

Decentralizing Generative Al (GenAl) Inference

On-Device Deployment of Lightweight Open Source GenAl Models, Including Large Language Models (LLMs), Can Improve Accessibility and Latency

Author

Bharath G. Srivats

Product Marketing Engineer

Credits

Sanjay Addicam Ilia Efimov Pinkesh Shah Balaji Srinivasan Michael Campbell Helena Kloosterman Michael Hansen

Abstract

Bringing Al inference closer to the data source offers significant advantages in cost, privacy and performance. Recent advancements in light weight GenAl models (i.e., 1-8B parameters) provide a disruptive opportunity to shift GenAl deployment from the cloud to the edge, but alternatives to cloud-based GenAl need to be practical and efficient. This white paper outlines a strategic approach to shift GenAl deployments from cloud-native (i.e., GPU based) solutions to edge (i.e., hardware based) solutions using the built-in compute acceleration of CPU-GPU-NPU (e.g., Intel® Core™ Ultra processors, Intel® Arc™ GPUs) and open source GenAl models. On-device deployment offers low total cost of ownership (TCO), offline capabilities, data sovereignty and reduced latency, making powerful GenAl models accessible across regions and sectors that may previously have faced barriers to deployment.

Table of Contents

Abstract1
Introduction1
Why it Matters2
Technical Performance Benefits of On-Device GenAl
Benchmarking GenAl on OpenVINO3
Cost, Privacy and Accessibility Benefits of Edge-Based GenAl4
The Impact of Decentralized AI on Education in Emerging Markets: Use Case Deep Dive5
Conclusion 5

Introduction

Large, sophisticated GenAl models have immense computational needs, and as a result, cloud-based GenAl solutions have dominated. However, these models require substantial infrastructure investment or recurring end-user API costs, which limits access to their advanced capabilities.

Recent advances in GenAl open ecosystem development are changing the landscape; models are optimized to be smaller, faster and less power hungry, enabling edge-based, decentralized deployments with the ability to run GenAl models locally. This shift from cloud-based Al towards offline, edge-based Al presents a compelling deployment alternative — one that democratizes GenAl access away from centralized infrastructure and offers benefits like reduced costs and increased data privacy by leveraging hardware (CPU-GPU-NPU).

On-device GenAl could have wide-reaching positive impacts, increasing access to GenAl tools in circumstances, locations or sectors that face challenges like lack of internet, low bandwidth, shortage of skilled people, etc.

1

Why it Matters

Low Cost

Localized deployments of on-device GenAI applications do not require incremental subscriptions or fees. By optimizing lightweight models for existing underlying hardware configurations (i.e., leveraging sunk costs), organizations can achieve competitive GenAI performance at a fraction of the cost associated with cloud-based instances.

Education:

UNICEF estimates that nearly 1 billion children may face unreliable internet access and a shortage of qualified teachers. Low cost, offline access to GenAI-powered educational tools presents an opportunity to bridge this gap.

Data Privacy

On-device GenAl is complementary to agentic Al workflows², offering a local inference tier that can function offline at low latency and act upon private data enclaves for sovereign data intelligence.

Healthcare:

In crisis response situations or remote locations without reliable internet, healthcare providers could access GenAl tools that support diagnostics, triage or secure access to critical patient records.

Latency

While cloud-based Gen Al scales for high-end workloads, it often entails network latency. Lightweight models with low latency and on-device compute offer a hybrid approach; Al workloads can be allocated between the cloud and ondevice depending on user requirements, resource utilization and desired quality or user experience.

Automotive:

Low-latency, on-device GenAI tools could automate certain workflows at the point of decision-making or help streamline repairs at the point of service.

Technical Performance Benefits of On-Device GenAl

Edge solutions, in contrast to compute-intensive advanced cloud-based models, offer a balance between performance and computational efficiency. The following table demonstrates how recent advancements in lightweight, open-source GenAl models and Intel OpenVINOTM toolkit-based optimizations enable efficient deployment and orchestration on Intel® CoreTM Ultra processors³.

Example Solution Ingredients:

To demonstrate how on-device GenAI deployment can replicate cloud-based deployment and functionality, we ran the following scenarios on OpenVINO: chatbot, image generation and live captioning/audio transcription.

- LLM (Large Language Model): Llama, Phi, Qwen, SmolLM, Mistral, etc.
- LVM** (Large Vision Model): MiniCPM-V, Qwen-VL, Llava, SmolVLM, etc.
- Image Generation**: Stable Diffusion, Flux
- Audio Transcription**: Whisper
- Console**: OpenVINO Test Drive

Benchmarking GenAI on OpenVINO

As the following table illustrates, we found satisfactory performance (i.e., $20-30\,\text{TPS}$) on real-life use cases such as conversing with a chatbot.

