Out-of-Distribution Detection with Attention Head Masking for Multimodal Document Classification

Anonymous ACL submission

Abstract

Detecting out-of-distribution (OOD) data is crucial in machine learning applications to mitigate the risk of model overconfidence, thereby 004 enhancing the reliability and safety of deployed systems. The majority of existing OOD detection methods predominantly address uni-modal 007 inputs, such as images or texts. In the context of multi-modal documents, there is a notable lack of extensive research on the performance of these methods, which have primarily been developed with a focus on computer vision tasks. We propose a novel methodology termed as attention head masking (AHM) for multi-modal OOD tasks in document classification systems. Our empirical results demonstrate that the pro-015 posed AHM method outperforms all state-of-017 the-art approaches and significantly decreases the false positive rate (FPR) compared to existing solutions up to 7.5%. This methodology 019 generalizes well to multi-modal data, such as documents, where visual and textual information are modeled under the same Transformer architecture. To address the scarcity of highquality publicly available document datasets and encourage further research on OOD detection for documents, we introduce FinanceDocs, 027 a new document AI dataset. Our code¹ and dataset² are publicly available.

1 Introduction

029

034

037

Out-of-distribution (OOD) detection presents a significant challenge in the field of document classification. When a classifier is deployed, it may encounter types of documents that were not included in the training dataset. This can lead to mishandling of such documents, causing additional complications in a production environment. Effective OOD detection facilitates the identification

²https://drive.google.com/drive/folders/ 1dV9obe_3hTsDoWJyYuNLBAXEiwOPwCw7



Figure 1: Visual demonstration of AHM on a transformer-based model: For each attention layer, we utilize the corresponding attention head mask from the AHM matrix. Following query-key multiplication and the subsequent softmax operation, the resulting attention scores undergo element-wise multiplication with the relevant attention head mask. This process effectively reduces the attention scores of certain heads to zero, thereby inhibiting the propagation of their respective information through the value matrix.

of unfamiliar documents, enabling the system to manage them appropriately which allows the classifier to maintain its reliability and accuracy in real-world applications. This has heightened the focus on OOD detection, where the primary objective is to determine if a new document belongs to a known in-distribution (ID) class or an OOD class. A significant challenge lies in the lack of supervisory signals from the unknown OOD data, which can encompass any content outside the ID classes. The complexity of this problem increases with the semantic similarity between the OOD and ID data (Fort et al., 2021).

A number of approaches have been developed to differentiate OOD data from ID data, broadly classified into three categories: (i) confidence-based methods, which focus on softmax confidence scores (Liu et al., 2020; Hendrycks and Gimpel, 2016;

¹https://anonymous.4open.science/r/ OOD-AHM-FE25/README.md

Hendrycks et al., 2019; Huang et al., 2021; Liang 056 et al., 2017), (ii) features/logits-based methods, 057 which emphasize logit outputs Sun and Li (2021); Sun et al. (2021); Wang et al. (2022); Djurisic et al. (2023), and (iii) distance/density-based methods, which concentrate on dense embeddings from the 061 final layers (Ming et al., 2023; Lee et al., 2018; 062 Sun et al., 2022). Recent research also investigates domain-invariant representations, such as HYPO (Ming et al., 2024), and introduces new OOD metrics like NECO (Ammar et al., 2024), which leverage neural collapse properties (Papyan et al., 2020). 067 Confidence-based methods can be unreliable as they often yield overconfident scores for OOD data. Features/logits-based methods attempt to combine class-agnostic scores from the feature space with the ID class-dependent logits. Our approach focuses on identifying more robust class-agnostic scores from the feature space, and as such, we conduct our experiments using distance/density-based methods.

076

077

084

880

090

094

096

100

101

102

103

Many OOD detection techniques have been developed, but most have been evaluated only in unimodal systems, such as text or images, and not extensively tested in the document domain (Gu et al., 2023). This may be due to the lack of highquality public document datasets, mostly based on IIT-CDIP (et al., 2006). To address the lack of comprehensive research in the document domain, we introduce a new document AI dataset, Finance-Docs. Additionally, we propose a novel technique called attention head masking (AHM) to effectively improve feature representations for distinguishing between ID and OOD data. Our method is illustrated in Figure 1. Our contributions can be summarized as follows: (1) FinanceDocs Dataset: We introduce FinanceDocs, the first high-quality digital document dataset for OOD detection with multimodal documents, offering digital PDFs instead of low-quality scans. (2) AHM: We propose a multihead attention masking mechanism for transformerbased models applied post-fine-tuning. By identifying masks that enhance similarity between ID training and evaluation features, we generate robust representations that improve the separation of ID and OOD data using distance/density-based OOD techniques. Our AHM method surpasses existing OOD solutions on key metrics.

2 Related Work

Learning embedding representations that general-ize effectively and facilitate better differentiation

between ID and OOD data is a well-recognized challenge in the field of machine learning (Zhou et al., 2023). To tackle this challenge, various studies have focused on specialized learning frameworks aimed at optimizing intra-class compactness and inter-class separation (Ye et al., 2021). Building on the principles of contrastive representation learning, researchers such as Chen et al. (2020) and Li et al. (2021) introduced prototypical learning (PL). This approach leverages prototypes derived from offline clustering algorithms to enhance unsupervised representation learning. Furthermore, Ming et al. (2024) integrated PL into their OOD learning framework, HYPO, achieving effective separation between ID and OOD data. This line of research was further advanced by Lu et al. (2024), who introduced the concept of multiple prototypes per cluster and employed a maximum likelihood estimation (MLE) loss to ensure that sample embeddings closely align with their corresponding prototypes. Additionally, approaches such as VOS (Du et al., 2022) and NPOS (Tao et al., 2023) have focused on regularizing the decision boundary between ID and OOD data by generating synthetic OOD samples, while Lin and Gu (2023) utilized open-source data as an OOD signal.

