Multimedia Event Extraction with LLM Knowledge Editing

Anonymous EMNLP submission

Abstract

Multimodal event extraction task aims to identify event types and arguments from visual and 003 textual representations related to events. Due to the high cost of multimedia training data, previous methods mainly focused on weakly alignment of excellent unimodal encoders. However, 007 they ignore the conflict between event understanding and image recognition, resulting in redundant feature perception affecting the understanding of multimodal events. In this paper, we propose a multimodal event extraction strategy with a multi-level redundant feature selection mechanism, which enhances the event understanding ability of multimodal large lan-014 guage models by leveraging knowledge editing 015 techniques, and requires no additional parame-017 ter optimization work. Extensive experiments show that our method outperforms the state-ofthe-art (SOTA) baselines on the M2E2 benchmark. Compared with the highest baseline, we 021 achieve a 34% improvement of precision on event extraction and a 11% improvement of F1 on argument extraction.

1 Introduction

024

034

040

In the real world, the representation of events frequently encompasses a multitude of expression modalities such as images and text. This catalyzes the evolution of Multimedia Event Extraction (MEE) task, which aims to extract the event type and argument information from the event's text description and associated image.

Due to the high cost of building multimedia data sets, existing multimedia event extraction dataset M2E2 only provides data for testing, which restricts the fine-tuning on multimodal large language model (MLLM) in multimedia event extraction task. Recent MEE methods mainly focus on the weak alignment of features obtained from wellpretrained unimodal encoders. According to their cross-modal alignment strategy, they can be divided into graph-based methods, which align the



Figure 1: Example of multimedia data on M2E2, the overunderstanding of visual features led to the incorrect prediction of MLLM on image-only task, which further effect the prediction on multimedia task.

visual objects and textual entities with a knowledge graph structure (Li et al., 2020, 2022; Liu et al., 2024); template-based methods, which align candidate items from different modalities with the template uniformly (Seeberger et al., 2024); and fine-tune-based methods, which align the features under different modal by parameter-efficient finetuning on multimodal large language models (Du et al., 2023; Sun et al., 2024).

043

044

045

046

057

060

061

062

063

064

065

066

067

These methods have achieved good performance, however, they ignore the differences of visual encoding requirements between image recognition and event extraction task. In fact, there is a tradeoff between the event extraction task and the image recognition task: the image recognition task desires to perceive the image features as much as possible, while the event extraction task prefers to focus on those important features and neglect redundant content. Due to the fact that the visual encoder utilized by previous methods are trained from image recognition tasks, there may exist redundant feature perception capabilities which affect event understanding of MLLMs. Taking Fig. 1 as an example, the image data from M2E2 are originally intended to show a person running away with a computer. However, due to the forward-leaning posture of the

body which resembles being detained, it is wrongly
predicted as an arrest when only inputting visual
features, which further misleads the judgment of
multimodal large language models.

075

077

079

087

090

100

101

102

104

105

To address this problem, we propose a multimedia event extraction framework that enhances the event understanding ability of MLLMs by editing the knowledge of the visual layer. Specifically, we first design a multilevel redundant neuron selection mechanism with information entropy and L1-normalization, and then conduct a mask matrix based knowledge editing strategy. The evaluations on M2E2 show that our method significantly outperforms all SOTA methods for multimedia event extraction. The main contributions of this paper are as follows:

> • We propose a multimedia event extraction framework based on LLM knowledge editing which enhances the understanding ability of MLLMs on event structure features, avoids the reliance on the quality of training data or the consistency of data distribution brought by strategies such as fine-tuning.

> • We propose a multi-level knowledge editing strategy which can effectively mitigate the influence of redundant neurons. Compared with the existing popular knowledge editing strategies, our method does not introduce any additional parameters and has extremely low computational and storage costs.

 The experiments on M2E2 benchmark show that our model achieves SOTA performance.
 We outperform the SOTA baselines by 34% precision on multimedia event extraction and by 11% F1 on multimedia argument extraction.

