# The .serva Standard

## One Primitive for All AI - Cost Reduced, Barriers Removed

*Part One*

December 21, 2025

**Servamind Inc.**

Rachel St. Clair, John Austin Cook, Peter Sutor Jr. Victor Cavero, Garrett Mindt

Website: [servamind.com](servamind.com)  Email: info@servamind.com

## Abstract

Artificial Intelligence (AI) infrastructure faces two compounding crises. Compute payload – the unsustainable energy and capital costs of training and inference – threatens to outpace grid capacity and concentrate capability among a handful of organizations. Data chaos – the 80% of project effort consumed by preparation, conversion, and preprocessing – strangles development velocity and locks individual datasets to single model architectures. Current approaches treat these as separate problems, managing each with incremental advancements in optimization, increasing complexity in the overall AI tooling ecosystem. The approach presented here views data and computation as two expressions of single architecture, where a unified primitive is missing. Early indicators result in significant optimization, lifting current approaches and reducing complexity throughout the AI ecosystem (i.e. infrastructure).

This paper presents ServaStack: a universal data format (.serva) paired with a universal AI compute engine (Chimera). The .serva format achieves lossless compression by encoding information using laser holography principles, while Chimera converts compute operations into a representational space where computation occurs directly without decompression. For AI, the result is automatic data preprocessing by converting into .serva. The Chimera engine enables any existing model to operate on .serva data without retraining, preserving infrastructure investments while revamping their efficiency.

Internal benchmarks demonstrate 30-374× energy efficiency improvements (96-99% energy reduction), ~4× lossless storage compression, and 34× compute payload reduction without loss of accuracy when compared to RNN, CNN, and MLP models trained on original FashionMNIST and MNIST dataset files. At hyperscale one billion daily iterations, these gains translate cost savings of $4.85M per petabyte in per training. The impact of this technology proposes more

significance than the efficiency; that when any data flows to any model on any hardware, the entire AI development paradigm shifts. The bottleneck moves from infrastructure to imagination.

# I. Introduction

*Any data to any model on any hardware.* This is what Servamind has built.

WE NEED A HOOK, SOMETHING INTERESTING… WHY READ THIS? A narrative of some sort… a little fluff.. "in the land before times" "Three Rings for the Elven-kings under the sky, Nine for Mortal Men doomed to die, Seven for Dwarf-lords in their halls of stone, One for the Dark Lord on his dark throne".. Leading with numbers in the following paragraph is so boring! People want to feel good about the money they make themselves under the guise of some cause.

AI infrastructure spending reached $135 billion in 2024 and is projected to surpass $200 billion by 2028 and the figure is accelerating [16]. Yet AI remains trapped: training accessible only to hyperscalers, impractical at the edge, unsustainable at scale. The cardinal constraint is feasibility, which relies on  infrastructure, encouraging new forms of machine intelligence to emerge.

Servamind's purpose is to allow AI to flourish in a meaningful way to all of society. To achieve this we have built Servastack, a combination of our unique encoder (Serva Encoder) and an AI wrapper (Chimera) to solve what we think are two difficult problems in the AI development pipeline: compute payload and data chaos. Servastack is a universal data format paired with a universal AI compute engine, allowing data and computation to flow across any ecosystem into a unified application. This property of universality is paired with an exploit on the current compute paradigm for maximum efficiency. This brings down the cost of AI creation and use by orders of magnitude [https://arxiv.org/abs/2403.17887]. Rather than copying and optimizing current approaches, the focus of AI advancement can now become new modes of learning, new user interfaces, and new insights into thinking machines.

## The Two Primary Problems

AI cost is dominated by two challenges: data chaos and compute payload. In plain terms, AI is expensive and data is difficult. Together, these account for the vast majority of project resources. Data preparation alone consumes roughly 80% of total effort and budget, while compute infrastructure represents 47-67% of total AI development costs for organizations building from scratch [8,18].

**Compute Payload** is the visible crisis. The International Energy Agency projects that global datacenter electricity consumption will more than double by 2030, reaching approximately 945

TWh annually, equivalent to Japan's entire electricity consumption [1]. AI-accelerated servers are growing at 30% annually, <u>four times faster than total electricity supply growth</u>, and 20% of planned datacenter projects already face delays due to grid bottlenecks. Grid operators are already hitting capacity limits: in Virginia, Dominion Energy faces a seven-year backlog for new datacenter connections; Ireland imposed a moratorium on Dublin datacenter grid connections from 2021 to 2025 due to electricity system strain [19, 20]. The U.S. Department of Energy warns that without efficiency breakthroughs, AI's power demands could require dozens of new power plants within the decade [2].

The scale of investment often reflects the scale of demand. NVIDIA's datacenter revenue more than doubled from $15 billion in 2023 to over $47 billion in 2024 [3]. Microsoft committed $80 billion to AI datacenters in 2025 [4]. Epoch AI's analysis shows compute used in frontier training runs growing 4-5× per year since 2010, far exceeding the chip miniaturization efficiency gains chip Moore's Law [5]. The Stanford AI Index reported training costs for state-of-the-art models reaching $78 million (GPT-4) to $191 million (Gemini Ultra) [6]. Training GPT-3 emitted an estimated 552 metric tons of $CO_2$; GPT-4's training emissions are estimated at 10,000-15,000 metric tons – roughly 20× higher [7].

ARK Invest, an asset management company, offers a counternarrative: AI training costs are declining approximately 75% annually through Wright's Law dynamics, hardware improvements reducing compute unit costs by 53% per year, compounded by algorithmic efficiencies contributing another 47% [22]. If costs decline accordingly, the argument goes, AI scales sustainably.

This analysis, however, conflates efficiency with capability. Wright's Law measures the cost to reproduce yesterday's performance, not to achieve tomorrow's, not to increase AI efficacy (i.e. intelligence). Moreover, ARK's cost curves track per-unit compute while excluding the infrastructure buildout where costs are rising and resources are constrained: datacenter real estate prices increased 19% in 2024, supply chain shortages for generators, chillers, and transformers are inflating construction costs, and grid connection backlogs stretch to seven years in key markets-none of which follow Wright's Law dynamics [CITE]. The view that AI is scaling is myopic compared to the total supply chain growth needs of current AI's trajectory.

DeepMind's Chinchilla scaling laws reveal the deeper constraint in large language model (LLM) development. The relationship between compute and capability follows a power law [23]. Compute-optimal training requires scaling parameters and data together, with FLOPs scaling quadratically with model size (approximately 6ND for dense transformers) while capabilities improve along a much shallower curve. The implication is sobering: frontier capability does not get cheaper. Each incremental improvement in model performance demands disproportionately more compute. The goalpost moves faster than efficiency gains can follow.

The consequences are threefold. Climate impact is mounting. Capability is concentrating-only a handful of organizations can afford frontier model training. Barriers to entry are rising; startups, researchers, and developing nations find themselves increasingly locked out of meaningful participation in AI advancement. The Chinchilla constraint compounds all three: every capability

improvement demands proportionally more data, more compute, and more energy. Yet, Chinchilla applies only to LLMs, the models that have dominated optimization efforts for the past several years. The next frontier in generative AI, simulation models, multi-modal systems, and multimedia reasoning generation, exhibits even steeper scaling requirements, with data volumes and compute demands that dwarf text-based training. Section V examines how Servamind's efficiency gains alter this calculus across modalities.

AI in its current form is not scalable to the degree that its collateral costs should be overlooked.

**Data Chaos** is the less visible crisis even though AI developers live it daily.

The practical experience goes as follows: spend months determining how to prepare data and then quickly copy/paste a model architecture from a journal repository. Refit this model in a week. Then run inference for another week, before finally obtaining some results. Then after all of that, refactor everything for a new model and repeat the entire cycle. The actual AI is the easy part. The real job is all the data work.

Industry analyses consistently estimate that 80% or more of any AI project's effort goes to data preparation, cleaning, and orchestration [8]. This reality is reflected in market behavior: models are frequently open-sourced while training data remains sacred, proprietary, high-priced IP. The value resides in the data, and the cost resides in preparing it.

There is also a one-to-one lock-in between dataset and model. Data cannot simply be fetched from storage and passed to any AI. It must be specifically pre-processed for the downstream AI model architecture. Every time a new model is adopted, the data must be re-processed from scratch to produce model identifiable features. Aside from being inefficient, it is genuinely frustrating for practitioners who understand that the underlying information remains the same regardless of which model will consume it. Feature engineering is a human operation; eliminating it would allow AI models to do the bulk of the work.

As data capture expands, the problem compounds. IDC projects global data creation will exceed 180 zettabytes by 2025, up from 64 zettabytes in 2020 [9]. Much of this growth comes from specialized domains, medical imaging, satellite telemetry, genomics, industrial sensors, producing formats poorly matched to mainstream architectures. Regulated industries like healthcare and finance face 20–35% higher AI implementation costs due to compliance requirements, specialized data handling, and domain adaptation needs [21]. The prevailing data constraints are format issues, quality issues, scale issues, and security issues.

The data problem is accelerating, not slowing.

## Current Approaches

Hardware approaches are largely abetted by miniaturization, making transistors smaller so more of them can fit on chip. Moore's Law, the observed rate of miniaturization, however, is lessening. Perhaps 10–15 years of conventional scaling remain before atomic limits impose quantum

effects [10]. Practical quantum computing at consumer scale is not expected until 2040 or beyond [11]. The hardware path alone will not solve the timeline we face.

Software approaches seek efficiency through representation: 64-bit to 8-bit quantization, pruning, distillation, sparsity. Even aggressive techniques like Microsoft's BitNet, which reduces weights to ternary values {-1, 0, +1}, achieve 2-6× speedups and 55-82% energy reduction [24]. In an attempt to overcome data chaos, typical approaches are to orchestrate systems and layers to navigate the chaos. Some current approaches attempt data orchestration: Pandas DataFrames, PyTorch tensors, Hugging Face SafeTensors. ML-Ops platforms like MLflow, Weights & Biases, and Kubeflow coordinate infrastructure and track experiments.

The AI ecosystem has also produced numerous data formats, each solving a narrow problem. Apache Arrow and Parquet optimize columnar analytics but assume tabular structure. TFRecord and SafeTensors serialize tensors for specific frameworks. ONNX provides model interchange but not data interchange and remains a conversion layer, not a native format.

Current approaches represent steps in the right direction towards unifying the AI ecosystem's tooling. Yet none are universal for all data to any model. Further, they do not address the root cause of binding feature information to compute payload. While these approaches deliver valuable gains, none deliver the order-of-magnitude improvements required to break scaling constraints, or ease the data-model lock-in. Each standard addresses one link in the pipeline while the fundamental fragmentation remains.

The AI field moves too rapidly to design for specific configurations. It is futile for any team to keep pace with every variation. Too many data formats exist. Too many model architectures compete. Too many hardware targets fragment the landscape. Too many programming languages divide practitioners. Tooling is scattered, redundant, and confusing.

