

Workload	Container Orchestrator	Distributed compute engine	Training & inference framework	Outcomes
Last-mile data processing +model training	<b>Kubernetes</b> -based PinCompute platform	Previously Spark, now <b>Ray</b> due to streaming execution and heterogeneous CPU + GPU compute	<b>PyTorch</b>	Dataset-iteration wall-clock cut by 6x (from 90h to 15h). Dev cycles cut from days to hours. GPU utilization over 90%. Training throughput up 45% while per-job cost down 25%. ✧ ( <a href="#">blog</a> )
Offline / batch inference	<b>Kubernetes</b> -based PinCompute platform	Previously Spark, now <b>Ray</b>	<b>PyTorch, vLLM</b>	30x decrease in cost for search-quality jobs. 4.5x throughput increase. 4x reduction in job runtime for GPU inference jobs. Running around 1800 jobs per month. ✧ ( <a href="#">blog</a> , <a href="#">blog</a> )
Large-scale training Especially recommender models	<b>Kubernetes</b> -based PinCompute platform	<b>Ray</b> on heterogeneous clusters	<b>PyTorch</b>	Running 5000+ training jobs per month. Transparent scaling. Greater developer velocity and interactive development. Enable heterogeneous resources to keep GPUs saturated. ✧ ( <a href="#">blog</a> , <a href="#">talk</a> )