


<div>Flexible Machine Learning Algorithm Lab</div>		<div>■ Contact information</div> <div>Professor : Insu Han TEL : 7477</div> <div>Lab. : Flexible ML Algorithm Lab</div> <div>Website : flexml.kaist.ac.kr / insuhan.github.io</div>
<div>■ Current state of the Lab. (in 2025 Spring Semester)</div> <div>2 incoming Master's Student, 8 undergraduate students</div>		
<div>■ Research Areas</div> <div>Our research group focuses on developing scalable algorithms for large-scale machine learning problems, combining practical efficiency with rigorous theoretical analysis. We tackle challenges at the forefront of AI and machine learning, aiming to improve the performance and applicability of cutting-edge models and techniques.</div> <div>Success of AI</div> <div>1. Big Data 2. Hardware 3. Software</div> <div></div> <div>Several research directions are</div> <div>1. Development of large-scale machine learning approximation algorithms</div> <div>2. Development of quantization algorithms for inference foundation models</div> <div>3. Development efficient algorithms for genomic foundation models for cancer classification</div> <div>3. Acceleration of core operations in foundation models through matrix sparsity approximation</div> <div><div>Quantization on query embeddings</div><div>$S \in \mathbb{R}^{m \times d} \sim \mathcal{N}(0, 1)$</div><div>$q \in \mathbb{R}^d \xrightarrow{\text{JL transform}} Sq \in \mathbb{R}^m$</div></div> <div><div>Quantization on key embeddings</div><div>$Sk \in \mathbb{R}^m$</div><div>$k \in \mathbb{R}^d \xrightarrow{\text{JL transform}} \text{binary projection} \rightarrow \text{sign}(Sk) \in \{-1, 1\}^m$</div></div> <div><div>Our KV Quantization</div><div>$K \xrightarrow{\text{QJL}} \text{cache} \xrightarrow{\text{uniform quantization}} V$</div><div>Lemma 3.3</div><div>$\langle Sq, \text{sign}(Sk) \rangle \approx_\epsilon \langle q, k \rangle$</div></div>		
<div>■ Recommended courses & Career after graduation</div> <div>- Programming structures (EE209, EE309, CS101, CS109, etc)</div> <div>- Data structure and algorithm (EE205, CS206, CS300, etc)</div> <div>- Machine learning (EE331, EE412, EE424, CS475, etc)</div> <div>- Mathematics (EE210, MS250, MS350, MS355, MS365, etc)</div> <div>- Industrial experiences before/after graduate school</div>		<div>■ Introduction to other activities besides research</div> <div>- Physical activities that achieve progressive goals, such as hiking, cycling, swimming</div>
<div>■ Introduction to the Lab.</div> <div>Established in September 2024, the lab offers you mentorship from highly motivated young professor working on cutting-edge researches, and the unique opportunity to become one of the first alumni of the group.</div>		
<div>■ Recent research achievements ('22~'25)</div> <div>- PolarQuant: Quantizing KV Caches with Polar Transformation, Under Review 2025</div> <div>- BalanceKV: KV Cache Compression through Discrepancy Theory, Under Review 2025</div> <div>- CalibQuant: 1-Bit KV Cache Quantization for Multimodal LLMs, ICML 2025 LCFM workshop</div> <div>- QJL: 1-Bit Quantized JL Transform for KV Cache Quantization with Zero Overhead, AAAI, 2025</div> <div>- Cell2Sentence: Teaching Large Language Models the Language of Biology, ICML 2024</div> <div>- HyperAttention: Long-context Attention in Near-Linear Time, ICLR 2024</div> <div>- Near Optimal Reconstruction of Spherical Harmonic Expansions, NeurIPS 2023</div> <div>- KDEformer: Accelerating Transformers via Kernel Density Estimation, ICML 2023</div> <div>- Fast Neural Kernel Embeddings for General Activations, NeurIPS 2022</div> <div>- Scalable MCMC Sampling for Nonsymmetric Determinantal Point Processes, ICML 2022</div> <div>- Random Gegenbauer Features for Scalable Kernel Methods, ICML 2022</div>		