KAI CHEN

HKUST, Clear Water Bay, New Territories, Hong Kong SAR

Email: kai.chen@connect.ust.hk \lor Homepage: www.cse.ust.hk/kchenbf

RESEARCH OVERVIEW

My research aims at constructing reliable **Multi-modal** AI systems from a **data-centric** perspective. Recently, we have witnessed the interim success of training foundational models on massive **human data**, which, however, is believed to come to an end. Towards the second half of AI, a scalable **synthetic data** generation and training pipeline is necessary, for which **reinforcement learning** serves as a well-formulated solution. Thus, **Scaling Reinforcement Learning for (M)LLMs** has become indispensable for achieving Artificial Super-Intelligence (ASI), which has obtained remarkable success for text-only LLMs, while remaining open challenges for (native) **Multi-modal LLMs**. Specifically, I'm currently trying to answer the following questions,

- Policy: How to build end-to-end Multi-modal LLMs with frontier visual, textual, and speech abilities?
- World: How to construct 3D visual world models in a controllable and scalable manner?
- Feedback: How to get effective intrinsic feedback from (M)LLMs themselves without a reward model?
- Optimization: Does more data always result in better performance?

Research Areas: Omni-modal LLMs, Visual World Modeling, Mixture-of-Experts (MoE)

EDUCATION

Hong Kong University of Science and Technology, Hong Kong SAR Sep 2020 - June 2026 (Expect) Ph.D. in Computer Science and Engineering

GPA: 4.10/4.0

Advisor: Prof. Dit-Yan Yeung

Fudan University(FDU), Shanghai, China

Sep 2016 - June 2020

B.S. in Computer Science, Minor in Economics (Outstanding Graduates of Shanghai)

Overall GPA: 3.70/4.0, Major GPA: 3.90/4.0, Ranking: 3/32

Advisor: Prof. Yanwei Fu

University of Manchester, Manchester, UK

Sep 2018 - Jan 2019

Exchange student in the **Department of Computer Science**

Advisor: Dr. Tingting Mu

EXPERIENCE

Mobile Intelligence Group (MIG), SenseTime

Oct 2019 - April 2020

Research Intern

Advisor:Dr. Wenxiu Sun

• Research on real-time (portrait) instance segmentation deployable on mobile devices.

Computer Vision Lab, Indiana University Bloomington (IUB)

June 2019 - Sep 2019

Global Talent Attraction Program (GTAP) Visiting Scholar

Advisor:Prof. David Crandall

• Research on semi-supervised semantic segmentation and indoor 3D reconstruction.

SELECTED HONORS

CVPR 2025 Doctoral Consortium Awards	May 2025
HKUST Research Travel Grant	2023-2025
HKUST Postgraduate Scholarship	Sep 2020
Outstanding Graduates of Shanghai [Wechat Post] (5%, by Shanghai Government)	April 2020
Scholarship for Outstanding Graduates (5%, by Fudan University)	April 2020