106

107

108

109

110

111

112

113

114

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

In our proposed methodology, we similarly aim to enhance the distinction between ID and OOD data through improved embedding representations. However, unlike previous studies that explore customized learning frameworks diverging from the standard cross-entropy loss, we concentrate on feature regularization during inference using our proposed attention head masking methodology. Our approach deliberately avoids altering the network's training procedure, thereby mitigating potential negative impacts on performance and preventing increased training costs. By focusing on inference rather than training modifications, our method ensures robust and cost-effective OOD detection.

Other inference-based methods, such as Avg-Avg (Chen et al., 2022) and Gnome (Chen et al., 2023), have also sought to enhance OOD detection through innovative techniques. Avg-Avg operates by averaging embeddings across both sequence length and different layers of a fine-tuned model, while Gnome combines embeddings from both a pre-trained and a fine-tuned model. These approaches, like our own, emphasize the importance of embedding manipulation during inference to achieve improved OOD detection without modifying the underlying training framework.

3 Method

158

159

160

161

162

163

164

165

166

168

169

170

171

The proposed AHM method, focuses on the feature extraction mechanisms inherent in transformer models, specifically the self-attention mechanism (Vaswani et al., 2017). Based on the premise that OOD data exhibit less semantic similarity to ID data, our goal is to generate embedding features that enhance the separation between ID and OOD data. The embeddings are then used in distance or density-based OOD detection methods, such as the Mahalanobis (Lee et al., 2018) or kNN+ (Sun et al., 2022). Our method is provided in Algorithm 1 (cf. Appendix A.1 for the theoretical framework) and the masking step is summarised in Figure 1.

> Algorithm 1 Optimization of Transformer-based Model using Attention Head Masking for OOD Detection – cf. Appendix A.2 for more details

- 1: **Input:** Budget T, model weights $W_{\text{pretrained}}$, percentage masking p, neighbors K, layers N, attention heads H, top attention head matrices to select F
- 2: Output: Optimal ensemble embedding
- 3: **1. Fine-tune Model** $W_{\text{pretrained}} \rightarrow W_{\text{finetuned}}$
- 4: for trial = 1 to T do
- 5: 2. InitializeAttention Head Matrix
- 6: Create $N \times H$ matrix A, A[i, j] = 1
- 7: **3. Mask Attention Heads**
- 8: Randomly set elements of A[i, j] to 0
- 9: **4. Extract Embeddings**
- 10: Extract embed_{train} $\in \mathbb{R}^{O \times Hid}$ and embed_{eval} $\in \mathbb{R}^{Q \times Hid}$
- 11: **5. Compute Similarity Scores**
- 12: For $e_i \in \text{embed}_{\text{eval}}$, get K nearest neighbors in embed_{train} and compute mean score S_i
- 13: 6. Assign and Collect Scores
- 14: Average similarity score: $\frac{1}{Q} \sum_{i=1}^{Q} S_i$. Collect scores S_i and their respective A[i, j]
- 15: end for
- 16: 7. Select Top Scores
- 17: Sort scores S_i , select top F masks A[i, j]
- 18: 8. Ensemble Embedding Generation
- 19: Use top F masks A[i, j] to generate and average embeddings for OOD detection

4 Results and Discussion

173 4.1 Datasets

We utilized two datasets in our experiments: Tobacco3482 and FinanceDocs. The Tobacco3482 dataset (Kumar et al., 2014) comprises 10 classes: Memo (619), Email (593), Letter (565), Form (372), Report (261), Scientific (255), Note (189), News (169), Advertisement (162), and Resume (120). As a subset of IIT-CDIP (et al., 2006), it was further processed to remove blank and rotated pages, preserving the rich textual and image modalities essential for a multi-modal system. Despite these efforts, some instances exhibit poor OCR quality due to the low-quality scans. 176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

224

225

226

We present FinanceDocs (cf. Appendix A.5 for per-category details and A.6 for dataset samples), a newly created dataset comprising 10 classes derived from open-source financial documents, including SEC Form 13 (663), Financial Information (360), Resumes (287), Scientific AI Papers (267), Shareholder Letters (256), List of Directors (188), Company 10-K Forms (181), Articles of Association (176), SEC Letters (141), and SEC Forms (121). Unlike Tobacco3482, FinanceDocs consists of high-quality digital PDFs (Annual Reports; SEC EDGAR Database; Companies House Service; ACL Anthology; Resume Dataset). The FinanceDocs dataset was labeled through the following process: a PDF parsing package (PyPDF2) was used to extract content from the original PDF documents. Each page was then visualized individually by a human annotator, who determined the relevance of the page to the collected classes and assigned the appropriate class label (cf. Appendix A.4 for annotator training and validation).

4.2 Experimental Setup

We employ two widely recognized OOD metrics to assess the performance of our proposed AHM method in comparison to other OOD benchmarks (Yang et al., 2024): AUROC, which measures the area under the ROC curve (higher values indicate better performance), and FPR, the false positive rate at a 95% true positive rate. A higher AUROC signifies better discrimination, while a lower FPR indicates greater robustness in rejecting OOD data.

For our experiments, we utilize LayoutLMv3 (Huang et al., 2022), a transformer-based multimodal model with 125.92 million parameters. We conduct both cross-dataset and intra-dataset OOD experiments. In cross-dataset OOD, the model is trained on the classes of one dataset and evaluated on the entirety of the other dataset as OOD. In intradataset OOD, one of the 10 classes is designated as OOD, and the model is trained on the remaining 9 classes, with the ID data split into training and