Our benchmarking demonstrates that Lunar Lake (LNL) is approximately 1.6x more power efficient at equivalent* iGPU AI inferencing compared to Meteor Lake (MTL); and LNL is approximately 1.4x faster (i.e., throughput) at equivalent* iGPU AI inferencing compared to MTL. Power draw and other utilization results are shown below.

iGPU*^	Models, HW:	Throughput TPS (2nd tkn+)	Latency (ms) (2nd token)	At Full iGPU Utilization*, Measured:	Power Draw*	
Llama	Llama 3.1- 8B -int4, MTL	14.20*/16.48^	70.50*/60.67^	16.4% CPU utilization 59000 MB/s Mem BW Idle Mem: 12GB Load Mem: 21GB	Idle: 5.5W Load: 29W	
	Llama-3.1- 8B -int4, LNL	20.38*/20.91^	49.07*/47.80^	21.8% CPU utilization 78401 MB/s Mem BW Idle Mem: 6GB Load Mem: 15GB	Idle: 2.4W Load: 25W	
	Llama 3.2- 1B -int4, MTL	66.06*	15.13*	16.8% CPU utilization 49698 MB/s Mem BW Idle Mem: 11GB Load Mem: 13.5GB	Idle: 5.1W Load: 29W	
	Llama 3.2- 1B -int4, LNL	79.33*	12.59*	22% CPU utilization 64676 MB/s Mem BW Idle Mem: 6.2GB Load Mem: 8GB	ldle: 2.4W Load: 23.5W	
iGPU^	Models & Generations:	Throughput [^] TPS (2nd tkn+)	Latency^ (ms) (2nd token)			
Llama	Llama-3.2-3B-instruct-int4, MTL	31.76	31.47			
	Llama-3.2-3B-instruct-int4, LNL	40.69	24.57			
Phi	Phi-3-mini-4k-instruct-int4, MTL	24.59	33.95			
	Phi-3-mini-4k-instruct-int4, LNL	35.21	28.4			
Qwen	Qwen2.5- 7b -int4, MTL	17.68	56.53			
	Qwen2.5-7b-int4, LNL	21.62	46.24			
DeepSeek	Deepseek-rl-distill-qwen-7b int4, MTL	17.82	56.09			
	Deepseek-rl-distill-qwen-7b int4, LNL	21.68	46.12			
NPU*	Model, HW:	Throughput	Latency* (ms) (2nd token)	At Full NPU Utilization*, Measured:	Power Draw*	
TF3 (ZHUKH)						
NPU data is preliminary and not fully optimized for						
Llama	Llama 3.2- 1B- int4, MTL *	22.91*	43.5*	2.9% CPU utilization 13375 MB/s Mem BW Idle Mem:11GB Load Mem: 13.4GB	ldle: 5.1W Load: 16.45W	
	Llama 3.2- 1B -int4, LNL	25.97*	38.5*	3% CPU utilization 27030 MB/s Mem BW Idle Mem: 6.2GB Load Mem: 8.9GB	Idle: 2.4W Load: 15.0W	

Sources & Notes:

 $Open-source\ instruct-tuned\ models, OV-quantized/optimized, please\ see\ sources\ for\ additional\ details$

Figure 1. Benchmarking GenAl on Intel MTL, LNL and OpenVINO

^{*}Config 1: SKUs tested: Meteor Lake Core Ultra 7 MTL 165H, Lunar Lake Core Ultra 7 LNL 268V, OpenVINO 2024.6.0.0

^{*}Memory and power usage metrics were measured under a custom Q&A inference workload implemented with OpenVINO GenAl, reflecting measured system resource performance defined test parameters

^{*}Input token size = <25, output token size = 512

[^]Config 2: SKUs tested: Meteor Lake Core Ultra 7 MTL 155H, Lunar Core Ultra 7 LNL 288V, OpenVINO 2025.0

[^]input token size = 32, output token size = 1024

[^]https://docs.openvino.ai/2025/index.html

[^]https://docs.openvino.ai/2024/openvino-workflow/model-preparation.html

[^]https://docs.openvino.ai/2024/openvino-workflow/model-optimization.html

^{**}Coming soon - Whisper audio transcription, image time-to-generate, multimodal models benchmarking, Openvino Test Drive console (beta)

^{**}https://github.com/openvinotoolkit/openvino_testdrive (beta)

Cost, Privacy and Accessibility Benefits of Edge-Based GenAl

With Intel® Core™ Ultra processors, Intel® Arc™ GPUs and OpenVINO, enterprises can meet their edge and on-prem GenAI application needs.