2 RELATED WORK

2.1 Multimedia Event Extraction

Multimedia event extraction task (MEE) mainly 106 includes event extraction and event argument ex-107 traction. Existing MEE methods can be divided 108 into graph-based methods, template-based methods, and fine-tune-based methods according to 110 their cross-modal alignment strategy. Graph-based 111 methods typically supervise visual objects with tex-112 tual entities with the association algorithm based 113 on knowledge graph structure. Following this 114

idea, WASE (Li et al., 2020) applies an attentionbased graph encoder to link the entities and objects extracted in different modals. Clip-event (Li et al., 2022) proposes an event graph alignment via optimal transport and evaluates their similarity at different granularities. MGIM (Liu et al., 2024) represents the multimedia information in a graph structure and performs coarse-grained alignment. Template-based methods identify the arguments according to the similarity calculated between candidates and roles in templates. For example, MMUTF (Seeberger et al., 2024) designs a unified template filling framework that encodes the corresponding event roles in templates to match the candidate textual or visual arguments. Fine-tunebased methods achieve cross-modal alignment by updating some parameters. CAMEL (Du et al., 2023) proposes an incremental training strategy that complements missing images and text with a bidirectional cross-modality data augmentation. UMIE (Sun et al., 2024) adaptively learns the correlation between textual embedding and visual objects with a gated attention mechanism.

115

116

117

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

Compared to these methods, our method is the first to apply knowledge editing technique to alleviate the redundancy feature understanding problem on MEE, obtaining more accurate event structure understanding with less cost overhead.

2.2 Instruction following

Instruction following aims to help LLMs to understand various instructions and produce the corresponding responses, so that the model can be suitable for downstream tasks. Existing instruction following technology can be divided into three types: natural-language-inference-oriented instructions (Yin et al., 2019; Xu et al., 2023; Zhong et al., 2021), which unifies various classification problems into an inference task by constructing hypotheses to explain the labels; LLM-oriented instructions (Gao et al., 2020; Bach et al., 2022a,b), which utilizes templates to convert the origin inputs into machine-friendly fill-in-blank questions; and human-oriented instructions (Mishra et al., 2021; Gupta et al., 2022; Yin et al., 2022), which splits the task into definitions, demonstrations, instances, and other information with obvious descriptive, and paragraph-style based on the way people understand the problem. These techniques can significantly improve the model's ability to understand downstream tasks through different instruction de-



Figure 2: The framework of our method, we first select the candidate layer with entropy and Gaussian smoothing, then identify the redundant neurons on visual layers, and finally conduct knowledge editing for MLLMs.

sign strategies, and have been widely used in fewshot and zero-shot classification scenarios.

2.3 Knowledge editing

165

166

167

168

171

172

173

174

175

176

177

178

179

181

184

185

188

189

190

193

194

195

196

Knowledge editing aims to update, correct, optimize, or supplement existing knowledge. Common LLM knowledge editing methods can be divided into external-memorization-based methods, reparameterization-based methods, and localmodification-based methods. In order to obtain an accurate representation of knowledge in downstream tasks, external-memorization-based methods update knowledge by maintaining a memory which retrieves the most relevant cases (Zheng et al., 2023; Mitchell et al., 2022; Zhong et al., 2023), or introducing additional trainable parameters while maintaining the pre-trained backbone unchanged (Pfeiffer et al., 2020; Lei et al., 2023; Houlsby et al., 2019); reparameterization-based methods construct low-dimensional reparameterization of original parameters for training and equivalently transforms it back for inference (Lin et al., 2024; Hu et al., 2022; Wu et al., 2024); localmodification-based methods aim to locate and optimize relevant parameters for specific knowledge learning in LLM (Guo et al., 2020; Zhang et al., 2024; He et al., 2023). Our method is designed based on the idea of local modification, but compared with the previous local-modification-based work, we propose a knowledge editing technique that does not require back propagation, which improves the performance on downstream tasks with extremely small computational overhead.

3 Model

Let d = (t, v) represent a multimedia document consisting of a text-image pair that describes a specific event, where t is the text and v is the corresponding image. The multimedia event extraction task contains the following two components. 197

198

200

201

202

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

221

222

223

225

Event Extraction: Given a multimedia document d, Event Extraction (EE) requires predicting the event type y_e from the candidates C.

Event Argument Extraction: Given a multimedia document d and event type y_e , Event Argument Extraction (EAE) aims to identify event arguments (entities or objects) with pre-defined event roles (all participants and attributes) associated with y_e .

