

Serverless Everywhere: A Comparative Analysis of WebAssembly Workflows Across Browser, Edge, and Cloud

Mario Colosi
MIFT Department, University of
Messina
Messina, Italy

Reza Farahani
Institute of Information Technology,
University of Klagenfurt
Klagenfurt, Austria

Lauri Lovén
Center for Ubiquitous Computing,
University of Oulu
Oulu, Finland

Radu Prodan
Department of Computer Science,
University of Innsbruck
Innsbruck, Austria

Massimo Villari
MIFT Department, University of
Messina
Messina, Italy

Abstract

WebAssembly (Wasm) is a binary instruction format that enables portable, sandboxed, and near-native execution across heterogeneous platforms, making it well-suited for serverless workflow execution on browsers, edge nodes, and cloud servers. However, its performance and stability depend heavily on factors such as startup overhead, runtime execution model (e.g., Ahead-of-Time (AOT) and Just-in-Time (JIT) compilation), and resource variability across deployment contexts. This paper evaluates a Wasm-based serverless workflow executed consistently from the browser to edge and cloud instances. The setup uses `wasm32-wasi` modules: in the browser, execution occurs within a web worker, while on Edge and Cloud, an HTTP shim streams frames to the Wasm runtime. We measure cold- and warm-start latency, per-step delays, workflow makespan, throughput, and CPU/memory utilization to capture the end-to-end behavior across environments. Results show that AOT compilation and instance warming substantially reduce startup latency. For workflows with small payloads, the browser achieves competitive performance owing to fully in-memory data exchanges. In contrast, as payloads grow, the workflow transitions into a compute- and memory-intensive phase where AOT execution on edge and cloud nodes distinctly surpasses browser performance.

CCS Concepts

• **Computing methodologies** → **Distributed algorithms.**

Keywords

Serverless Computing, Workflow, WebAssembly, Edge Computing, Browser-Edge-Cloud Continuum.

1 Introduction

In the serverless paradigm, providers manage the scalability, placement, and lifecycle of short-lived, event-driven functions. While this abstraction significantly reduces the operational overhead for developers, it exposes drawbacks, including cold starts, queue latency, performance variability, and cost unpredictability [6]. Optimizing these aspects requires consistent measurement methodologies across the computing continuum, i.e., browsers, edge instances, and cloud servers [7]. Recent works have analyzed the causes of cold

starts, offering a deeper understanding of their cost implications and underscoring the need for rigorous comparative evaluations [9].

WebAssembly (Wasm), with its compact bytecode, fast validation and instantiation, and strong sandboxed isolation, serves as a unifying execution substrate across the computing continuum [5, 20]. Recent advancements, such as the maturation of runtimes and the standardization of the *WebAssembly System Interface* (WASI) ¹, have brought the “*compile once, run anywhere*” paradigm closer to reality. In practice, the same Wasm artifact (i.e., a compiled and portable binary module) can be executed across diverse contexts without modification, offering two key benefits: (i) methodological comparability and (ii) a simplified software supply chain with uniform packaging and minimal system dependencies.

Recent studies confirm the competitiveness of Wasm compared to solutions like containers and Function-as-a-Service (FaaS) platforms, particularly under resource-constrained conditions [3]. However, trade-offs remain in terms of startup latency and memory footprint, indicating that Wasm’s advantages are not uniform across all deployment contexts. Comparisons between x86 and ARM architectures further demonstrate the growing maturity of this technology while revealing runtime variations with practical implications for system design and development choices [11]. Nevertheless, existing evaluations are restricted to a single environment or micro-benchmark, offering limited insight into the comparison of end-to-end latency (i.e., function cold/warm behavior, function execution and communication time, and workflow makespan) for realistic application workflows.

To our knowledge, no unified methodological work has yet executed the same Wasm-based workflow across browsers, edge nodes, and cloud platforms using consistent and homogeneous metrics. Thus, this paper presents a comparative analysis of serverless workflows implemented in Wasm and executed across three distinct environments: the browser, the edge, and the cloud. The entire evaluation is performed under a unified experimental setup, ensuring methodological consistency and enabling a fair, cross-environment comparison of performance and resource efficiency. The results of this analysis provide practical guidance for making informed orchestration decisions (e.g., function placement) in serverless workflows. The key contributions of this paper are:

 This work is licensed under a Creative Commons Attribution International 4.0 License.