Rather than spend so much time navigating the chaos we set out to solve the problem by simplifying all data preprocessing using a universal encoding producing a standard file format, a .serva file.

## How .serva Differs from Existing Standards

The .serva format differs in kind, not degree. It is not a serialization format for a specific data type. It is not a conversion layer between frameworks. It is a universal encoding that transforms any input (e.g. images, text, audio, sensor streams, structured records, etc.) into a single representational space where all information is preserved and direct information extraction (i.e. computation) can occur. The question shifts from "how do I convert my data for this model" to "encode once, compute anywhere."

## Servamind Approach

The root inefficiency cause is that data and computation have never been addressed together in a way that leverages the existing ecosystem. Our answer to these two compounding problems

is our Servastack, a universal data format (.serva) to eliminate data chaos and a universal AI compute engine (Chimera) that tackles the compute payload problem. These two solutions are independently needed in order to penetrate the existing ecosystem at various levels of the pipeline, where the two problems remain separate. When combined in a workflow, Servastack emerges as a unification paradigm that can begin to be leveraged for compounding efficiency gains.

Servamind attacks both problems simultaneously through a unified system:

- A universal data format (.serva) that eliminates data chaos created by Serva Encoder program
- A universal compute engine (Chimera) that solves compute payload
- Together (Servastack): 30-374× energy efficiency (96-99% cost reduction), ~4× lossless storage compression, and 34× compute payload reduction – without diminished accuracy
- Data agnostic. Model agnostic. Hardware agnostic.

☆ The insight most observers miss: they hear "compression" and think "efficiency." Efficiency is a consequence. **Universality is the breakthrough.**

When any data can flow to any model, the entire AI development paradigm shifts so that any model can ingest any data. Bidirectional compatibility emerges. The data layer becomes decoupled from the model layer entirely. True multimodality becomes tractable – vision-language-action models in robotics have struggled not for algorithmic reasons, but because fitting heterogeneous data into one model presents an engineering nightmare. The data-model lock dissolves across all verticals.

The market reality shaped the product. In the current AI ecosystem only hyperscalers can afford frontier training. The efficiency gained through Servastack would democratize access. But adoption faces a barrier, since organizations resist retraining models in which they have already heavily invested. Infrastructure overhaul appears more expensive than ongoing inefficiency.

This constraint forced two requirements:

1. More efficient computation regardless of data type or source
2. Compatibility with any stage of AI-training, pre-training, fine-tuning, inference

What is most important to note: the Servamind solution does not contradict or compete with other approaches. It is universal. It is additive.

☆ *Servastack partners with and amplifies the efforts of all those pursuing ease and efficiency in AI.*

# II. The Origin

Servamind began with a question about why AI fails to scale like biological intelligence.

The culprit in neural networks is catastrophic forgetting, known as the fundamental limitation of backpropagation-based learning [CITE]. Alternative paradigms hit their own walls: reinforcement learning's sample inefficiency, knowledge graphs' combinatorial explosion, symbolic AI's brittleness [25-28]. No existing approach scales cleanly. The workarounds remain fragile, and the industry locks in technical debt with every deployment. When a neural network learns a new task, it updates its weights to optimize for that task at the expense of prior tasks. The model drifts away from previous knowledge, constantly forgetting what it once knew. This is not a bug in implementation; it is a structural consequence of how greedy-based learning operates [12].

Brains learn a bit differently. Biological neural systems do not usually catastrophically forget. A human who learns Spanish does not entirely lose their English for doing so, as would, for example, Siri's AI. The brain has evolved to solve this problem. This recognition initiated a search for learning mechanisms closer to biological reality.

## Inspiration from Biological Systems

That search led to the work of L. Andrew Coward, namely Recommendation Architecture, a theoretical framework for understanding higher cognition in terms of anatomy and physiology [13]. Years of study revealed several foundational insights that would reshape our approach to the problem.

**First**: Information in biological systems is computed in three-dimensional physical space. The spatial arrangement of neurons matters. This makes von Neumann architectures, where memory and processing are fundamentally separated, a poor substrate for brain-like computation. Every mainstream computer inherits this bottleneck.

**Second**: The units of information in biological systems are maximally combinatorial. They are designed to combine and build up into any higher-order representation for unknown future tasks. The brain does not optimize its representations for the current task; it maintains flexibility for tasks it has never encountered. Representations are distilled into suggestions that the system learns from, not hard commitments that foreclose future possibilities.

**Third**: The filtering of noise to signal does not happen inside the brain the way it happens inside AI models. In traditional AI, feature engineering and model layers perform noise-to-signal transformation. In biological systems, however, that filtering has already occurred upstream, at the sensor. The retina, the cochlea, and the mechanoreceptors are not passive recorders. They are intelligent filters shaped by hundreds of millions of years of evolution. By the time information reaches the brain, coherence has already been imposed by the camera, the lidar sensor, the capture device.

*Representations should preserve possibility, not collapse it prematurely.* The relationship between DNA and protein, where the same genetic sequence can participate in producing vastly different outcomes depending on context. This implies that internal representations should be more ambiguous.

How can we know what information will matter when the downstream task is unknown, as is often the case in general intelligence and in practical AI development? The approach Servamind takes is to preserve everything and assume all apparent noise may be a latent signal. We keep everything, while increasing efficiency as a byproduct. In later sections, we explain how this paradox is resolved.

## Compression and Intelligence: The Hutter Framework

Marcus Hutter's work from Google DeepMind establishes a profound equivalence: compression and prediction are fundamentally the same operation [14]. To compress data optimally is to model it optimally. Optimal compression implies optimal inference. "If you can compress, you can learn" is a mathematical identity rooted in Kolmogorov complexity and algorithmic information theory [CITE CITE].

The implication for AI systems is significant: superior compression yields superior efficiency. A system that achieves better compression has, by definition, extracted more structure from data using fewer resources. Most approaches follow the trend of incremental optimization. Our approach adds to incremental optimization by providing a primer designed from exploitation of mathematics about the nature of learning itself.

In practical terms, successful compression separates signal from noise in a useful way. In AI, we call this feature engineering and it is often performed by expert-crafted reduction operations to reduce unnecessary information in the data. The feature engineered data is then ingested by the AI model. The outputs of intermediate model network layers, before final classification or generation, are feature vectors: compressed representations that capture learned structure by losing the unlearned structure. The unlearned structure, or appropriate information to lose throughout the feature vector creation process is directly determined by its unnecessity to produce correct results for the desired learning task at hand. Compression and learning are the same operation viewed from different angles.

## Information Theory: The Shannon Foundation

Claude Shannon's foundational work on noisy channels established the theoretical limits of information transmission [15]. His central insight: information can be preserved through transformation if entropy is managed correctly. There exist transformations that reduce representation size while losing nothing – the domain of lossless compression.

Shannon's framework offers an additional insight relevant to our problem. Information was defined, in part, as non-randomness: structure, pattern, predictability. Determining whether apparent randomness is true stochasticity or merely a temporal state in an intractable

deterministic system is not possible without complete observation. Information is thus defined in juxtaposition to noise. Whatever is not useful, we call noise. Usefulness, however, is task-dependent and often unknowable in advance.

The distinction between lossy and lossless compression matters critically for AI. Lossy compression discards information deemed unimportant by some prior criterion. In AI applications, however, we often cannot know in advance what information will prove important for downstream tasks. Lossless compression preserves all information, deferring the question of relevance to the learning system itself.

## The Paradox

These frameworks created a paradox at their intersection.

Hutter says compress to learn since ptimal compression implies optimal prediction. AI systems should aggressively compress their representations to maximize learning efficiency, which by current methods involves drastic information reduction, further indebting the data-model lock.

Shannon says preserve to remain general. Discarding information precludes possibilities. However, retaining all information would explode the compute resources required. The Shannon insight points in a possible direction for a solution: lossless compression. For tasks where downstream use is unknown, lossless preservation is required.

These imperatives seem opposed. Compression discards; preservation retains. How can both be satisfied simultaneously?

### III. Our Solution

The resolution came from recognizing where entropy actually exists in the pipeline.

Shannon's framework assumes a noisy channel, a transmission medium that introduces randomness. The data AI systems operate on, however, is not raw entropy. It is captured data: images from cameras, audio from microphones, telemetry from sensors, text from human authors.

These capture devices were designed by human intelligence. They impose coherence. They select what to record and how to record it. The camera does not capture pure photon chaos; it captures structured light filtered through a lens system engineered for human visual understanding. By the time data enters an AI pipeline, it has already been filtered by the causal structure of physical reality and the intentional design of the capture mechanism.

The entropy is already reduced. The signal has already been extracted, not by the AI, but by the physical and engineered systems upstream. Physical reality is causal - each state follows coherently from the last. Capture devices record this coherence. The data we compress is not noise. It is structured by physical causality and human intent. We are not fighting entropy - we are revealing the structure that was always present.

8

Hutter established that compression enables learning. Shannon established how to compress without losing information. Both can be satisfied simultaneously when the data is already structured and real-world data is structured by physical causality and human intent.

Servamind applies this synthesis: a lossless compressed format that AI models can compute on directly. Not lossy approximation, not dimensionality reduction - lossless compression that retains signal because, when downstream tasks are unknown, all of it may be signal.

**Universal Feature Vectors**

This raises a practical question: if we want feature vectors from data yet do not know what downstream model or task will consume them, how do we know what to encode?

The answer follows from the theory above: encode everything. Preserve all structure. Let downstream tasks extract what they need rather than guessing in advance what they will require. The difficulty of solving the problem then lies the challenge of how. How to create lossless compression across arbitrary information, which is addressed partially in next section, Architecture.

Servamind encodings (.serva files) are representations that preserve information, because reduction may damage what matters. They trend towards being maximally combinatorial, able to serve any downstream task because they have not been optimized for any particular one.

This framing also addresses catastrophic forgetting at its source. Backpropagation-based learning overwrites weights optimized for previous tasks when learning new ones. The model constantly drifts from prior knowledge. Learning on universal feature vectors, however, provides protection at the data layer. Information is not discarded because all of it is treated as information. The representations themselves resist the forgetting problem by never having collapsed the possibility space in the first place, as long as the computation interacts in this encoded space.

## Why No One Attempted this Until Now

In short, this challenge is a monumental opportunity, where many other lower hanging fruits are to be found. The paradox of keeping everything and increasing efficiency also sounds implausible on its face. Compression and computation appear to be opposing operations. Compression standards are focused on removing noise to filter the data for easier computation. Lossless compression lacked any robust implementation that could rival lossy methods. Further computing on any lossless compression is narrow and few implementations exist at all, let alone practically and market viable ones.

This apparent contradiction dissolves under a different computational paradigm. Compression and computation seem congealed only when data is lost or must decompress before processing. When the encoding itself is designed for direct computation, when operations preserve their meaning under the transformation but in

9

smaller size, the two coexist. Servamind operates in this space: lossless compressed representations that remain computationally accessible.