Oversea Visiting Student Stipend of (15,000 CNY, Fudan University) Joel & Ruth Spira Scholarship (1%, by Lutron Electronics) National Scholarship (1%, by Ministry of Education of P.R.China) Scholarship for Outstanding Undergraduate Students (5%, by Fudan University)	Dec 2019 Mar 2019 Sep 2018 Oct 2017
PUBLICATIONS	
Full publication list on my Google Scholar. (* denotes equal contribution)	
I. Multi-modal Foundation Models - Omni-modality and Reasoning RQ: How to construct multi-modal LLMs with visual, textual, and speech reasoning abilities s	imultaneously?
[C23] EMOVA: Empowering Language Models to See, Hear and Speak with Vivid Emotions Kai Chen*, Yunhao Gou*, Runhui Huang*, Zhili Liu*, Daxin Tan*, and other 26 authors	CVPR 2025 [link]
[C22] Perceptual Decoupling for Scalable Multi-modal Reasoning via Reward-	Arxiv 2025
Optimized Captioning Yunhao Gou*, <u>Kai Chen*</u> , Zhili Liu*, Lanqing Hong, Xin Jin, Zhenguo Li, James T. Kwok, Yu Zhang.	$[\underline{ ext{link}}]$
II. Multi-modal Foundation Models - Mixture of Cluster-conditional Experts (MRQ: Does more data always result in better performance during model pre-training and fine-training and fine-trainin	,
[C21] Mixture of Cluster-conditional LoRA Experts for Vision-language Instruction Tuning Yunhao Gou*, Zhili Liu*, <u>Kai Chen*</u> , Lanqing Hong, Hang Xu, Aoxue Li, Dit-Yan Yeung, James Kwok, Yu Zhang.	Arxiv 2023 [link]
[C20] Task-customized Masked Autoencoder via Mixture of Cluster-conditional Experts Zhili Liu*, Kai Chen*, Jianhua Han, Lanqing Hong, Hang Xu, Zhenguo Li, James Kwok.	2023 Spotlight
[C19] Task-Customized Self-Supervised Pre-training with Scalable Dynamic Routing Zhili Liu, Jianhua Han, <u>Kai Chen</u> , Lanqing Hong, Hang Xu, Chunjing Xu, Zhenguo Li.	AAAI 2022 [link]
III. Multi-modal Foundation Models - Scalable Oversight for (M)LLM Self-align RQ: Are there any intrinsic scalable oversight from (M)LLMs to supervise their own capability	
[C18] Corrupted but Not Broken: Rethinking the Impact of Corrupted Data EMN in Visual Instruction Tuning Yunhao Gou, Hansi Yang, Zhili Liu, <u>Kai Chen</u> , Yihan Zeng, Lanqing Hong, Zhenguo Li,	NLP 2025 Oral [link]
Qun Liu, James T Kwok, Yu Zhang.	TMID 2025
[J2] Unified Triplet-Level Hallucination Evaluation for Large Vision-Language Models Junjie Wu*, Tsz Ting Chung*, <u>Kai Chen*</u> , Dit-Yan Yeung.	[link]
[C17] Mixture of insightful Experts (MoTE): The Synergy of Thought Chains and Expert Mixtures in Self-Alignment Zhili Liu*, Yunhao Gou*, <u>Kai Chen*</u> , Lanqing Hong, Jiahui Gao, Fei Mi, Yu Zhang, Zhenguo Li, Xin Jiang, Qun Liu, James T. Kwok.	ACL 2025 [link]
[C16] Eyes Closed, Safety On: Protecting Multimodal LLMs via Image-to-Text	ECCV 2024
Transformation Yunhao Gou*, <u>Kai Chen*</u> , Zhili Liu*, Lanqing Hong, Hang Xu, Zhenguo Li, Dit-Yan Yeung, James Kwok, Yu Zhang.	[<u>link</u>]