¹<https://wasi.dev>

- (i) we implement a serverless workflow in Wasm, exposing WASI interfaces compatible with three environments, browser, edge, and cloud;
- (ii) we isolate the effects of the execution environment and quantify workflow performance on cold and warm starts, function makespan and workflow makespan, throughput, and resource utilization (CPU, memory);
- (iii) we compare the performance of ahead-of-time (AOT) and just-in-time (JIT) as compilation strategies, along with a pre-warming method, to evaluate their impact within each environment.

This paper has six sections: Section 2 summarizes the background and related work on Wasm, serverless computing, workflow processing, and cold start latency reduction. Section 3 describes the strategy adopted for browsers, edge instances and cloud servers. Section 4 explains our evaluation setup before describing the experimental results in Section 5. Section 6 finally concludes the paper.

2 Related Work

This section reviews related work in three main areas: Wasm as an execution substrate for serverless platforms, serverless workflow orchestration strategies, and cold-start mitigation and cost modeling in serverless systems.

2.1 Wasm as a serverless execution substrate

Recent measurements show that Wasm can serve as a competitive alternative to, or an integration layer with, containers [17]. However, its maturity and compatibility still vary across runtime implementations and their integration with the underlying execution environment [12]. Systematic comparison of Wasm runtimes ranks standalone and browser-integrated executors across execution models, interpreter, just-in-time (JIT), and ahead-of-time (AOT) compilation, as well as system interfaces and security aspects, highlighting significant differences in initialization latency and throughput [21]. These findings guide the selection of runtimes and configuration options, such as AOT and JIT compilation modes, as well as caching strategies evaluated in our experiments.

On the edge side, benchmarks show that Wasm-based serverless platforms can achieve reduced startup times and competitive efficiency profiles compared to widely used container-based solutions [14]. However, methodological differences in workloads, metrics, and network configurations across studies highlight the need for a homogeneous measurement protocol [3]. In parallel, recent analyses show the maturation of client-side, in-browser execution [2, 15]. While Wasm offers advantages in portability and time-to-use for scientific computing, it remains constrained by limited tooling and restricted system access [16]. Furthermore, in-browser large language model (LLM) inference systems demonstrate that compute-intensive operations can be executed directly on the client using WebGPU for acceleration and Wasm for CPU-bound computation, achieving near-native performance on the same hardware [1, 19].

2.2 Serverless workflow orchestration strategies

On the other hand, existing research on serverless workflow orchestration investigates the effects of placement strategies on workflow

makespan and resource consumption [7, 8, 18]. Edge resource allocation mechanisms using serverless platforms typically employ formalized policies derived from observable runtime metrics [13], making placement decisions repeatable and reproducible by third parties under identical inputs and configurations. The comparative work of centralized and distributed orchestration approaches has demonstrated that the underlying architecture significantly influences workflow completion times and performance variability [4]. A recurring observation across these works is the necessity of empirical characterization of behaviors, such as cold and warm start patterns, queuing latency, and resource footprint—supported by homogeneous benchmarking. Therefore, without such standardized measurements, transferring placement strategies across heterogeneous environments becomes unreliable. These insights motivate our focus on comparable measurements of the same workflow executed in the browser, at the edge, and in the cloud.

2.3 Cold-start and cost management

One of the key factors in serverless computing is the *cold start*, referring to the additional latency that occurs when a function is invoked without an already warm instance available [6]. This forces the platform to provision a new sandbox environment (e.g., microVM or container), load and initialize the runtime along with its dependencies, and bring the function code to a ready-to-execute state. For Wasm, this process includes module validation, compilation, and instantiation. However, recent work [9] shows that although cold start can sometimes dominate end-to-end latency in certain production workloads, many functions are not latency-sensitive. Consequently, the authors argue that serverless systems should be evaluated and designed with a broader perspective, considering not only startup latency but also predictability, throughput, and resource efficiency. With respect to cost, recent work [10] investigates how pricing models, billing granularity, and data egress fees affect the total cost of serverless applications. These works introduce parameterized evaluation methodologies that facilitate normalized comparisons across providers, yet require workload and runtime measurements for proper calibration.