Method	Tobacco3482 AUROC	(ADVE OOD) FPR	Tobacco3482 (Cross-dataset OOD) AUROC FPR		FinanceDocs (Resume OOD) AUROC FPR		FinanceDocs (C AUROC	ross-dataset OOD) FPR
energy	0.951 ± 0.012	0.267 ± 0.057	0.944 ± 0.014	0.157 ± 0.042	0.848 ± 0.093	0.413 ± 0.218	0.846 ± 0.016	0.567 ± 0.039
gradNorm	0.940 ± 0.025	0.330 ± 0.116	0.824 ± 0.040	0.410 ± 0.094	0.742 ± 0.153	0.664 ± 0.251	0.724 ± 0.128	0.817 ± 0.145
kl	0.914 ± 0.016	0.448 ± 0.099	0.970 ± 0.014	0.071 ± 0.035	0.902 ± 0.040	0.295 ± 0.106	0.840 ± 0.025	0.630 ± 0.047
knn	0.958 ± 0.011	0.269 ± 0.074	0.991 ± 0.004	0.030 ± 0.018	0.965 ± 0.023	0.172 ± 0.127	0.891 ± 0.017	0.589 ± 0.067
Mahalanobis	0.976 ± 0.009	0.155 ± 0.053	0.996 ± 0.002	0.010 ± 0.009	0.977 ± 0.013	0.122 ± 0.100	0.898 ± 0.017	0.541 ± 0.090
mah _{AvgAvg}	0.942 ± 0.008	0.375 ± 0.054	0.997 ± 0.001	0.0004 ± 0.0005	0.996 ± 0.003	0.006 ± 0.005	0.949 ± 0.015	0.353 ± 0.196
mahGnome	0.971 ± 0.009	0.155 ± 0.054	0.992 ± 0.003	0.037 ± 0.016	0.938 ± 0.035	0.314 ± 0.165	0.822 ± 0.024	0.646 ± 0.114
maxLogit	0.946 ± 0.012	0.311 ± 0.063	0.945 ± 0.013	0.151 ± 0.033	0.851 ± 0.086	0.410 ± 0.203	0.846 ± 0.017	0.584 ± 0.037
msp	0.929 ± 0.009	0.471 ± 0.103	0.952 ± 0.016	0.140 ± 0.050	0.883 ± 0.041	0.400 ± 0.142	0.846 ± 0.032	0.612 ± 0.048
neco	0.971 ± 0.012	0.164 ± 0.046	0.995 ± 0.002	0.013 ± 0.011	0.975 ± 0.012	0.132 ± 0.096	0.888 ± 0.020	0.546 ± 0.114
residual	0.976 ± 0.008	0.149 ± 0.051	0.996 ± 0.002	0.011 ± 0.009	0.976 ± 0.014	0.130 ± 0.106	0.896 ± 0.016	0.541 ± 0.089
vim	0.976 ± 0.008	0.147 ± 0.044	0.996 ± 0.002	0.011 ± 0.009	0.976 ± 0.014	0.125 ± 0.101	0.899 ± 0.015	0.537 ± 0.086
knn _{AHM} mah _{AHM} mah _{AvgAvg_AHM}	$\begin{array}{c} \textbf{0.969} \pm \textbf{0.009} \\ \textbf{0.985} \pm \textbf{0.005} \\ \textbf{0.956} \pm \textbf{0.007} \end{array}$	$\begin{array}{c} 0.182 \pm 0.039 \\ 0.071 \pm 0.041 \\ 0.267 \pm 0.007 \end{array}$	$\begin{array}{c} \textbf{0.991} \pm \textbf{0.003} \\ \textbf{0.997} \pm \textbf{0.002} \\ \textbf{0.998} \pm \textbf{0.001} \end{array}$	$\begin{array}{c} \textbf{0.024} \pm \textbf{0.013} \\ \textbf{0.006} \pm \textbf{0.006} \\ \textbf{0.0001} \pm \textbf{0.0009} \end{array}$		$\begin{array}{c} 0.114 \pm 0.088 \\ 0.099 \pm 0.086 \\ 0.004 \pm 0.003 \end{array}$	$\begin{array}{c} 0.885 \pm 0.011 \\ 0.892 \pm 0.013 \\ 0.951 \pm 0.012 \end{array}$	$\begin{array}{c} 0.562 \pm 0.096 \\ 0.522 \pm 0.126 \\ 0.302 \pm 0.012 \end{array}$

Table 1: Performance metrics (arithmetic mean and standard deviation) for different methods across two datasets with intra-dataset and cross-dataset experiments configurations per dataset using AUROC (higher is better) and FPR (lower is better) – (cf. Appendix A.3 for hyperparameter tuning details).

evaluation sets. We select Advertisement (ADVE) and Resumes as the OOD classes for Tobacco3482 and FinanceDocs, respectively.

The models are trained over 5 random runs, with checkpoints saved at high ID classification metrics. Checkpoints with low silhouette scores $s(i) = \frac{b(i)-a(i)}{\max(a(i),b(i))}$ are filtered out to optimize intra-class similarity and inter-class separation. Our experiments were conducted using a single NVIDIA A100 GPU (80GB) for 72 GPU compute hours. We trained the models for a maximum of 15 epochs with an initial learning rate of 5×10^{-5} .

4.3 Current Benchmarks

We evaluated the peformance of various OOD detection methods, comparing them with our proposed methods, **knn**_{AHM}, **mah**_{AHM}, and **mah**_{AvgAvg_AHM}, which apply k-Nearest Neighbor (kNN) and Mahalanobis methods to dense embeddings generated by AHM. **mah**_{AvgAvg_AHM} is similar to **mah**_{AHM} but uses the AvgAvg embedding aggregation method (Chen et al., 2022).

As shown in Table 1, for the Tobacco3482 dataset with ADVE as the OOD class, our proposed mah_{AHM} outperformed other methods, achieving an AUROC of 0.985 and an FPR of 0.071. The high AUROC indicates that our method significantly enhances the Mahalanobis distance-based approach in distinguishing between ID and OOD samples. The notably lower FPR compared to previous methods like *vim* and *residual* (FPRs of 0.147 and 0.149, respectively) demonstrates the robustness of **mah**_{AHM} in correctly rejecting OOD samples.

For the FinanceDocs dataset, with Resumes as the OOD class, both \mathbf{knn}_{AHM} and \mathbf{mah}_{AHM} achieved superior performance, with AUROCs of 0.975 and 0.978, and FPRs of 0.114 and 0.099, respectively. Our **mah**_{AvgAvg_AHM} method also improved performance over mah_{AvgAvg}, highlighting the effectiveness of our approach in creating more separable embeddings between ID and OOD data. This is further evidenced by cross-dataset results in Table 1, where **mah**_{AvgAvg_AHM} consistently outperformed mah_{AvgAvg}, notably reducing the FPR by 5% on FinanceDocs and achieving an AUROC of 0.99 with an FPR of 0.0001 on Tobacco3482. This performance surpasses the respective method mah_{AvgAvg} without AHM applied. In fact, across all methods tested mah_{AvgAvg}, Mahalanobis and knn, the application of our AHM technique consistently resulted in improved performance.

