT}(8E,A2)$	26.22	34.04	<u>51.57</u>	28.95	<u>39.84</u>	13.57	33.86	53.22	66.49	38.64
$Upcycle_{PT}(8E,A4)$	25.61	<u>34.27</u>	50.59	<u>30.12</u>	40.63	13.68	34.98	53.22	<u>67.02</u>	<u>38.90</u>
UpIT(4E,A1)	19.51	28.73	46.81	<u>29.56</u>	37.69	14.40	34.56	48.81	66.31	36.26
UpIT(4E,A2)	31.34	33.81	48.97	29.53	40.84	14.71	36.99	47.80	65.96	38.88
UpIT(4E,A4)	<u>33.43</u>	37.62	48.02	29.25	<u>40.76</u>	15.28	40.19	47.46	69.19	<u>40.13</u>
UpIT(8E,A1)	20.12	30.40	42.76	29.44	37.98	13.99	35.71	50.17	66.31	36.32
UpIT(8E,A2)	32.19	35.64	49.15	30.23	40.38	14.57	<u>37.93</u>	49.10	<u>68.43</u>	39.74
UpIT(8E,A4)	34.56	36.73	50.27	29.17	40.76	<u>14.79</u>	37.49	51.26	68.25	40.36

Table 10: Detailed results of different numbers of experts and activated parameters, where (xE,Ay) indicates that y out of x experts are activated. LoRA-based UpIT(16E,A2) is expanded from LoRA-based UpIT(8E,A2). Bold text and underlined text denote the best and second-best results in each group.

	HumanE.	GSM8K	HellaS.	BBH	MMLU	NQ	TriviaQA	ARC-c	ARC-e	Avg.
dense LLMs										
LoRA(1)	15.85	40.26	50.54	40.40	49.12	12.60	60.50	64.07	82.19	46.17
LoRA(2)	22.56	43.90	58.03	39.25	45.61	13.30	<u>59.15</u>	<u>63.73</u>	81.31	47.43
LoRA(3)	17.68	46.10	58.79	38.88	47.97	14.68	57.84	64.07	81.66	47.52
LoRA(4)	22.56	45.72	65.36	37.14	49.33	14.99	50.15	61.36	81.13	<u>47.53</u>
LoRA(5)	21.62	<u>47.71</u>	63.29	38.88	48.15	13.30	53.64	61.36	82.01	47.77
LoRA(6)	23.17	47.61	<u>64.89</u>	<u>40.46</u>	47.49	15.07	50.80	60.68	74.60	47.20
LoRA(7)	26.22	45.11	50.84	38.80	48.03	15.43	54.72	64.07	82.19	47.27
LoRA(8)	<u>25.73</u>	47.96	65.30	37.26	47.35	14.07	46.19	61.36	82.01	47.47
LoRA-based upcyc	ling models									
$LoRAMoE_{PT}(1)$	25.61	43.90	52.49	39.21	51.17	24.71	56.29	61.36	77.95	48.08
$LoRAMoE_{PT}(2)$	31.71	48.60	55.33	<u>38.69</u>	47.92	22.83	<u>51.18</u>	62.03	76.37	48.30
$LoRAMoE_{PT}(3)$	35.56	49.66	54.78	38.25	50.60	20.80	50.15	60.36	77.10	48.58
$LoRAMoE_{PT}(4)$	34.15	47.61	60.89	37.40	46.61	17.62	46.33	60.68	74.60	47.32
$LoRAMoE_{PT}(5)$	<u>35.98</u>	<u>48.82</u>	63.10	35.84	48.15	18.20	44.52	<u>62.71</u>	<u>78.48</u>	<u>48.42</u>
$LoRAMoE_{PT}(6)$	<u>35.98</u>	46.47	<u>63.20</u>	37.60	47.40	16.26	42.45	62.03	76.37	47.53
$LoRAMoE_{PT}(7)$	33.54	<u>48.82</u>	62.71	36.99	47.17	16.48	41.60	65.08	78.66	47.89
$LoRAMoE_{PT}(8)$	37.80	46.47	63.71	37.29	47.88	15.96	38.84	65.08	76.37	47.71
LoRA-based UpIT	with 2 epochs	s expert prej	paration							
UpIT(2,1)	32.93	49.13	64.89	39.73	49.39	25.48	54.05	64.07	80.76	51.16
UpIT(2,2)	35.37	49.51	66.00	40.27	50.31	24.52	55.27	65.08	83.60	52.21
UpIT(2,3)	36.48	50.29	67.19	38.80	49.88	25.12	54.72	68.81	82.19	52.61
UpIT(2,4)	37.69	50.49	65.85	41.08	50.96	<u>25.69</u>	56.95	67.13	<u>84.91</u>	53.42
UpIT(2,5)	<u>39.12</u>	<u>50.95</u>	<u>66.70</u>	40.27	<u>51.13</u>	25.48	58.54	<u>68.03</u>	85.19	53.93
UpIT(2,6)	41.46	51.93	66.14	<u>41.08</u>	51.33	26.59	<u>57.61</u>	68.81	84.66	54.40
LoRA-based UpIT	with 4 epochs	s expert prej	paration							
UpIT(4,1)	34.39	48.22	70.10	<u>39.43</u>	51.45	25.58	53.62	63.43	78.48	51.63
UpIT(4,2)	37.82	48.98	70.24	39.04	<u>51.52</u>	<u>25.72</u>	54.72	64.56	80.32	52.55
UpIT(4,3)	38.62	<u>49.58</u>	70.63	39.32	51.10	25.37	57.29	<u>66.79</u>	83.42	<u>53.57</u>
UpIT(4,4)	40.12	50.82	70.98	40.27	51.93	26.19	57.37	67.74	84.69	54.46