2.4 Contributions beyond the state-of-the-art

Prior works evaluated Wasm runtimes, serverless orchestration, and cold-start behaviors largely in isolation, focusing on a single environment (browser, edge, or cloud) or on micro-benchmarks. To our knowledge, no unified, cross-environment work has executed the same Wasm-based workflow across these three layers under a common experimental setup and set of metrics. This paper instead provides an end-to-end comparison of Wasm-based serverless workflows executed in the browser, at the edge, and in the cloud, offering new insights into how runtime configurations (AOT, JIT) and environmental factors jointly affect performance and resource efficiency across the computing continuum.

3 Experimental Design

This section explains how we evaluate the impact of the execution environment (browser, edge, and cloud instances), on the performance of Wasm-based serverless workflows under a consistent measurement setup.

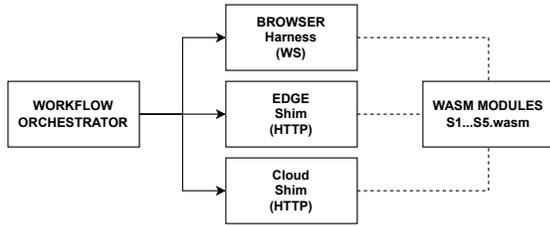


Figure 1: High-level execution pipeline of the serverless workflow across browser, edge, and cloud environments.

Workflow artifacts: We model the application as a directed acyclic graph (DAG), where each node represents a Wasm module with explicit input/output interfaces. The same artifacts are executed across all environments, allowing us to attribute observed performance differences solely to environmental factors rather than to code or packaging variations.

Workflow orchestrator: We developed an external component that serves as a workflow orchestrator, interpreting the DAG and triggering each task according to its structure (e.g., sequence, fan-out, and fan-in). Each task’s output is passed as input to its successors, ensuring that tasks remain isolated and do not invoke one another directly. This design preserves identical flow control across all three environments. The orchestrator also records timing and resource utilization metrics, enabling consistent and comparable measurements across experiments. The orchestrator runs in a control environment isolated from the workflow Wasm executors to prevent interference. In the browser setup, it is co-located on the same machine but runs in a separate process, communicating via a loopback WebSocket with a harness page that forwards messages to the Web Worker. The Web Worker executes the workflow tasks in a separate thread. At the edge, the orchestrator runs on a different machine within the same LAN as the executor node and invokes tasks through a lightweight HTTP shim co-located with the Wasm runtime. In the cloud configuration, it runs on a virtual machine (VM) in the same region, using the same HTTP shim interface to trigger executions. Fig 1 illustrates the resulting high-level software architecture.

Metrics and data collection We organize the collected metrics into two complementary perspectives for evaluating application performance. The first captures end-to-end workflow behavior and results, while the second focuses on resource utilization and runtime-level effects within the platform:

(i) *User-level metrics:*

- *Cold start:* invocation without an existing warm instance, requiring sandbox provisioning (e.g., container), artifact loading, runtime initialization, and dependency setup;
- *Warm start:* served by an already initialized instance without reloading or reinitialization;
- *Function Latency:* Start/end timestamps per step;
- *Workflow makespan:* total elapsed time from acquiring the first step’s input to producing the final output, including internal queuing and inter-step data transfers.

(ii) *System-level metrics:*

- *Throughput:* number of completed workflow invocations per unit time under steady-state conditions with fixed concurrency;
- *Overhead and resources:* CPU (average and peak utilization), memory (RSS and peak), I/O activity, and artifact size (Wasm binary and AOT cache), with the sampling period and measurement tools specified.

Analysis and reproducibility: We adopt a robust statistical analysis, reporting median, 25th/75th percentiles (p25/p75), interquartile range (IQR), and, where applicable, 95 % confidence intervals for key performance indicators. To ensure fair comparisons, we explicitly specify the varying factors between measurements (e.g., execution environment or ablation parameters) while keeping all other parameters constant. In addition, reproducibility is guaranteed through strict artifact versioning, fixed runtime configurations, and deterministic random seeds.