Different from the currently popular methods of introducing additional parameters to fine-tune the pre-trained model for downstream tasks (Pfeiffer et al., 2020; Houlsby et al., 2019; Hu et al., 2022), we propose a large model knowledge editing approach that does not require fine-tuning. As mentioned above, we focus on the visual perception layer. Our approach consists of two stages: redundant neuron selection and knowledge editing. We first propose a multilevel redundant neuron selection mechanism for MLLMs, and then use a mask-matrix-based editing strategy to edit redundant visual perception knowledge.

3.1 Multilevel redundant neuron selection

Relevant studies show that the pre-trained trunk exhibits different feature patterns at different parameter positions. In order to improve the ability of large models to understand event features, a redundant neuron recognition strategy is proposed in this paper. Specifically, given the multimodal large model $\theta = \{W^1, W^2, ..., W^n\}$, where $W^k = \{w_{ij}^k | i \in I, j \in O\}$ represents the network parameter set of layer k of the large language model, I and O represent the number of input feature and output feature dimensions of layer k respectively. We first convert the network weight parameter into a probability distribution, and the formula is expressed as:

$$P_{ij}^{k} = \frac{\exp(w_{ij}^{k})}{\sum_{b=1}^{O} \sum_{a=1}^{I} \exp(w_{ab}^{k})}$$
(1)

We use entropy to measure the importance of neural network layers, which can be expressed as:

$$H_i = -\sum_{b=1}^{O} \sum_{a=1}^{I} P_{ij}^k \log_2(P_{ij}^k)$$
(2)

240

241

242

259

260

261

265

269

$$H_i^{avg} = \sum_{k=-K}^{K} \mathcal{G}(k;\sigma) \cdot H_{i+k}$$
(3)

where $\mathcal{G}(k; \sigma)$ is a Gaussian kernel with bandwidth 244 σ, K is the kernel size. The consideration behind 245 246 this is that if the probability distribution of weights within a neural network layer is very different, and 247 such difference is more obvious than the weight 248 distribution of the neighboring layer, it means that there is a certain preference from input features to outputs in that layer, resulting in the layer being sensitive to feature learning. We select the layer of neural network with the most average weight distribution according to the entropy value, which is represented as θ^r , and consider that redundant information transfer in this layer is the cause of incorrect knowledge learning. 257

> Next, we evaluated the differences in the information value of neurons within θ^r . We think that neurons represent the ability to perceive different feature dimensions, and we want to identify those neurons that are used to capture redundant features. For each weight belonging to θ^r , we use the L1norm to identify redundant neurons and information transfer pathways. Neuron c_i importance can be expressed as:

$$c_i^r = \sum_{j=1}^{O} |w_{ij}^r|$$
 (4)

We use δ to represent the threshold for determining whether a neuron is redundant, here δ is defined in [0,1] in order to avoid the case of layer collapse. It is used to represent the importance of neurons at the decomposition position represented by the threshold value. Neurons that are determined to be redundant can be expressed as:

$$C^r = \{c_i^r \mid c_i^r < s\} \tag{5}$$

270

271

272

273

275

276

277

278

279

285

287

291

292

293

294

3.2 Mask matrix based knowledge editing strategy

In order to modify the understanding of event knowledge of the MLLMs, we apply a mask matrix to mask the redundant neurons. The mask matrix can be expressed as:

$$M_{ij} = \begin{cases} 0, c_j^r \in C^r \\ 1, c_j^r \notin C^r \end{cases} \forall i \in I, j \in O \qquad (6) \end{cases}$$

The network weight parameters after the mask can be expressed as:

$$\tilde{W^r} = w_{ij}^r * M_{ij} \tag{7}$$

 M_{ij} indicates the state information of the mask matrix at (i, j). For each parameter in \tilde{W}^r , it is masked if it represents the weight connected to the neuron in C^r , otherwise it is retained as original.

