

White Paper

# Small But Mighty: The Strategic Advantage of Small Language Models

Sponsored by: Dell Technologies

Arnal Dayaratna

February 2025

## SITUATION OVERVIEW

---

The maturation of generative AI has driven widespread adoption of large language models (LLMs) that have catalyzed a profound transformation of contemporary application development practices. For starters, foundation models such as LLMs have become central to application development because of their ability to accelerate development and expand the universe of use cases that are amenable to digitization. As LLMs become increasingly central to contemporary application development, organizations face a growing need to evaluate a broad and diverse range of options. Large language models, while powerful, often come with significant resource demands and operational complexities. This has prompted questions about how organizations can select models that align with their technical capabilities and strategic objectives. These challenges have opened the door for the proliferation of alternative approaches, particularly the adoption of small language models (SLMs) as a practical and scalable solution.

The proliferation of small language models underscores their growing relevance in addressing the challenges posed by LLMs. With their reduced computational requirements, SLMs offer organizations a path to AI adoption that balances performance with efficiency. For technology suppliers, this means tailoring solutions to meet the unique needs of enterprises seeking cost-effective and flexible AI options. The adaptability of SLMs to cloud, edge, and mobile environments positions them as ideal for enterprises aiming to deploy AI in diverse operational contexts. Moreover, their ease of customization allows organizations to accelerate time to market and optimize AI solutions for specific use cases. As enterprises continue to navigate the complexities of model selection, technology suppliers would do well to focus on providing solutions that emphasize the accessibility, privacy, and efficiency of small language models.

## DEFINITION

---

### **Defining Small Language Models**

A small language model is an AI model designed to process, generate, or analyze language, characterized by having fewer than 10 billion parameters. This parameter limit distinguishes SLMs from large-scale models such as GPT-4, Llama 3, or Anthropic Claude 3.5, which often contain tens or hundreds of billions of parameters. The streamlined architecture of SLMs enables them to operate efficiently, making them ideal for applications where computational resources, energy efficiency, or privacy concerns are critical.

Optimized for efficiency and lightweight deployments, SLMs are well suited for constrained environments such as edge devices, mobile applications, and IoT sensors. In these settings, limited computational capacity and energy restrictions often render large models impractical. By focusing on delivering solid performance in targeted use cases, SLMs offer organizations the ability to implement AI solutions that align closely with their specific operational needs and resource constraints.

## BENEFITS

---

### **Benefits of Small Language Models**

The unique attributes of small language models confer a range of benefits that make them increasingly valuable in modern AI development. The key advantages are discussed in the sections that follow.

#### **Lower Resource Requirements**

Small language models require significantly fewer computational resources than their large language model counterparts. This reduced requirement for computational resources allows organizations to deploy AI solutions on hardware that would otherwise be incapable of supporting large models. For start-ups and smaller enterprises, this means they can leverage advanced AI capabilities without the need for substantial investment in high-end infrastructure. The reduced resource demand also facilitates smoother integration into existing systems in ways that minimize disruption and accelerate implementation timelines. In addition, lower resource requirements reduce operational complexities, enabling organizations in emerging markets or those with limited technical infrastructure to adopt AI technologies more readily. This inclusivity fosters innovation across a wider spectrum of industries and regions.

## **Cost Efficiency**

The diminished computational needs of SLMs enable lower infrastructure and maintenance expenses. Organizations can achieve substantial cost savings on hardware, energy consumption, and operational overhead. This cost efficiency is appealing to businesses operating under budget constraints or those seeking to maximize return on investment. By enabling affordable access to AI technologies, SLMs democratize innovation and allow a broader range of companies to compete and innovate in the AI space. Moreover, the reduced financial burden permits organizations to reallocate resources toward other strategic initiatives such as research and development, marketing efforts, or acquiring top talent, thereby strengthening their competitive position.

## **Faster Inference Times**

The compact size of small language models results in faster response times, which is crucial for real-time applications. Industries such as finance, healthcare, and customer service rely on instantaneous data processing to make timely decisions and enhance user experiences. Faster inference times enable applications such as chatbots, recommendation engines, and predictive analytics to operate with minimal latency, thereby increasing performance and user satisfaction. This responsiveness can be a significant differentiator in competitive markets where speed is essential. In addition, faster processing allows systems to handle higher volumes of requests efficiently that enable scalability and reliability during peak usage periods. In time-sensitive scenarios, such as emergency response or fraud detection, quicker inference can have critical impacts on outcomes and safety.