4 Evaluation Setup

This section explains how we executed experiments across browsers, edge, and cloud instances under identical conditions, i.e., Wasm artifact, invocation Application Binary Interface (ABI), and workflow orchestrator. To accommodate browser constraints, the workflow is designed to operate without file-system access. The ABI remains identical across all environments; only the transport mechanism differs: WebSocket with `postMessage` in the browser, and HTTP `multipart/form-data` at the edge and in the cloud instances.

Each workflow task is implemented in Rust and compiled into Wasm modules using the same toolchain (Rust \rightarrow wasm32-wasi), allowing AOT and JIT compilation to be toggled as ablation parameters. All builds and corresponding module images are version-controlled via commit hashes or container digests for reproducibility. Compilation flags and enabled Wasm features are kept identical across the environments for consistent execution behavior.

4.1 Workflow

We evaluate a serverless workflow, modeled as a DAG with ($N = 5$) functions and fan-out/fan-in $\langle f_{out}, f_{in} \rangle = \langle 4, 4 \rangle$, as shown in Fig. 2. As mentioned, each step is implemented in Rust and compiled to wasm32-wasi modules that exchange data frames encoded in the Concise Binary Object Representation (CBOR) format entirely in memory. The workflow processes deterministic synthetic payloads (generated with a fixed seed) in three input sizes: $s_{small} = 16$ KiB, $s_{medium} = 1$ MiB, $s_{large} = 4$ MiB. The workflow functions are:

- (S1) *Ingest&Deserialize (serialization-bound):* decodes the CBOR frame in memory buffer, validates the schema, and calculates a CRC32 of the input for a quick integrity check. Complexity cost: $O(n)$.
- (S2) *Preprocess (memory-bound):* applies normalization and a scan over the buffer (sliding window, clamp, type conversion $u8 \rightarrow f32$); the output is a buffer of the same size. Complexity cost: $O(n)$;
- (S3) *Map (CPU-bound, fan-out=4):* divides the buffer into four blocks and processes them in parallel with a blocked *dense matrix multiplication*. Complexity cost: $O(d^3)$, with $d = \sqrt{s/4}$;

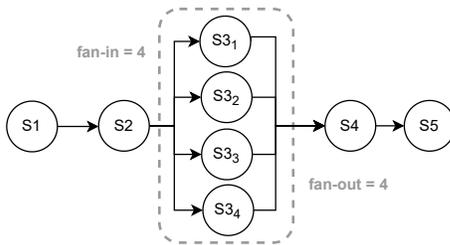


Figure 2: Workflow DAG.

- (S4) *Reduce (serialization-bound, fan-in=4)*: aggregates the four results (sum/concatenation according to the scheme) into a single block and calculates a partial digest (*BLAKE3*) for verification. Complexity cost: $O(n)$;
- (S5) *Serialize&Finalize (serialization-bound)*: recodes the final output in CBOR and produces the final digest (*BLAKE3*) to be compared with the expected value. Complexity cost: $O(n)$.

CPU kernel sizing (S3). The size of the matrix multiplication is derived from the payload:

$$d(s) = \begin{cases} 64 & \text{if } s = s_{\text{small}} (\approx 64^2 \cdot 4 \text{ B} \approx 16 \text{ KiB}), \\ 512 & \text{if } s = s_{\text{medium}} (\approx 512^2 \cdot 4 \text{ B} \approx 1 \text{ MiB}), \\ 1024 & \text{if } s = s_{\text{large}} (\approx 1024^2 \cdot 4 \text{ B} \approx 4 \text{ MiB}). \end{cases}$$

Each $S3[k]$ processes a separate block in parallel (fan-out = 4); S4 combines the results.

4.2 Platforms, Runtimes, and Tooling

To ensure consistency across browsers, edges, and clouds, we use the same Wasm artifact (target wasm32-wasi, identical toolchain and flags) and fix runtime versions and configurations. The invocation details are described in Section 4.3; here we summarize the platforms, runtimes, and tools used. Table 1 lists, for each environment, hardware/OS, runtime/browser versions, Wasm flags, and resource limits.