4 Experiment

In this section, we extensively evaluate our framework by comparing against existing SOTA approaches and variants of MLLM with different redundant neuron selection strategies.

| Event Type | Argument Role | | | | | |
|-----------------------|---------------------------|--|--|--|--|--|
| Movement:Transport | Agent, Artifact, Vehicle, | | | | | |
| | Destination, Origin | | | | | |
| Conflict:Attack | Attacker, Target, Instru- | | | | | |
| | ment, Place | | | | | |
| Conflict.Demonstrate | Entity, Police, Instru- | | | | | |
| | ment, Place | | | | | |
| Justice:ArrestJail | Agent, Person, Instru- | | | | | |
| | ment, Place | | | | | |
| Contact:PhoneWrite | Entity, Instrument, Place | | | | | |
| Contact:Meet | Participant, Place | | | | | |
| Life:Die | Agent, Instrument, Vic- | | | | | |
| | tim, Place | | | | | |
| Transaction:Transfer- | Giver, Recipient, Money | | | | | |
| Money | | | | | | |

Table 1: Event types and corresponding argument roles in multimedia dataset M2E2.

| Metrics Task | Event Extraction Task | | | | | | | | |
|-------------------|-----------------------|-------|-------|------------|-------|-------|------------|-------|-------|
| | text-only | | | image-only | | | multimedia | | |
| Baselines | Р | R | F1 | P | R | F1 | P | R | F1 |
| WASE | 42.8 | 61.9 | 50.6 | 43.1 | 59.2 | 49.9 | 43.0 | 62.1 | 50.8 |
| CLIP-EVENT | - | - | - | 41.3 | 72.8 | 52.7 | - | - | - |
| CAMEL | 45.1 | 71.8 | 55.4 | 52.1 | 66.8 | 58.5 | 55.6 | 59.5 | 57.5 |
| UMIE | - | - | - | - | - | - | - | - | 62.1 |
| MMUTF | 45.1 | 71.8 | 55.4 | 55.1 | 59.1 | 57.0 | 47.9 | 63.4 | 54.6 |
| Qwen2vl-7b | 80.05 | 73.94 | 74.49 | 69.11 | 65.38 | 64.83 | 83.77 | 73.14 | 70.64 |
| Variant-1 | 81.36 | 71.81 | 73.06 | 68.55 | 63.09 | 64.34 | 83.88 | 73.46 | 70.99 |
| Variant-2 | 82.07 | 70.74 | 72.70 | 69.03 | 64.15 | 63.46 | 85.41 | 79.93 | 82.66 |
| Variant-3 | 81.07 | 70.89 | 72.26 | 69.15 | 65.10 | 64.89 | 89.13 | 87.06 | 86.59 |
| Ours | 82.10 | 76.06 | 76.01 | 69.32 | 65.77 | 65.30 | 89.82 | 87.70 | 87.39 |

Table 2: Results of Event Extraction on M2E2, we compare our framework against 5 baselines and 3 variants, the bold numbers denote the best results of all methods.

Table 3: Results of Argument Extraction on M2E2.

| Metrics Task | Argument Extraction Task | | | | | | | | |
|-------------------|--------------------------|-------|-------|------------|-------|-------|------------|-------|-------|
| | text-only | | | image-only | | | multimedia | | |
| Baselines | Р | R | F1 | P | R | F1 | P | R | F1 |
| WASE | 23.5 | 30.3 | 26.4 | 14.5 | 10.1 | 11.9 | 19.5 | 18.9 | 19.2 |
| CLIP-EVENT | - | - | - | 21.1 | 13.1 | 17.1 | - | - | - |
| CAMEL | 24.8 | 41.8 | 31.1 | 21.4 | 28.4 | 24.4 | 31.4 | 35.1 | 33.2 |
| UMIE | - | - | - | - | - | - | - | - | 24.5 |
| MMUTF | 33.6 | 44.2 | 38.2 | 23.6 | 18.8 | 20.9 | 39.9 | 20.8 | 27.4 |
| Qwen2vl-7b | 50.10 | 36.19 | 37.83 | 27.30 | 27.63 | 27.31 | 47.13 | 42.70 | 43.04 |
| Variant-1 | 48.07 | 36.25 | 36.81 | 25.42 | 26.95 | 25.67 | 46.45 | 40.63 | 42.77 |
| Variant-2 | 50.69 | 36.19 | 37.95 | 25.08 | 24.58 | 25.88 | 48.02 | 42.25 | 44.08 |
| Variant-3 | 51.45 | 37.05 | 38.22 | 26.64 | 27.63 | 25.90 | 47.47 | 41.03 | 43.39 |
| Ours | 51.92 | 37.48 | 39.34 | 27.30 | 27.63 | 27.31 | 48.11 | 43.89 | 44.25 |

297

4.1

- 299
- 301 304

305

types of arguments, as shown in Table. Specifically, M2E2 contains 245 multimedia documents with 6167 sentences and 1014 images. There are 1,297 text events and 391 visual events, of which 192 text event mentions and 203 visual event mentions are arranged into 309 multimedia events.

