

# DeepSeek: Rethinking AI Efficiency and Performance

DeepSeek's success exemplifies the potential of efficiency-driven AI, reinforcing Kaiser's prior insights on decentralized model training and the development of smaller, use-case-specific models to overcome power and chip constraints. This shift challenges traditional scaling and necessitates innovation from chip manufacturers, AI model developers, and enterprise software companies to optimize hardware, refine architectures, and advance AI-driven solution development.



## GenAI Application Implications

Enterprises integrating AI with core product suites must assess whether efficiency-driven models align with customer needs, costs, and differentiation, while exploring use-case-driven models as a scalable alternative

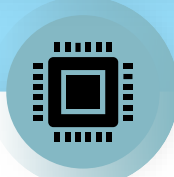
- + Technology providers bringing to market new AI features should explore **how use-case-driven, efficiency-focused models or agents can be tailored to meet customer needs** while balancing performance expectations, responsiveness, and computational efficiency
- + As these companies look to differentiate, they will need to prove that **efficiency-driven models can match or exceed the performance of larger, high-powered alternatives** while maintaining cost-effectiveness
- + The **lower COGS of efficiency-driven AI models enables broader AI deployment across product suites**, making AI integration more feasible while reshaping AI pricing and cost structures



## AI Model Developer Implications

AI model developers must explore new architectures that optimize for efficiency without sacrificing performance, ensuring models remain competitive in both capability and scalability

- + Balancing power efficiency with model performance will require the **integration of RLHF and CoT techniques used by DeepSeek and new techniques such as enhanced MoE and quantization** to maintain reasoning ability, accuracy, and generalization while operating with lower resource demands
- + Efficiency-driven models challenge the notion that scaling large models is the primary way forward, pushing model developers to **rethink parameter optimization and architectural design for smaller, use-case driven models**
- + Increased **focus on efficiency should lead to a shift toward specialized, domain-specific models** that optimize for industries or applications rather than broad, general-purpose AI



## Chip Manufacturer Implications

Chip manufacturers must adapt by optimizing AI hardware for efficiency-first architectures, refining AI accelerator strategies, and aligning with the evolving economics of AI-as-a-Service

- + Increased efficiency in AI models will drive demand for **AI-specific chips optimized for power savings**, requiring hardware manufacturers to focus on performance-per-watt rather than raw computational power
- + **Lower power requirements could reshape AI-as-a-Service offerings**, leading cloud providers to rethink infrastructure investments, prioritize AI accelerators for more efficient compute, and introduce new pricing structures based on power consumption rather than sheer processing capability
- + Dominant **GPU programming languages like CUDA may need optimizations** to better support efficiency-focused models, ensuring they fully utilize AI accelerators without unnecessary power overhead

**For over 40 years, Kaiser Associates has employed robust methodologies to craft innovative go-to-market strategies. Our TMT practice is experienced in empowering clients to navigate and seize strategic opportunities in the rapidly expanding field of GenAI.**

**Connect with Kaiser's Technology Practice experts to learn more about how Kaiser can support your business:**



**Joe Kestel**  
VICE PRESIDENT  
JKESTEL@KAISERASSOCIATES.COM



**Azulina Green**  
VICE PRESIDENT  
AGREEN@KAISERASSOCIATES.COM



**Rodger Heidgerken**  
SR. MANAGER  
RHEIDGERKEN@KAISERASSOCIATES.COM