## **Greater Privacy Control**

The ability of SLMs to run locally on edge devices enhances data privacy and security by eliminating the need to transmit sensitive information to remote servers. This feature is invaluable in sectors that handle confidential data, such as healthcare, finance, and legal services. By processing data on device, organizations can mitigate risks associated with data breaches and comply more easily with stringent privacy regulations. This localized processing also builds trust with customers and stakeholders by demonstrating a commitment to safeguarding personal information. Furthermore, keeping data within local systems can simplify compliance with international data protection laws by avoiding complexities related to cross-border data transfers. Enhanced privacy controls not only protect the organization but also contribute to higher customer satisfaction and loyalty.

## **Customization and Fine-Tuning**

Small language models require fewer computational resources for fine-tuning, making them a cost-effective choice for developing specialized AI solutions. Their reduced size leads to shorter training times and less demand on hardware, which lowers expenses associated with model development. Organizations can tailor SLMs to address specific industry challenges, niche markets, or unique operational needs without incurring significant costs. For instance, a company operating in the legal sector might fine-tune an SLM to efficiently process legal documents and contracts, enhancing productivity without substantial investment.

The reduced complexity of SLMs simplifies the fine-tuning process, enabling faster iteration and deployment cycles. Developers can experiment with different configurations more readily, refining models swiftly to achieve optimal performance. This agility is particularly beneficial for start-ups and small businesses looking to adapt quickly to market demands and validate their offerings in a competitive environment. By accelerating development timelines, these organizations can respond to customer feedback promptly and correspondingly improve their products continuously and gain a competitive edge.

In addition, the ease of customization allows organizations to update models more frequently to keep pace with evolving industry standards, regulations, or customer preferences. In fast-moving industries, the ability to adjust AI models swiftly ensures that solutions remain relevant, compliant, and effective. The capability to rapidly prototype and test models not only accelerates innovation cycles but also fosters a culture of continuous improvement within the organization. This approach helps businesses stay competitive by consistently enhancing their AI applications to better meet user needs and expectations.

Furthermore, the lower barriers to fine-tuning mean that domain experts can collaborate more closely with AI developers. Subject matter experts can contribute their specialized knowledge during the customization process, improving the accuracy and relevance of the models. This collaborative effort enhances the overall quality of AI solutions and ensures they are closely aligned with organizational goals and industry best practices.

## **Deployment Flexibility**

The adaptability of SLMs allows them to function effectively in low-bandwidth environments or on devices with limited processing power. This flexibility extends the reach of AI applications beyond traditional datacenters and into everyday consumer products, remote locations, and emerging markets. By overcoming barriers associated with infrastructure limitations, SLMs enable organizations to tap into new customer

segments and create innovative solutions that were previously unattainable due to technological constraints. In addition, SLMs can be deployed in environments with intermittent or unreliable connectivity, supporting applications in remote or underdeveloped areas. Flexible deployment options also facilitate compliance with local regulations regarding data storage and processing, enabling organizations to operate globally with greater ease.

## **Reduced Risk of Hallucinations**

Because SLMs tend to be rigorously fine-tuned or leverage implementations of retrieval-augmented generation (RAG), they tend to produce more predictable and reliable outputs that reduce the likelihood of hallucinations or unpredictable responses. This reliability is crucial for applications where accuracy and consistency are paramount, such as legal document analysis, medical diagnostics, or financial forecasting. By minimizing errors and enhancing output quality, SLMs help organizations maintain credibility and trust with their users. Improved reliability also supports compliance with industry regulations that require consistent and accurate data processing. Moreover, the reduced need for human oversight in monitoring AI outputs can lower operational costs and allow staff to focus on more strategic tasks.

