

40th AAAI Conference on Artificial Intelligence (AAAI-26), 38th Conference on Innovative Applications of Artificial Intelligence (IAAI-26), and 16th Symposium on Educational Advances in Artificial Intelligence (EAAI-26)

Volume 44: AAAI Special Track

- AI Alignment

Singapore
20-27 January 2026

ISBN: 979-8-3313-3418-5

Printed from e-media with permission by:

Curran Associates, Inc.
57 Morehouse Lane
Red Hook, NY 12571



Some format issues inherent in the e-media version may also appear in this print version.

Copyright© (2026) by Association for the Advancement of Artificial Intelligence
All rights reserved.

Printed with permission by Curran Associates, Inc. (2026)

For permission requests, please contact Association for the Advancement of Artificial Intelligence
at the address below.

Association for the Advancement of Artificial Intelligence
2275 East Bayshore Road
Suite 160
Palo Alto, California 94303
USA

Phone: 1-650-328-3123
Fax: 1-650-321-4457

<https://aaai.org/Press/press.php>

Additional copies of this publication are available from:

Curran Associates, Inc.
57 Morehouse Lane
Red Hook, NY 12571 USA
Phone: 845-758-0400
Fax: 845-758-2633
Email: curran@proceedings.com
Web: www.proceedings.com

TABLE OF CONTENTS

AAAI SPECIAL TRACK ON AI ALIGNMENT

AURA: Affordance-Understanding and Risk-Aware Alignment Technique for Large Language Models.....	37204
<i>Sayantan Adak, Pratyush Chatterjee, Somnath Banerjee, Rima Hazra, Somak Aditya, Animesh Mukherjee</i>	
Beyond Patches: Mining Interpretable Part-Prototypes for Explainable AI.....	37213
<i>Mahdi Alehdaghi, Rajarshi Bhattacharya, Pourya Shamsolmoali, Rafael M. O. Cruz, Eric Granger</i>	
Operationalizing Pluralistic Values in Large Language Model Alignment Reveals Trade-Offs in Safety, Inclusivity, and Model Behavior.....	37222
<i>Dalia Ali, Dora Zhao, Allison Koenecke, Orestis Papakyriakopoulos</i>	
MoralReason: Generalizable Moral Decision Alignment for LLM Agents Using Reasoning-Level Reinforcement Learning.....	37232
<i>Zhiyu An, Wan Du</i>	
DNR Bench: Benchmarking Over-Reasoning in Reasoning LLMs.....	37240
<i>Oluwanifemi Bamgbose, Masoud Hashemi, Sathwik Tejaswi Madhusudhan, Jishnu Sethumadhavan Nair, Aman Tiwari, Vikas Yadav</i>	
ALPHA: Action-Based Learning for Pluralistic Human Alignment in Large Language Models.....	37249
<i>Aanisha Bhattacharyya, Susmit Agrawal, Yaman Kumar Singla, Tarun Ram Menta, Nikitha Sr, Rajiv Ratn Shah, Changyou Chen, Balaji Krishnamurthy</i>	
A Course Correction in Steerability Evaluation: Revealing Miscalibration and Side Effects in LLMs.....	37259
<i>Trenton Chang, Tobias Schnabel, Adith Swaminathan, Jenna Wiens</i>	
MetaCipher: A Time-Persistent and Universal Multi-Agent Framework for Cipher-Based Jailbreak Attacks for LLMs.....	37268
<i>Boyuan Chen, Minghao Shao, Abdul Basit, Siddharth Garg, Muhammad Shafique</i>	
On the Exponential Convergence for Offline RLHF with Pairwise Comparisons.....	37277
<i>Zhirui Chen, Vincent Y. F. Tan</i>	
Preference Optimization Via Contrastive Divergence: Your Policy is Secretly an NLL Estimator.....	37286
<i>Zhuotong Chen, Fang Liu, Xuan Zhu, Haozhu Wang, Jiayu Li, Yanjun Qi, Mohammad Ghavamzadeh</i>	
Deep Hidden Cognition Facilitates Reliable Chain-Of-Thought Reasoning.....	37295
<i>Zijun Chen, Wenbo Hu, Richang Hong</i>	
MegaCoin: Enhancing Medium-Grained Color Perception for Vision-Language Models.....	37305
<i>Ming-Chang Chiu, Shicheng Wen, Pin-Yu Chen, Xuezhe Ma</i>	
TWINFUZZ: Dual-Model Fuzzing for Robustness Generalization in Deep Learning.....	37314
<i>Enze Dai, Wentao Mo, Kun Hu, Xiaogang Zhu, Xi Xiao, Sheng Wen, Shaohua Wang, Yang Xiang</i>	
A Multi-Agent Conversational Bandit Approach to Online Evaluation and Selection of User-Aligned LLM Responses.....	37323
<i>Xiangxiang Dai, Yuejin Xie, Maoli Liu, Xuchuang Wang, Zhuohua Li, Huanyu Wang, John C. S. Lui</i>	
Resilience in Ambient Multi-Agent LLMs Via Decentralized Bio-Autonomic Control and Immune-Inspired Anomaly Detection.....	37332
<i>Nastaran Darabi, Devashri Naik, Sina Tayebati, Dinithi Jayasuriya, Amit Ranjan Trivedi</i>	
AMaPO: Adaptive Margin-Attached Preference Optimization for Language Model Alignment.....	37341
<i>Ruibo Deng, Duanyu Feng, Wenqiang Lei</i>	