Experimental settings

Baselines We compare our proposed method against a wide range of SOTA models: WASE (Li

Datasets We evaluated our framework on the

M2E2 benchmark, a large-scale multimedia event

extraction dataset with 8 types of events and 15

et al., 2020) uses weakly aligned structured embedding to encode the multimodal events, CLIP-Event (Li et al., 2022) performs visual event extraction with a pre-trained CLIP network, CAMEL (Du et al., 2023) learns more accurate text alignment with image generator and image captioner, UMIE (Sun et al., 2024) constructs a series of instruction following templates, and MMUTF (Seeberger et al., 2024) addresses the MEE with candidate-query matching.

307

308

309

310

311

312

313

314

315

316

317

318

Experimental Setup In this work, we use Qwen2-VL-7B as the foundational model. The



Figure 3: Cases on M2E2, the above two shows the error correction cases in the MEE task, and the following two correct errors for text argument and image argument in the MAE task.

experiments are conducted on an NVIDIA A100 GPU with a memory capacity of 40GB, and due to GPU memory limitations, the size of images in M2E2 is set to 512*512 resolution. In addition, we also build several MLLM-based variants with other redundant neuron selection strategies which expands the scope of knowledge editing to all fully connected layers, uses L1-normalization of neurons for identification, and focuses on the redundancy of weights rather than the neurons.

4.2 Results

319

320

321

326

327

328

Tab. 2 and Tab. 3 show the comparison in detail. 330 Compared with previous SOTA methods, the original Qwen2-VL-7b has a significantly outperformance, demonstrating the strong semantic understanding ability of the pre-trained multimodal large model for text and images. Our method inherits such ability, and enhances the event structure understanding of MLLMs with knowledge editing that makes a comprehensive performance improvement. Especially for multimedia event extraction, the selection of redundant feature makes a 6.15% Pre-340 cision improvement, 14.56% Recall improvement, 341 16.75% F1 improvement on multimedia event ex-342 traction. Compared with Variant-1, the focus on the visual encoding layer makes a 5.94% Precision improvement, which validates our proposed hypothesis that the bottleneck of multimodal event 346 extraction lies in the capability of image event understanding. The comparison against Variant-2

demonstrates the superiority of the entropy-based redundant identification mechanism over directly judging based on L1 values. And the comparison against Variant-3 may due to the unstable weight masking causes the broken of structured sparsity, which may causes incomplete input received by downstream neurons, affecting feature expression ability. For multimedia argument extraction task, our method still achieve a 8.21% Precision improvement, 8.79% Recall improvement, 11.05% F1 improvement compared with the nearest SOTA performance. We attribute this as our appropriate feature extraction granularity which revises the understanding of roles in events by MLLM. Note that we do not utilize any additional image-text paired datasets or cross-modal data augmentation for fine-tuning in this work. Incorporating crossmodal information from the document's context and external datasets, as well as customized Prompt Engineering for each task, might further improve our method's performance.

349

350

351

352

353

354

355

356

357

358

360

361

362

364

365

366

367

368

370

371

372

373

374

375

376

378

4.3 Case Study

We selected four representative cases to illustrate the effect of multimodal large language models after knowledge editing, as shown in Fig. 3.

In case-1, the main reason for the original MLLM's error lies in the presence of multiple triggers strongly associated with different candidate event types in the text: the simultaneous existence of "Protestor" and "detained" leads to a

misunderstanding of the event from a textual per-379 spective. Our method does not focus on achieving 380 more accurate textual identification but corrects the erroneous event understanding by optimizing the main subject's actions in visual modality. Another type of error is shown in case-2: the event "Confilct: Attack" is not reflected in the corresponding image. We attribute this to the temporal misalignment between the text and image information. Despite the painstaking efforts of M2E2 to clean the mismatched text-image data, such temporal misalignment still affects model's understanding of the event content. In contrast, our method alleviates the over-interpretation of image information by filtering redundant image features, thereby achieving an appropriate fusion of event text and image features.