## **Adaptability to Low-Power Devices**

Small language models are well suited for deployment on low-power devices such as mobile phones, tablets, and IoT sensors. This adaptability enables AI functionalities in contexts where power consumption is a critical concern, such as wearable technology, remote monitoring systems, and environmental sensors. By facilitating AI integration into these devices, SLMs expand the potential for data collection, analysis, and decision-making in real time, even in resource-constrained environments. This capability promotes the development of new products and services that leverage AI in innovative ways, opening up opportunities in sectors such as agriculture, energy, and transportation. Integrating AI into low-power devices also provides access to valuable data and insights that were previously inaccessible, driving informed decision-making and strategic planning.

## **Scalable Maintenance**

The reduced complexity of SLMs simplifies maintenance and updates, allowing businesses to scale and manage their AI infrastructure without significant operational overhead. Organizations can implement updates, security patches, and performance enhancements more efficiently, ensuring that their AI solutions remain current and effective. This scalability is essential for companies looking to grow their AI capabilities alongside their business without incurring prohibitive costs or resource demands. Simplified maintenance also reduces system downtime, enhancing availability and

reliability for end users. In addition, the ease of updating SLMs facilitates the incorporation of the latest advancements in AI technology, ensuring that organizations remain at the forefront of innovation.

## **Sustainability**

The lower energy consumption of SLMs supports sustainable computing practices that reduce the environmental impact of deploying AI solutions. As organizations increasingly prioritize sustainability and corporate social responsibility, adopting SLMs aligns with environmental goals and demonstrates a commitment to reducing carbon footprints. This sustainability not only benefits the planet but also enhances the organization's reputation among eco-conscious consumers and partners. Energy-efficient AI solutions can also lead to long-term cost savings in energy expenditures, contributing to the organization's financial sustainability. Moreover, sustainable practices can be differentiators in the market, attracting customers and partners that value environmental responsibility and ethical operations.

## **TRENDS**

---

### **Trends Shaping Small Language Models**

The evolving landscape of AI has given rise to several trends that underscore the growing importance and utility of small language models. These trends reflect the shifting priorities of organizations and the innovative ways in which SLMs are being integrated into various applications.

#### **Edge Computing Integration**

The integration of small language models into edge computing environments continues to expand, enabling AI capabilities without the need for constant internet connectivity. By processing data locally on edge devices, organizations can achieve real-time analytics, reduce latency, and enhance the responsiveness of applications. This trend is significant for industries such as manufacturing, where immediate data processing can improve operational efficiency, and for consumer products that require seamless user experiences.

#### **Model Quantization for Efficient AI Deployment**

Model quantization reduces the memory footprint and computational requirements of language models by lowering a model's weights and activations from higher-precision formats such as 32-bit floating-point numbers to 8-bit integers. This optimization enhances efficiency by allowing small language models to run on resource-constrained hardware such as edge devices and mobile processors. Quantization also improves

latency and power consumption while maintaining accuracy for most applications. As AI adoption expands beyond cloud environments, quantization plays a critical role in making language models more accessible, cost effective, and practical for real-time, on-device, and privacy-focused deployments. Its growing adoption supports AI scalability across industries that require efficient and responsive machine learning protocols.

## **Hybrid Architectures**

There is a growing adoption of hybrid foundation model architectures that combine small language models with larger models. In this scenario, the smaller model typically handles more specialized tasks, while generalist computations are managed by larger systems. This approach optimizes resource utilization, balances performance with efficiency, and allows organizations to leverage the strengths of both model types. Hybrid architectures enable scalable solutions that can adapt to varying workload demands, improving overall system robustness.

## **Specialization for Specific Domains**

An increasing focus on training small language models for specific domains or industries allows for highly efficient and accurate performance in targeted use cases. By honing models on domain-specific data, organizations can achieve superior results in areas like medical diagnostics, legal analysis, or financial forecasting. This specialization enhances the relevance and value of AI solutions, providing organizations with tools that address their unique challenges more effectively than generalized models.

## **Privacy-Focused Solutions**

The emphasis on privacy-preserving AI has intensified interest in small language models that operate offline or locally. Organizations are increasingly seeking solutions that minimize data exposure and comply with stringent privacy regulations such as GDPR and CCPA. By utilizing SLMs in privacy-focused applications, companies can offer AI-driven services that protect user data, build consumer trust, and avoid potential legal repercussions associated with data breaches or noncompliance.