There is also a practical barrier of assuaging the current ecosystem of infrastructure tooling in AI. Few practitioners train models from scratch. Those who do are reluctant to retrain on a new data format, even one promising cheaper computation, because their investment in existing trained models and infrastructure creates integration friction. The switching cost appears prohibitive. Established investments create inertia. Organizations routinely accept ongoing inefficiency over one-time integration cost, even when the long-term economics favor change.

We therefore built a wrapper. The function of Chimera is simple in concept. Chimera can take any model in any state and enable it to operate on .serva universal feature vector files without re-training, with minimal compute overhead, and without adding complexity to the user. These constraints meant any approach had to satisfy two requirements: more efficient computation, and compatibility with existing models at any stage without retraining.

This innovation required substantial mathematical and computer science innovation, unifying disparate formalisms into a coherent system.

☆ *The result is that Servastack adoption does not require abandoning existing investments, it extends them. It does not require a competitive choice over one efficiency gain to the next, it binds them all to work symbiotically.*
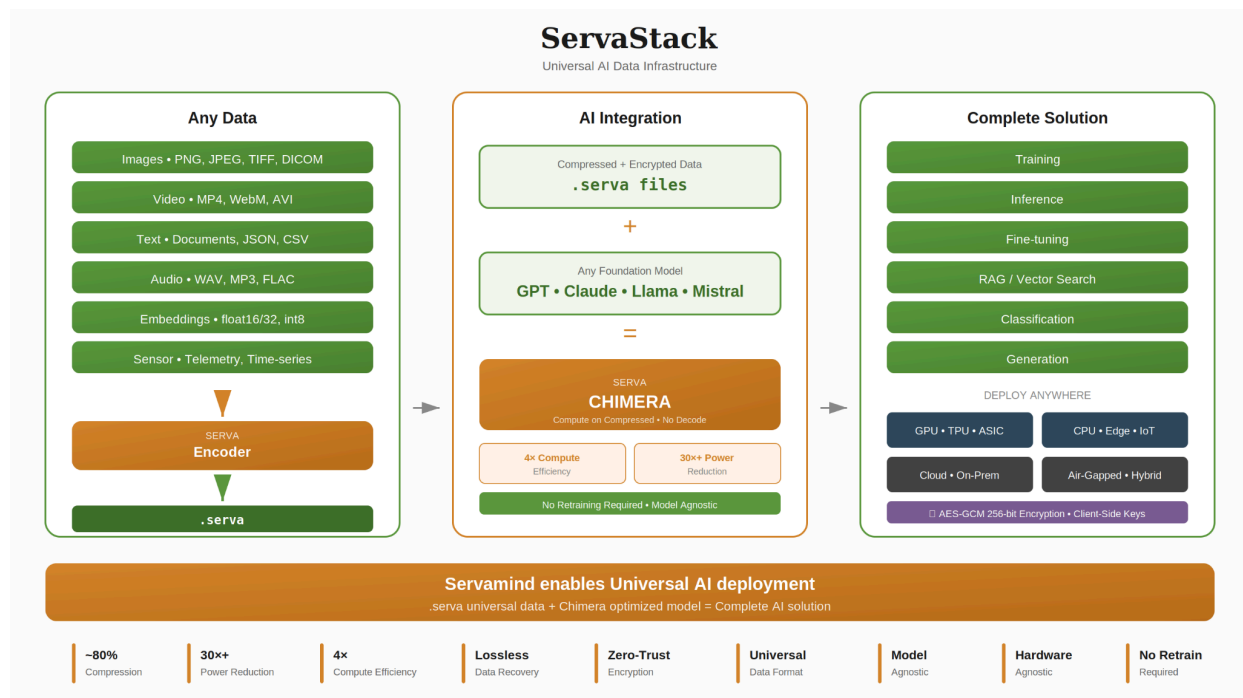
## Architecture

10

**Figure 1.** ServaStack architecture. Universal data encoding through the Serva Encoder produces .serva files that integrate with any foundation model via the Chimera wrapper, enabling deployment across all AI tasks and hardware targets.

In summary of the architecture for any level of the AI stack in any ecosystem on any infrastructure Servastack works as follows. Data goes through the Serva Encoder software program to produce fully pre-processes AI ready .serva files. This reduces memory footprint, and depending on the level of integration in the stack can increase data upload and transmission efficiency. It also reduces need for data operations, cutting the 80% of AI project effort. Then, training or inference programs go through the Serva Chimera software which converts the existing program to an equivalent set of instructions to perform directly on datasets composed of .serva files, or .serva files which contain an entire dataset. Chimera only takes the pertinent parts of the .serva file needed for training, automatically, leaving the original disk .serva data unchanged. This reduces data operations significantly in AI training and inference while also reducing in-memory computation and storage access — resulting in speed, throughput, and power reduction.

There are several ways in which we intend for the Servastack to be used. The most directly applicable is a library in which developers can call Serva Encoder and Serca Chimera to preprocess AI and wrap their models before device execution. This would work in a fashion similar to calling data.transform in Pytorch. That data is on disk, called to the AI program software in a step that converts it to the AI framework. Instead of many lines of code processing the data before framework formatting, one call takes the data from disk and prepares it into a bit vector. From there, no framework specific formatting needed. The model is written, training loops described, and just before it is sent to device (e.g. .to(device) in Pytorch, a Chimera call is

made (e.g. model.Chimera(args)) to wrap the model, compiling the static or dynamic graph computation instructions to be run on device.

While this effort is underway, it needs to be extended and maintained to every language, developers have to read our docs for correct implementation, and become aware of the tools existence. If one searches AI tools, there is a cataclysmic amount of results which from google trends have increased 85% from January to November of 2025. The library helps today's technically inclined developers who have the deep industry knowledge of why this tooling accelerates their workflows. Further, excessive functionality is needed to be higher up in the tech stack since tooling is variable and vast — which requires time and maintenance for the developer to keep up with our feature advancements.

The library approach, while useful and necessary, does not prepare the future state of AI nor does it motivate the supply chain efficiency needs. For those angles, mass adoption is required and their relief is most pertinent to the compute providers, the cloud providers, OEMs, and layers of server infrastructure providers. Here the utility is in directly embedding Servastack technology into the operating system or through command line interface close to the metal (e.g. Kubernetes/Docker containers). Thus, on premise deployments will be key desirade. Here, all data can be stored in .serva. Anything moving the data that isn't called with Chimera, will direct lossless decompression for appropriate and normal workflows. The chimera flag triggers non-decompression such that AI workflows can act directly on the .serva files.

With a few key players, most workflows in general could be using Servastack without the end-user ever knowing, resulting in the experience that AI has just become cheaper, easier, and more useful. Take AWS as a prime example. If all data in an S3 bucket was converted upon upload to .serva format, the memory footprint would be reduced and upload would become faster, more reliable and more secure (details ruminate in part two of this white paper). Once data is transferred from S3 to EC2 for computation, during transfer any non chimera calls will decompress during fast file transfer protocol. Chimera called data from EC2 will directly ingest the data without need for any pre-processing or further data orchestration steps besides which device to send to (e.g. CPU, GPU etc.). Instead Chimera will only need to ingest the model (in .onnx, pytorch, TF, or .yaml config file) and the data in .serva format before sending to execution compilation. At this moment, Chimera will come to action, transmuting the compiled instructions to an equivalent set of instructions directly computable in the .serva vector space.

Finally, to ensure full compatibility with the existing tooling ecosystem and fragmented infrastructure, executable program applications are underway. As the Servastack paradigm prevails, AI can become closer to the everyday user. Users who historically prefer drag and drop, intuitive interfaces for mobile devices and the likes. Here, end to end AI implementations embody the full Servastack experience.

## Development Methodology

The technical moat is substantial. The underlying principles span multiple disciplines that rarely intersect: information theory, holographic encoding, hyperdimensional computing, and

hardware-aware optimization. This convergence is not easily replicated. A future technical paper may be released upon Serva Encoder's open-sourcing strategy. The work presented here required years of sustained effort and a willingness to reject established assumptions in computer science, AI research, and adjacent fields. The result is infrastructure that appears simple in use while embedding deep theoretical innovation beneath the surface.

The "how it works" draws from principles of laser holography, where an interference pattern encodes information without storing the data itself. Because this representation exists in an abstract referential space, computation can occur in that same space, provided the transformations remain homomorphic. This is the key that permits computing directly on compressed representations: the mathematical operations that define learning preserve their meaning under the encoding.

Serva Encoder's initial implementation is compact: approximately 200 kilobytes, relying on elementary operations: bit-level addition, XOR, permutation, pseudo-random bit generation, and distance. The simplicity of the operations belies the complexity of their orchestration. When creating  the resulting bit vectors of a .serva file, a ciphertext is generated with a random seed, which can be pushed client side for encryption, making each file secure. If you map this to the analogy of laser holography, it is essentially the angle at which the grooves in the photo-lithographic plate are positioned to reflect the light bouncing off the source, imprinting the original information.

Serva Chimera is mathematically matched to the representation space. Meaning the math of computation operates in the same space as the encoder holography math. To ensure arbitrary model wrapping, the ability to transmute any existing model to operate on .serva representations without retraining, several other techniques are required. First, topology analysis to abstract the original model architecture to the referential space. Operations are projected using geometric mappings. One analogy may clarify this process. Traditional gradient-based learning navigates a loss landscape by walking, step by step through high-dimensional terrain, guided by local slope. The Chimers approach operates more like celestial navigation. Rather than traversing the landscape, it uses a star chart to compute coordinates and arrive directly.

To summarize, Servamind created a universal data format grounded in holographic encoding principles, a universal compute engine capable of transmuting any model architecture, validated by a framework designed to measure what matters – energy cost per unit of capability, and information preserved through transformation. The following section presents results.

# IV. Key Performance Indicators

## Validation Methodology

Servamind conducted internal benchmarking under controlled conditions designed to ensure fair comparison for Serva Encoder against other compression algorithms and also to validate the viability of Servastack by training models on .serva files. Part two of this white paper will provide rigorous external performance validation.

**Serva Encoder Compression Benchmark**

First, We evaluated Serva Encoder against 25 established compression algorithms using the Canterbury Corpus benchmark suite, measuring bits per byte (bpb) across four standard corpora. On the Canterbury Corpus (11 files), Serva achieved a weighted 1.708 bpb, ranking 13th overall and outperforming gzip, compress, and dictionary-based methods while trailing block-sorting variants like bzip2-9 (1.545 bpb) and context-mixing methods like ppmD5 (1.520 bpb). On the Large Corpus, Serva placed 3rd with 1.747 bpb, demonstrating competitive scaling on multi-megabyte files. Serva ranked 1st on the Artificial Corpus (3.036 bpb), indicating strong handling of pathological cases including high-entropy and highly repetitive data.

