

# ELASWAVE: An Elastic-Native System for Scalable Hybrid-Parallel Training

Xueze Kang<sup>1†</sup>, Guangyu Xiang<sup>1†</sup>, Yuxin Wang<sup>4\*</sup>, Hao Zhang<sup>4</sup>, Yuchu Fang<sup>4</sup>,  
Yuhang Zhou<sup>3</sup>, Zhenheng Tang<sup>2</sup>, Youhui Lv<sup>5</sup>, Eliran Maman<sup>7</sup>, Mark Wasserman<sup>7</sup>,  
Alon Zameret<sup>7</sup>, Zhipeng Bian<sup>6</sup>, Shushu Chen<sup>5</sup>, Zhiyou Yu<sup>5</sup>, Jin Wang<sup>5</sup>, Xiaoyu Wu<sup>4</sup>,  
Yang Zheng<sup>4</sup>, Chen Tian<sup>3</sup>, Xiaowen Chu<sup>1\*</sup>

<sup>1</sup>*HKUST(GZ)* <sup>2</sup>*HKUST* <sup>3</sup>*NJU*  
*Huawei* <sup>4</sup>*TTE Lab* <sup>5</sup>*ICT BG* <sup>6</sup>*Cloud* <sup>7</sup>*TRC Team*

Large-scale LLM pretraining is rapidly moving to the  $10^5$ – $10^6$  accelerator scale, where failures are no longer exceptional events but part of normal operation. In this setting, efficient and elastic recovery is not just a backup feature; it is becoming a core systems requirement at production-scale. In this talk, we will briefly discuss the path toward an industrially deployable elastic-native training stack. We will focus on the design intuition behind ELASWAVE, the key mechanisms that make per-step elastic recovery practical, and the current results from our prototype.

We start building ELASWAVE from a simple observation: practical elastic training must jointly preserve parameter consistency and computation consistency, while also keeping recovery time (MTTR) low and sustaining high throughput after topology changes. Existing systems typically optimize only a subset of these goals. ELASWAVE addresses these challenges with a per-step fault-tolerant design that treats elasticity as multi-dimensional scheduling across graph, dataflow, DVFS, and RNG. It preserves training semantics while adapting execution online through micro-batch/pipeline resharding, asynchronous parameter migration, optimizer-partitioned recovery, DVFS-based bubble absorption, and RNG resharding for computation consistency, supported by a dynamic communicator and per-step in-memory snapshots for in-place group edits, verification, and redistribution.

Preliminary results on 96 NPUs show that ELASWAVE is promising relative to state-of-the-art baselines, with improvements in throughput, recovery efficiency, and training stability. As a work in progress, we are improving the system towards an industrially deployable stack.

---

<sup>†</sup>Equal contribution.

\*Corresponding authors: [yuxin.wang11@huawei.com](mailto:yuxin.wang11@huawei.com), [xwchu@hkust-gz.edu.cn](mailto:xwchu@hkust-gz.edu.cn).