

# Kaleem Ullah Qasim

PhD in Artificial Intelligence | AI Researcher & LLM Engineer  
Chengdu, China | +86-13111895637 | kaleem@my.swjtu.edu.cn  
[LinkedIn](#) | [UpWork](#) | [GitHub](#) | [Google Scholar](#) | [HuggingFace](#)

## Research Areas

LLM & SLM Reasoning | Recursive Task Decomposition | Agentic Systems & Multi-Agent Coordination | RL Alignment (RLHF, DPO, GRPO, Verifiable Rewards) | Context Engineering

## Technical Skills

**Programming & Development:** Python | TypeScript | JavaScript | SQL | FastAPI | Django | Flask | Git | REST APIs | GraphQL  
**LLM & Generative AI:** HuggingFace (Transformers, TRL, PEFT, Datasets, Accelerate) | Fine-tuning (SFT, LoRA, QLoRA) | vLLM | DeepSpeed | FSDP | RAG Architecture | Context Engineering | In-context Learning  
**Agentic Systems:** MCP (Model Context Protocol) | OpenAI Agents SDK | Pydantic AI | smolagents | Google ADK | LangGraph | LangChain | LlamaIndex | AutoGen | CrewAI | ReAct Pattern | Tool Use & Function Calling  
**Alignment & Agentic RL:** RLHF (PPO, REINFORCE) | Preference Optimization (DPO, CPO, ORPO) | GRPO | Verifiable-Reward RL | TRL | verl | OpenRLHF | Reward Modeling | Multi-Agent RL  
**Machine Learning & Evaluation:** PyTorch | scikit-learn | XGBoost | LightGBM | BERT | Transformers | NLP | Transfer Learning | lm-evaluation-harness | OpenCompass  
**MLOps & Cloud:** AWS (SageMaker, Lambda, S3) | Huawei Ascend | Docker | Kubernetes | MLflow | Weights & Biases | Vector Databases (Pinecone, Weaviate, ChromaDB)

## Professional Experience

### Research Contract – Recursive Reasoning Agent System, Huawei Technologies (Lead Researcher) 2025 – Present

- Won Huawei-commissioned research contract on the strength of RDoLT paper (JAIR Q1, 16 citations) as lead researcher, building a recursive decomposition agent framework for complex reasoning
- Designed core recursive decomposition logic and memory modules, validated on AIME, U-MATH, and custom algorithm benchmarks – demonstrating measurable gains over pure CoT baselines
- Evaluated heterogeneous large/small model combinations to optimize performance-cost tradeoffs; researched and selected RL training strategies for agent improvement
- Building RL training loop on Huawei Ascend platform targeting top rankings on AIME/U-MATH benchmarks, with goal of CCF-A publication or patent

### AI Engineer & LLM Specialist, Upwork (Freelance) 2023 – Present

- Top Rated (top 10%) with 100% job success rate and 5-star ratings from 20+ clients across AI/ML projects
- Built production RAG chatbots using LangChain/LlamaIndex/CrewAI with Pinecone and Weaviate vector databases, reducing task completion time by 20% with 95% accuracy on domain-specific queries
- Optimized local LLMs (Llama 2/3, Mistral) using LoRA/QLoRA fine-tuning, improving task accuracy by 25% while ensuring GDPR compliance for privacy-focused enterprise deployments
- Developed multi-agent AI systems using LangGraph and AutoGen for workflow orchestration, boosting semantic search accuracy by 35% and reducing API latency by 40%
- Architected multimodal AI pipelines (Deepseek OC) for document analysis, processing 10K+ documents monthly with 92% accuracy using custom prompt engineering

### Research Contractor – Traffic AI, University of Jeddah (Dr. Tariq Alsaifi) 2024 – Present

- Co-authored 2 papers on LLM-based traffic analysis published in Alexandria Engineering Journal and arXiv, focusing on spatio-temporal reasoning and accident severity prediction using hybrid LLM architectures
- Developed TrafficCoT-R framework combining Chain-of-Thought prompting with multi-agent coordination for traffic prediction, integrating GPT-4 with GIS data and graph neural networks (PyTorch Geometric)
- Built SAFE (Semantic-Augmented Fusion Ensemble) hybrid system using Random Forest and Qwen3-4B LLM for accident severity classification, achieving 85.7% recall on severe accidents (from 0% with traditional ML) and 53.1% overall accuracy on naturally imbalanced data

### Research Contractor – AI Security, Zhejiang University (Dr. Haitao Xu) 2022 – 2024

- Co-authored paper on online affiliate marketing published in IEEE INFOCOM 2025, conducting empirical study of deceptive marketing practices and affiliate network analysis using NLP and network graph analysis techniques
- Developed ADsFlow Chrome extension for detecting dynamic web advertisements, analyzing ad positioning and loading patterns on webpages using DOM analysis and computer vision (OpenCV) for real-time ad classification
- Built webpage classification system using fine-tuned RoBERTa and BERT embeddings to identify webpage intent (cybersecurity threats, marketing, phishing), processing webpage code and content for multi-class classification (patent application filed)
- Designed ML pipeline using XGBoost for deceptive ad detection, achieving 89% precision on 50K+ samples with automated feature extraction from URL structure, content, and affiliate link patterns