Table 11: Detailed results of performance upper bound with LoRA-based upcycling. All models are under(8E,A2) settings and(1) represents totally 1 training epoch. Bold text and underlined text denote the best and second-best results in each group.

	HumanE.	GSM8K	HellaS.	BBH	MMLU	NQ	TriviaQA	ARC-c	ARC-e	Avg.
dense LLMs										
SFT(1)	17.07	24.94	30.10	28.29	29.77	8.01	20.48	40.68	46.56	27.32
SFT(2)	19.51	27.98	27.59	26.69	32.41	8.53	20.35	38.98	55.91	28.66
SFT(3)	20.73	<u>31.77</u>	40.28	26.23	34.54	9.22	<u>20.81</u>	42.03	58.73	<u>31.59</u>
SFT(4)	26.22	29.19	38.01	26.46	33.93	8.42	18.61	42.37	<u>58.55</u>	31.31
SFT(5)	<u>24.39</u>	25.32	39.29	27.28	32.99	8.53	21.15	<u>43.39</u>	53.44	30.64
SFT(6)	22.56	32.07	41.67	27.82	<u>34.35</u>	<u>9.00</u>	17.20	44.75	55.03	31.61
SFT(7)	21.34	30.10	37.60	27.14	33.25	8.59	19.55	42.71	58.73	31.00
SFT(8)	22.56	29.42	<u>41.16</u>	26.46	32.26	8.48	18.18	34.58	51.32	29.38
FFN-based upc	ycling models	5								
$Upcycle_{PT}(1)$	18.29	31.54	42.58	29.42	35.54	<u>14.71</u>	35.18	49.49	<u>65.96</u>	35.86
$Upcycle_{PT}(2)$	20.73	37.23	51.79	28.29	40.33	14.24	34.12	47.80	64.02	37.62
$Upcycle_{PT}(3)$	26.22	34.80	45.66	30.12	38.52	14.79	34.59	49.49	66.31	<u>37.83</u>
$Upcycle_{PT}(4)$	31.71	35.10	43.40	30.23	37.93	13.74	<u>34.72</u>	45.08	58.73	36.74
$Upcycle_{PT}(5)$	28.05	35.03	<u>48.59</u>	30.22	37.08	13.91	34.57	50.85	65.26	38.17
$Upcycle_{PT}(6)$	28.05	35.41	43.40	29.93	39.27	13.93	33.07	<u>50.51</u>	61.39	37.22
$Upcycle_{PT}(7)$	29.27	35.03	46.26	<u>30.30</u>	38.18	14.24	33.04	49.49	62.26	37.56
$Upcycle_{PT}(8)$	<u>30.62</u>	<u>35.71</u>	44.19	30.71	35.99	10.75	33.02	47.80	62.43	36.80
FFN-based Upl	T with 2 epoc	chs expert p	reparation							
UpIT(2,1)	17.07	29.11	45.77	31.03	36.07	<u>15.24</u>	38.34	43.05	<u>68.43</u>	36.01
UpIT(2,2)	31.34	33.81	48.97	29.53	40.84	14.71	36.99	47.80	65.96	38.88
UpIT(2,3)	33.78	35.03	51.43	29.75	41.22	14.88	35.09	<u>55.59</u>	64.02	40.09
UpIT(2,4)	35.69	<u>37.62</u>	<u>51.06</u>	28.06	40.88	14.16	34.39	56.08	67.90	40.65
UpIT(2,5)	<u>36.75</u>	38.61	50.59	28.62	38.52	15.24	<u>37.12</u>	<u>55.59</u>	68.96	<u>41.11</u>
UpIT(2,6)	39.12	<u>37.62</u>	51.95	<u>30.12</u>	39.60	15.96	36.18	<u>55.59</u>	66.67	41.42
FFN-based Upl	T with 4 epoc	chs expert p	reparation							
UpIT(4,1)	17.07	27.37	43.90	29.84	38.92	15.40	38.30	50.17	64.55	36.17
UpIT(4,2)	30.62	31.08	48.74	<u>30.62</u>	<u>39.69</u>	13.71	<u>37.38</u>	49.83	70.90	39.17
UpIT(4,3)	<u>34.69</u>	<u>37.23</u>	51.73	30.84	39.64	13.38	35.18	<u>53.22</u>	<u>70.19</u>	<u>40.68</u>
UpIT(4,4)	38.92	40.26	<u>50.56</u>	30.05	39.81	<u>14.79</u>	36.18	53.61	69.49	41.52