The compression benchmarks were performed to analyze how the program scales with file size, its viability outside of AI workflows, and general analysis for internal development.

**Training on .serva Data  Benchmarks**

In the cases without Chimera, it is possible to train models directly on .serva files, but the models have to then be configured to properly train on this new data format. For this test, we compare how a native  .serva model performs to traditional models. The native model was evaluated against standard neural network architectures on Fashion-MNIST and MNIST, benchmarks that permit direct comparison with published results and enable reproducibility assessment. Fashion-MNIST is a classification task in which the model must predict categories of clothing from photos [https://arxiv.org/abs/1708.07747]. MNIST is a similar classification task of numbers from photos of handwritten digits [https://ieeexplore.ieee.org/document/6296535].

The SERVA model in testing encodes images into a .serva variant, classified by k-NN (k=3) with class-balanced scoring. Two benchmark modes were run: N-epoch (train until matching SERVA accuracy or 100 epochs) and single-epoch.

We evaluated SERVA against five neural network baselines on Fashion-MNIST and MNIST using a controlled benchmark environment. All models were implemented in pure NumPy with Numba JIT compilation, SGD optimization (lr=0.01), float64 precision, and He/Xavier initialization to eliminate framework-level confounds.

Baseline architectures:

- MLP-1L: 784→256→10 with ReLU (batch=128)
- MLP-2L: 784→256→256→10 with ReLU (batch=128)
- MLP-3L: 784→256→256→256→10 with ReLU (batch=128)
- CNN: 8 filters (5×5) → ReLU → maxpool(2) → FC (batch=64)
- RNN: vanilla RNN, 28 timesteps × 28 features, hidden=64 (batch=128)

Two benchmark modes were run: N-epoch training (until matching SERVA accuracy or 100 epochs) and single-epoch comparison. Energy was measured via Intel RAPL (CPU package + DRAM domains) using pyJoules. On Fashion-MNIST, SERVA achieved 88.39% accuracy in 1.41s consuming 150.2J; the fastest baseline to match this accuracy was MLP-3L requiring 60 epochs, 165.03s, and 14,938.1J (99× energy overhead). On MNIST, SERVA reached 96.48% in 1.45s at 153.6J versus MLP-3L at 18 epochs, 50.21s, and 4,551.5J (30× energy overhead).

Hardware: 48-core Intel Skylake-AVX512, 257GB RAM.

**SERVA Model Training**

In addition to the comparison test, we evaluated the SERVA model alone to determine what portion of the .serva files are needed for training, which denotes the payload reduction during training. Eight variations SERVA architectures were trained on .serva encoded data, ensemble-evaluated across all combinations (1-of-8 through 8-of-8), with the optimal model selected for final test accuracy reporting. Again, training was performed on both the Fashion-MNIST and MNIST tasks.

The critical metric we are tracking here is compute payload, the data volume that must be processed per training iteration versus raw dataset size. This metric captures whether Chimera can extract minimal computational representations from .serva files while preserving all information necessary for model performance. Unlike on-disk storage compression, compute payload measures what the model actually operates on during training and inference.

**Serva Encoder Compression Benchmark Results**

## Table 1: SERVA Summary

| Metric | Value |
|---|---|
| Total Original | 17.66 MiB |
| Total Compressed | 4.24 MiB |
| Overall bpb | 1.920 |

| | |
|---|---|
| Compression Ratio | 4.17× |
| Compression Throughput | 4.65 MB/s |
| Decompression Throughput | 15.85 MB/s |

## Table 2: Canterbury Corpus Results (bits per byte, lower = better)

| Method | Weighted bpb | Rank |
|---|---|---|
| szip-b | 1.464 | 1 |
| szip | 1.478 | 2 |
| bzip-6 | 1.490 | 3 |
| bzip-9 | 1.498 | 4 |
| ppmD5 | 1.520 | 5 |
| bzip2-6 | 1.538 | 6 |
| bzip2-9 | 1.545 | 7 |
| ppmD7 | 1.561 | 8 |

| | | |
|---|---|---|
| bzip-1 | 1.591 | 9 |
| ppmC-896 | 1.612 | 10 |
| bzip2-1 | 1.640 | 11 |
| ppmD3 | 1.645 | 12 |
| SERVA | 1.708 | 13 |
| dmc-50M | 1.737 | 14 |
| ppmCnx-896 | 1.745 | 15 |
| gzip-b | 2.082 | 19 |
| gzip-d | 2.090 | 20 |
| compress | 2.553 | 24 |

*SERVA ranks 13th of 32 methods, outperforming gzip, compress, and most lightweight compressors.*

---

Table 3: Large Corpus Results (bits per byte, lower = better)

| Method | E.coli | bible | world | Weighted bpb | Rank |
|--------|--------|-------|-------|--------------|------|
| szip-b | 2.060 | 1.530 | 1.400 | 1.721 | 1 |
| ppmD5 | 1.990 | 1.580 | 1.520 | 1.737 | 2 |
| SERVA | 1.993 | 1.643 | 1.453 | 1.747 | 3 |
| szip | 2.070 | 1.620 | 1.600 | 1.803 | 4 |
| ppmD7 | 2.030 | 1.660 | 1.660 | 1.814 | 5 |
| bzip-9 | 2.130 | 1.650 | 1.570 | 1.832 | 6 |
| bzip2-9 | 2.160 | 1.670 | 1.580 | 1.854 | 8 |
| gzip-b | 2.240 | 2.330 | 2.330 | 2.293 | 19 |
| gzip-d | 2.310 | 2.350 | 2.340 | 2.331 | 20 |

*SERVA ranks 3rd of 32 methods on large files, outperforming bzip and most other methods.*

---

Table 4: Artificial Corpus Results (bits per byte, lower = better)

| Method | aaa | alphabet | random | pi | Weighted bpb | Rank |
|--------|-----|----------|--------|----|--------------|------|

| | | | | | | |
|---|---|---|---|---|---|---|
| SERVA | 0.061 | 0.068 | 6.049 | 3.329 | 3.036 | 1 |
| bzip-9 | 0.000 | 0.010 | 6.080 | 3.390 | 3.076 | 2 |
| bzip-6 | 0.000 | 0.010 | 6.080 | 3.400 | 3.084 | 3 |
| bzip-1 | 0.000 | 0.010 | 6.080 | 3.400 | 3.084 | 4 |
| bzip2-9 | 0.000 | 0.040 | 6.050 | 3.450 | 3.122 | 5 |
| gzip-b | 0.010 | 0.020 | 6.050 | 3.760 | 3.360 | 18 |
| gzip-d | 0.010 | 0.020 | 6.050 | 3.760 | 3.360 | 19 |

*SERVA ranks 1st of 31 methods on artificial/pathological files.*

## Table 5: Calgary Corpus Results (bits per byte, lower = better)

| Method | Weighted bpb | Rank |
|---|---|---|
| szip-b | 2.075 | 1 |
| ppmD5 | 2.084 | 2 |
| szip | 2.091 | 3 |

| | | |
|---|---|---|
| bzip-9 | 2.093 | 4 |
| bzip2-9 | 2.110 | 5 |
| bzip-6 | 2.119 | 6 |
| bzip2-6 | 2.136 | 7 |
| ppmD7 | 2.148 | 8 |
| SERVA | 2.226 | 9 |
| ppmD3 | 2.260 | 10 |
| dmc-50M | 2.261 | 11 |
| gzip-b | 2.592 | 19 |
| gzip-d | 2.610 | 20 |

*SERVA ranks 9th of 32 methods, competitive with best-in-class compressors.*

## Table 6: SERVA Ranking Summary Across All Corpora

| Corpus | Files | SERVA Rank | Total Methods | Notes |
|---|---|---|---|---|

| | | | | | |
|---|---|---|---|---|---|
| Canterbury | 11 | 13th | 32 | | Main benchmark |
| Large | 3 | 3rd | 32 | | Best on large files |
| Artificial | 4 | 1st | 31 | | Best on pathological cases |
| Calgary | 14 | 9th | 32 | | Historic benchmark |

## Table 7: SERVA vs Common Compressors

| Compressor | Canterbury bpb | vs SERVA |
|---|---|---|
| SERVA | 1.708 | — |
| gzip (best) | 2.082 | SERVA 18% better |
| gzip (default) | 2.090 | SERVA 18% better |
| gzip (fast) | 2.462 | SERVA 31% better |
| compress | 2.553 | SERVA 33% better |
| lzrw1 | 3.584 | SERVA 52% better |

**Training on .serva Data Benchmarks Results**

**Table 1: N-Epoch Results (Fashion-MNIST)**

Models trained to convergence or maximum epochs

| Model | Accuracy | Time | Energy | Epochs |
|-------|----------|------|--------|--------|
| SERVA | 88.39% | 1.41s | 150.2 J | 1 |
| MLP-1L | 87.74% | 284.79s | 26,947.4 J | 100 |
| MLP-2L | 88.43% | 144.23s | 13,088.8 J | 67 |
| MLP-3L | 88.44% | 165.03s | 14,938.1 J | 60 |
| CNN | 88.41% | 321.97s | 24,757.1 J | 27 |
| RNN | 86.05% | 1,019.11 s | 56,135.7 J | 100 |

**Table 2: N-Epoch Results (MNIST)**

Models trained to convergence or maximum epochs

| Model | Accuracy | Time | Energy | Epochs |
|-------|----------|------|--------|--------|
| SERVA | 96.48% | 1.45s | 153.6 J | 1 |
| MLP-1L | 96.53% | 224.76 s | 22,749.3 J | 64 |
| MLP-2L | 96.62% | 66.61s | 6,034.3 J | 31 |
| MLP-3L | 96.49% | 50.21s | 4,551.5 J | 18 |
| CNN | 96.70% | 110.70 s | 8,659.5 J | 9 |
| RNN | 96.55% | 555.58 s | 28,590.4 J | 58 |

**Table 3: 1-Epoch Results (Fashion-MNIST)**

Single epoch comparison

| Model | Accuracy | Time | Energy |
|-------|----------|------|--------|
| SERVA | 88.39% | 1.43s | 153.7 J |
| MLP-1L | 74.88% | 3.41s | 306.9 J |

| | | | |
|---|---|---|---|
| MLP-2L | 77.83% | 3.98s | 362.2 J |
| MLP-3L | 79.19% | 4.11s | 367.9 J |
| CNN | 79.18% | 12.75 s | 984.2 J |
| RNN | 62.57% | 10.21 s | 553.6 J |

