

# Chi Xing

chi.xing2025@gmail.com | github.com/MartinRepo | openchi.life

## Education

|  |                   |
|--|-------------------|
| <b>University of Edinburgh</b>   | Edinburgh, UK     |
| M.Sc Artificial Intelligence   | 2024/09 – 2025/09 |
| <ul style="list-style-type: none"><li>• Focused on various machine learning architectures, ranging from basic neural networks to advanced modern architectures (Transformer, Diffusion Model, Visual Auto-Regressive Model, etc.)</li><li>• Thesis is focused on Carbon-Aware Geospatial Shifting of LLM Training Workloads, which is supervised by <a href="#">Prof. Luo Mai</a>.</li></ul> |                   |
| <b>University of Liverpool &amp; Xi'an Jiaotong-Liverpool University</b>   | Liverpool, UK     |
| B.Sc Computer Science (1st class with honors)  | 2020/09 – 2024/07 |
| <ul style="list-style-type: none"><li>• Research Interest Points: Algorithm Design, C++/C/C#, Machine Learning, Trustworthy AI, Web Development.</li><li>• Thesis is focused on scheduling algorithms for modern smart grid. This project is supervised by <a href="#">Prof. Prudence Wong</a>.</li></ul>  |                   |

## Open Source Projects

|   |                    |
|---|--------------------|
| <b>ServerlessLLM</b>  | 500+ GitHub Stars  |
| Core Contributor, Code Reviewer   | Nov 2024 – Present |
| <ul style="list-style-type: none"><li>• Proficient in building large-scale distributed inference systems using Hugging Face Transformers and vLLM.</li><li>• Designed and implemented an <b>end-to-end serverless PEFT LoRA fine-tuning solution</b> within the ServerlessLLM ecosystem to provide on-demand, cost-effective model customization services (<a href="#">#251</a>, <a href="#">#189</a>).</li><li>• Developed a <b>multi-tenant serverless serving solution for LoRA adapters</b> using Ray, achieving up to a 4.4x faster loading speed compared to the safetensors format by leveraging a multi-tier checkpoint loading mechanism (<a href="#">#248</a>, <a href="#">#221</a>).</li></ul>   |                    |
| <b>Casibase</b>   | 4k+ GitHub Stars   |
| Core Contributor, OSPP (Open Source Promotion Plan) 2024 Mentor   | Jan 2024 – Present |
| <ul style="list-style-type: none"><li>• <b>Enhanced the platform’s core multi-modal capabilities:</b> Deeply integrated various large multi-modal models to enable end-to-end functionalities for image understanding, generation, and mixed-media dialogue. Optimized user experience with features like drag-and-drop uploads and URL parsing (<a href="#">#925</a>, <a href="#">#895</a>, <a href="#">#717</a>, <a href="#">#716</a>).</li><li>• <b>Expanded and optimized LM support:</b> Integrated multiple industry-leading models and engineered a model provider multiplexing mechanism, allowing the system to dynamically select models based on load and cost (<a href="#">#785</a>, <a href="#">#783</a>, <a href="#">#703</a>).</li><li>• <b>Improved the core RAG workflow:</b> Significantly boosted the quality of knowledge base vectorization and retrieval relevance by designing novel text-splitting strategies (<a href="#">#778</a>, <a href="#">#727</a>).</li><li>• <b>Led full-stack development and performance optimization:</b> Utilized Go (BeeGo) and React.js to independently deliver features including real-time billing &amp; usage statistics (<a href="#">#898</a>, <a href="#">#735</a>), rich text rendering (LaTeX, code highlighting) (<a href="#">#775</a>, <a href="#">#776</a>), and front-end optimizations that enhanced message rendering speed and system stability (<a href="#">#777</a>, <a href="#">#954</a>).</li></ul> |                    |

## Work Experience

|   |                      |
|---|----------------------|
| <b>N8 CIR</b>   | Liverpool & York, UK |
| Research Intern@Computational Biology Facility  | 2024/06 – 2024/09    |
| <ul style="list-style-type: none"><li>• <b>Focused on benchmarking various LLMs</b> for reading biomedical literature, utilizing Llama.cpp to quantize open-source models such as Llama3.1-70B, Llama3.1-405B, DBRX, and Mixtral-8x22B.</li><li>• <b>Developed an objective scoring system</b> that extracts key information from model outputs and evaluates their similarity to manually extracted data for performance benchmarking.</li><li>• <b>Designed a summarization method</b> to reduce input size, enabling the use of models with smaller context windows.</li><li>• <b>This work also involved comparing model performance</b> across different hardware platforms, including NVIDIA GH200, A100, and CPU/GPU references, and deploying LLMs on high-performance computing (HPC) architectures.</li></ul> |                      |
| <b>IFLYTEK</b>  | Suzhou, China        |
| SDE@R&D   | 2022/06 – 2022/09    |
| <ul style="list-style-type: none"><li>• Enhanced the accuracy of location-based NLP tasks within the IFLYTEK “Police Super Brain” system by <b>conducting meticulous data annotation</b> and quality assurance, <b>correcting machine-labeled address POIs</b>, and applying foundational knowledge of <b>entity relationship extraction</b>.</li></ul>   |                      |

## Publication

---

Preference Alignment on Diffusion Model: A Comprehensive Survey for Image Generation and Editing Feb, 2025

- Preference Alignment on DMs Application section, investigated and summarised a set of application paradigms.

## Hackathon

---

2023 BMW Hackathon Shenyang, China

2nd Place, HVB Reuse for Energy Saving in Production Channel 2023.08

- **Designed a power scheduling algorithm:** formulated a dynamic programming model based on electricity price fluctuations, photovoltaic generation, and solar radiation intensity; derived the dynamic transition equations and successfully solved for the optimal scheduling strategy.
- **Developed a battery dispatch strategy:** proposed a scheduling method based on greedy algorithms for energy storage cabinets utilizing retired automotive batteries, effectively mitigating battery degradation; theoretically validated the strategy by proving its greedy choice and optimal substructure properties.
- **Engineered and deployed the system:** containerized the scheduling solution using Docker, completed final system submission, and gained proficiency in basic Docker operations.

## Skills

---

- Programming Language & Tools: C/C++, Python, Go, JavaScript, Rust, Java, Git, Linux, Shell
- Deeplearning Framework : Pytorch, Huggingface-Transformers, Deepspeed, Megatron-LM
- Distributed System/Computing : Docker, MPI, Ray