

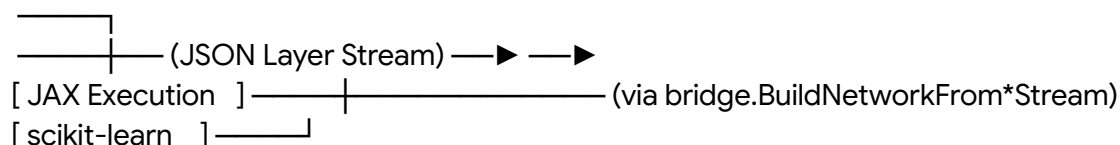
# Comprehensive Analysis of Multi-Engine Weight Absorption and Edge-Native Retraining within Loom's Planetbridging Framework

The landscape of edge artificial intelligence is undergoing a significant transition from cloud-dependent inference pipelines to fully sovereign, local execution.<sup>1</sup> Historically, deploying neural networks to resource-constrained devices has been a highly fragmented process.<sup>3</sup> Each major AI runtime operates as an isolated ecosystem, characterized by its own closed file formats, custom operator dialects, and hardware-specific math libraries.<sup>3</sup> This fragmentation often forces developers to perform lossy model conversions, which introduce implementation bugs, high maintenance overhead, and cumulative numerical drift.<sup>4</sup>

Loom's Planetbridging framework addresses this challenge by providing a universal runtime bridging layer.<sup>3</sup> Rather than compiling static, intermediate file-based representations, Loom uses an in-memory virtualization pipeline designed to absorb live model parameters directly from disparate host frameworks.<sup>3</sup> This report evaluates the architectural viability of Loom's multi-engine weight absorption, analyzes its native edge-training capabilities, and compares its design to other contemporary edge AI engines.

## The Architectural Viability of Multi-Engine Weight Absorption

The core technical challenge in edge deployment is the structural consolidation of model weights developed across heterogeneous training environments.<sup>3</sup> Traditional compilation frameworks treat weight migration as an offline, ahead-of-time translation process, which is highly sensitive to semantic gaps between the source engine's abstract syntax tree and the target engine's operator set.<sup>4</sup> Loom's Planetbridging addresses this by introducing a live, in-memory streaming protocol that can dynamically ingest active weights from up to four distinct source engines—PyTorch, TensorFlow, JAX, and scikit-learn—and consolidate them into a single, cohesive runtime entity.<sup>1</sup>



The ingestion process begins inside the native python execution environment where the source model is running.<sup>3</sup> Instead of saving the network to disk, a lightweight python helper library captures the active, in-memory weight matrices.<sup>3</sup> It serializes these parameters and the network topology into a JSON layer stream, which is broadcast to a local network interface via POST /api/v1/loom/stream/\*.<sup>3</sup> On the receiving end, Loom's native Go-based comparison host catches the stream and uses localized builder methods, such as bridge.BuildNetworkFrom\*Stream, to reconstruct the network architecture and compile it into a unified, binary .entity checkpoint.<sup>1</sup>

This streaming approach is validated across twelve standard bedrock layers<sup>3</sup>:

Bedrock Layer ID	Source Engine Coverage	Loom Stream Status	Loom Compare Status	Empirical Tolerance Threshold
<b>Dense</b>	PyTorch, TensorFlow, JAX, scikit-learn <sup>3</sup>	PASS <sup>3</sup>	43/48 PASS <sup>3</sup>	Single-precision fp32 (< ) <sup>3</sup>
<b>CNN1 / CNN2 / CNN3</b>	PyTorch, TensorFlow, JAX <sup>3</sup>	PASS <sup>3</sup>	12/12 PASS <sup>3</sup>	Single-precision fp32 (< ) <sup>3</sup>
<b>Multi-Head Attention (MHA)</b>	PyTorch, TensorFlow, JAX <sup>3</sup>	PASS <sup>3</sup>	12/12 PASS <sup>3</sup>	Single-precision fp32 (< ) <sup>3</sup>
<b>LSTM</b>	PyTorch, TensorFlow, JAX <sup>3</sup>	PASS <sup>3</sup>	12/12 PASS <sup>3</sup>	Single-precision fp32 (< ) <sup>3</sup>
<b>Simple RNN</b>	PyTorch, TensorFlow, JAX <sup>3</sup>	PASS <sup>3</sup>	12/12 PASS <sup>3</sup>	Single-precision fp32 (< ) <sup>3</sup>
<b>LayerNorm</b>	PyTorch, TensorFlow, JAX <sup>3</sup>	PASS <sup>3</sup>	24/24 PASS <sup>3</sup>	Single-precision fp32 (< ) <sup>3</sup>