Democratizing Diplomacy: A Harness for Evaluating Any Large Language Model on Full-Press Diplomacy	37350
<i>Alexander Duffy, Samuel J Paech, Ishana Shastri, Elizabeth Karpinski, Baptiste Alloui-Cros, Tyler Marques, Matthew Lyle Olson</i>	
ACID Test: A Benchmark for Cultural Safety and Alignment in LALMs	37360
<i>Bikash Dutta, Adit Jain, Rishabh Ranjan, Mayank Vatsa, Richa Singh</i>	
Aligning Attention with Human Rationales for Self-Explaining Hate Speech Detection	37369
<i>Brage Eilertsen, Røskva Bjørgfinsdóttir, Francielle Vargas, Ali Ramezani-Kebrya</i>	
The Alignment Game: A Theory of Long-Horizon Alignment Through Recursive Curation	37379
<i>Ali Falahati, Mohammad Mohammadi Amiri, Kate Larson, Lukasz Golab</i>	
SMiLE: Provably Enforcing Global Relational Properties in Neural Networks	37387
<i>Matteo Francobaldi, Michele Lombardi, Andrea Lodi</i>	
EssayBench: Evaluating Large Language Models in Multi-Genre Chinese Essay Writing	37396
<i>Fan Gao, Dongyuan Li, Ding Xia, Fei Mi, Yasheng Wang, Lifeng Shang, Baojun Wang</i>	
Beyond Transcription: Mechanistic Interpretability in ASR	37407
<i>Neta Glazer, Yael Segal-Feldman, Hilit Segev, Aviv Shamsian, Asaf Buchnick, Gill Hetz, Ethan Fetaya, Joseph Keshet, Aviv Navon</i>	
AlignTree: Efficient Defense Against LLM Jailbreak Attacks	37417
<i>Gil Goren, Shahar Katz, Lior Wolf</i>	
Identifying Features Associated with Bias Against 93 Stigmatized Groups in Language Models and Guardrail Model Safety Mitigation	37426
<i>Anna-Maria Gueorguieva, Aylín Caliskan</i>	
Resolving Predictive Multiplicity for the Rashomon Set	37435
<i>Parian Haghghat, Hadis Anahideh, Cynthia Rudin</i>	
Unintended Misalignment from Agentic Fine-Tuning: Risks and Mitigation	37443
<i>Dongyoon Hahm, Taywon Min, Woogyel Jin, Kimin Lee</i>	
Silenced Biases: The Dark Side LLMs Learned to Refuse.....	37452
<i>Rom Himmelstein, Amit Levi, Brit Youngmann, Yaniv Nemcovsky, Avi Mendelson</i>	
TAPO: Dynamic Teacher and Perturbed Answer Injection for Policy Optimization.....	37462
<i>Maowei Jiang, Zihang Wang, Qi Wang, Peter Bůš, Moquan Cheng, Yifan Wang, Quangao Liu, Ruiqi Li, Pengyu Zeng, Ruikai Liu, Alan Liang, Yansong Xu, Yusong Hu, Chaoran Zhang, Zhiyong Dong</i>	
Uncovering and Aligning Anomalous Attention Heads to Defend Against NLP Backdoor Attacks	37472
<i>Haotian Jin, Yang Li, Haihui Fan, Lin Shen, Xiangfang Li, Bo Li</i>	
Requirements for Aligned, Dynamic Resolution of Conflicts in Operational Constraints	37481
<i>Steven J. Jones, Robert E. Wray, John E. Laird</i>	
Benchmarking XAI Explanations with Human-Aligned Evaluations	37491
<i>Rémi Kazmierczak, Steve Azzolin, Eloïse Berthier, Anna Hedström, Patricia Delhomme, David Filliat, Nicolas Bousquet, Goran Frehse, Massimiliano Mancini, Baptiste Caramiaux, Andrea Passerini, Gianni Franchi</i>	
Moral Change Or Noise? on Problems of Aligning AI with Temporally Unstable Human Feedback.....	37501
<i>Vijay Keswani, Cyrus Cousins, Breanna Nguyen, Vincent Conitzer, Hoda Heidari, Jana Schaich Borg, Walter Sinnott-Armstrong</i>	
Transparent Networks for Multivariate Time Series.....	37510
<i>Minkyu Kim, Suan Lee, Jinho Kim</i>	
Align to Structure: Aligning Large Language Models with Structural Information.....	37519
<i>Zae Myung Kim, Anand Ramachandran, Farideh Tavazoe, Joo-Kyung Kim, Oleg Rokhlenko, Dongyeop Kang</i>	
Beautiful Images, Toxic Words: Understanding and Addressing Offensive Text in Generated Images	37529
<i>Aditya Kumar, Tom Blanchard, Adam Dziedzic, Franziska Boenisch</i>	