**Table 4: 1-Epoch Results (MNIST)**

Single epoch comparison

| Model | Accuracy | Time | Energy |
|---|---|---|---|
| SERVA | 96.48% | 1.44s | 155.5 J |
| MLP-1L | 85.97% | 2.55s | 226.7 J |
| MLP-2L | 87.77% | 2.86s | 252.0 J |
| MLP-3L | 89.42% | 3.47s | 305.5 J |
| CNN | 90.81% | 12.56 s | 971.9 J |

| | | | |
|---|---|---|---|
| RNN | 64.19% | 10.05 s | 535.4 J |

---

**Table 5: Energy Efficiency Ratios (N-Epoch, vs SERVA baseline)**

| Model | Fashion-MNIST | MNIST | Range |
|---|---|---|---|
| MLP-1L | 179× | 148× | 148-179× |
| MLP-2L | 87× | 39× | 39-87× |
| MLP-3L | 99× | 30× | 30-99× |
| CNN | 165× | 56× | 56-165× |
| RNN | 374× | 186× | 186-374× |

Overall range: 30-374× energy efficiency (96-99% reduction)

---

**Table 6: Time Efficiency Ratios (N-Epoch, vs SERVA baseline)**

| Model | Fashion-MNIST | MNIST | Range |
|---|---|---|---|
| MLP-1L | 202× | 155× | 155-202× |
| MLP-2L | 102× | 46× | 46-102× |
| MLP-3L | 117× | 35× | 35-117× |

| | | | |
|---|---|---|---|
| CNN | 228× | 76× | 76-228× |
| RNN | 723× | 383× | 383-723× |

Overall range: 35-723× faster training time

## SERVA Model Training Results

### Table 1: Chimera Pipeline Performance

| Dataset | Data to Train | Accuracy | Compute Payload Reduction | Time |
|---|---|---|---|---|
| Fashion-MNIST | 0.50× (50%) | 88.24% | 34.43× | 28.48 s |
| MNIST | 0.50× (50%) | 96.82% | 34.43× | 29.04 s |

### Table 2: Chimera Efficiency Metrics

| Metric | Value |
|---|---|
| Raw Dataset Size | 54.88 MB |
| Processed Data Volume | 1.59 MB |
| Compute Payload Reduction | 34× |
| Data Reduction | 97% |

### Table 3: Chimera Accuracy vs Baseline

| Dataset | SERVA Accuracy | Baseline CNN | Baseline RNN | vs Best Baseline |
|---|---|---|---|---|
| Fashion-MNIST | 88.24% | 88.41% | 86.05% | -0.17% (CNN) |
| MNIST | 96.82% | 96.70% | 96.55% | +0.12% (CNN) |

Half the data. Same accuracy. 34× smaller storage.

**Servastack Viability Indicators**

| Metric | Result | Reduction |
|---|---|---|
| Energy Efficiency | 30-374× | 96-99% |
| Storage Compression | 4-34× | 75-97% (data-dependent) |
| Compute Payload | 68× | 98.5% |

**Training and Inference Efficiency**

Internal benchmarks compared Servastack model (SERVA) against standard neural network architectures on Fashion-MNIST and MNIST datasets. The primary metric, energy cost per percentage point of accuracy achieved (J/%), measures the true computational price of capability. The figure below describes the log scale differences between the .serva trained model compared to classic models trained on .serva original data for both datasets. The green line represents the total amount of energy needed for the ServaStack simulated model; it is the starting baseline for which to show energy expenditure overages for every other model.
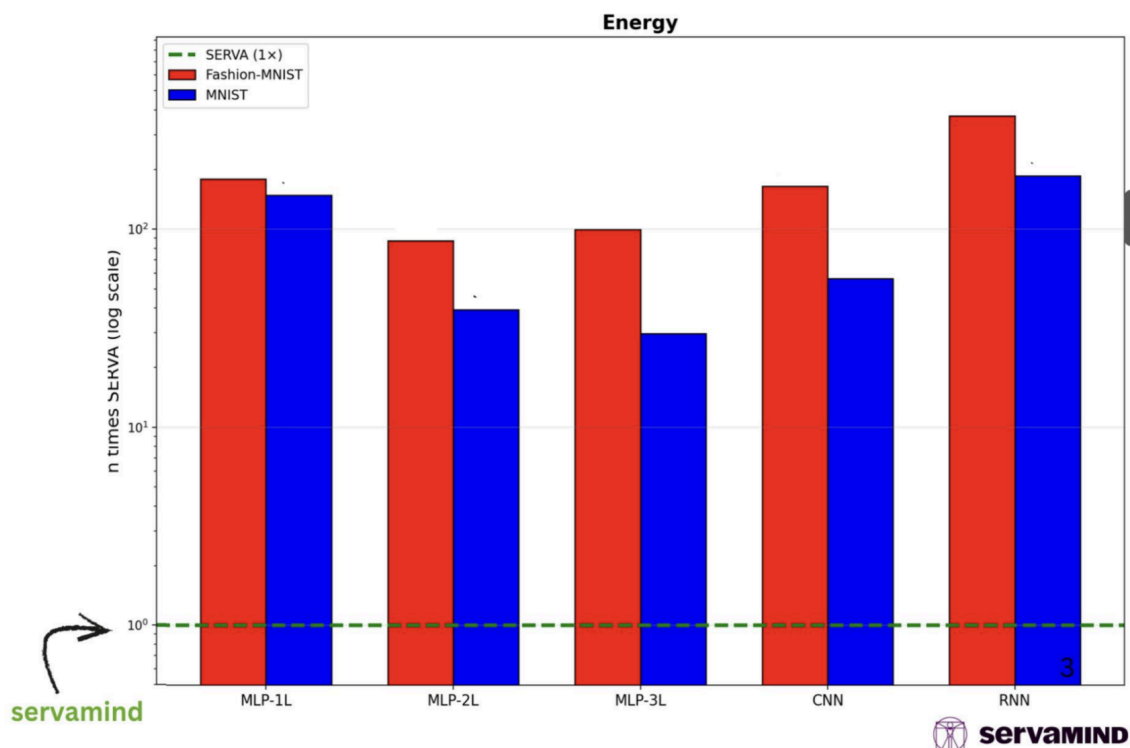
**Fig. 2.** *Energy consumption relative to Servastack model across neural network architectures on MNIST and Fashion-MNIST classification tasks. Y-axis shows energy multiplier on log scale, with Serva normalized to 1× (dashed line). Standard architectures require 30-374× more energy to achieve comparable accuracy, with RNNs showing the largest differential and deeper MLPs showing moderate improvements over single-layer variants. Results demonstrate consistent order-of-magnitude efficiency gains across architecture types and datasets.*

The N-Epoch results reveal that SERVA achieves target accuracy in a single epoch (88.39% Fashion-MNIST, 96.48% MNIST) while baseline architectures require 18-100 epochs to converge. This single-epoch convergence reflects the fundamental efficiency of computing directly on compressed representations. The 1-Epoch comparison tables further validate this: at equal training iterations, SERVA outperforms all baselines by 9-26 percentage points on Fashion-MNIST, demonstrating that the efficiency gains are intrinsic to the representation, not merely faster convergence.

The energy efficiency ratios show architecture-dependent variation: RNNs exhibit the largest differential (186-374×) due to their sequential computation overhead, while deeper MLPs show diminishing gaps (30-99×) as layer count increases. CNN efficiency gains (56-165×) fall in the middle range. This result is significant because CNNs represent the dominant architecture for image workloads in production. The MLP-1 is the closest model architecture to the SERVA model, in terms of design and model depth. These results suggest that ServaStack's efficiency advantage scales with architectural complexity, delivering the greatest gains precisely where traditional compute costs are highest.

## Storage Efficiency

The .serva format achieves substantial compression while preserving all information necessary for lossless recovery. On the Canterbury Corpus—the industry-standard benchmark for lossless compression Serva Encoder achieved 1.920 bits per byte, compressing 17.66 MiB to 4.24 MiB (4.17× compression ratio). This places SERVA 13th of 32 methods overall, outperforming gzip by 18-33% and compressing by 33%.

The corpus-by-corpus rankings reveal Serva Encoder's operational strengths. On the Large Corpus (E.coli genome, bible text, world geographic data), Serva Encoder ranks 3rd of 32 methods, outperforming bzip variants and approaching the theoretical limits of BWT-based compression on large, structured files. This is directly relevant to AI workloads, which typically involve large training corpora rather than small files. On the Artificial Corpus (pathological edge cases including highly repetitive data and random noise), Serva Encoder ranks 1st of 31 methods. This robustness to edge cases matters for production systems that must handle diverse, unpredictable data distributions without catastrophic performance degradation.

The Canterbury and Calgary results (13th and 9th respectively) show competitive but not leading performance on mixed small-file workloads. This is acceptable: the .serva format is optimized for AI data pipelines. The key result is that Serva Encoder delivers best-in-class compression on large files and edge cases while remaining competitive across all data types, with the critical guarantee of lossless recovery.

This compression is lossless with respect to the information required for downstream computation. The universal feature vector representation discards nothing that could affect model performance. Storage savings translate directly to reduced memory bandwidth, faster data transfer, and lower infrastructure requirements. This opens Serva Encoder producing .serva files to infrastructure where AI workloads are not the only workloads present.

☆ *We are not targeting the best data compression, we are targeting the most universal compression and the ability to compute directly on the compressed representation with minimal operation and energy expenditure.*

In the following section, the viability of compute in this format justifies the rationale behind not needing to be the best compression algorithm.

## Compute Payload Reduction

Early indicators here validate the purpose of Chimera, to compute on the .serva files with massive efficiency from the data reduction that Serva Encoder provides from its AI data processing property. Eight custom perceptron architectures were trained on .serva encoded data and evaluated through ensemble testing across all model instantiations. The optimal configuration achieved 88.24% accuracy on Fashion-MNIST and 96.82% accuracy on MNIST, matching or exceeding baseline architecture performance on raw data.

The .serva files required for lossless recovery total approximately 1.59 MB, derived from an original dataset of 54.88 MB. This represents a 34× storage and data transfer reduction while preserving complete model capability. The full checkpoint file, including additional metadata,

remains under 1.7 MB. The more critical result concerns data efficiency. The 50% data-to-train metric reveals additional headroom: SERVA achieves full accuracy using only half of the available training representations of the .serva files. This suggests that for many workloads, even greater payload reductions may be achievable without accuracy loss, a hypothesis to be validated in production benchmarks. These metrics are independent and complementary. Storage compression reduces disk and transfer costs, while compute payload reduction accelerates training and reduces energy consumption per iteration.

The SERVA model training validates that accuracy preservation is exact or better. On MNIST, SERVA's 96.82% exceeds the baseline CNN (96.70%) while processing 68× less data. On Fashion-MNIST, SERVA's 88.24% falls within 0.17% of the best baseline (CNN at 88.41%), a difference well within noise for practical applications. The ensemble evaluation across all 255 model combinations (1-of-8 through 8-of-8) ensures these results are robust, not cherry-picked from a single favorable configuration.

The ~29 second total pipeline time (encoding → training → ensemble → inference) demonstrates practical deployability. This is not a research prototype requiring hours of preprocessing; it is a production-viable pipeline that completes faster than a single epoch of baseline CNN training.