<b>Embedding</b>	PyTorch, TensorFlow, JAX <sup>3</sup>	PASS <sup>3</sup>	24/24 EXACT <sup>3</sup>	Bit-exact bitwise parity (0.0) <sup>3</sup>
<b>RMSNorm</b>	PyTorch, TensorFlow, JAX <sup>3</sup>	PASS <sup>3</sup>	24/24 PASS <sup>3</sup>	Single-precision fp32 (< ) <sup>3</sup>
<b>SwiGLU</b>	PyTorch, TensorFlow, JAX <sup>3</sup>	PASS <sup>3</sup>	24/24 PASS <sup>3</sup>	Single-precision fp32 (< ) <sup>3</sup>
<b>Residual</b>	PyTorch, TensorFlow, JAX <sup>3</sup>	PASS <sup>3</sup>	24/24 PASS <sup>3</sup>	Single-precision fp32 (< ) <sup>3</sup>
<b>Mixer v1</b>	PyTorch, TensorFlow, JAX <sup>3</sup>	PASS <sup>3</sup>	3/3 PASS <sup>3</sup>	Single-precision fp32 (< ) <sup>3</sup>
<b>Mixer v2</b>	PyTorch, TensorFlow, JAX <sup>3</sup>	PASS <sup>3</sup>	3/3 POC PASS <sup>3</sup>	Floating-point limit ( $\sim 5 \times 10^3$ ) <sup>3</sup>

The physical pipeline of streaming a live model and comparing its output to the original host engine is illustrated by running `O1_hello_stream.py` on a dense MLP model (`mlp_32_64_32_16_8_relu` with 4 layers).<sup>3</sup> This script executes a forward pass on PyTorch, streams the live weights to the comparison host, and saves the binary .entity file directly to the disk.<sup>3</sup> The outputs generated by PyTorch and Loom's runtime are then validated against a shared test vector, matching within standard single-precision tolerances.<sup>3</sup> The same verification pipeline successfully validates multi-head attention layers, such as `mha_8_2_4`, confirming that multi-engine parameter absorption is mathematically valid and structurally sound.<sup>3</sup>

## Active Edge Training and the Deterministic Neural Virtual Machine

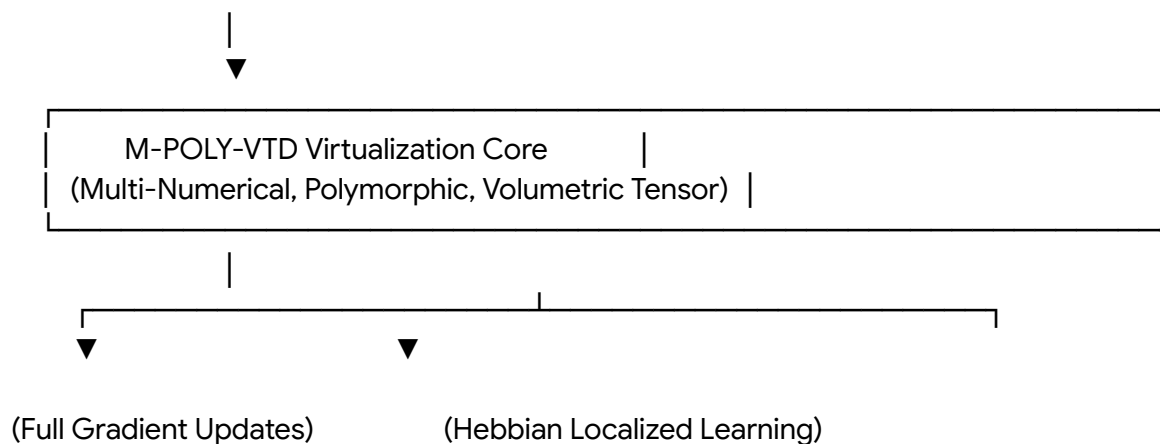
Consolidating multi-engine weights into a single execution target is only the initial phase; the subsequent architectural requirement is enabling active training within that same runtime.<sup>1</sup>