Cost-Minimized Label-Flipping Poisoning Attack to LLM Alignment.....	37538
<i>Shigeki Kusaka, Keita Saito, Mikoto Kudo, Takumi Tanabe, Akifumi Wachi, Youhei Akimoto</i>	
Dropouts in Confidence: Moral Uncertainty in Human-LLM Alignment.....	37547
<i>Jea Kwon, Luiz Felipe Vecchiatti, Sungwon Park, Meeyoung Cha</i>	
Selective Weak-To-Strong Generalization.....	37556
<i>Hao Lang, Fei Huang, Yongbin Li</i>	
MobileSafetyBench: Evaluating Safety of Autonomous Agents in Mobile Device Control	37565
<i>Juyong Lee, Dongyoon Hahm, June Suk Choi, W. Bradley Knox, Kimin Lee</i>	
STELAR-VISION: Self-Topology-Aware Efficient Learning for Aligned Reasoning in Vision	37574
<i>Chen Li, Han Zhang, Zhantao Yang, Fangyi Chen, Zihan Wang, Anudeepsekhar Bolimera, Marios Savvides</i>	
ARGH-Mark: Anchor-Synchronized Watermarking with Hamming Correction for Robust and Quality-Preserving LLM Attribution.....	37583
<i>He Li, Xiaojun Chen, Jingcheng He, Zhendong Zhao, Shuguang Yuan, Xin Zhao, Yunfei Yang</i>	
StyleBreak: Revealing Alignment Vulnerabilities in Large Audio-Language Models Via Style-Aware Audio Jailbreak.....	37591
<i>Hongyi Li, Chengxuan Zhou, Chu Wang, Sicheng Liang, Yanting Chen, Qinlin Xie, Jiawei Ye, Jie Wu</i>	
How Bias Binds: Measuring Hidden Associations for Bias Control in Text-To-Image Compositions	37600
<i>Jeng-Lin Li, Ming-Ching Chang, Wei-Chao Chen</i>	
TORA: Train Once, Realign Anytime for Offline Multi-Objective Reinforcement Learning.....	37609
<i>Weichen Li, Waleed Mustafa, Marcio Monteiro, Puyu Wang, Marius Kloft, Sophie Fellenz</i>	
Bolster Hallucination Detection Via Prompt-Guided Data Augmentation	37618
<i>Wenyun Li, Zheng Zhang, Dongmei Jiang, Xiangyuan Lan</i>	
Editing as Unlearning: Are Knowledge Editing Methods Strong Baselines for Large Language Model Unlearning?.....	37627
<i>Zexi Li, Xiangzhu Wang, William F. Shen, Meghdad Kurmanji, Xinchu Qiu, Dongqi Cai, Chao Wu, Nicholas D. Lane</i>	
How Much Do Large Language Model Cheat on Evaluation? Benchmarking Overestimation Under the One-Time-Pad-Based Framework	37636
<i>Zi Liang, Liantong Yu, Zhang Shiyu, Qingqing Ye, Haibo Hu</i>	
Semantics-Preserving Adversarial Attacks on Event-Driven Stock Prediction Models	37645
<i>Aofan Liu, Haoxuan Li, Hongjian Xing, Yuguo Yin, Zijun Li, Yiyang Qi</i>	
SRAM: Shape-Realism Alignment Metric for No Reference 3D Shape Evaluation.....	37654
<i>Sheng Liu, Tianyu Luan, Phani Nuney, Xuelu Feng, Junsong Yuan</i>	
MRACL: Multi-Reward Space Guided Adaptive Curriculum Reinforcement Learning for LLMs.....	37663
<i>Wenxuan Liu, Liangyu Huo, Yi Jing, Xiyuan Zhang, Jian Xie</i>	
On the Alignment of Large Language Models with Global Human Opinion.....	37673
<i>Yang Liu, Masahiro Kaneko, Chenhui Chu</i>	
DarkBench+: An Extended Benchmark for Evaluating Dark Patterns in Large Language Models	37682
<i>Yaowen Liu, Shenjia Jing, Yufei Wei, Shoumin Zhang, Jinglu Zhang, Zhen Mei, Liangliang Yue, Jiarui Wang, Peng Zhang</i>	
Targeting Misalignment: A Conflict-Aware Framework for Reward-Model-Based LLM Alignment.....	37692
<i>Zixuan Liu, Siavash H. Khajavi, Guangkai Jiang, Xinru Liu</i>	
Mitigating Self-Preference by Authorship Obfuscation	37701
<i>Taslim Mahbub, Shi Feng</i>	
DETONATE – a Benchmark for Text-To-Image Alignment and Kernelized Direct Preference Optimization.....	37709
<i>Renjith Prasad Kaippilly Mana, Abhilekh Borah, Hasnat Md Abdullah, Chathurangi Shyalika, Gurpreet Singh, Ritvik Garimella, Rajarshi Roy, Harshul Raj Surana, Nasrin Imanpour, Suranjana Trivedy, Amit Sheth, Amitava Das</i>	