Browser host. The machine running the browser is an Ubuntu 24.04 LTS workstation with an Intel® Core i7-8700K (6 cores, 12 threads, 3.70 GHz) and 16 GiB RAM. We use Mozilla Firefox 143.0.4 (64-bit, Snap for Ubuntu), with WebAssembly SIMD enabled. The workflow is executed in a Web Worker, while the orchestrator runs in a separate process on the same machine and communicates via loopback WebSocket with a harness page.

Edge Node. The edge runs the same module inside a minimal OCI/Docker container with WasmEdge v0.13.5. The host is Ubuntu 24.04 LTS (linux/arm64) with Docker Engine 28.5.1 (cgroups v2); resources are capped via container limits (`--cpus 4, --memory 4`). Wasm features (SIMD) are aligned with the browser; warm state is achieved through a pool of pre-initialized instances, while cold state is forced by resetting the pool. No disk I/O in the workflow; messages pass through memory to the executor.

Cloud VM. We replicate the edge setup on a linux/amd64 VM co-located with the orchestrator (same region): 4 vCPU, 24 GiB RAM, Ubuntu 24.04; container runtime Docker Engine 28.5.1 (containerd v1.7.28, runc 1.3.0). The Wasm executor is WasmEdge v0.13.5, with the same toolchain configuration as on the edge. Resource limits and cold/warm policies are identical to edge to make the results comparable. Again, the workflow does not access the file system; orchestration uses in-memory messages.

4.3 Execution Environments

We adopt a unified invocation ABI for all environments. Each task receives a request frame using CBOR serialization for metadata fields. On edge/cloud, large binary payloads are transmitted separately via HTTP multipart/form-data to avoid a base64 or similar encoding overhead within the CBOR envelope. This hybrid approach maintains protocol type safety while supporting payloads up to several megabytes without buffer overflow issues. The minimum structure of exchanged messages is as follows:

- request: {op: "invoke", step_id, run_id, payload}
- response: {op: "result", status, payload[, error]}

The orchestrator sends/waits for these frames and correlates the measurements via `run_id`. Since we do not use the file system in the workflow, input/output travels as buffers in memory.

Browser. The orchestrator communicates with a harness page via WebSocket. The page forwards the frames to the Web Worker with transferable `ArrayBuffer` `postMessage` and sends the Web Worker's response back to the orchestrator. Cold start: new Web Worker instance and cache/Service Worker invalidation; warm start: Web Worker and module already compiled remain alive. Timestamps are collected with `performance.now()`.

Edge. The orchestrator invokes a minimal HTTP endpoint on the edge node: `POST /invoke` (CBOR binary body). A shim in the container delivers the frame to the Wasm executor (WasmEdge) via pipe/IPC and returns the response frame. Cold start: one-shot mode; warm start: pool of pre-initialized instances.

Cloud. Same edge pipeline on VM/container in the same region as the orchestrator: `POST /invoke` (CBOR binary body, HTTP/1.1) and same cold/warm behavior as above. Timestamps are recorded with a monotonic clock (`Rust Instant`) and the shim returns a CBOR response including the phase breakdown (load, compile, instantiate, init, execute), total latency, and best-effort CPU/RSS samples collected from `/proc`.

4.4 Measurement Protocol

In the experiments, we use a mixed set of micro and macro workloads. The micro workloads include: (i) CPU-bound (e.g., `matmul`) with synthetic payloads of size $\langle s_{\text{small}}, s_{\text{medium}}, s_{\text{large}} \rangle$; (ii) memory-bound (allocation/scan); (iii) serialization-bound (CBOR \leftrightarrow in-memory). The macro workflow is a DAG with $\langle N \rangle$ steps, fan-out/fan-in $\langle f_{\text{out}}, f_{\text{in}} \rangle$, and explicit I/O scheme; for each step, we specify the format and expected size of input/output.

The measurements of latency, cold/warm behavior, and throughput are performed in a closed-loop configuration. Each configuration (environment \times payload \times mode ablation) is repeated k

Table 1: Platforms and runtimes.