Case-3 shows the mistake of identifying the entity of "Confilct: Attack" as "Islamist terroism". In fact, it illustrates a typical problem of multimedia textual argument extraction with original LLM. Due to the lack of understanding of the event structure, the pre-trained model interprets the object of 400 the event as the subject of the event. Although our model only edits the visual layer, it still significantly enhances the ability to extract textual argu-403 ments, which are also observed on text-only MAE. 404 We infer that it may be because the moderate visual 405 sparsity can induce text-side capability compen-406 sation, resulting in an overall performance boost. Case-4 gives an example of correcting faulty visual 408 argument extractions. Compared with the origi-409 nal MLLM, our method provides more accurate 410 visual event object identification with appropriate size boundaries. We attribute this to our knowledge editing strategy which focus on eliminating redun-413 dant visual features. It not only solves the feature 414 dilution caused by excessive visual layer channels, 415 but also reduces the interference of noise on sub-416 sequent layers, making the model more inclined to focus on some key features such as lines and 418 colors, ultimately alleviating the boundary blurring 419 problem of the original MLLM. 420

5 Conclusion

401

402

407

411

412

417

421

In this study, we propose a multimedia event ex-422 traction method based on the knowledge editing of 423 large language models. We first analyze the reasons 424 for the poor performance of multimodal large lan-425 guage models in multimedia event extraction tasks, 426 and then design a multi-level redundant neuron 427 selection mechanism. Finally, based on the identi-428

fied redundancy, we edited the knowledge of LLM 429 to enhance the understanding of event structure. 430 Compared to previous work, our method does not 431 require fine-tuning and has extremely low compu-432 tational and storage costs. Substantial experiments 433 on the M2E2 benchmark show that our method 434 significantly improves the LLM's understanding 435 of multimedia event structures, and has become a 436 new SOTA with a 34% precision improvement on 437 MEE and a 11% F1 performance improvement on 438 MAE. Future work mainly focus on the research 439 of domain adaptation capabilities for multimodal 440 large language models. 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

469

470

471

472

473

474

Limitations

Although our work achieved a significant improvement in accuracy in multimodal event extraction and argument extraction tasks through knowledge editing of the visual layer, this editing is not as sensitive to image only argument extraction as other subtasks. This may be due to our multimodal large model's tendency towards object level boundingboxes rather than instance level, and we leave more validation work to future experiments.

In addition, we found in the experiment that the M2E2 dataset itself has some missing argument labels. In fact, the extracted events can be further subdivided into events with clear subject object relationships (e.g. Justice:Arrest-Jail) and events with unclear subject object relationships (e.g. Contact:Meet). Normally, they should have different mining difficulties, but due to the limitations of multimodal datasets, we have to leave these tasks for further research.

- Stephen H Bach, Victor Sanh, Zheng-Xin Yong, Albert Webson, Colin Raffel, Nihal V Nayak, Abheesht Sharma, Taewoon Kim, M Saiful Bari, Thibault Fevry, et al. 2022a. Promptsource: An integrated development environment and repository for natural language prompts. arXiv preprint arXiv:2202.01279.
- Stephen H Bach, Victor Sanh, Zheng-Xin Yong, Albert Webson, Colin Raffel, Nihal V Nayak, Abheesht Sharma, Taewoon Kim, M Saiful Bari, Thibault Fevry, et al. 2022b. Promptsource: An integrated development environment and repository for natural language prompts. arXiv preprint arXiv:2202.01279.
- Zilin Du, Yunxin Li, Xu Guo, Yidan Sun, and Boyang 475 Li. 2023. Training multimedia event extraction with 476 generated images and captions. In Proceedings of the 477