[C15	Gaining Wisdom from Setbacks: Aligning Large Language Models via Mistake Analysis	ICLR	2024
	<u>Kai Chen*</u> , Chunwei Wang*, Kuo Yang, Jianhua Han, Lanqing Hong, Fei Mi, Hang Xu, Zhengying Liu, Wenyong Huang, Zhenguo Li, Dit-Yan Yeung, Lifeng Shang, Xin Jiang, Qun L	iu.	[link]
	IV. Visual World Models - Corner Cases for Autonomous Driving (CODA) RQ: How to enhance the robustness of self-driving agents towards road corner cases? A: 1) corner case collection, 2) corner case generation, and 3) multi-modal reasoning		
[C14	ECCV 2024 W-CODA: 1st Workshop on Multimodal Perception and Comprehension of Corner Cases in Autonomous Driving	ECCV	
	<u>Kai Chen*</u> , Ruiyuan Gao*, Lanqing Hong*, Hang Xu, Xu Jia, Holger Caesar, Dengxin Dai, Bingbing Liu, Dzmitry Tsishkou, Songcen Xu, Chunjing Xu, Qiang Xu, Huchuan Lu, Dit-Yan	Yeung.	[<u>link</u>]
[C13	CODA-LM: Automated Evaluation of Large Vision-Language Models on Self-driving Corner Cases Kai Chen*, Yanze Li*, Wenhua Zhang*, Yanxin Liu, Pengxiang Li, Ruiyuan Gao,	WACV	2025
	Lanqing Hong, Meng Tian, Xinhai Zhao, Zhenguo Li, Dit-Yan Yeung, Huchuan Lu, Xu Jia.		[<u>IIIIK]</u>
[C12	CODA: A Real-World Road Corner Case Dataset for Object Detection in Autonomous Driving	ECCV	2022
	Kaican Li*, <u>Kai Chen*</u> , Haoyu Wang*, Lanqing Hong, Chaoqiang Ye, Jianhua Han, Yukuai Chen, Wei Zhang, Chunjing Xu, Dit-Yan Yeung, Xiaodan Liang, Zhenguo Li, Hang Xu	1.	[link]
	V. Visual World Models - Geometric-controllable Visual Generation RQ: How to generate the 3D visual world in a controllable and scalable manner?		
[C11	MagicDrive3D: Controllable 3D Generation for Any-View Rendering in Street Scenes Ruiyuan Gao, Kai Chen, Zhihao Li, Lanqing Hong, Zhenguo Li, Qiang Xu.	Arxiv	2024 [link]
[C10	MagicDrive-V2: High-Resolution Long Video Generation for Autonomous Driving with Adaptive Control Ruiyuan Gao, Kai Chen, Bo Xiao, Lanqing Hong, Zhenguo Li, Qiang Xu.	ICCV	
[C9]	Implicit Concept Removal of Diffusion Models Zhili Liu*, <u>Kai Chen*</u> , Yifan Zhang, Jianhua Han, Lanqing Hong, Hang Xu, Zhenguo Li, Dit-Yan Yeung, James Kwok.	ECCV	2024 [<u>link</u>]
[C8]	DetDiffusion: Synergizing Generative and Perceptive Models for Enhanced Data Generation and Perception Yibo Wang*, Ruiyuan Gao*, <u>Kai Chen*</u> , Kaiqiang Zhou, Yingjie Cai, Lanqing Hong, Zhenguo Li, Lihui Jiang, Dit-Yan Yeung, Qiang Xu, Kai Zhang.	CVPR	2024 [<u>link</u>]
[C7]	MagicDrive: Street View Generation with Diverse 3D Geometry Control Ruiyuan Gao*, <u>Kai Chen*</u> , Enze Xie, Lanqing Hong, Zhenguo Li, Dit-Yan Yeung, Qiang Xu	ICLR	2024 [<u>link</u>]
[C6]	TrackDiffusion: Tracklet-Conditioned Video Generation via Diffusion Models Pengxiang Li*, <u>Kai Chen*</u> , Zhili Liu*, Ruiyuan Gao, Lanqing Hong, Dit-Yan Yeung, Huchuan Lu, Xu Jia.	WACV	2025 [<u>link</u>]
[C5]	GeoDiffusion: Text-Prompted Geometric Control for Object Detection Data Generation Kai Chen*, Enze Xie*, Zhe Chen, Yibo Wang, Lanqing Hong, Zhenguo Li, Dit-Yan Yeung.	ICLR	2024 [<u>link</u>]
	VI. Representation Learning - Object-level Self-supervised Learning (SSL) RQ: How to perform object-level SSL for better transferability on downstream dense perception	tasks?	