Env.	Host / Instance	CPU	RAM	OS	Browser/Container	Wasm Runtime
Browser	Workstation (linux/amd64)	6 cores	16 GiB	Ubuntu 24.04	Firefox 143.0.4	Firefox WebAssembly engine (SpiderMonkey)
Edge	Raspberry Pi 4 (linux/arm64)	4 cores	4 GiB	Ubuntu 24.04	Docker 28.5.1 (cgroups v2)	WasmEdge v0.13.5
Cloud	VM (linux/amd64)	4 vCPU	24 GiB	Ubuntu 24.04	Docker 28.5.1	WasmEdge v0.13.5

Wasm features: WASI preview1 (edge/cloud; JS shim, subset in browser); SIMD on; Threads on.
Resource limits: Edge/Cloud via cgroups (--cpus, --memory); Browser via Web Worker isolation.

Table 2: Artifact sizes comparison between WASM modules and their AOT compiled outputs.

Step	Size WASM Module (MB)	Size AOT Artifact (MB)	Size Increase
S1	0.152	0.342	+ 125%
S2	0.135	0.317	+ 135%
S3	0.141	0.322	+ 128%
S4	0.151	0.354	+ 134%
S5	0.153	0.360	+ 135%

times ($k = 20$ with randomized order, separate warm-up, and fixed seed. Timestamps are collected with monotonic clocks (browser: `performance.now()`; edge/cloud: high-resolution clock) and correlated via `trace-id/span-id`; we record: cold-start breakdown (load \rightarrow compile/AOT-load \rightarrow instantiate \rightarrow init \rightarrow run), per-step and end-to-end latencies, workflow makespan, throughput, CPU (average/peak), and memory (RSS/peak).

Outlier, retry, and exclusions follow predefined rules: we discard a run if the monitor detects abnormal conditions; for valid samples, we apply an outlier filter based on $1.5 \times \text{IQR}$.

5 Experimental Results

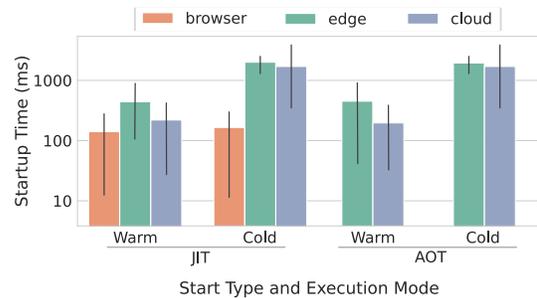
In this section, we present the experimental results obtained by running the same WebAssembly workflow separately in the browser, edge, and cloud. The analyses are organized along two complementary perspectives: *User-Perceived Performance* and *System-Level Metrics*. The results are aggregated over k repetitions.

Before discussing performance, we quantify the size of the artifacts. Table 2 compares, for each workflow module, the size of the Wasm bytecode and that of the AOT (precompiled) output, also reporting the AOT/Wasm percentage increase; this provides the context for interpreting the benefits observed in cold tests with AOT and for estimating the impact of deployment in the three environments.

5.1 User-Perceived Performance

We analyze the application’s performance that directly impacts the user experience, using three predefined metrics: (i) cold/warm startup times, (ii) latencies per workflow step, and (iii) end-to-end makespan.

Cold and warm startup. As shown in Figure 3, warm starts are faster than cold starts in all environments, although this gap is less pronounced in browsers. The use of AOT, which is only applicable


Figure 3: Cold vs Warm startup times.

on edge/cloud, significantly reduces the amount of time spent on compilation/validation in cold tests. The browser shows particularly low warm start times, benefiting from compilation/instance already residing in the Worker; edge shows the highest cost in cold JIT, while cloud ranks in the middle with slightly lower variability.

Function latency. We run the workflow with a warm profile and 1MiB as payload size and collect the start/end timestamps for each step (S1–S5). In Figure 4 we represent the distributions with box-plots on a logarithmic scale: box = IQR (p25–p75), center line = median; whisker = $1.5 \times \text{IQR}$; outliers shown explicitly. The browser performs better on S3–S4–S5 (lower medians and narrower IQRs), thanks to in-memory execution and lower orchestration/IPC overhead. On S1–S2, it ranks between edge and cloud: the overhead of (de)serialization and the Worker-main thread transition (postMessage/buffer transfer) make it less advantageous for the browser, while edge/cloud runtimes with optimized containers and I/O pipes mitigate the cost of these steps. In summary, when computation or aggregation prevails (S3–S4–S5), the browser is competitive; when (de)serialization dominates (S1–S2), the advantage diminishes.