## **Open Source and Democratization**

The availability of open source small language models is fostering innovation and experimentation among developers and smaller organizations that may lack significant resources. Open source models lower barriers to entry and allow a wider range of contributors to improve and adapt AI technologies. This democratization accelerates the advancement of AI by encouraging collaborative development, knowledge sharing, and the proliferation of customized solutions tailored to diverse needs.

## **Growing Interest in Decentralized AI Systems**

The trend toward decentralized AI has increased the popularity of small language models capable of operating independently without relying on central servers. Decentralized systems enhance resilience, reduce single points of failure, and allow for more localized and personalized AI experiences. This approach is valuable in scenarios where connectivity is unreliable or where autonomy is desired, such as in remote monitoring, autonomous vehicles, or distributed sensor networks.

## **KEY CONSIDERATIONS**

---

### **Challenges in Small Language Model Adoption**

Despite the numerous advantages, organizations must navigate certain challenges when adopting small language models to ensure they align with strategic objectives and operational capabilities.

#### **Limited Generalization Ability**

Small language models may lack the capacity to generalize effectively across a wide range of topics, limiting their applicability in tasks that require broad language comprehension. This constraint necessitates careful consideration of the specific use cases where SLMs are appropriate. Organizations may need to supplement SLMs with additional models or systems to address tasks outside the SLM's specialized domain, potentially increasing complexity.

#### **Performance Trade-Offs**

The reduction in size and complexity of SLMs can lead to sacrifices in accuracy and sophistication. In applications that demand nuanced understanding, contextual awareness, or high levels of creativity, larger models may outperform SLMs. Organizations must evaluate whether the performance trade-offs are acceptable for their specific needs and consider hybrid solutions or complementary technologies to mitigate these limitations.

#### **Dependency on High-Quality Training Data**

The dependency of small language models on high-quality training data poses a challenge to their adoption because they require extensive, diverse, and accurate data sets to perform effectively. Inadequate or biased data can lead to poor model performance, limited generalization, and potential safety issues that make it difficult for organizations to trust and implement these models in real-world applications.

## DELL PRO MAX FOR SMALL LANGUAGE MODELS

---

To fully maximize the benefits of small language models, robust and reliable computing infrastructure is essential. Dell Pro Max offers a wide range of configurations and form factors that provide high-performance solutions that are well suited for developing, deploying, and fine-tuning SLMs. These workstations leverage the scalability derived from being configurable with one to four NVIDIA RTX professional graphics cards. These graphics cards deliver the requisite computational power and efficiency to optimize model performance. While SLMs require less compute resources from the GPU, the GPU remains crucial for running AI workloads efficiently. By investing in a solid hardware foundation, businesses can support the demanding computational tasks associated with AI development, leading to smoother workflows and more effective implementation of small language models.

### Fixed Versus Mobile Workstations

Dell Pro Max fixed (desktop, tower, and rack) workstations provide exceptional computational capabilities that surpass those of mobile alternatives. This makes them ideal for intensive AI development tasks that require significant processing power. Fixed workstations enable faster training cycles, which allows developers to iterate more quickly and increase model performance. They also support more complex model customization, giving organizations the flexibility to tailor SLMs to their specific needs. In addition, fixed workstations can handle large data sets efficiently, which is essential for training and fine-tuning models with substantial amounts of data. This level of performance is vital for organizations aiming to refine their small language models and achieve the best possible results. The stability and scalability offered by fixed workstations make them a strong choice for enterprises that prioritize high-performance computing in their AI strategies.

Mobile workstations, although offering slightly less raw computational power, provide significant advantages in terms of portability and flexibility. They are suitable for scenarios where mobility is essential, such as field deployments, onsite client presentations, or remote work environments. Mobile workstations allow developers to continue working on AI projects without being confined to a specific location. This enhances collaboration among team members and enables quick responses to new opportunities or challenges that may arise. Mobility supports agile development practices and ensures that teams can remain productive regardless of their physical location. For organizations that value flexibility and need to adapt to dynamic work settings, mobile workstations offer a practical solution that balances performance with convenience.