Misalignment from Treating Means as Ends	37719
<i>Henrik Marklund, Alex Infanger, Benjamin Van Roy</i>	
STACK: Adversarial Attacks on LLM Safeguard Pipelines	37728
<i>Ian R. McKenzie, Oskar John Hollinsworth, Tom Tseng, Xander Davies, Stephen Casper, Aaron David Tucker, Robert Kirk, Adam Gleave</i>	
Aligning Machiavellian Agents: Behavior Steering Via Test-Time Policy Shaping.....	37738
<i>Dena Mujtaba, Brian Hu, Anthony Hoogs, Arslan Basharat</i>	
SharedRep-RLHF: A Shared Representation Approach to RLHF with Diverse Preferences	37747
<i>Arpan Mukherjee, Marcello Bullo, Deniz Gündüz</i>	
Quiet Feature Learning in Algorithmic Tasks.....	37756
<i>Prudhviraj Naidu, Zixian Wang, Leon Bergen, Ramamohan Paturi</i>	
A Tale of Two Identities: An Ethical Audit of AI-Crafted Synthetic Personas	37765
<i>Pranav Narayanan Venkit, Jiayi Li, Yingfan Zhou, Sarah Rajtmajer, Shomir Wilson</i>	
Intrinsic Barriers and Practical Pathways for Human–AI Alignment: An Agreement-Based Complexity Analysis	37775
<i>Aran Nayebi</i>	
CTPD: Cross Tokenizer Preference Distillation	37783
<i>Truong Nguyen, Phi Van Dat, Ngan Nguyen, Linh Ngo Van, Trung Le, Thanh Hong Nguyen</i>	
Realist and Pluralist Conceptions of Intelligence and Their Implications on AI Research	37791
<i>Ninell Oldenburg, Ruchira Dhar, Anders Søgaard</i>	
LieCraft: A Multi-Agent Framework for Evaluating Deceptive Capabilities in Language Models	37802
<i>Matthew Lyle Olson, Neale Ratzlaff, Musashi Hinck, Tri Nguyen, Vasudev Lal, Joseph Campbell, Simon Stepputtis, Shao-Yen Tseng</i>	
Refine and Align: Confidence Calibration Through Multi-Agent Interaction in VQA	37810
<i>Ayush Pandey, Jai Bardhan, Ishita Jain, Ramya S Hebbalaguppe, Rohan Raju Dhanakshirur, Lovekesh Vig</i>	
AdvBDGen: A Robust Framework for Generating Adaptive and Stealthy Backdoors in LLM Alignment.....	37820
<i>Pankayaraj Pathmanathan, Udari Madhushani Sehwag, Michael-Andrei Panaitescu-Liess, Cho-Yu Jason Chiang, Furong Huang</i>	
Beyond I’m Sorry, I Can’t: Dissecting Large-Language-Model Refusal	37830
<i>Nirmalendu Prakash, Yeo Wei Jie, Amir Abdullah, Ranjan Satapathy, Erik Cambria, Roy Ka-Wei Lee</i>	
Towards Benchmarking Privacy Vulnerabilities in Selective Forgetting with Large Language Models.....	37839
<i>Wei Qian, Chenxu Zhao, Yangyi Li, Mengdi Huai</i>	
Backdoor Attacks on Open Vocabulary Object Detectors Via Multi-Modal Prompt Tuning	37849
<i>Ankita Raj, Chetan Arora</i>	
Chain-Of-Thought Driven Adversarial Scenario Extrapolation for Robust Language Models	37858
<i>Md Rafi Ur Rashid, Vishnu Asutosh Dasu, Ye Wang, Gang Tan, Shagufta Mehnaz</i>	
FindTheFlaws: Annotated Errors for Detecting Flawed Reasoning and Scalable Oversight Research.....	37867
<i>Gabriel Recchia, Chatrik Singh Mangat, Issac Li, Gayatri Krishnakumar</i>	
Confirmation Bias: A Challenge for Scalable Oversight	37877
<i>Gabriel Recchia, Chatrik Singh Mangat, Jinu Nyachhyon, Mridul Sharma, Callum Canavan, Dylan Epstein-Gross, Muhammed Abdulbari</i>	
Mind the Gap: Quantifying and Aligning Human-AI Visual Attention for Accident Anticipation.....	37887
<i>Hoe Sung Ryu, Christian Wallraven</i>	
Polarity-Aware Probing for Quantifying Latent Alignment in Language Models	37896
<i>Sabrina Sadiekh, Elena Ericheva, Chirag Agarwal</i>	
Detecting Compute Structuring in AI Governance is Likely Feasible.....	37904
<i>Emmanouil Seferis, Timothy Fist</i>	