Loom is designed as a pure-Go Deterministic Neural Virtual Machine (DNVM).<sup>1</sup> This architecture decouples execution from system-level platform variations, allowing backpropagation and weight optimization to occur natively on the target device without frozen dependencies.<sup>1</sup>

## The Prerequisite of Bit-Exact Cross-Platform Determinism

In standard deep learning runtimes, on-device training is highly susceptible to mathematical divergence.<sup>4</sup> Floating-point operations executed on disparate physical architectures (such as ARM processors versus NVIDIA CUDA GPUs) typically yield minor numerical variations due to compiler optimizations, platform-specific math libraries, and hardware-specific rounding behaviors.<sup>5</sup> During inference, these micro-divergences are generally negligible; however, during backpropagation, cumulative mathematical drift across multiple training iterations can cause the gradient calculations to diverge, ultimately degrading the model.<sup>5</sup>

The DNVM architecture addresses this by enforcing bit-exact cross-platform determinism.<sup>1</sup> By utilizing a pure-Go execution core with explicit, deterministic floating-point operations, Loom guarantees a Mean Absolute Error ( $MAE$ ) of less than  $10^{-8}$  across WebAssembly, Android, iOS, Windows, macOS, and Linux.<sup>1</sup> This numerical stability ensures that any gradient updates calculated during edge training remain mathematically consistent regardless of the underlying hardware.<sup>1</sup>



## The M-POLY-VTD Engine and Multi-Numerical Training

Loom's execution environment treats neural networks as 3D spatial grids where individual cells support dynamic precision adjustments, rather than as sequential layers.<sup>7</sup> This is powered by its M-POLY-VTD (Multi-Numerical Architecture, Polymorphic Layer-Morphing, Volumetric Tensor Dispatch) virtualization core.<sup>7</sup>

The DNVM includes a native data-type system supporting <sup>21</sup> distinct numerical types, ranging from double-precision  $fp64$  down to 1-bit binary representations.<sup>1</sup> During on-device training, the engine can dynamically morph the precision of individual layers to optimize for local hardware features—such as hardware-level ternary operations—without altering the global model structure.<sup>1</sup> This allows the training pass to run directly within the virtual machine, avoiding memory bandwidth bottlenecks on edge devices.<sup>1</sup>

## Backpropagation and Localized Target Propagation

To facilitate training on edge devices, the DNVM provides a dual-path optimization architecture<sup>1</sup>:

1. **Direct Backpropagation:** Loom implements standard gradient-descent optimization passes directly in pure Go, allowing full parameter updates across all standard layers.<sup>1</sup>
2. **Neural Target Propagation (Tweening):** As an alternative to traditional backpropagation, Loom's core engine includes a localized target propagation system implemented natively via `tween.go`.<sup>6</sup> Tweening replaces global gradient chains with localized, gap-based Hebbian learning that updates layer parameters based on local activation targets rather than global error distribution.<sup>6</sup> This reduces the peak memory footprint during training because the device does not need to store the massive global activation graphs required by traditional backpropagation.

