(54) Title of the invention : "Next-Generation Cloud Architectures for Large-Scale Generative AI Systems"

| | |
|---|---|
| (51) International classification :G06F0009500000, G06F0021620000, G06N0003063000, G06N0003045000, G06F0009455000 | (71)**Name of Applicant :**<br>  1)MURUGALAKSHMI k<br>    Address of Applicant :Assistant Professor, Department of Artificial Intelligence and Data Science, VSB College of Engineering Technical Campus, Kinathukadavu. ----------- -----------<br>  2)Dr. R. Murugadoss<br>  3)Mrs. D.S. Jayakumari<br>  4)Ms. D.Jeevitha<br>  5)Mrs. B. Jebaranjani<br>  6)Mr. S. Soundhar<br>  7)Mr. F. Theophilus<br>  8)Mrs. G. Nithya<br>  9)Ms. G. Yogapriya<br>**Name of Applicant : NA**<br>**Address of Applicant : NA** |
| (86) International Application No :NA<br>   Filing Date :NA | |
| (87) International Publication No : NA | (72)**Name of Inventor :**<br>  1)MURUGALAKSHMI k<br>Address of Applicant :Assistant Professor, Department of Artificial Intelligence and Data Science, VSB College of Engineering Technical Campus, Kinathukadavu. ----------- -----------<br>  2)Dr. R. Murugadoss<br>Address of Applicant :Professor, Department of AI & DS, V.S.B College of Engineering and Technical Campus, Coimbatore Coimbatore ----------- -----------<br>  3)Mrs. D.S. Jayakumari<br>Address of Applicant :Assistant Professor, Department of AI & DS, V.S.B College of Engineering and Technical Campus, Coimbatore coimbatore ----------- -----------<br>  4)Ms. D.Jeevitha<br>Address of Applicant :Assistant Professor, Department of AI & DS, V.S.B College of Engineering and Technical Campus, Coimbatore coimbatore ----------- -----------<br>  5)Mrs. B. Jebaranjani<br>Address of Applicant :Assistant Professor, Department of AI & DS, V.S.B College of Engineering and Technical Campus, Coimbatore coimbatore ----------- -----------<br>  6)Mr. S. Soundhar<br>Address of Applicant :Assistant Professor, Department of AI & DS, V.S.B College of Engineering and Technical Campus, Coimbatore coimbatore ----------- -----------<br>  7)Mr. F. Theophilus<br>Address of Applicant :Assistant Professor, Department of AI & DS, V.S.B College of Engineering and Technical Campus, Coimbatore coimbatore ----------- -----------<br>  8)Mrs. G. Nithya<br>Address of Applicant :Assistant Professor, Department of AI & DS, V.S.B College of Engineering and Technical Campus, Coimbatore coimbatore ----------- -----------<br>  9)Ms. G. Yogapriya<br>Address of Applicant :Assistant Professor, Department of AI & DS, V.S.B College of Engineering and Technical Campus, Coimbatore coimbatore ----------- ----------- |
| (61) Patent of Addition to Application Number :NA<br>   Filing Date :NA | |
| (62) Divisional to Application Number :NA<br>   Filing Date :NA | |

(57) Abstract :

Abstract The rapid evolution of generative AI systems—such as large language models (LLMs), diffusion models, and multi-modal transformers—demands a fundamental rethinking of cloud architectures to meet the scale, latency, security, and efficiency requirements of modern workloads. This paper explores next-generation cloud architectures designed to support large-scale generative AI systems. It presents a layered framework that integrates heterogeneous computing (GPUs, TPUs, FPGAs), high-throughput networking, advanced memory hierarchies, and disaggregated storage. The architecture leverages innovations in distributed training, model parallelism, and elastic inference serving, orchestrated via AI-native platforms such as Kubernetes with custom scheduling and autoscaling policies. Moreover, it addresses critical challenges in energy efficiency, cost optimization, and compliance with data privacy regulations through federated and edge-cloud hybrid deployments. We also examine the role of AI accelerators, service meshes, serverless AI pipelines, and zero-trust security models in enabling reliable, real-time AI model delivery. The proposed architectures aim to provide scalable, resilient, and sustainable infrastructure foundations for the future of generative AI in industries ranging from healthcare and finance to autonomous systems and creative content generation.

No. of Pages : 9  No. of Claims : 10