[C4]	Mixed Autoencoder for Self-supervised Visual Representation Learning Kai Chen*, Zhili Liu*, Lanqing Hong, Hang Xu, Zhenguo Li, Dit-Yan Yeung.	$\begin{array}{c} \text{CVPR 2023} \\ [\underline{\text{link}}] \end{array}$
[C3]	MultiSiam: Self-supervised Multi-instance Siamese Representation Learning for Autonomous Driving	ICCV 2021
	Kai Chen, Lanqing Hong, Hang Xu, Zhenguo Li, Dit-Yan Yeung.	$[\underline{ ext{link}}]$
[C2]	SODA10M: A Large-Scale 2D Self/Semi-Supervised Object Detection Dataset for Autonomous Driving Jianhua Han, Xiwen Liang, Hang Xu, <u>Kai Chen</u> , Lanqing Hong, Jiageng Mao, Chaoqiang Ye, Wei Zhang, Zhenguo Li, Xiaodan Liang, Chunjing Xu.	NeurIPS 2021 [<u>link</u>]
	Early Works	
[J1]	Automatic Dense Annotation for Monocular 3D Scene Understanding Md. Alimoor Reza, <u>Kai Chen</u> , Akshay Naik, David Crandall, Soon-Heung Jung.	EE Access 2020 [<u>link</u>]
[C1]	Automatic Annotation for Semantic Segmentation in Indoor Scenes Md Alimoor Reza, Akshay Naik, <u>Kai Chen</u> , David Crandall.	IROS 2019 [<u>link</u>]
A	CADEMIC SERVICES	
	Program Committee / Organizer	
	The 1st W-CODA Workshop at ECCV 2024 on Multimodal Perception and Comprehension Corner Cases in Autonomous Driving.	
	The 2nd SSLAD workshop at ECCV 2022. The 1st SSLAD workshop at ICCV 2021 on Self-supervised Learning for Next-generation Industry-level Autonomous Driving.	2022 2021
	Area Chair	
•	International Joint Conferences on Artificial Intelligence (IJCAI)	2025
	Conference Reviewer	
•	IEEE Conference on Computer Vision and Pattern Recognition (CVPR) IEEE International Conference on Computer Vision (ICCV) European Conference on Computer Vision (ECCV) International Conference on Learning Representations (ICLR) International Conference on Machine Learning (ICML) Neural Information Processing Systems (NeurIPS) International Joint Conferences on Artificial Intelligence (IJCAI) AAAI Conference on Artificial Intelligence (AAAI) International Conference on Robotics and Automation (ICRA) ACM International Conference on Multimedia (ACM MM) IEEE Winter Conference on Applications of Computer Vision (WACV) Asian Conference on Computer Vision (ACCV)	2022-2025 2023-2025 2022-2024 2023-2026 2025 2021-2025 2023-2025 2022 2022 2025 2026 2024
	Journal Reviewer	
•	IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) IEEE Transactions on Circuits and Systems for Video Technology (TCSVT) IEEE Transactions on Image Processing (TIP) IEEE Access	

PATENTS

- [CN116665219A] GeoDiffusion: Text-Prompted Geometric Control for Object Detection Data Generation. Enze Xie, <u>Kai Chen</u>, Lanqing Hong, Zhenguo Li. *Published in May 26th*, 2023.
- [CN115731530A] MultiSiam: Self-supervised Multi-instance Siamese Representation Learning for Autonomous Driving. Kai Chen, Lanqing Hong, Hang Xu, Zhenguo Li. Published in Aug. 24th, 2021.

TEACHING

- HKUST COMP 4211 Machine Learning, Teaching Assistant, Fall 2025.
- HKUST COMP 2012 Object-Oriented Programming and Data Structures, Teaching Assistant, Fall 2024.
- HKUST COMP 2012 Object-Oriented Programming and Data Structures, Teaching Assistant, Fall 2021.
- HKUST COMP 2012 Object-Oriented Programming and Data Structures, Teaching Assistant, Spring 2021.

INVITED TALKS

- [AI TIME Online] EMOVA: Empowering Language Models to See, Hear and Speak with Vivid Emotions. [Recording]
- [VALSE Webinar] Geometric-controllable Visual Generation: A Systemetric Solution. [Recording]
- [AIDriver Online] Controllable Corner Case Generation for Autonomous Driving. [Recording]
- [AI TIME Online] Gaining Wisdom from Setbacks: Aligning Large Language Models via Mistake Analysis. [Recording]
- [TechBeat Online] Gaining Wisdom from Setbacks: Aligning Large Language Models via Mistake Analysis. [Recording]
- [VALSE 2023@Wuxi] Mixed Autoencoder for Self-supervised Visual Representation Learning. [Recording]
- [VALSE 2023@Wuxi] CODA: A Real-World Road Corner Case Dataset for Object Detection in Autonomous Driving. [Recording]

TECHNICAL SKILLS

Program Languages Python, Matlab, C/C++/C#, SQL, LATEX

Framework Pytorch, Tensorflow

Language Native in Mandarin, Fluent in English and Japanese

CET-4(649), CET-6(619), TOEFL-iBT(101)