Low Total Cost of Ownership (TCO)

Offline edge and enterprise GenAl deployments can supplement reliance on expensive cloud infrastructure, thereby reducing long-term API token costs. Once the hardware is deployed (i.e. sunk cost), open-source GenAl models can run offline, significantly reducing costs over an estimated five year depreciation period. Based on assumptions regarding task-specific performance and precision requirements, enterprises can optimize workload distribution between edge and cloud to best match the demands of each application. In the long run, enterprises may be able to realize significant savings by adopting a hybrid strategy to leverage more on-device functionality relative to cloud-based. As lightweight models' portability and quality continues to rapidly improve, they can be customized to fit in a hybrid framework and utilized for local intelligence.

The decentralization of GenAI through edge devices creates a more accessible entry point for developers and enterprises alike, lowering both capital expenditure (CAPEX) and operational expenditure (OPEX).

No Internet Needed

On-device GenAl models do not require internet to work. This is especially important in emerging markets where intermittent connectivity and poor bandwidth impact operations. For example, an on-device chatbot can be initialized in a matter of seconds from the time the system is powered on and remain available thereafter with nearinstantaneous response times and full offline functionality. Furthermore, the ability to deploy multimodal* GenAl models — text, audio, vision inputs and text, audio, vision outputs — on-device expands sophisticated offline use cases and could include image recognition*, image generation* or language translation, to name a few.

Optimized Performance

OpenVINO's sophisticated optimization framework maximizes performance through Intel GPU XMX acceleration and NPU capabilities, allowing maximum flexibility for developers through intelligent workload distribution across architectures. Independent Software Vendors (ISV) can create solutions with models across the GPU+NPU mix. Advanced quantization techniques reduce memory requirements while accelerating inference speeds, making 1-8B parameter models viable for on-device deployments that prioritize privacy protection, performance optimization and cost efficiency.

Improved Latency

Computation (e.g., TOPS) residing next to the data source or device offers fundamental latency benefits for running GenAI. Processing GenAI models directly on premand local edge deployments reduces latency to milliseconds—a critical advantage over cloud-based GenAI systems where latency can be affected by network bandwidth and remote server loads.

Examples (see Figure 1): Under 100ms* iGPU latency on lightweight LLM enables near-instantaneous intelligence, such as a yes/no that can be derived as alerts or decisions (i.e., agentic). Through sophisticated pipeline engineering (e.g., preprocessing, context loading), a response can be extracted in a few milliseconds, fundamentally accelerating inference for specialized scenarios. This low-latency verdict is ideal for simple yes/no or pass/fail classifications, providing an on-device gating superpower that outpaces cloud latency. This initial low-latency assessment can revolutionize system responsiveness in carefully selected, high-impact applications where immediate categorical decisions drive operational success, such as in healthcare and automation.

Data Sovereignty

On-premise and local GenAl application deployments ensure that sensitive data remains on the device, preserving privacy and adhering to regulations like the General Data Protection Regulation (GDPR) or Health Insurance Portability and Accountability Act (HIPAA). Because they don't need an internet connection to process GenAl in real-time, organizations gain enhanced control over their data while avoiding the risks of data breaches and leaks associated with cloud-based systems.

Cost-Effectiveness

On-premise and edge GenAl solutions present compelling economic and privacy advantages compared to cloud-based solutions. Organizations can significantly reduce or eliminate subscription and API expenses and achieve precise operational cost management with local inference capabilities.

Low Maintenance

Edge and on-premise GenAI models deployed on local infrastructure require only periodic manual updates to stay current. Given that lightweight models evolve rapidly — often in six-month cycles — an annual maintenance schedule is recommended for edge-deployed GenAI. This allows for critical updates to models and optimizations through OpenVINO release updates, ensuring the models remain efficient and effective for their intended use cases.

Across multi-year deployment cycles, model maintenance and updates become increasingly predictable, simplifying maintenance cycles and reducing total expenditure compared to cloud-centric approaches that demand continuous network connectivity and consumption-based pricing models.

The Impact of Decentralized AI on Education in Emerging Markets

Use Case Deep Dive

The deployment of offline LLMs through low-cost on-device GenAl solutions can have a transformative impact on education in emerging markets. With an estimated 1 billion 1 students in regions such as South America, India and Africa facing unreliable internet access and a shortage of qualified teachers, GenAl-powered educational tools present an opportunity to bridge this gap. On-device LLMs can contextualize student requirements, understand their strengths and weaknesses and provide individually tailored responses — all in a matter of seconds while preserving student privacy. By combining local and cloud-based models for a hybrid approach, educational institutions could:



Address the Teacher Shortage with Offline GenAI

In regions where access to quality education is limited, on-device GenAl can be a proxy to learning resources that anyone can access, any time, for free. LLMs running on low-cost edge devices can offer personalized tutoring, generate content, answer student queries and facilitate self-paced learning, complementing the role of teachers in areas with inadequate educational infrastructure and democratizing access to high-quality learning.