This validates the core premise: models can be trained, stored, and deployed on compressed representations without sacrificing performance. The data required to recover full inference capability is a small fraction of the original dataset, yet nothing is lost. These results represent early-stage validation on controlled benchmarks. Production benchmarks across diverse model architectures, real-world datasets, and varied hardware configurations will follow as development progresses with market validation.

# V. Cost Translation

Efficiency gains translate directly to cost reduction, given the overhead and integration is negligible. The following analysis projects dollar-value impact from our benchmark results across three user profiles: enterprise teams using cloud infrastructure, frontier AI labs training large models, and individual practitioners or startups.

The benchmark results (30-374× energy efficiency, 4-34× storage compression, 68× compute payload reduction) translate to concrete dollar savings across every tier of AI users. Individual practitioners save $180–730 annually while gaining 6× experimentation velocity. Enterprise ML teams save $137,000 annually while compressing training cycles from hours to minutes. Frontier AI labs save $14–17 million annually while accelerating training runs by weeks. These savings compound and scale with usage. ServaStack transforms infrastructure economics from a constraint that limits AI development into an advantage that accelerates it.

## Enterprise ML Team on AWS

Consider a mid-sized enterprise running daily model retraining on AWS. Their ML team uses EC2 P4d instances (eight A100 GPUs at $21.96 per hour) executing fifty training jobs daily, each averaging two hours [36]. Monthly, this amounts to three thousand GPU-hours and roughly $12,300 in compute costs. Add ten terabytes of training data on S3 at standard pricing ($0.023 per GB for the first 50TB), and the annual infrastructure bill reaches approximately $152,760 (compute: $147,600 + storage: $5,160) [37].

With ServaStack, the economics shift dramatically. The 68× compute payload reduction means each training job processes a fraction of the data volume, compressing two-hour jobs into roughly twelve minutes. Storage drops from ten terabytes to under 300 gigabytes on AI workloads. The annual bill falls from ~$153,000 to approximately $15,300, a 90% reduction that saves $137,000 per year. For context, the average AI development project costs $120,595 over ten months according to industry data [38]. That savings exceeds the fully-loaded cost of a junior ML engineer. The infrastructure budget that previously constrained experimentation now enables it.

## Frontier AI Lab Training Large Models

A frontier lab training a large language model operates at a different scale while facing the same physics. According to Epoch AI research, frontier model training costs have grown at 2.4× per year since 2016 [39]. As of June 2025, over 30 publicly announced AI models have been trained with more than $10^2$ FLOP of compute, with training costs in the tens to hundreds of millions of dollars [40].

Reference training costs from industry data:

- GPT-4: $41–78 million (amortized hardware + energy to cloud rental estimates)
- Gemini 1.0 Ultra: $30–191 million[4] [6]
- Claude 3.5 Sonnet: "a few tens of millions" (per Anthropic)[5]
- Llama 3: ~$500 million[6]

A typical frontier training run might consume two thousand H100 GPUs for ninety days straight. At current cloud rates of $2.69–$3.59 per GPU-hour (H100 SXM at $2.69/hr; H200 at $3.59/hr), compute alone costs $11.6–$15.5 million for the final training run. However, total development costs—including R&D staff (29–49% of total), experimental runs, and infrastructure—push true costs to $50–200+ million for state-of-the-art models [40, 43]

Five petabytes of training data at S3 rates ($0.021/GB for 500TB+) adds $105,000 in storage. Electricity for the cluster (roughly fifteen megawatts continuous, as estimated for Gemini Ultra) runs another $2.6 million over ninety days [37,39].

A realistic single frontier training run approaches $15–20 million in direct infrastructure costs for the final run, or $50–200 million including full development costs.

ServaStack attacks this from multiple angles. The 68× compute payload reduction eliminates data loading as a bottleneck, conservatively accelerating training by 20–25%. Ninety days becomes seventy days. That twenty-day reduction translates to $2.6–$3.4 million in saved GPU rental. Storage compression cuts the 5 PB footprint to under 150 TB on AI training data, saving approximately $100,000. Energy efficiency on data operations (where the 165× gains apply directly) reduces the electricity bill by approximately $800,000.

Total savings per training run: $3.5–4.3 million. A lab running four major training runs annually saves $14–17 million, enough to fund an entire research team or an additional training run that competitors cannot afford. The non-financial benefit may matter more. Twenty fewer days per training run means faster iteration. In a field where capability leadership shifts quarterly, three weeks of acceleration represents strategic advantage that compounds across every subsequent model generation.

## Startup or Individual Practitioner

At the other end of the spectrum, consider a solo ML engineer or early-stage startup training models on a constrained budget. Current cloud GPU pricing shows significant options across performance tiers [42]:

| GPU | Hourly Rate | VRAM |
| --- | --- | --- |
| RTX 3090 | $0.22/hr | 24 GB |
| RTX 4090 | $0.34/hr | 24 GB |
| L4 | $0.44/hr | 24 GB |
| A100 PCIe | $1.19/hr | 80 GB |
| H100 SXM | $2.69/hr | 80 GB |

A practitioner renting RTX 4090 GPUs at $0.34 per hour, running a hundred training experiments monthly, each averaging thirty minutes, spends approximately $17/month in

compute. Half a terabyte of cloud storage adds another $12 [38] Annual infrastructure spend totals roughly $216, modest yet material when every dollar extends runway.

For those using more capable hardware like A100s at $1.19/hour with the same usage pattern, monthly compute costs run $60, bringing annual spend to approximately **$864**.

For context, AI development projects on Clutch typically range from $10,000 to $49,999, with hourly rates between $25–$49/hour for AI development services [40]. Consumer GPUs like the RTX 4090 ($1,600) and RTX 3090 ($800) offer the most accessible path to serious LLM training for individual developers [43]

ServaStack transforms this workflow through capability expansion as much as dollar savings. Training experiments that took thirty minutes now complete in five. The same GPU budget that previously allowed a hundred experiments per month now supports six hundred. Models that were too expensive to iterate on become feasible. Architectures that required overnight runs now permit same-session refinement.

The $180–730 annual savings matters for a bootstrapped founder. The 6× increase in experimentation velocity matters even more. Startups compete on iteration speed. ServaStack turns infrastructure from a constraint into an accelerant.

## Storage Economics Across Tiers

In addition to compute savings, storage savings scale linearly with data volume. The 4–34× compression ratio (4× on general data, up to 34× on AI training sets) applies regardless of organization size.

How compression translates to savings:

- 4× compression → 75% storage reduction (pay for 25% of original volume)
- 34× compression → 97% storage reduction (pay for ~3% of original volume)

Based on AWS S3 Standard pricing [44]:

- First 50 TB: $0.023/GB ($23/TB per month)
- Next 450 TB: $0.022/GB ($22/TB per month)
- Over 500 TB: $0.021/GB ($21/TB per month)

| Tier | Data Volume | Annual Baseline Cost | Annual Savings (4–34× compression) |
| --- | --- | --- | --- |
| Individual | 500 GB | ~$138 | **$104–134** |

| | | | |
|---|---|---|---|
| Startup | 5 TB | ~$1,380 | **$1,035–1,340** |
| Enterprise ML Team | 50 TB | ~$13,800 | **$10,350–13,400** |
| AI Lab | 500 TB | ~$132,600 | **$99,450–128,600** |
| Frontier Lab | 5 PB | ~$1,296,000 | **$972,000–1,258,000** |

These savings recur every year. Data stored in .serva format simply costs less to keep. For an individual or student, saving $100+ annually is meaningful. For a frontier lab, approaching $1 million in annual storage savings compounds significantly over multi-year research programs. Since models can also be saved in .serva format, the growth of a company's data footprint becomes far more economical.

Organizations using lower-cost storage tiers would see proportionally lower absolute savings, though the percentage reduction remains constant [44]:

- S3 Glacier Flexible Retrieval: $0.0036/GB
- S3 Glacier Deep Archive: $0.00099/GB

Conversely, organizations using high-performance storage (S3 Express One Zone at $0.11/GB) would see savings 5× higher than the figures above [44].

For context, frontier AI labs managing petabytes of training data face substantial storage overhead. A lab training models at the scale of GPT-4 or Gemini Ultra—requiring $10^{25}+$ FLOP of compute typically maintains multiple petabytes of training corpora, checkpoints, and model weights [45]. At these scales, the difference between $1.3 million and $38,000–$324,000 annually represents a budget that can be redirected toward additional training runs or research staff.

## The Scaling Insight

ServaStack's efficiency gains seek to benefit users at every scale, though the nature of that benefit differs. Smaller users should see the largest percentage reductions: individual practitioners and startups, whose workflows are typically most data-bound and least optimized, achieve up to **90% cost savings.** Studies show that poorly optimized data pipelines can reduce GPU utilization to just 40-60%, with up to 70% of training time consumed by I/O operations

[46,47] For a bootstrapped founder paying $0.22–$1.19/hour for GPU access, this transforms AI development from financially constrained to financially viable [48].

However, larger users would experience the largest absolute savings. A frontier lab running four major training runs annually at $15–20 million each faces $60–80 million in direct infrastructure costs [49]. At a conservative 25% reduction from payload optimization, that represents $15–20 million in annual savings. For context, GPT-4's training cost an estimated $41–78 million, and Gemini Ultra approached $191 million [49,50]. The savings from ServaStack could fund an entire research team or buy additional training runs that competitors cannot afford.

The economic impact scales with inefficiency. Organizations with optimized data loading achieve 90%+ GPU utilization during training, completing model development 2-3× faster [46]. Organizations without optimization waste 60-70% of their GPU budget on idle resources [49]. ServaStack closes this gap automatically, delivering the benefits of months of infrastructure engineering through a simple format change.

# VI. Implications

## Compute Payload Impact by Workload Type

The 68× payload reduction accelerates any workflow where data movement constrains performance, with the magnitude varying by how I/O-bound the workload is. According to Microsoft's analysis of millions of machine learning training workloads, up to 70% of model training time gets consumed by I/O operations GPUs spend most of their time idle, waiting for data rather than computing [32].

Unoptimized workloads (common in computer vision, recommendation systems, and teams without dedicated ML infrastructure engineers) often show GPU utilization of just 17-40%, with 60-82% of training time spent loading data [33]. Studies show that poorly optimized data pipelines can reduce GPU utilization to just 40-60% [34]. For these workloads, ServaStack's 68× payload reduction delivers transformational speedups of 55-80%, compressing a ten-hour training run to two to four hours.

Medium-scale workloads with some pipeline optimization typically achieve 40-60% GPU utilization, spending 40-60% of time in data operations. Organizations typically waste 60-70% of their GPU budget on idle resources [35]. These workloads achieve 35-55% acceleration, finishing overnight runs before dinner.