Tianyu Gao, Adam Fisch, and Danqi Chen. 2020. Making pre-trained language models better few-shot learners. arXiv preprint arXiv:2012.15723. Demi Guo, Alexander M Rush, and Yoon Kim. 2020. Parameter-efficient transfer learning with diff pruning. arXiv preprint arXiv:2012.07463. Prakhar Gupta, Cathy Jiao, Yi-Ting Yeh, Shikib Mehri, Maxine Eskenazi, and Jeffrey P Bigham. 2022. Instructdial: Improving zero and few-shot generalization in dialogue through instruction tuning. arXiv preprint arXiv:2205.12673. Haoyu He, Jianfei Cai, Jing Zhang, Dacheng Tao, and Bohan Zhuang. 2023. Sensitivity-aware visual parameter-efficient fine-tuning. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 11825–11835. Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In International conference on machine learning, pages 2790-2799. PMLR. Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2022. Lora: Low-rank adaptation of large language models. ICLR, 1(2):3. Tao Lei, Junwen Bai, Siddhartha Brahma, Joshua Ainslie, Kenton Lee, Yanqi Zhou, Nan Du, Vincent Zhao, Yuexin Wu, Bo Li, et al. 2023. Conditional adapters: Parameter-efficient transfer learning with fast inference. Advances in Neural Information Processing Systems, 36:8152-8172. Manling Li, Ruochen Xu, Shuohang Wang, Luowei Zhou, Xudong Lin, Chenguang Zhu, Michael Zeng, Heng Ji, and Shih-Fu Chang. 2022. Clip-event: Connecting text and images with event structures. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 16420-16429. Manling Li, Alireza Zareian, Qi Zeng, Spencer Whitehead, Di Lu, Heng Ji, and Shih-Fu Chang. 2020. Cross-media structured common space for multimedia event extraction. arXiv preprint arXiv:2005.02472. Yang Lin, Xinyu Ma, Xu Chu, Yujie Jin, Zhibang Yang, Yasha Wang, and Hong Mei. 2024. Lora dropout as a sparsity regularizer for overfitting control. arXiv preprint arXiv:2404.09610. Yang Liu, Fang Liu, Licheng Jiao, Qianyue Bao, Long Sun, Shuo Li, Lingling Li, and Xu Liu. 2024. Multigrained gradual inference model for multimedia event extraction. IEEE Transactions on Circuits and Systems for Video Technology. 8

31st ACM International Conference on Multimedia,

pages 5504–5513.

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

506

507

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

527

- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2021. Cross-task generalization via natural language crowdsourcing instructions. *arXiv preprint arXiv:2104.08773*.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D Manning, and Chelsea Finn. 2022. Memorybased model editing at scale. In *International Conference on Machine Learning*, pages 15817–15831. PMLR.
- Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. 2020. Adapterfusion: Non-destructive task composition for transfer learning. *arXiv preprint arXiv:2005.00247*.
- Philipp Seeberger, Dominik Wagner, and Korbinian Riedhammer. 2024. Mmutf: Multimodal multimedia event argument extraction with unified template filling. *arXiv preprint arXiv:2406.12420*.
- Lin Sun, Kai Zhang, Qingyuan Li, and Renze Lou. 2024. Umie: Unified multimodal information extraction with instruction tuning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 19062–19070.
- Xun Wu, Shaohan Huang, and Furu Wei. 2024. Mixture of lora experts. *arXiv preprint arXiv:2404.13628*.
- 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. *arXiv preprint arXiv:2304.12244*.
- Wenpeng Yin, Jamaal Hay, and Dan Roth. 2019. Benchmarking zero-shot text classification: Datasets, evaluation and entailment approach. *arXiv preprint arXiv:1909.00161*.
- Wenpeng Yin, Jia Li, and Caiming Xiong. 2022. Contintin: Continual learning from task instructions. *arXiv preprint arXiv:2203.08512*.
- Zhi Zhang, Qizhe Zhang, Zijun Gao, Renrui Zhang, Ekaterina Shutova, Shiji Zhou, and Shanghang Zhang. 2024. Gradient-based parameter selection for efficient fine-tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 28566–28577.
- Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. 2023. Can we edit factual knowledge by in-context learning? *arXiv preprint arXiv:2305.12740*.
- Ruiqi Zhong, Kristy Lee, Zheng Zhang, and Dan Klein. 2021. Adapting language models for zero-shot learning by meta-tuning on dataset and prompt collections. *arXiv preprint arXiv:2104.04670*.
- Zexuan Zhong, Zhengxuan Wu, Christopher D Manning, Christopher Potts, and Danqi Chen. 2023. Mquake: Assessing knowledge editing in language models via multi-hop questions. *arXiv preprint arXiv:2305.14795*.