263

264

265

266

268

269

271

272

273

274

275

276

277

278

279

281

282

283

284

285

286

287

289

290

291

292

293

294

296

298

Overall, the AHM technique significantly enhances the performance of kNN, Mahalanobis, and mah_{AvgAvg}, resulting in superior outcomes for **knn_{AHM}**, **mah_{AHM}**, and **mah_{AvgAvg_AHM}**, as evidenced by higher AUROCs and lower FPRs across intra-dataset and cross-dataset experiments, demonstrating strong generalizability across diverse datasets and methods.

5 Conclusion

In this study, we present the AHM technique for OOD detection in transformer-based document classification. Our methods, **knn**_{AHM}, **mah**_{AHM} and **mah**_{AvgAvg_AHM}, demonstrated significant improvements in AUROC and FPR metrics across various datasets. These results underscore the effectiveness of optimizing attention mechanisms to enhance feature separation between ID and OOD data. Additionally, we introduce the FinanceDocs dataset, contributing valuable resources to OOD detection research. Our findings highlight AHM as a promising approach for achieving robust and accurate OOD detection in document classification.

227

229

235

6 Limitations

299

311

314

316

319

320

321

322

323

324

325

329

331

332

333

334

336

341

343

344

345

347

While AHM techniques significantly reduced FPR in most cases, the improvements were marginal in cross-dataset scenarios where the Tobacco dataset served as the OOD data. This suggests a potential dependency on specific datasets. Additionally, AHM is a technique limited to attention-based DNN architectures that employ multi-head selfattention. Future research should aim to broaden the range of datasets explored.

References

- ACL Anthology. ACL Anthology. https:// aclanthology.org/. Accessed: 15 June 2024.
- Mouïn Ben Ammar, Nacim Belkhir, Sebastian Popescu, Antoine Manzanera, and Gianni Franchi. 2024. Neco: Neural collapse based out-of-distribution detection. *Preprint*, arXiv:2310.06823.
- Annual Reports. Annual Reports. https:// www.annualreports.com/Browse/Industry. Accessed: 15 June 2024.
- Sishuo Chen, Xiaohan Bi, Rundong Gao, and Xu Sun. 2022. Holistic sentence embeddings for better out-ofdistribution detection. *Preprint*, arXiv:2210.07485.
- Sishuo Chen, Wenkai Yang, Xiaohan Bi, and Xu Sun. 2023. Fine-tuning deteriorates general textual outof-distribution detection by distorting task-agnostic features. *Preprint*, arXiv:2301.12715.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 1597–1607. PMLR.
- Companies House Service. Companies House Service. https://find-and-update. company-information.service.gov.uk/. Accessed: 15 June 2024.
- Andrija Djurisic, Nebojsa Bozanic, Arjun Ashok, and Rosanne Liu. 2023. Extremely simple activation shaping for out-of-distribution detection.
- Xuefeng Du, Zhaoning Wang, Mu Cai, and Yixuan Li. 2022. Vos: Learning what you don't know by virtual outlier synthesis. *Preprint*, arXiv:2202.01197.
- D. Lewis et al. 2006. Building a test collection for complex document information processing.
- Stanislav Fort, Jie Ren, and Balaji Lakshminarayanan. 2021. Exploring the limits of out-of-distribution detection. *Preprint*, arXiv:2106.03004.

Jiuxiang Gu, Yifei Ming, Yi Zhou, Jason Kuen, Vlad Morariu, Handong Zhao, Ruiyi Zhang, Nikolaos Barmpalios, Anqi Liu, Yixuan Li, Tong Sun, and Ani Nenkova. 2023. A critical analysis of document out-of-distribution detection. In *Findings of the Association for Computational Linguistics: EMNLP* 2023, pages 4973–4999, Singapore. Association for Computational Linguistics.

348

351

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

373

374

375

376

378

379

381

382

384

386

390

391

392

393

394

396

397

398

399

400

- Dan Hendrycks, Steven Basart, Mantas Mazeika, Mohammadreza Mostajabi, Jacob Steinhardt, and Dawn Song. 2019. A benchmark for anomaly segmentation. *CoRR*, abs/1911.11132.
- Dan Hendrycks and Kevin Gimpel. 2016. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *CoRR*, abs/1610.02136.
- Rui Huang, Andrew Geng, and Yixuan Li. 2021. On the importance of gradients for detecting distributional shifts in the wild. *CoRR*, abs/2110.00218.
- Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and Furu Wei. 2022. Layoutlmv3: Pre-training for document ai with unified text and image masking. *Preprint*, arXiv:2204.08387.
- Jayant Kumar, Peng Ye, and David Doermann. 2014. Structural similarity for document image classification and retrieval. *Pattern Recognition Letters*, 43:119–126. ICPR2012 Awarded Papers.
- Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. 2018. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. *Preprint*, arXiv:1807.03888.
- Junnan Li, Pan Zhou, Caiming Xiong, and Steven C. H. Hoi. 2021. Prototypical contrastive learning of unsupervised representations. *Preprint*, arXiv:2005.04966.
- Shiyu Liang, Yixuan Li, and R. Srikant. 2017. Principled detection of out-of-distribution examples in neural networks. *CoRR*, abs/1706.02690.
- Haowei Lin and Yuntian Gu. 2023. Flats: Principled out-of-distribution detection with feature-based like-lihood ratio score. *Preprint*, arXiv:2310.05083.
- Weitang Liu, Xiaoyun Wang, John D. Owens, and Yixuan Li. 2020. Energy-based out-of-distribution detection. *CoRR*, abs/2010.03759.
- Haodong Lu, Dong Gong, Shuo Wang, Jason Xue, Lina Yao, and Kristen Moore. 2024. Learning with mixture of prototypes for out-of-distribution detection. *Preprint*, arXiv:2402.02653.
- Yifei Ming, Haoyue Bai, Julian Katz-Samuels, and Yixuan Li. 2024. Hypo: Hyperspherical out-of-distribution generalization. *Preprint*, arXiv:2402.07785.
- Yifei Ming, Yiyou Sun, Ousmane Dia, and Yixuan Li. 2023. How to exploit hyperspherical embeddings for out-of-distribution detection?