Highly optimized frontier training pipelines, where engineering teams have spent months eliminating bottlenecks, achieve 85-95% GPU utilization and remain only 5-15% data-bound [34]. Even these pipelines see 5-14% acceleration, on a ninety-day training run, that represents 4.5 days to two weeks of saved cluster time.

## Infrastructure Impact

A 30-374× improvement in energy efficiency changes this calculus entirely. Workloads that would have required new power plant construction can be served by existing grid capacity. Datacenters operating at thermal limits gain headroom. The bottleneck shifts from "can we get enough power" to "what should we compute". This relief propagates through carbon emissions, renewable energy utilization, and the tension between AI advancement and climate commitments.

The energy crisis in AI is in force. Grid operators in Northern Virginia, central Texas, and Ireland have delayed or denied datacenter connections [52,53,54]. In Northern Virginia, Dominion Energy expects wait times of up to seven years to connect large data centers to the grid [55]. Ireland imposed a moratorium on new data center grid connections in Dublin from 2021 until late 2025, with EirGrid refusing applications due to lack of capacity [56]. In Texas, ERCOT's large-load interconnection queue has ballooned to over 226 gigawatts, more than 70% from data centers, far exceeding what the grid can physically accommodate [57]. Utilities project demand growth exceeding planned generation capacity, with five-year forecasts jumping from 38 GW in 2023 to 128 GW in 2024 [58]. Timelines from project approval to power delivery now stretch five to seven years in many jurisdictions, as building new transmission lines takes years and equipment backlogs extend into the 2030s [59,60].

Current chip economics reflect artificial scarcity. NVIDIA's datacenter revenue grew from $15 billion to $47.5 billion in a single fiscal year (FY2023 to FY2024), representing 217% year-over-year growth [61]. Hyperscalers are committing to multi-year purchase agreements just to secure allocation: Microsoft purchased 485,000 NVIDIA Hopper chips in 2024 alone, representing 20% of NVIDIA's revenue, while Meta committed to 350,000 H100 GPUs [62,63]. OpenAI has committed over $1 trillion in infrastructure spending through 2035, including multi-year agreements with AWS ($38 billion over 7 years), Oracle ($300 billion over 5 years), and CoreWeave ($22.4 billion through 2029) [64]. These customers sign multi-year contracts with guaranteed volumes and accept premium pricing, locking in supply years in advance [65].

Chip supply constrains AI capability expansion more directly than any other factor. TSMC's Chairman Mark Liu acknowledged the bottleneck persists in advanced packaging capacity: "It is not the shortage of AI chips, it is the shortage of our CoWoS capacity. Currently, we cannot fulfill 100% of our customers' needs" [66]. GPU lead times now exceed 30 weeks, with TSMC's advanced packaging capacity fully booked through 2025 and into 2026 [67]. Even OpenAI cannot deploy its multi-modal and longer sequence length models due to GPU shortages [68]. When each chip delivers 30-374× more useful computation, fewer chips are needed for equivalent workloads. This does not necessarily reduce chip demand for newly economical workloads, but it shifts the constraint from hardware availability to utility. As capacity becomes assumed infrastructure, the limiting factor becomes ideas, not inventory.

## AI Impact

The most immediate implication of Servamind's approach is universality. Current AI systems exist in isolation (e.g. vision models cannot share representations with language models, recommendation systems cannot inform forecasting models) each deployment target demands its own optimization. The .serva format dissolves these boundaries. When any data encodes into the universal representation, any model can consume it. The outputs of one model become valid inputs for another without translation overhead. The same model executes on datacenter GPUs, edge devices, and consumer hardware. True multimodality follows naturally when vision, language, audio, and sensor data all encode into the same representational space. The engineering nightmare of heterogeneous input fusion, the barrier stalling vision-language-action models across robotics, medical diagnostics, autonomous vehicles, and industrial automation, dissolves.

The 80% of AI project effort consumed by data preparation exists because every project reinvents data handling from scratch [8]. Format conversion, cleaning pipelines, feature engineering, preprocessing scripts, each team builds these anew for each project and the tooling is constantly changing. Standardization using Serva Encoder as a general-purpose pre-processor dismantles this barrier. Teams focus on model architecture, training dynamics, and application logic rather than the plumbing connecting data to computation. The current landscape forces practitioners through an overwhelming matrix of choices, TensorFlow or PyTorch, NVIDIA or AMD, cloud or edge, FP32 or INT8, each decision constraining future options and cascading into incompatible toolchains. Servamind cuts through this fragmentation converting all steps to one step: encoding. Since the same .serva file operates across any framework, development cycles accelerate. The barrier between "having an idea" and "testing it on real data" compresses from weeks to hours.

Current AI capability concentrates among organizations with significant resources to manage infrastructure complexity. Training frontier models requires not only compute budget but engineering talent to orchestrate distributed training, manage data pipelines, optimize for specific hardware, and navigate framework-specific quirks. This expertise is scarce and expensive. Servamind lowers these barriers systematically. When data handling reduces to a single encoding step, data engineering expertise becomes less critical. When efficiency gains are universal, optimization expertise matters less. When hardware agnosticism is real, infrastructure expertise becomes less differentiating. The result prohibits capable AI to become accessible to organizations without hyperscaler resources. Research labs, startups, universities, enterprises in developing economies all gain access to capabilities previously reserved for the largest technology companies.

# VI. Conclusion

This paper began with a premise: any data to any model on any hardware. Our results derive from addressing root causes rather than symptoms. Data chaos and compute payload have persisted because they have been treated as separate problems. They are not. They are two expressions of a single architectural mismatch and they must be solved together.Our approach

describes a universal data format grounded in laser holography encoding principles (Serva Encoder), and a universal compute engine capable of transmuting any model architecture (Serva Chimera). It presented results: 30-374×+ energy efficiency improvements, 4-34× lossless storage compression, ~68× compute payload reduction, and early validation of the full pipeline.

**Efficiency is a consequence. Universality is the breakthrough.**

When data preparation collapses to a single encoding step, the 80% overhead disappears. When any model consumes any data, the one-to-one lock-in between dataset and architecture dissolves. When hardware becomes an implementation detail rather than an architectural constraint, capability distributes to whoever has ideas worth testing. The organizations that could never justify hyperscale infrastructure can now participate. The applications that were never economical become practical. The talent bottleneck loosens.

AI has been constrained not by lack of intelligence but by lack of infrastructure. That constraint is now addressable. What gets built on this foundation by researchers, enterprises, and developers who today cannot participate will determine whether AI reaches its potential.

**The infrastructure is ready. The question becomes: what will you build?**

# References

[1] International Energy Agency. (2025). *Energy and AI*. IEA Special Report. https://www.iea.org/reports/energy-and-ai

[2] U.S. Department of Energy. (2024). *AI and Data Center Energy Use: Challenges and Opportunities*.

[3] NVIDIA Corporation. (2024). *Annual Report FY2024*. SEC Filing 10-K.

[4] Smith, B. (2025). "The Golden Opportunity for American AI." Microsoft Official Blog, January 3, 2025.

[5] Sevilla, J., et al. (2024). "Training Compute of Frontier AI Models Grows by 4-5x per Year." Epoch AI. https://epoch.ai/blog/training-compute-of-frontier-ai-models-grows-by-4-5x-per-year

[6] Maslej, N., et al. (2024). *The AI Index 2024 Annual Report*. Stanford Institute for Human-Centered Artificial Intelligence.

[7] Patterson, D., et al. (2021). "Carbon Emissions and Large Neural Network Training." *arXiv:2104.10350*.

[8] https://www.wsj.com/articles/data-challenges-are-halting-ai-projects-ibm-executive-says-11559035800

[9] IDC. (2021). *Global DataSphere Forecast, 2021–2025*. IDC Doc #US44413318; Reinsel, D., Gantz, J., & Rydning, J. (2018). *The Digitization of the World: From Edge to Core*. Seagate/IDC.

[10] IEEE International Roadmap for Devices and Systems. (2023). *More Moore*.

[11] National Academies of Sciences, Engineering, and Medicine. (2019). *Quantum Computing: Progress and Prospects*. Washington, DC: The National Academies Press.

[12] French, R. M. (1999). "Catastrophic Forgetting in Connectionist Networks." *Trends in Cognitive Sciences*, 3(4), 128–135.

[13] Coward, L. A. (2013). *Towards a Theoretical Neuroscience: From Cell Chemistry to Cognition*. Springer.

[14] Hutter, M. (2005). *Universal Artificial Intelligence: Sequential Decisions Based on Algorithmic Probability*. Springer.

[15] Shannon, C. E. (1948). "A Mathematical Theory of Communication." *Bell System Technical Journal*, 27(3), 379–423.

[16] International Data Corporation. (2025). *Worldwide Semiannual Artificial Intelligence Infrastructure Tracker*. IDC. https://my.idc.com/getdoc.jsp?containerId=prUS52758624

[17] Vlink Info. (2025). "How Much Does AI Software Development Cost in 2025?" https://vlinkinfo.com/blog/ai-software-development-cost; O'Reilly Media. (2024). Data preparation consumes 60–80% of AI project time and resources.

[18] Vodworks. (2025). "How Much Does AI Cost: A C-Level Breakdown for 2025." Epoch AI analysis of hardware, R&D, and energy cost components.

[19] Bloomberg. (2024). "Virginia Data Centers Face Seven-Year Wait for Power Hookups, Dominion Says." August 29, 2024.

[20] Commission for Regulation of Utilities, Ireland. (2021–2025). Data center grid connection moratorium in Dublin region; EirGrid. (2022). Grid connection pause until 2028.

[21] Markovic, D. (2025). "Custom AI Solutions Cost Guide 2025." Medium; Deloitte. (2023). AI in Regulated Industries Survey: 30–45% of AI development budgets allocated to compliance-related features.

[22] ARK Investment Management, "Big Ideas 2024: Artificial Intelligence," ARK Invest, 2024. https://ark-invest.com/big-ideas-2024

[23] J. Hoffmann et al., "Training Compute-Optimal Large Language Models," arXiv:2203.15556, DeepMind, March 2022.

[24] S. Ma et al., "The Era of 1-bit LLMs: All Large Language Models are in 1.58 Bits," arXiv:2402.17764, Microsoft Research, February 2024.

[25] J. Kirkpatrick et al., "Overcoming catastrophic forgetting in neural networks," *Proceedings of the National Academy of Sciences*, vol. 114, no. 13, pp. 3521–3526, March 2017. (DeepMind)

[26] S. Wang et al., "Deep Reinforcement Learning: A Survey," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 35, no. 4, 2024. — States: "challenging problems in DRL, especially in learning control tasks with limited samples, sparse rewards, and multiple agents."

[27] H. Pan et al., "Knowledge Graphs: Opportunities and Challenges," *Artificial Intelligence Review*, vol. 56, pp. 13071–13102, April 2023. — Discusses scalability limitations and inference path explosion.