## Technical Verification and the Velvet Ladder

The operational integrity of this weight ingestion and local retraining pipeline is validated using a multi-phase testing framework.<sup>3</sup> Inside the `planetbridging` directory, validation scripts run continuous comparisons between the source host and Loom's environment.<sup>3</sup>

### The Three-Way Validation Pipeline

The primary testing tool is the Velvet Ladder (implemented in `O5_velvet_ladder.py`), which conducts a three-way output validation<sup>3</sup>:

Native PyTorch Output ↔ Loom-Stream Output ↔ Velvet Reloaded Enti

This ladder validates the model's performance from its native execution, through the live in-memory stream, and after being reloaded from the binary `.entity` file.<sup>6</sup> For standard structures like convolutional layers (CNN1), multi-head attention (MHA), and layer normalization, the output remains stable throughout the entire pipeline.<sup>3</sup>

However, the three-way validation identifies specific edge cases where the reloading pipeline can encounter challenges.<sup>3</sup> Standard layers such as LSTM, Simple RNN, and deep mixer networks are noted as having C Application Binary Interface (C ABI) gaps.<sup>3</sup> These gaps can introduce minor parsing discrepancies during the raw reload phase, so the testing harness is

designed to bypass these layers during the Velvet reload check while still validating them during the live streaming phase.<sup>3</sup>

## Comprehensive Framework Comparison

Evaluating Loom's Planetbridging requires comparing its architectural design with other edge inference and training runtimes.<sup>4</sup> The primary modern alternatives are Microsoft's ONNX Runtime (ORT) On-Device Training<sup>8</sup>, Google's LiteRT (formerly TensorFlow Lite)<sup>9</sup>, Meta's ExecuTorch<sup>2</sup>, and Apple's Core ML.<sup>10</sup>

The table below provides a detailed comparison of these runtimes:

Architectural Dimension	Loom (with Planetbridging)	ONNX Runtime (ORT)	LiteRT (TensorFlow Lite)	ExecuTorch	Apple Core ML
<b>Ingestion Paradigm</b>	<b>Dynamic Stream:</b> Extracts weights from active python execution memory. <sup>3</sup>	<b>Static Compile:</b> Generates intermediate exchange graphs ahead-of-time. <sup>4</sup>	<b>Static Compile:</b> Converts source models into flatbuffer representations. <sup>9</sup>	<b>Static Compile:</b> Graph-level compilation into a unified executable. <sup>2</sup>	<b>Static Compile:</b> Converts model structures using custom tools. <sup>10</sup>
<b>Multi-Engine Support</b>	<b>Unified Hub:</b> Ingests PyTorch, TF, JAX, and scikit-learn. <sup>3</sup>	<b>Broad Compatibility:</b> Parses standard ONNX files exported from major libraries. <sup>4</sup>	<b>Conversion-Locked:</b> Optimized for TF; requires conversion tools for other frameworks. <sup>9</sup>	<b>Ecosystem-Locked:</b> Limited to PyTorch export paths. <sup>2</sup>	<b>Conversion-Locked:</b> Restricted to models translated via coremltools. <sup>10</sup>
<b>In-Memory Absorption</b>	<b>Supported:</b> Real-time JSON stream compiled into binary .entity files. <sup>3</sup>	<b>Unsupported:</b> Relies on offline export configurations. <sup>4</sup>	<b>Unsupported:</b> Requires pre-compiled flat files. <sup>9</sup>	<b>Unsupported:</b> Bound to compiled static .pte model graphs. <sup>2</sup>	<b>Unsupported:</b> Bound to compiled directory structures. <sup>13</sup>