Vardan Papyan, X. Y. Han, and David L. Donoho. 2020. Prevalence of neural collapse during the terminal phase of deep learning training. Proceedings of the National Academy of Sciences, 117(40):24652–24663.

402

403

404 405

406

407

408

409

410

411

412

413

414 415

416

417

418

419

420

421

422 423

424

425

426

427

428

429

430

431

432 433

434

435

436

437

- Resume Dataset. Resume Dataset. https: //www.kaggle.com/datasets/snehaanbhawal/ resume-dataset. Accessed: 15 June 2024.
- SEC EDGAR Database. SEC EDGAR Database. https://www.sec.gov/edgar/search/. Accessed: 15 June 2024.
 - Yiyou Sun, Chuan Guo, and Yixuan Li. 2021. React: Out-of-distribution detection with rectified activations. *CoRR*, abs/2111.12797.
 - Yiyou Sun and Yixuan Li. 2021. On the effectiveness of sparsification for detecting the deep unknowns. *CoRR*, abs/2111.09805.
 - Yiyou Sun, Yifei Ming, Xiaojin Zhu, and Yixuan Li. 2022. Out-of-distribution detection with deep nearest neighbors. *Preprint*, arXiv:2204.06507.
 - Leitian Tao, Xuefeng Du, Xiaojin Zhu, and Yixuan Li. 2023. Non-parametric outlier synthesis. *Preprint*, arXiv:2303.02966.
 - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *CoRR*, abs/1706.03762.
 - Haoqi Wang, Zhizhong Li, Litong Feng, and Wayne Zhang. 2022. Vim: Out-of-distribution with virtuallogit matching.
 - Jingkang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. 2024. Generalized out-of-distribution detection: A survey. *Preprint*, arXiv:2110.11334.
- Haotian Ye, Chuanlong Xie, Tianle Cai, Ruichen Li, Zhenguo Li, and Liwei Wang. 2021. Towards a theoretical framework of out-of-distribution generalization. *Preprint*, arXiv:2106.04496.
- Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. 2023. Domain generalization: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4):4396–4415.

443

Α

A.1

Appendix

concepts presented in this work.

Framework

form of dataset examples, implementation details,

etc. to bolster the reader's understanding of the

The central hypothesis underlying the proposed

solution is predicated on the assumption that ID

data should exhibit greater similarity in their fea-

ture representations when compared to OOD data.

Consequently, we posit that when considering a

pair of data points from two similar ID classes (de-

noted as Pair A) and a pair consisting of one ID and

one OOD data point (denoted as Pair B), the ap-

plication of a masking procedure on input features

(whether textual or visual) would result in a more

pronounced divergence in the feature space for Pair

B as compared to Pair A. Initial experiments were

conducted with random masking of input features.

For textual data, this involved replacing tokens ran-

domly with the '[MASK]' token. For visual data,

random image patches were set to zero, effectively

splitting the image into patches and nullifying se-

lected segments. These preliminary experiments

revealed two critical factors influencing the final

feature embeddings used in distance-based OOD

detection methods, such as the Mahalanobis dis-

tance: (a) the input tensors provided to the model,

and (b) the feature extraction mechanism employed

by the model, specifically the attention mechanism.

cused on input masking, achieving a consistent

masking strategy proved challenging. While a con-

sistent mask could be established for visual data by

dividing images into uniformly sized chunks and

consistently masking specific segments, such con-

sistency was elusive for textual features. The vari-

ability in sequence length across different tokens

complicated the masking process, often leading to

strategies that involved masking padding tokens

In light of these challenges, our focus shifted

from input masking to the feature extraction pro-

cess itself, particularly the attention mechanism

within the model. We discovered that consistent

masking could be achieved by selectively mask-

ing attention heads within different layers of the

encoder. These heads are responsible for learning

different representations and capturing different as-

rather than meaningful data.

Although the early experiments primarily fo-

Proposed Methodology and Theoretical

- 444 445
- 446 447

448

449

- 450 451
- 452

454

455

459

453

456 457 458

468

472

473

462 463

460 461

464

466

467

469 470 471

> 474 475

476 477

479 480 481

478

482 483

484 485 486

487 488 489

> 490 491 492

This section provides supplementary material in the

pects of the input sequence. Hence by shutting down heads we are effectively deactivating certain pattern-extracting mechanisms within the attention architecture.

A.2 Description of Algorithm 1

As detailed in Algorithm 1, we begin with a finetuned model and proceed by randomly initializing various attention head masks based on a masking hyperparameter p. This hyperparameter represents the percentage of attention heads H set to zero within each attention layer N of the model. For each random mask, we extract dense hidden representations from both the training and evaluation datasets. The objective is to identify which of these randomly generated attention head masks minimizes the divergence between the representations of the evaluation and training data in the feature space. This is accomplished by calculating the average similarity score among the top Knearest neighbors for each evaluation data point.

The attention head masks are then ranked based on these aggregated similarity scores. Finally, we select the top F masks with the highest similarity scores between the evaluation and training data and use them to generate new feature representations. These features are then ensembled (i.e., averaged) and subsequently utilized in a distance-based OOD detection method, such as the Mahalanobis distance.

A.3 Hyperparameter Tuning

Table 2 summarizes the hyperparameters for model training. The model was trained using a carefully selected set of hyperparameters to optimize its performance. The training batch size per device was set to 32, while the evaluation batch size was configured at 8, ensuring efficient computation throughout the process. To stabilize updates, gradient accumulation was performed over 8 steps. The learning rate was set at 5×10^{-5} , with no weight decay applied, to prevent the risk of overfitting.