## Optimizing Small Language Models

Fixed workstations offer significant benefits for organizations deeply involved in extensive AI development workflows. They are equipped to handle the demanding computational processes required for fine-tuning small language models, which often involve iterative adjustments and evaluations. Fixed workstations facilitate the execution of complex simulations that can test model performance under various conditions. They also support the deployment of AI applications at scale, ensuring that models can handle production-level workloads effectively. By utilizing the high-performance capabilities of fixed workstations, businesses can accelerate development timelines, reduce time to market for their AI solutions, and enhance the overall accuracy and reliability of their models. This leads to improved efficiency and can provide a competitive advantage in the rapidly evolving AI industry. Organizations that invest in robust fixed workstations are better positioned to optimize their small language models and fully realize the potential benefits these models offer.

Mobile workstations complement fixed systems by facilitating ongoing development and deployment activities in a variety of settings. They enable organizations to extend their AI development capabilities beyond traditional office environments, supporting more agile operations. For example, developers can work on refining models while traveling, during client site visits, or in collaborative workspaces. This flexibility allows for rapid iteration and quick adjustments in response to market changes or customer feedback. For start-ups and small businesses that need to be highly responsive and adaptable, mobile workstations provide the tools necessary to maintain momentum and stay competitive. Supporting development activities wherever they are needed, mobile workstations contribute to more dynamic and efficient AI workflows. The ability to work seamlessly across different locations can enhance collaboration and innovation within teams, further optimizing the development and deployment of small language models.

## CONCLUSION

---

Small language models are redefining the future of AI development and represent a significant evolution in the realm of foundation models. Their ability to deliver high performance with enhanced efficiency positions them as a pivotal force in the next generation of AI technologies. As computational constraints, sustainability concerns, and the need for privacy-preserving solutions become more critical, SLMs offer practical answers that align with these emerging priorities.

By providing greater accessibility and adaptability, SLMs empower organizations of all sizes to integrate advanced AI capabilities into their operations. They enable the creation of specialized applications that address specific industry challenges and niche

markets, supporting innovation across a diverse range of sectors. This shift reflects a broader movement toward AI systems that are not only powerful but also efficient, scalable, and tailored to real-world constraints.

Embracing small language models allows organizations to stay at the forefront of technological advancements. These models facilitate a move away from reliance on massive, resource-intensive systems toward more sustainable and deployable solutions. Their development encourages a future where AI is more accessible, environmentally friendly, and aligned with the nuanced needs of various industries.

By leveraging robust computing infrastructure to support the deployment and performance of small language models, businesses can optimize their AI initiatives. This approach accelerates development cycles, enhances model effectiveness, and positions organizations to meet market demands with agility and precision. In an increasingly competitive and rapidly evolving technological environment, small language models offer a pathway to harnessing AI's full potential, driving progress, and achieving strategic objectives.

Overall, small language models are not merely an alternative to their larger counterparts; they represent a fundamental shift toward more intelligent, efficient, and practical AI solutions. Their growing importance signals a future where AI development prioritizes efficiency without compromising performance. By adopting small language models, organizations can confidently navigate the complexities of modern AI, ensuring they remain competitive and innovative in the years to come.

## ABOUT IDC

---

International Data Corporation (IDC) is the premier global provider of market intelligence, advisory services, and events for the information technology, telecommunications, and consumer technology markets. With more than 1,300 analysts worldwide, IDC offers global, regional, and local expertise on technology, IT benchmarking and sourcing, and industry opportunities and trends in over 110 countries. IDC's analysis and insight helps IT professionals, business executives, and the investment community to make fact-based technology decisions and to achieve their key business objectives. Founded in 1964, IDC is a wholly owned subsidiary of International Data Group (IDG, Inc.).

### Global Headquarters

140 Kendrick Street  
Building B  
Needham, MA 02494  
USA  
508.872.8200  
Twitter: @IDC  
blogs.idc.com  
www.idc.com

---

#### Copyright Notice

External Publication of IDC Information and Data — Any IDC information that is to be used in advertising, press releases, or promotional materials requires prior written approval from the appropriate IDC Vice President or Country Manager. A draft of the proposed document should accompany any such request. IDC reserves the right to deny approval of external usage for any reason.

Copyright 2025 IDC. Reproduction without written permission is completely forbidden.