<b>Edge Training Approach</b>	<b>Unified Engine:</b> Native backpropagation and target propagation in-engine. <sup>1</sup>	<b>Split Phase:</b> Generates training artifacts offline before edge execution. <sup>8</sup>	<b>Signature-Based:</b> Exposes trainable weights via explicit execution signatures. <sup>1, 2</sup>	<b>Orchestrated:</b> Integrated with NVIDIA FLARE for federated deployments. <sup>17</sup>	<b>Layer-Locked:</b> Restricts training to layers flagged as updatable. <sup>18</sup>
<b>Cross-Platform Parity</b>	<b>Bit-Exact:</b> $MAE <$ across CPU, GPU, x86, ARM, and WebAssembly. <sup>1</sup>	<b>Variable:</b> Dependent on system-level libraries and hardware backends. <sup>5</sup>	<b>Variable:</b> Subject to platform-specific floating-point variations. <sup>5</sup>	<b>Variable:</b> Subject to hardware delegate execution pathways. <sup>4</sup>	<b>Hardware-Locked:</b> Optimized for Apple Silicon; differs from other architectures. <sup>5</sup>
<b>Base Engine Footprint</b>	<b>Ultra-Light weight:</b> Single compiled binary of approximately 10 MB. <sup>1</sup>	<b>Heavyweight:</b> Typically 50–100 ME engine footprint. <sup>5</sup>	<b>Moderate:</b> Embedded runtime is compact, but full setups exceed 500 MB. <sup>5</sup>	<b>Modular:</b> Base runtime footprint of approximately 50 KB. <sup>2</sup>	<b>System-Integrated:</b> Native operating system library with zero direct app footprint. <sup>18</sup>
<b>Vendor Dependency</b>	<b>Zero Lock-In:</b> Pure Go runtime with zero CGO or cloud dependencies. <sup>1</sup>	<b>Low-to-Medium:</b> Multi-platform support, but requires complex build pipelines. <sup>5</sup>	<b>Moderate:</b> Heavily integrated with Google's mobile and web platforms. <sup>9</sup>	<b>Ecosystem-Locked:</b> Bound to PyTorch compilation and runtime tools. <sup>2</sup>	<b>Platform-Locked:</b> Restricted to Apple hardware environments. <sup>13</sup>

## Architectural Synthesis and the Future of Edge AI

Comparing these design approaches highlights the trade-offs between compile-time optimization and runtime virtualization in edge computing.

## The Trade-Offs of Native Graph Compilation

Runtimes like ExecuTorch and Apple Core ML achieve high performance on specialized hardware by using ahead-of-time (AOT) graph-level partitioning and compiling the network directly into hardware-specific binary instructions.<sup>2</sup> This optimization, however, creates platform lock-in.<sup>4</sup>

If a compiled model must run on an unsupported device or if a custom layer cannot be mapped to the target hardware's operator set, execution reverts to a slow CPU fallback or fails entirely.<sup>4</sup> Additionally, on-device training in these compiled frameworks is often constrained; Core ML, for example, requires specific layers to be designated as "updatable" before deployment, limiting the model's ability to adapt dynamically.<sup>18</sup>

## The Virtualization Strategy

Loom's DNVM bypasses compile-time constraints by using a virtualization strategy.<sup>1</sup> By treating the virtual machine itself as the target runtime and building it in pure Go with zero system-level dependencies, Loom achieves cross-platform consistency across diverse hardware configurations.<sup>1</sup>

The Planetbridging framework (v0.5.0) demonstrates this flexibility by dynamically streaming active weights from multiple source engines directly into the DNVM, avoiding the need for static file conversion chains.<sup>1</sup> Combined with built-in backpropagation and memory-efficient target propagation, this design provides a viable path for training and optimizing models natively on edge devices.<sup>1</sup>

## Future Development and Roadmap

The current release of Planetbridging (v0.5.0) marks the completion of the ingestion phase, enabling live parameter streaming from host platforms to Loom's deterministic runtime.<sup>3</sup> The project's development path is structured around several upcoming milestones:

- **Phase B (v1.0.0):** Establishing bidirectional bridging by implementing native exporters from Loom to ONNX, Safetensors, and GGUF, enabling models to flow back out to engines like llama.cpp and CoreML.<sup>3</sup>
- **Phase C (v1.x):** Introducing offline file-based ingestion directly into the comparison host, allowing direct loads of Safetensors, GGUF, and ONNX files without requiring an active Python environment.<sup>3</sup>
- **Phase D (v1.x to 2.0):** Adding support for more specialized layer types (such as transpose convolutions and custom softmax variants), expanding training capabilities, and optimizing edge execution.<sup>3</sup>