The Adam optimizer was configured with parameters $\beta_1 = 0.9$, $\beta_2 = 0.999$, and an epsilon value of 1×10^{-8} to ensure effective convergence. To maintain stability during training, the maximum gradient norm was capped at 1.0. The model underwent training for 65 epochs, with evaluations delayed by 5 steps to monitor progress at appropriate intervals, allowing for a well-tuned and stable learning process.

The hyperparameters chosen for the proposed

494 495 496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

493

519 520 521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

568

569

543AHM method are presented in Table 3. Follow-544ing the procedure outlined in Algorithm 1, an ex-545ploration budget of 25 was allocated for potential546AHM configurations. To assess the effectiveness547of different configurations, masking percentages of5480.1 and 0.2 were applied during the process.

549

550

551

554

555 556

557

558

559

561

563

564

567

To ensure robust performance, similarity scores between ID validation data and ID training data were computed. These scores were determined by averaging the similarity of the top 10 nearest neighbors for each validation data point. Using these similarity scores, the top five AHM heads were selected to generate the final representation embeddings, which were then combined through an ensemble approach to enhance the overall model performance.

Table 2: Hyperpara	neters for	model	training.
--------------------	------------	-------	-----------

Hyperparameter	Value
per_device_train_batch_size	32
per_device_eval_batch_size	8
gradient_accumulation_steps	8
eval_delay	5
learning_rate	5e-05
weight_decay	0.0
adam_beta1	0.9
adam_beta2	0.999
adam_epsilon	1e-08
max_grad_norm	1.0
num_train_epochs	65

Table 3: Hyperparameters for A	HM	I.
--------------------------------	----	----

Hyperparameter	Value
Exploration budget (T)	25
Percentage masking (p)	[0.1, 0.2]
Neighbors (K)	10
Top AHM matrices select (F)	5

A.4 Annotator Training and Validation

To maintain high-quality annotation in line with ethical standards, we enlisted three postgraduate students fluent in English. They received instruction and participated in sessions with finance professionals to address any task-related questions. The annotation process spanned about four months, involving 90 training sessions, with breaks scheduled every 45 minutes. The students were compensated through gift vouchers and honorariums per minimum wage requirements³.

A.5 Dataset description of FinanceDocs

The FinanceDocs dataset comprises a diverse collection of financial and legal documents sourced from various reliable platforms, offering a comprehensive view of corporate disclosures, shareholder communications, and regulatory filings. Each document type serves a distinct purpose, providing insights into different aspects of corporate governance, financial performance, and regulatory compliance, as detailed below:

- **SEC form documents:** These documents were collected from the Securities Exchange Commission (SEC) website. These forms are statements of changes in beneficial ownership.
- Shareholder letter documents: These documents were collected from annual reports. A shareholder letter in an annual report provides a summary of the company's financial performance, highlighting key achievements, strategic initiatives, and market conditions over the past year. It offers leadership's perspective on successes and challenges while outlining future goals and potential risks. The letter also emphasizes the company's commitment to corporate governance, social responsibility, and long-term growth.
- **SEC letter documents:** These documents were collected from the SEC website. These are letters from companies to the SEC about various company disclosures.
- SEC-13 form documents: These documents were collected from the SEC website. These forms disclose significant information about an entity's ownership or control over securities, typically required for investors with large holdings.
- **10k form documents:** These documents were collected from annual reports. These represent the 10k forms of an annual report
- Financial info documents: These documents were collected from annualreports (Annual Reports). They consist of various financial information, including the income statement,

³https://www.minimum-wage.org/international/ united-states

613balance sheet, and cash flow statement, which614detail the company's revenue, expenses, as-615sets, liabilities, and cash movements. It also616includes financial ratios and metrics to assess617profitability, liquidity, and leverage.

- Articles of scientific paper documents:
 These documents were collected from ACL
 Anthology⁴. It is a comprehensive digital
 archive of research papers in computational
 linguistics and natural language processing,
 published by the Association for Computational Linguistics.
 - Articles of resume documents: These documents were collected from Kaggle. They represent resumes from different occupations.

625

626

627

628 629

630

633

634

- Articles of Association documents: These documents were collected from Companies House Services UK. They represent documents relating to articles of association of a company. These involve information such as directors powers and responsibilities, interpretation and limitation of liability as well as distribution of shares.
- Director documents: These documents were
 collected from annual reports and Companies
 House Services UK⁵. It involves information
 about the directors of a company.

⁴https://aclanthology.org/

⁵https://www.gov.uk/government/organisations/ companies-house

A.6 Dataset examples of FinanceDocs

Presented below are examples from each document category included in FinanceDocs, providing the reader with a comprehensive visual overview of the dataset.



Figure 2: Examples of SEC form documents.



Figure 3: Examples of shareholder letter documents.



Figure 4: Examples of SEC letter documents.

