

# Heejun Lee

Seoul, South Korea | ainl@kaist.ac.kr | +82 10-7757-5176 | <https://github.com/gmlwns2000>

## Professional Summary

Ph.D. student at KAIST specializing in efficient deep learning and large language models. Proven track record of developing novel sparse attention mechanisms with multiple first-author publications at top-tier conferences (ICLR). Passionate about inventing cutting-edge research and translating it into cost-effective, high-performance AI products.

## Education

Korea Advanced Institute of Science and Technology, Combined M.S./Ph.D. in Artificial Intelligence	Sep 2024 – Feb 2029 (Expected)
Korea Advanced Institute of Science and Technology, BS in Computer Science	Mar 2020 – Aug 2024
<ul style="list-style-type: none"><li>GPA: 3.97/4.3</li><li>College of Engineering Dean's List (Spring 2022)</li><li>College of Engineering Leadership Award on Research Excellence (Spring 2022, Spring 2023)</li></ul>	

## Experience

AI Research Engineer, DeepAuto.ai – Seoul, South Korea	Dec 2023 – Present
<ul style="list-style-type: none"><li>Developed <b>ScaleServe</b>, a cost-efficient LLM serving framework that <b>reduces end-to-end serving costs by approximately 52%</b> by integrating novel, training-free attention mechanisms.</li><li>Invented <b>HiP Attention (ICLR 2025)</b>, a training-free attention algorithm that <b>speeds up long-context inference by 50%</b> and enables serving million-token contexts on a single GPU via KV cache offloading.</li><li>Designed <b>Delta Attention</b>, a novel correction algorithm that <b>boosts sparse attention accuracy by 20-30%</b> on the RULER benchmark with only a marginal (<math>&lt;10\%</math>) latency overhead.</li><li>Engineered and integrated custom attention modules into serving frameworks like vLLM and SGLang, reducing computational complexity for long contexts from quadratic (<math>O(n^2)</math>) to near-linear (<math>O(n)</math>).</li></ul>	

## Publications

*\* Denotes equal contribution*

<b>Delta Attention: Fast and Accurate Sparse Attention Inference by Delta Correction</b>	arXiv Preprint
Jeffery Willette, <i>Heejun Lee</i> , Sung Ju Hwang (Github)	
<b>InfiniteHiP: Extending Language Model Context Up to 3 Million Tokens on a Single GPU</b>	arXiv Preprint
<i>Heejun Lee*</i> , Geon Park*, Jaduk Suh*, Sung Ju Hwang (Github)	
<b>A Training-free Sub-quadratic Cost Transformer Model Serving Framework With Hierarchically Pruned Attention</b>	ICLR 2025
<i>Heejun Lee*</i> , Geon Park*, Youngwan Lee*, Jaduk Suh*, et al. (Github)	
<b>Training-Free Exponential Extension of Sliding Window Context with Cascading KV Cache</b>	ICLR 2025
Jeffrey Willette, <i>Heejun Lee</i> , Youngwan Lee, Myeongjae Jeon, Sung Ju Hwang (Github)	
<b>SEA: Sparse Linear Attention with Estimated Attention Mask</b>	ICLR 2024
<i>Heejun Lee</i> , Jina Kim, Jeffery Willette, Sung Ju Hwang	

(Github)

## **Sparse Token Transformer with Attention Back Tracking**

ICLR 2023

*Heejun Lee*, Minki Kang, Youngwan Lee, Sung Ju Hwang

(Github)

## **Skills**

---

**Languages:** Python, C++, C#

**Frameworks & Libraries:** PyTorch, Hugging Face, vLLM, SGLang, OpenAI Triton, .NET