Build Localized Solutions

 $Language and vision GenAl models can provide diverse educational tools from interactive learning experiences to creative content support while running efficiently on low-power edge devices. Independent Software Vendors (ISVs) and Systems Integrators (SIs) play a critical role in this ecosystem by providing the front-end UI/UX wrappers, ensuring they are user-friendly, intuitive and tailored to specific educational needs. ISVs can develop custom interfaces that cater to both teachers and students, while SIs can focus on the setup, maintenance and activation of the edge GenAl systems. Through customized interface design and user experience optimization, these lightweight LLMs effectively serve diverse educational environments with adaptations such as the ability to translate learning materials into anyone's first language. And when integrated with retrieval-augmented generation (RAG)^4 or extended context workflows 5, these lightweight models align precisely with specific curricular requirements and moderation protocols like safety, ethics, age appropriateness for K-12 segments and others.$



Reduce Operational Costs

By using on-premise GenAl, schools can reduce or eliminate subscription or API fees and gain tighter control over operational costs. In other words, they can implement the above solutions in a cost-effective way.

In summary, the deployment of optimized lightweight LLMs on edge devices offers a scalable and cost-effective solution to improving education in underserved regions. By leveraging the rapid evolution of open-source GenAl models and relying on ISVs and SIs for customization and maintenance, emerging markets can take significant strides toward bridging the education gap, bringing high-quality, GenAl-driven learning to millions of students.

Conclusion

Decentralized, on-device GenAl architectures can transform both operational performance and cost efficiency in myriad real-world deployments, from enterprises to emerging markets. Together, Intel® $Core^{TM}$ Ultra processors, Intel® Arc^{TM} GPUs and OpenVINO offer an accessible route to democratizing GenAl by accelerating the development of optimized on-device GenAl solutions.

- Learn more about the <u>AIPC powered by Intel</u> and deployment with the <u>OpenVINO toolkit</u>
- Try <u>OpenVINO GenAl examples</u> on your hardware

Disclaimer

This whitepaper analyzes various AI models, including fully open-source models (e.g., Apache 2.0, MIT-licensed models) and openly available models (e.g., Meta's LLaMA). Each model remains the intellectual property of its respective creators and is governed by its original license. Fully open-source models such as Qwen, Mistral, and Phi allow unrestricted use, modification, and distribution under their respective licenses. Models like LLaMA are available under specific usage terms that may include restrictions. This document does not redistribute, modify, or alter any model weights, nor does it claim ownership over them. Users should consult the official licensing terms of each model before use. The inclusion of specific models in this whitepaper does not imply endorsement by their respective developers. The recommendations presented are based on independent research and do not constitute official guidance from Meta, OpenAI, Google, or any other model provider.

Performance varies by use, configuration and other factors. Learn more at www.Intel.com/PerformanceIndex. Performance results are based on testing as of dates shown in configurations and may not reflect all publicly available updates. See backup for configuration details. No product or component can be absolutely secure. Your costs and results may vary. Intel technologies may require enabled hardware, software or service activation. Current characterized errata are available on request.

© Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.

Statements in this document that refer to future plans or expectations are forward-looking and subject to change without notice. Actual results may vary.



Sources

- $^{1}\underline{How\,many\,children\,and\,young\,people\,have\,internet\,access\,at\,home?}\,UNICEF.$
- ² <u>Agentic Al and Confidential Computing.</u> Agentic Al is: "a paradigm that involves deploying agents or other programs that act autonomously on behalf of humans or other Al systems to make decisions and take actions to achieve specific goals. Unlike current agent models, which are rules-based, agentic Al uses sophisticated reasoning, knowledge, and intelligence to help drive process automation by making decisions and taking actions rather than just responding to questions."
- ³ Optimizing Large Language Models with the OpenVINO™ Toolkit. Intel. April 2024.
- 4 What is retrieval-augmented generation? RAG is an AI framework that "helps LLMs deliver more-accurate and -relevant AI responses. RAG supplements an LLM with data from an external knowledge base, ensuring LLMs can access the most-reliable and -current information. This additional data helps LLMs deliver up-to-date and contextually meaningful responses."
- ⁵ What is a context window? McKinsey & Company. December 2024. "A context window is how information is entered into LLMs. The larger the window, the more information an LLM is able to process at once."