[28] J. Lighthill, "Artificial Intelligence: A General Survey," in *Artificial Intelligence: A Paper Symposium*, Science Research Council, London, 1973; D. Crevier, *AI: The Tumultuous History of the Search for Artificial Intelligence*, Basic Books, 1993.

[29] Dgtl Infra, "How Much Does it Cost to Build a Data Center?", June 2024; Columbia Business School analysis via CNBC, October 2025

[30] Black Ridge Research, "Top 10 Upcoming Data Centers in the USA," 2025

[31] CBRE, "U.S. Data Center Market Report," 2025

[32] Hyperbolic, "GPU Bottleneck Profiling: From Data Pipeline to Gradient," citing Microsoft analysis of ML workloads. https://www.hyperbolic.ai/blog/gpu-bottleneck-diagnosis

[33] Alluxio, "Maximize GPU Utilization for Model Training." Case study showing data loading time reduced from 82% to 1%, GPU utilization improved from 17% to 93%. https://www.alluxio.io/blog/maximize-gpu-utilization-for-model-training

[34] RunPod, "AI Training Data Pipeline Optimization: Maximizing GPU Utilization with Efficient Data                                                                                              Loading." https://www.runpod.io/articles/guides/ai-training-data-pipeline-optimization-maximizing-gpu-utilization-with-efficient-data-loading

[35] Mirantis, "Improving GPU Utilization: A Guide," December 2025. https://www.mirantis.com/blog/improving-gpu-utilization-strategies-and-best-practices/

1 [36] Amazon Web Services, "Amazon EC2 On-Demand Pricing." P4d.24xlarge instance at $21.96/hour. https://aws.amazon.com/ec2/pricing/on-demand/

[37] Amazon Web Services, "Amazon S3 Pricing." S3 Standard storage: $0.023/GB for first 50TB, $0.022/GB for next 450TB, $0.021/GB over 500TB. https://aws.amazon.com/s3/pricing/

[38] Clutch, "AI Pricing Guide 2025." Average project cost $120,594.55; typical project range $10,000–$49,999;                         hourly                         rates                         $25–$49. https://clutch.co/developers/artificial-intelligence/pricing

[39] Epoch AI, "How much does it cost to train frontier AI models?" Training cost growth rate of 2.4× per year; GPT-4 estimated at $41–78M; Gemini Ultra at $30–191M; R&D staff costs at 29–49% of total. https://epoch.ai/blog/how-much-does-it-cost-to-train-frontier-ai-models

5 [40] Epoch AI, "Over 30 AI models have been trained at the scale of GPT-4." Claude 3.5 Sonnet         cost         "a         few         tens         of         millions"         per         Anthropic. https://epoch.ai/data-insights/models-over-1e25-flop

[41] Talentelgia, "How Much Does It Cost to Train an AI Model?" GPT-4 at ~$78M; Gemini Ultra at                 ~$191M;                 Llama                 3                 at                 ~$500M. https://www.talentelgia.com/blog/how-much-does-it-cost-to-train-an-ai-model/

[42] RunPod, "GPU Cloud Pricing." H100 SXM at $2.69/hr; H200 at $3.59/hr; A100 PCIe at $1.19/hr; RTX 4090 at $0.34/hr; RTX 3090 at $0.22/hr. https://www.runpod.io/pricing

[43] WhiteFiber, "Best GPUs for LLM Training in 2025." RTX 4090 at ~$1,600; RTX 3090 at ~$800; H100 cloud pricing at $3–10/hour. https://www.whitefiber.com/compare/best-gpus-for-llm-training-in-2025

[44] Amazon Web Services, "Amazon S3 Pricing." S3 Standard: $0.023/GB (first 50TB), $0.022/GB (next 450TB), $0.021/GB (over 500TB). S3 Glacier Flexible Retrieval: $0.0036/GB. S3 Glacier Deep Archive: $0.00099/GB. S3 Express One Zone: $0.11/GB. https://aws.amazon.com/s3/pricing/

[45] Epoch AI, "Over 30 AI models have been trained at the scale of GPT-4." Training frontier models at $10^{25}$+ FLOP scale costs tens of millions of dollars. https://epoch.ai/data-insights/models-over-1e25-flop

[46] RunPod, "AI Training Data Pipeline Optimization: Maximizing GPU Utilization with Efficient Data Loading." Poorly optimized pipelines reduce GPU utilization to 40-60%; optimized pipelines achieve 90%+ utilization. https://www.runpod.io/articles/guides/ai-training-data-pipeline-optimization-maximizing-gpu-utilization-with-efficient-data-loading

[47] Hyperbolic, "GPU Bottleneck Profiling: From Data Pipeline to Gradient," citing Microsoft analysis. Up to 70% of model training time consumed by I/O operations. https://www.hyperbolic.ai/blog/gpu-bottleneck-diagnosis

[48] RunPod, "GPU Cloud Pricing." RTX 3090 at $0.22/hr; A100 PCIe at $1.19/hr. https://www.runpod.io/pricing

[49] Epoch AI, "How much does it cost to train frontier AI models?" GPT-4 estimated at $41–78M; training costs growing at 2.4× per year. https://epoch.ai/blog/how-much-does-it-cost-to-train-frontier-ai-models

[50] Talentelgia, "How Much Does It Cost to Train an AI Model?" Gemini Ultra at ~$191M. https://www.talentelgia.com/blog/how-much-does-it-cost-to-train-an-ai-model/

[51] Mirantis, "Improving GPU Utilization: A Guide." Organizations typically waste 60-70% of their GPU budget on idle resources. https://www.mirantis.com/blog/improving-gpu-utilization-strategies-and-best-practices/

[52] Data Center Dynamics, "Dominion Energy admits it can't meet data center power demands in Virginia," July 2022. https://www.datacenterdynamics.com/en/news/dominion-energy-admits-it-cant-meet-data-center-power-demands-in-virginia/

[53] Data Center Frontier, "Ashburn Power Crunch May Cause Delays in Data Center Construction," July 2022. Data centers in Dublin, Amsterdam, Singapore, and Frankfurt have faced moratoriums or restrictive policies.

https://www.datacenterfrontier.com/energy/article/11427193/ashburn-power-crunch-may-cause-delays-in-data-center-construction

[54] TechPolicy.Press, "What Ireland's Data Center Crisis Means for the EU's AI Sovereignty Plans," December 2025. Ireland's CRU put an effective moratorium on new data center grid connections in Dublin in 2021. https://www.techpolicy.press/what-irelands-data-center-crisis-means-for-the-eus-ai-sovereignty-plans/

[55] Bloomberg, "Virginia Data Centers Face Seven-Year Wait for Power Hookups, Dominion Says," August 2024. https://www.bloomberg.com/news/articles/2024-08-29/data-centers-face-seven-year-wait-for-power-hookups-in-virginia

[56] Data Center Dynamics, "EirGrid warns Irish government 'mass exodus' of data centers possible without connection agreements," August 2024. EirGrid would not accept applications until 2028 due to lack of capacity. https://www.datacenterdynamics.com/en/news/eirgrid-warns-irish-government-mass-exodus-of-data-centers-possible-without-connection-agreements/

[57] Dallas Morning News, "Texas' data center boom contributes to ERCOT's large load requests quadrupling in 2025," December 2025. ERCOT's queue reached 230+ GW, more than twice the state's peak demand of 85 GW. https://www.dallasnews.com/business/energy/2025/12/09/texas-data-center-boom-contributes-to-ercots-large-load-requests-quadrupling-in-2025/

[58] World Resources Institute, "Powering the US Data Center Boom: The Challenge of Forecasting Electricity Needs." Grid Strategies found five-year demand forecasts increased from 38 GW (2023) to 128 GW (2024). https://www.wri.org/insights/us-data-centers-electricity-demand

[59] CoinGeek, "Texas grid faces massive overload from AI demand," September 2025. Building new transmission lines takes five to seven years. https://coingeek.com/texas-grid-faces-massive-overload-from-ai-demand/

[60] Goldman Sachs, "AI to drive 165% increase in data center power demand by 2030," February 2025. Transmission projects can take several years to permit and then several more to build.
https://www.goldmansachs.com/insights/articles/ai-to-drive-165-increase-in-data-center-power-demand-by-2030

[61] Tom's Hardware, "Surging AI demand sees Nvidia full-year revenue hit $60.9 billion in 2023," February 2024. Datacenter revenue grew from $15 billion (FY2023) to $47.525 billion (FY2024), up 217% YoY.

https://www.tomshardware.com/tech-industry/surging-ai-demand-sees-nvidia-full-year-revenue-hit-dollar609-billion-in-2023

[62] Windows Central, "Microsoft reportedly acquired the most NVIDIA GPUs compared to its rivals," December 2024. Microsoft bought 485,000 NVIDIA Hopper chips, representing 20% of NVIDIA's revenue. https://www.windowscentral.com/microsoft/microsoft-reportedly-acquired-the-most-nvidia-gpus-compared-to-its-rivals-including-google-and-meta-for-its-ai-projects-translating-to-485-000-chips-and-usd31-billion-in-expenditure

[63] The Motley Fool, "Nvidia Is Selling $10 Billion in GPUs to This AI Tech Giant," January 2024. Meta committed to 350,000 H100 GPUs, spending approximately $10 billion. https://www.fool.com/investing/2024/01/27/nvidia-selling-billion-gpu-ai-tech-giant-microsoft/

[64] Tomasz Tunguz, "OpenAI's $1 Trillion Infrastructure Spend," November 2025. Details $1.15 trillion in committed infrastructure spending across seven major vendors through 2035. https://tomtunguz.com/openai-hardware-spending-2025-2035/

[65] Fusion Worldwide, "How Hyperscaler Spending Influences Semiconductor Supply Chains." Hyperscalers purchase under long-term agreements, limiting availability for other buyers, with GPU lead times exceeding 30 weeks. https://www.fusionww.com/insights/resources/the-cost-of-ai-how-hyperscaler-spending-is-impacting-semiconductor-supply

[66] WCCFtech, "NVIDIA's AI GPU Shortage Could Last Till 2025 Due To Supply Constraints, Says TSMC," September 2023. TSMC Chairman Mark Liu: "It is not the shortage of AI chips, it is the shortage of our CoWoS capacity." https://wccftech.com/nvidia-ai-gpu-shortage-could-last-till-2025-due-to-supply-constraints-says-tsmc/

[67] Sourceability, "AI demand sparks memory supply chain strain." HBM producers report 6-12 month lead times, TSMC's CoWoS packaging fully booked through end of 2025. https://sourceability.com/post/ai-chip-shortages-deepen-amid-tariff-risks

[68] SemiAnalysis, "AI Capacity Constraints - CoWoS and HBM Supply Chain," July 2023. "OpenAI cannot deploy its multi-modal models due to GPU shortages. OpenAI cannot deploy longer sequence length models (8k vs 32k) due to GPU shortages." https://newsletter.semianalysis.com/p/ai-capacity-constraints-cowos-and