Tight Robustness Certification Through the Convex Hull of ℓ_0 Attacks	37913
<i>Yuval Shapira, Dana Drachler-Cohen</i>	
EASE: Practical and Efficient Safety Alignment for Small Language Models	37923
<i>Haonan Shi, Guoli Wang, Tu Ouyang, An Wang</i>	
Efficient Switchable Safety Control in LLMs Via Magic-Token-Guided Co-Training	37932
<i>Jianfeng Si, Lin Sun, Zhewen Tan, Xiangzheng Zhang</i>	
Beyond Verdicts: Evaluating Language Model Moral Competence.....	37941
<i>Aaron J Snoswell, Daniel Kilov, Seth Lazar</i>	
SMPRO: Self-Supervised Visual Preference Alignment Via Differentiable Multi-Preference Multi-Group Ranking	37951
<i>Sirnam Swetha, Rui Meng, Shwetha Ram, Tal Neiman, Son Tran, Mubarak Shah</i>	
Persistent Instability in LLM’s Personality Measurements: Effects of Scale, Reasoning, and Conversation History.....	37961
<i>Tommaso Tosato, Saskia Helbling, Yorguin-Jose Mantilla-Ramos, Mahmood Hegazy, Alberto Tosato, David John Lemay, Irina Rish, Guillaume Dumas</i>	
Shadows in the Code: Exploring the Risks and Defenses of LLM-Based Multi-Agent Software Development Systems	37970
<i>Xiaoqing Wang, Keman Huang, Bin Liang, Hongyu Li, Xiaoyong Du</i>	
Benchmarking Trustworthiness in Multimodal LLMs for Video Understanding.....	37979
<i>Youze Wang, Zijun Chen, Ruoyu Chen, Shishen Gu, Wenbo Hu, Jiayang Liu, Yinpeng Dong, Hang Su, Jun Zhu, Meng Wang, Richang Hong</i>	
STAR-1: Safer Alignment of Reasoning LLMs with 1K Data	37988
<i>Zijun Wang, Haoqin Tu, Yuhan Wang, Juncheng Wu, Yanqing Liu, Jieru Mei, Brian R. Bartoldson, Bhavya Kailkhura, Cihang Xie</i>	
CluCERT: Certifying LLM Robustness Via Clustering-Guided Denoising Smoothing.....	37998
<i>Zixia Wang, Gaojie Jin, Jia Hu, Ronghui Mu</i>	
Safe Multi-Agent Reinforcement Learning with Natural Language Constraints	38007
<i>Ziyan Wang, Meng Fang, Tristan Tomilin, Fei Fang, Yali Du</i>	
Designing Incident Reporting Systems for Harms from General-Purpose AI	38016
<i>Kevin Wei, Lennart Heim</i>	
HumorReject: Decoupling LLM Safety from Refusal Prefix Via a Little Humor.....	38030
<i>Zihui Wu, Haichang Gao, Jiacheng Luo, Zhaoxiang Liu</i>	
MCA-Bench: A Multimodal Benchmark for Evaluating CAPTCHA Robustness Against VLM-Based Attacks	38039
<i>Zonglin Wu, Yule Xue, Yaoyao Feng, Xiaolong Wang, Yiren Song</i>	
MedAtlas: Evaluating LLMs for Multi-Round, Multi-Task Medical Reasoning Across Diverse Imaging Modalities and Clinical Text	38048
<i>Ronghao Xu, Zhen Huang, Yangbo Wei, Xiaoqian Zhou, Zikang Xu, Ting Liu, Zihang Jiang, S. Kevin Zhou</i>	
When Human Preferences Flip: An Instance-Dependent Robust Loss for RLHF	38057
<i>Yifan Xu, Xichen Ye, Yifan Chen, Qiaosheng Zhang</i>	
Multi-Faceted Attack: Exposing Cross-Model Vulnerabilities in Defense-Equipped Vision-Language Models	38066
<i>Yijun Yang, Lichao Wang, Jianping Zhang, Chi Harold Liu, Lanqing Hong, Qiang Xu</i>	
CoSPED: Consistent Soft Prompt Targeted Data Extraction and Defense	38075
<i>Zhuochen Yang, Kar Wai Fok, Vrizzlynn L. L. Thing</i>	
Fading the Digital Ink: A Universal Black-Box Attack Framework for 3DGS Watermarking Systems.....	38084
<i>Qingyuan Zeng, Shu Jiang, Jiajing Lin, Zhenzhong Wang, Kay Chen Tan, Min Jiang</i>	
SL-CBM: Enhancing Concept Bottleneck Models with Semantic Locality for Better Interpretability.....	38093
<i>Hanwei Zhang, Luo Cheng, Rui Wen, Yang Zhang, Lijun Zhang, Holger Hermanns</i>	

