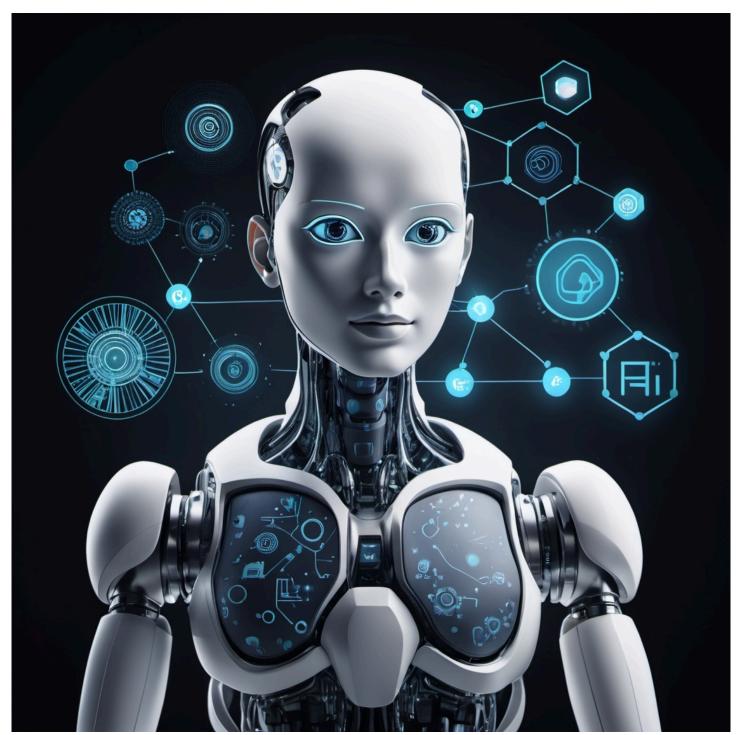
NeurochainAl's Solution for Cutting Al Compute Costs by Up to 68%

Streamlining GenAl Inference with Unique NeurochainAl's Inference Routing Solution that saves Al Compute costs drastically.



Florida City, Florida Nov 21, 2024 (<u>Issuewire.com</u>) - Today enterprises are racing to adopt GenAl to stay ahead of the competition and optimize processes to supercharge performance and slash costs of their business. However, the steep prices of Al compute infrastructure turn many of the GenAl solutions not feasible cost-wise.

Cutting Down Al Inference Costs by up to 68%

NeurochainAl's IRS: Transforming Al Compute Efficiency The Inference Routing Solution from NeurochainAl sets itself apart from conventional cloud services like AWS and Azure with its innovative routing and optimization features. By enabling dynamic distribution of Al workloads across a decentralized network, the IRS delivers unparalleled scalability and dependability. This cutting-edge system draws parallels with blockchain technology in terms of its operational agility and adaptability.

Key attributes of the IRS include:

- Al Model Optimization (Quantization): Reduces the precision of numbers in neural networks without sacrificing accuracy, leading to smaller (yet as accurate), faster models that use less computational power.
- IRS Software (P2P Load Balancer): A unique software combination of a load balancer and network orchestrator to monitor the availability of GPUs in relation to AI tasks and optimizes the assigning of all AI tasks to servers or virtual machines based on the optimal execution. It aso enables multiple AI models to run efficiently on a single GPU, reducing the total number of GPUs required, which dramatically lowers infrastructure costs.
- Infrastructure Set-Up Analysis and Consultation: Due to the unique nature of IRS Software and its abilities, we review, analyze, and suggest ways to optimize the existing infrastructure composition.

IRS software layer on top of the existing infrastructure can be implemented in 2 ways:

- **Fully managed service** where NeurochainAl performs the implementation of IRS on top of the infrastructure and provide maintenance and monitoring services.
- **Self-managed service** where NeurochainAl provides end-points and instructions on how to run compute nodes and the client's team does the implementation.

NeurochainAl's IRS: Versatile, Cost-Efficient Al Compute Platform For Web2 enterprises where Al plays a pivotal role, NeurochainAl's IRS presents a versatile and economical solution. With Al adoption already at 40% among global firms and another 42% in the exploration phase, NeurochainAl tackles one of the primary obstacles to Al implementation: the steep costs associated with compute infrastructure.

As mentioned above, the IRS can be deployed either as a fully managed or self-managed service, catering to diverse business requirements. Its compatibility with popular open-source AI frameworks like Llama and Mistral allows companies to leverage their existing AI investments while reaping the benefits of IRS optimizations.

By reducing infrastructure needs, businesses can redirect resources towards enhancing AI applications and expediting product development. As projections indicate the global AI market will surge to \$1.8 trillion by the end of the decade, NeurochainAI is well-positioned to play a crucial role in democratizing AI technology, making it more accessible, scalable, and efficient for businesses across the board.

Media Contact

NeurochainAl

hello@neurochain.ai

Source: NeurochainAl

See on IssueWire