**Data Scientist, Chengdu Ayurveda Biotechnology Co., Ltd****2020 – 2023**

- Contributed to company's achievement of #1 ranking on Alibaba marketplace through ensemble ML models (Random Forest, Gradient Boosting) for demand forecasting and dynamic pricing optimization, supporting 180% YoY revenue growth
- Implemented predictive SEO strategy using time-series forecasting (ARIMA, Prophet) and NLP-based keyword analysis, achieving 95% increase in search appearances and 65% improvement in organic traffic over 12-month period
- Built real-time analytics dashboard using Streamlit, Plotly, and PostgreSQL with automated ETL pipelines for sales, inventory, and customer data, improving data-driven decision response time by 60%

**Publications**

- **Kaleem Ullah Qasim**, J. Zhang, MK Shaheen, R. Alharith, H. Zhang (2026). *The Residual Stream Is All You Need: On the Redundancy of the KV Cache in Transformer Inference*. arXiv:2603.19664.
- **Kaleem Ullah Qasim**, J. Zhang, H. Li, MK Shaheen (2026). *VERIFY-RL: Verifiable Recursive Decomposition for Reinforcement Learning in Mathematical Reasoning*. arXiv:2602.07559.
- **Kaleem Ullah Qasim**, J. Zhang, T. Alsahfi, A. U. R. Butt (2025). *Recursive Decomposition of Logical Thoughts: Framework for Superior Reasoning and Knowledge Propagation in Large Language Models*. Journal of Artificial Intelligence Research (JAIR), 83. **16 citations**.
- T. Alsahfi, **Kaleem Ullah Qasim** (2025). *TrafficToT-R: A Framework for Advanced Spatio-Temporal Reasoning in Large Language Models*. Alexandria Engineering Journal, 128, 464-475. **5 citations**.
- H. Xu, Y. Sun, **Kaleem Ullah Qasim**, et al. (2025). *Understanding the Business of Online Affiliate Marketing: An Empirical Study*. IEEE INFOCOM 2025, 1-10. 2 citations.
- **Kaleem Ullah Qasim**, J. Zhang, HS Ur Rehman (2025). *Complexity Aware Recursive Decomposition for Math Reasoning*. Proceedings of the 2025 International Conference on Embodied Intelligence.
- R. Ali, J. Xu, M. H. Baig, H. S. U. Rehman, M. W. Aslam, **Kaleem Ullah Qasim** (2024). *From Data to Decisions: Enhancing Financial Forecasts with LSTM for AI Token Prices*. Journal of Economic Studies, 51(8), 1677-1693. **10 citations**.
- A. U. Rehman, M. Asif, **Kaleem Ullah Qasim**, et al. (2025). *AdvancedHybridNet: An AI-Powered Hybrid Ensemble for High-Accuracy Thyroid Disease Diagnosis Using Dynamic Feature Selection*. In Press.
- **Kaleem Ullah Qasim**, J. Zhang (2025). *MARBLE: A Multi-Agent Rule-Based LLM Reasoning Engine for Accident Severity Prediction*. arXiv:2507.04893. 1 citation.
- **Kaleem Ullah Qasim**, J. Zhang (2025). *Accelerating Training Speed of Tiny Recursive Models via Curriculum Guided Adaptive Recursion*. arXiv:2511.08653. 1 citation.
- M. W. Aslam, Z. Zhang, **Kaleem Ullah Qasim** (2025). *LLMFacility: A Large Language Model Driven Evolutionary Framework for Interpretable and Scalable Facility Layout Optimization*. SSRN 5692433.
- H. S. U. Rehman, L. Liu, **Kaleem Ullah Qasim** (2025). *ASTIF: Adaptive Semantic-Temporal Integration for Cryptocurrency Price Forecasting*. arXiv:2512.18661.

**Selected Research Projects****RDOLT**

JAIR (Q1) · 2025 · First Author

Recursive task decomposition with knowledge propagation, expanding reasoning into a verified sub-problem tree; significant gains over CoT/ToT on GSM8K, MATH, MMLU-STEM (**16 citations**; basis for Huawei research contract).

**VERIFY-RL**

arXiv 2026 · First Author

Process-level verifiable rewards combined with GRPO for end-to-end RL on 7B reasoning models over recursive decomposition trajectories; consistent gains over SFT and vanilla-GRPO on AIME and U-MATH.

**CGAR – Curriculum-Guided Adaptive Recursion**

arXiv 2025 · First Author

Curriculum-driven adaptive recursion-depth scheduling that adjusts inference steps per sample difficulty; accelerated training and lower inference cost while preserving accuracy.

**Education****Ph.D. in Artificial Intelligence****2022 – 2026**

Southwest Jiaotong University (SWJTU), School of Computing and Artificial Intelligence, Chengdu, China  
Research: LLM/SLM Reasoning, Recursive Decomposition, Agentic Systems, Reinforcement Learning Alignment

**Master in Computer Application Technology****2019 – 2022**

Southwestern University of Finance and Economics (SWUFE), Chengdu, China  
Focus: Natural Language Processing, Machine Learning, Deep Learning

**Certifications**

**Professional Certifications:** [Google Prompting Essentials](#) | [Generative AI: Prompt Engineering](#) | [Generative AI & LLMs](#) | [Introduction to AI](#) | [Claude Code in Action](#) | [Reinforcement Fine-Tuning LLMs With GRPO](#) | [Post-training of LLMs](#)

**Languages & Communication:** English (Fluent), Chinese (HSK 5 & HSKK Intermediate), Urdu (Native), Hindi (Fluent)