Table 12: Detailed results of performance upper bound with FFN-based upcycling. All models are under(4E,A2) settings and(1) represents totally 1 training epoch. Bold text and underlined text denote the best and second-best results in each group.

	HumanE.	GSM8K	HellaS.	BBH	MMLU	NQ	Tri.QA	ARC-c	ARC-e	Avg.
UpIT(Front.Half)	21.34	44.88	64.48	40.55	51.36	27.56	59.36	68.47	84.30	51.37
UpIT(Uniform)	27.44	43.90	<u>65.39</u>	40.73	50.93	<u>25.90</u>	<u>58.31</u>	<u>68.14</u>	84.66	<u>51.71</u>
UpIT(Back.Half)	35.37	49.51	66.00	40.27	50.31	24.52	55.27	65.08	83.60	52.21

Table 13: Detailed results of different checkpoint selection strategies during expert preparation. Front.Half represents selecting the first half of checkpoints, Uniform represents uniformly selecting checkpoints and Back.Half represents selecting the back half ones which is used in our paper. All models are LoRA-based models under (8E,A2)

	HumanE.	GSM8K	HellaS.	BBH	MMLU	NQ	Tri.QA	ARC-c	ARC-e	Avg.
UpIT(w/o EE)	40.62	48.82	65.58	40.60	51.59	25.24	57.00	67.12	83.25	<u>53.31</u>
UpIT(Random)	38.72	44.71	66.14	40.27	<u>50.96</u>	<u>25.48</u>	54.72	<u>67.13</u>	<u>83.60</u>	52.41
UpIT(Genetic)	40.62	<u>48.37</u>	66.62	39.43	50.70	25.62	<u>56.61</u>	67.46	84.66	53.34

Table 14: Detailed results of different expanding strategies during expert expansion. w/o EE represents directly using checkpoints for expert preparation without expert expansion. Random represents randomly selecting two experts to merge a new one during expert expansion and Genetic represents the selection approach shown in Algorithm 3. All models are LoRA-based models under (16E,A2)

	HumanE.	GSM8K	HellaS.	BBH	MMLU	NQ	Tri.QA	ARC-c	ARC-e	Avg.
UpIT(w/o Init)	<u>32.32</u>	49.81	66.86	<u>39.15</u>	48.88	21.08	<u>49.67</u>	62.71	79.19	<u>49.96</u>
UpIT(Random)	28.05	46.70	64.27	38.67	<u>49.63</u>	<u>21.11</u>	48.67	64.41	<u>82.19</u>	49.30
UpIT(Skilled)	35.37	<u>49.51</u>	66.00	40.27	50.31	24.52	55.27	65.08	83.60	52.21

Table 15: Detailed results of different data selection strategies during router initialization. w/o Init represents training models without our proposed router initialization. Random represents randomly construct the expert-specific data and Skilled represents our PPL-based data selection method as shown in Algorithm 3. All models are LoRA-based models under (8E,A2).