10

640 641

wight with an			0.000/2014.00:53 we protective suggestion VEPO/54000 See Conductor Sec. 24 (14), and Sec. 2010			1016203-010	9013905	036204.0.30	UNIXER IN REPORT OF THE OWNER AND A STOCK		
								CUSP No. 3496/E307			
CUSP No. 33302202	130	Page 8 of 9	CUSP No. 694103102	13.0	Page 7 of 12			1	NAME OF REPORTING PERSON		
	SIGNATURES		1. NAMES OF REPORTING FER	SONS LIES. IDENTIFICATION NOS. OF AIRAVE PERS	ONS (ENTITIES ONLY)				PITERA PER		
After resonable inquiry and to set forth in this Schedule 13D is true, o	the best of my low-vieige and belief, each of the undersigned certifies that to complete and convext.	the information	Value C. Kanan 2. CHEX THE APPROPRIATE DOX IF A MEMBER OF A GROUP SEE INSTRUCTIONS)			This Amendment No. 8 ("Amendment	EXPLANATORY NOTE No. F1 amends and supplements the Scholale 13D field with the Securities and Dathanae	2	CHECK THE APPROPRIATE BOX IF A MEMBER OF A GROUP DATE OF A GRO		
Data: May 10, 2024			61 C 3. SEC Une Only			Commission (the "Commission") on F No. 2 to Schedule 13D filed on August to Schedule 13D filed on Neurophyr 2	Verany I, 2022, Assentiment No. 1 in Schodule 13D filed on February I, 2022, Assentiment 7 23, 2022, Assentiment No. 3 to Schodule 13D filed on August 28, 2022, Assentiment No. 4 II. 2023, Assentiment No. 1 in Schodule 13D Filed on Dimension ID, 2023, Assentiment No. 6	,	SECUSE ONLY		
	MINURSET FAILURS INC.		4 SCHOL OF AND SELENTIALTION			to Schedule 13D filed on January 19 "Schedule 13D"). Each liters below Schedule 13D: Castadized sums use	2024 and Amendment No. 7 to Schedule 13D filed on March 11, 2024 toolocitoty, the amends and supplements the information disclored under the corresponding lism of the J her on defined herein shall have the reasoning antibuland to them in the Schedule 13D.	4	SOLICE OF FLNDS		
	By: <u>N. ELIAN FERNANDEZ NANCHEZ</u> Name: Elian Fernandez Sanchez	By: <u>N. ELLAR FERNANDEZ SANCHEZ</u> Name: Ellar Formadez Sanchez	3 CHECK BOX # DISCLOSURE 320	FOF LIGAL PROCEEDINGS IS REQUIRED PURSUAN	IT TO TITALS 201-OR D	Except is otherwise set firsh herein, this Amendment No. 8 does not modify any of the infer Reporting Persons in the Schedule 13D.	this Amendment No. 8 does not modify any of the information previously reported by the λ	5	CHECK BOX & DESCLOSURE OF LIGAL PROCEEDINGS IS REQUIRED FURSIONT : TO TEED SHE OF 2H		
	Title President		6 CITIZINSIPOR PLACE OF	OBGANIZATION		The purpose of this Schedule 110 ES WellEnterprises USA, LLC (collective	ig is to update the connecting by Damian Lancendein, HillClose Investment Fund, LLC, and ly the "Hisporting Persons") of the lower's Common Stock.		CITIZENSIEP OR PLACE OF ORGANIZATION		
	MINERSET INTERNATIONAL LTD		United Status 2 SOLES	LOTING POWER		Except as specifically arounded below,	all other provisions of the Schedule 13D romain in effect.		134		
	IN: A TUAS PERMANEZ SANCHEZ		NUMBER OF 2011			Box 2. Meetity and Background		NUMBER O SHARES	7 7 SOLE VOTING POWER		
	Name: Elias Forsandoz Sanchez Ticle: CDO	0.112	BENEFICIALLY DESERVICE A	D YOTING POWER		Item 2 of the Schohole 13D is benefy amended and supplemented by inverting the following paragraph before the first paragraph theored.			3 NUMED VOTING POWER		
	KOSMO INVESTMENTS PTE. LTD		REPORTING & SOLET PERSON WITH 204.01	DEPOSITIVE NYATR		On May 10, 2024, the Reporting Perspect of \$2,4899. The Reporting Person	on parchased 10;000 shares of Common Nitak of the leaser at a weighted average parchase a pold such consideration using personal funds.	REPORTIN PERSON WE	J T.MI.000 DI 9 SOLE DEPOSITIVE POWER		
	IV: 5: ELIAS FERNANDEZ SANCHEZ	2	IR. SILARS	D DEPOSITIVE POWER		On May 23, 2024, the Reporting Perso	n granted 600,000 shares of Rastriced Stock Units to be vested over the next 3 years.		11 SEASED DEPOSITIVE POASE		
	Name: Elius Fernander Sancher TRA: CDO		1. ANGENERAL CARACTERISTICATION OF A RAY OF DESCRIPTION STREAM OF THE ANGENERAL DESCRIPTION OF THE			Ben J. Source and Amount of Far	ab or Other Considerations		1,983,080		
						Item 3 of the Schedule 13D is hereby thereof	arceded and supplemented by inserting the following paragraph after the last paragraph	11	AGGREGATE AMOUNT RENEFICIALLY OWNED RY EACH REPORTING PERSON 7,995,008		
	8: ILLAS FERNANDEZ SANCHEZ Thin Fernandez Sancher, individually					The Stones and Amount of Funds or Other Consideration for the purchase of the shares of Common Stock on May 10, 2024 is set forth in lasm 2.			OBEX BEX IF THE AGGREGATE ANCENT IN ROW (11) EXCLUDES CHERAIN T SEARES		
			1.9%	al des la secolar d		Dem 5. Interest in Securities of the	honer	0	PERCENT OF CLASS REPRESENTED BY AMOUNT IN ROW OFF		
			14. TVPE OF BARKETSU PERCEN (see methodane) IN			liens Su) of the Schedule 13D are her	thy amended and supplemented by the following paragraphs:		5.0%		
						(i) Danies Lamendela		14	TYPE OF REPORTING PERSON		
			Percentage calculated based on 36, Trust, filed with the Securities and	503,158 units outstanding as of August 1, 2019, as reports Exchange Commission on August 1, 2019.	d in the 10-Q of Pacific Coast Of	As of May 13, 2024, Danie Issuer's Common Stock held which Mr. Lanundola is the 3 Issuer Common Dank (Id)	n Lamandula may be deemed to be the beneficial owner of: (1) 3,041,346 shares of the directly by BillCour Investment Fund, LLC (27,97) of the constanting Common Socie), of damages, and sees which is holds the voting and dispositive powers (29,931,847 shares of the tricking and sees which is holds the voting and dispositive powers (19,931,847 shares of the tricking and the solution of the tr		13		
						connect solvaidary of HERCon and he holds the voting and a of Commun. Block Restricted of which Nr. Laurandon is fl	where the solutioning operation is a second				
						bield directly by ERECON Evolutions Fluck LLC (2019) of the constanting Cosmon Reduct, of Hilds MM. Lamonds is the Hampung, and we will solve build have build over all and placeoptic spraces (LK32) shows of the have Cosmon Reduct models appends to activate of placeoptic sprace (LK32) shows of the have Cosmon Reduct models appends to activate of a placeoptic sprace (LK32) shows of the have Cosmon Reduct models appends to activate of a placeoptic sprace (LK32) shows of the have Cosmon Reduct models appends to activate of a placeoptic sprace (LK32) shows of the have Cosmon Reduct placeoptic sprace (LK32) shows the have Cosmon Reduct append and the have Cosmon Reduct appendix optime model to the Lamondol and an activation give of 46.45, and and 100 kg (2012) and cosmon Reduct appendix optime model to the Lamondol and an activation give of 46.45, and and 100 kg (2012) and cosmon Reduct appendix optime model to the Lamondol and an activation give of 46.45, and and 100 kg (2012) and cosmon Reduct appendix optime model to the Lamondol and an activation give of 46.45, and and 100 kg (2012) and cosmon Reduct appendix optime model to activate appendix optime appendix optime Reduct appendix optime Reduct appendix optime model to the Lamondol and an activation give of 46.45, and and 100 kg (2012) and cosmon Reduct appendix optime model to the Lamondol and an activation give of 46.45, and and 100 kg (2012) and cosmon Reduct appendix optime model to the Lamondol and an activation give of 46.45, and and 100 kg (2012) and cosmon Reduct appendix optime model to the Lamondol and an activate appendix optime Reduct appendix optime Reduct appendix optime Reduct appendix optime Reduct appendix optime Reduct appendix optime Reduct appendix optime Reduct appendix optime Reduct appendix optime Reduct appendi					