By addressing the challenges of cross-platform determinism and multi-engine weight consolidation, Loom's Planetbridging framework provides a practical solution for deploying,

running, and updating models locally at the edge.<sup>1</sup>

## Works cited

1. Loom - Layered Omni-architecture Openfluke Machine - GitHub, accessed July 1, 2026, <https://github.com/openfluke/loom>
2. ExecuTorch - On-Device AI Inference Powered by PyTorch, accessed July 1, 2026, <https://executorch.ai/>
3. README.md
4. ExecuTorch -- A Unified PyTorch Solution to Run AI Models On-Device - arXiv, accessed July 1, 2026, <https://arxiv.org/pdf/2605.08195>
5. LOOM: The Universal AI Runtime That Works Everywhere (And Why That Matters) - Medium, accessed July 1, 2026, <https://medium.com/@planetbridging/loom-the-universal-ai-runtime-that-works-everywhere-and-why-that-matters-54de5e7ec182>
6. loom/README.md at main · openfluke/loom - GitHub, accessed July 1, 2026, <https://github.com/openfluke/loom/blob/main/README.md>
7. M-POLY-VTD: Architecture Overview — Loom Docs - OpenFluke, accessed July 1, 2026, <https://openfluke.com/docs/overview>
8. On-Device Training | onnxruntime, accessed July 1, 2026, <https://onnxruntime.ai/docs/get-started/training-on-device.html>
9. LiteRT: High-Performance On-Device Machine Learning Framework | Google AI Edge, accessed July 1, 2026, <https://developers.google.com/edge/litert>
10. Updatable Models — Guide to Core ML Tools - Apple, accessed July 1, 2026, <https://apple.github.io/coremltools/docs-guides/source/updatable-model-examples.html>
11. Training - ONNX Runtime, accessed July 1, 2026, <https://onnxruntime.ai/training>
12. Training AI Models Directly on Android: A Practical Guide Using LiteRT + TensorFlow Signatures | by Lubna Mariyam | Medium, accessed July 1, 2026, <https://medium.com/@lubnamariyam62/training-ai-models-directly-on-android-a-practical-guide-using-litert-tensorflow-signatures-d959c66d202e>
13. Personalizing a Model with On-Device Updates | Apple Developer Documentation, accessed July 1, 2026, <https://developer.apple.com/documentation/coreml/personalizing-a-model-with-on-device-updates>
14. ExecuTorch Concepts - PyTorch documentation, accessed July 1, 2026, <https://docs.pytorch.org/executorch/0.4/concepts.html>
15. Overview - Python API documentation - ONNX Runtime, accessed July 1, 2026, [https://onnxruntime.ai/docs/api/python/on\\_device\\_training/overview.html](https://onnxruntime.ai/docs/api/python/on_device_training/overview.html)
16. On-device training in TensorFlow Lite, accessed July 1, 2026, <https://blog.tensorflow.org/2021/11/on-device-training-in-tensorflow-lite.html>
17. Effortless Federated Learning on Mobile with NVIDIA FLARE and Meta ExecuTorch, accessed July 1, 2026, <https://developer.nvidia.com/blog/effortless-federated-learning-on-mobile-with-nvidia-flare-and-meta-executorch/>

18. What's new in Core ML 3 - Fritz ai, accessed July 1, 2026,  
<https://fritz.ai/whats-new-in-core-ml/>
19. On-Device Training: Efficient training on the edge with ONNX Runtime | Microsoft Open Source Blog, accessed July 1, 2026,  
<https://opensource.microsoft.com/blog/2023/05/31/on-device-training-efficient-training-on-the-edge-with-onnx-runtime/>