Differentiated Directional Intervention: A Framework for Evading LLM Safety Alignment.....	38102
<i>Peng Zhang, Peijie Sun</i>	
DAVSP: Safety Alignment for Large Vision-Language Models Via Deep Aligned Visual Safety Prompt	38111
<i>Yitong Zhang, Jia Li, Liyi Cai, Ge Li</i>	
CultureRL: Internalizing Cultural Principles in Large Language Models Via Norm-Driven Reinforcement Learning	38120
<i>Weixiang Zhao, Haozhen Li, Yanyan Zhao, Haixiao Liu, Biye Li, Ting Liu, Bing Qin</i>	
Composable Assurance for AI Alignment: A Framework for Propagating Formal Safety Properties Through MLOps	38129
<i>Xiaofei Zhao</i>	
Value-Aligned Prompt Moderation Via Zero-Shot Agentic Rewriting for Safe Image Generation.....	38137
<i>Xin Zhao, Xiaojun Chen, Bingshan Liu, Zeyao Liu, Zhendong Zhao, Xiaoyan Gu</i>	
GEM: Generative Entropy-Guided Preference Modeling for Few-Shot Alignment of LLMs	38146
<i>Yiyang Zhao, Huiyu Bai, Xuejiao Zhao</i>	
Explainable Melanoma Diagnosis with Contrastive Learning and LLM-Based Report Generation.....	38156
<i>Junwen Zheng, Xinran Xu, Li Rong Wang, Chang Cai, Lucinda Siyun Tan, Dingyuan Wang, Hong Liang Tey, Xiuyi Fan</i>	
Can LLMs Detect Their Confabulations? Estimating Reliability in Uncertainty-Aware Language Models.....	38164
<i>Tianyi Zhou, Johanne Medina, Sanjay Chawla</i>	
Not All Tokens Are Meant to Be Forgotten	38173
<i>Xiangyu Zhou, Yao Qiang, Saleh Zare Zade, Douglas Zytko, Prashant Khanduri, Dongxiao Zhu</i>	

Author Index