Figure 5: Examples of SEC-13 form documents.



<form>





Figure 6: Examples of 10k form documents.



Figure 7: Examples of financial info documents.



Figure 8: Examples of scientific paper documents.



Figure 9: Examples of resume documents.



Figure 10: Examples of Articles of Association documents.

Company Directo	v 1		Proposed Officers						803.089	
Type: Pub Forename(s):	Person MR TARCISIUS ZIEDOU D. DUITON	Company Director	1		Pursuant to the requirements of this report to be signed on its be	SHENATURES Section 13 or 15(d) of the Securities Enchange Act of 1934, half by the undersigned, thereanto doly authorized.	the registrant has duly caused	Person to be reprinted of the Sociality Art adhesized, in the Troom, Ontario, Canada, Spell	C POL as anomald, the Regiment has duly caused this report to be signed or include 17, 2023. EXERCISE GAME TREEWORDER, COMP.	nil hy fie undesigned, therease-bily
Service Address:	DLASAGU 71-75 SHELFON STREET COYENT GARDEN LONDON UNITED INFORMA WC2H NIQ UNITED INFORMA WC2H NIQ	Type: Per Full Forename(s): MR Surname: TAL	rsee R BAHIM JIRAT		Date: February 23, 2024	ABC FEST CORPORATION By: (c) July B. Midleenolds July B. Midleenolds Chairwan, Provident and C (Principal Executive Offic	Thief Essentive Officer 21)		 N. Michael Connella Name Michael Connella File: Insuis Out Double Office POPER OF ATTOENTY 	
Country-State Unsafty Resident: Date of Resiz. **/12/20	BENIN BU Katawalar BENINESE	Service Address: 7 Cl LOS UNI Country/State Unsafty MO Posidae:	TORONATION ROAD, BEPRINA ROUSE, LAUN INDON SITTED KINGBOM NW16 TPQ OROCCO	CHESE 408	Personant to the requirements of pursons on bulaif of the negistry <u>Ngastury</u> A Judy R. McReywolds	the Securities Exchange Act of 1554, this report has been a at and in the capacities and on the data: indicated. — <u>The</u> Chairman, President and Chief Executive Officer	igned below by the following two Following Following 23, 2024	Bernin et al a factor en al carda arrange in opecifica, to sign are stal al andra arrange in opecifica, to sign are stal al anormalization to the Exchange Commention, passing areas call allows alone in connection fearereds, as fally to all term agenes, of the or the absolution or adultance, the 1 Personnel to the explorements of the Kanashine Acti-	C true are particle totale tiggioure approximation for the community and approximations. A final approx of all advances and a result totales of the first and approximation of the particle of advances and a state of the first and approximation of a state particle of a state particle of the particle and approximation of advances and and approximation of advances and and approximation of advances and approximation of advances and approximation of advances and and approximation of advances and advances an	ments, place and could, in any and all dots frequeld, with the Societies and and fiting superior and microscopy in feet using all that said anteneys in fact and time and on the date or dates indicated.
Occupation: DIREC The suffscribers confir	TOR m that the person named has consented to act as a director.	Date of Birth: **/#22806 Occupation: CEO	Scienceley: MOROCCAN		July R. McRaynolds 5: J. Mathew Beasley J. Mathew Boosley	(Principal Executive Officer) Chief Financial Officer (Principal Financial Officer)	Fahrany 23, 2824	Signature 3. Match Cleanelle	The Inviso Out Exceptive Officer	Ber
		The subscribers confirm that	it the person named has consented to act as a c	l to act as a director:	N Jason T. Parks Jacon T. Parks N Subvatore A. Abbate Subvatore A. Abbate	Vier Periolest - Controller and Chief Accounting Officer (Principal Accounting Officer) Disorter	February 23, 2024	Mahda Caravita 'a Cato Apdi Cata Insi	Principal December 2010 or Index 10 Principal Transition Officer Principal Transition Officer and Principal Sciencesing Officery	April 17, 2028
					-sc Eduardo F. Consulte Eduardo F. Consulte	Discur	February 23, 2024	y Anter Good Anter Ander Maddi Manris y Paul Salvanar Fad Salvanar	Disor	April 11, 2021
					Norman A. Lanton Fradrik J. Eliason Michael P. Hogan Michael P. Hegan	Director	Edmany 22, 2021	v Storer A. Shiftons Saven A. Shiftons	bedar	April 17, 2025
					A Kathleen D. McFilgert Kathleen D. McFilgert	_ Director	February 23, 2024			
					Caug E. Platp A Steven L. Spinner Steven L. Spinner	Distar	February 23, 2024			
					. & Amice E. Scipp Junice E. Stipp	Discur	Fubmary 23, 2024			
Electronically filed docum	erd for Company Number: 14854689	Electronically filed document for C	Company Number: 180	27188		123				

Figure 11: Examples of list of director documents.