Workflow makespan. We measure the time S1 \rightarrow S5 (including queues and transfers) in *warm* mode by varying the *payload size* and the *JIT/AOT* mode. In Figure 5, we observe a clear reversal as the payload increases: with small sizes, the browser (JIT) is often faster because the startup cost is minimal, everything happens in-memory, and Worker orchestration has little impact. Moving to medium/large sizes, the makespan becomes compute/memory-bound: the startup cost is amortized and memory bandwidth, cache, and code quality prevail; here, AOT configurations on edge/cloud are faster and more stable. In particular, edge AOT tends to emerge

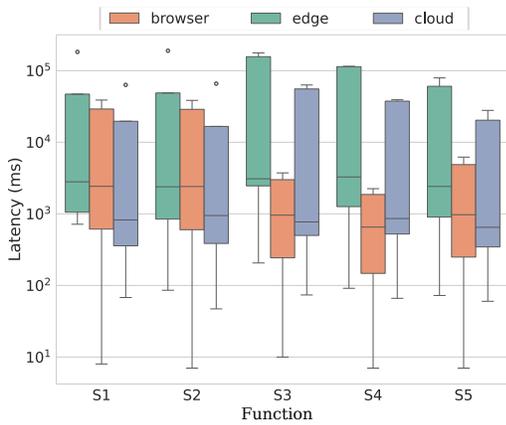


Figure 4: Function latency.

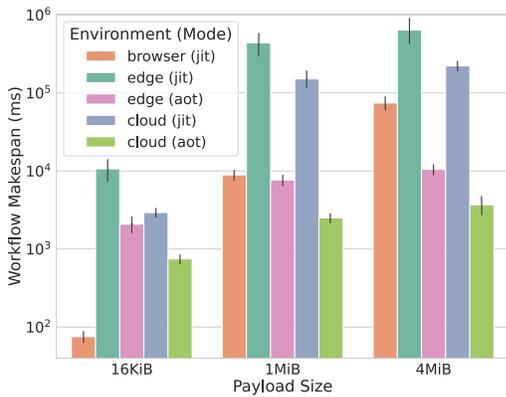


Figure 5: Workflow makespan.

when the load is markedly numerical, while the browser’s advantage diminishes until it reverses. In all cases, warm and AOT improve predictability and end-to-end times, with more pronounced effects as size increases.

5.2 System-Level Metrics

This section observes the platform’s behavior under load, complementing user-facing metrics. We consider sustainable capacity (throughput) and resource usage profile (CPU/RSS), measured with the same experimental setup.

Throughput. A new invocation is issued upon completion of the previous one, so that the reported value corresponds to the inverse of the average completion time. Figure 6 shows results consistent with those for makespan. With 16 KiB, the browser (JIT) prevails: lightweight invocations, in-memory exchanges, and reduced startup costs allow for tens of requests per second. At 1 MiB and 4 MiB, the compute/memory-bound effect emerges: the initial overhead is amortized, and AOT configurations on edge/cloud outperform the browser due to the elimination of JIT compilation and more favorable CPU/memory-bandwidth/cache.

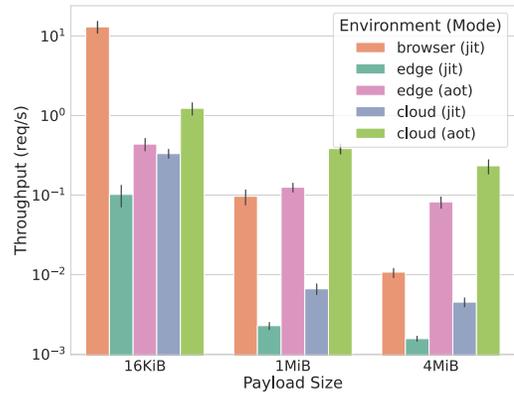


Figure 6: Workflow throughput.

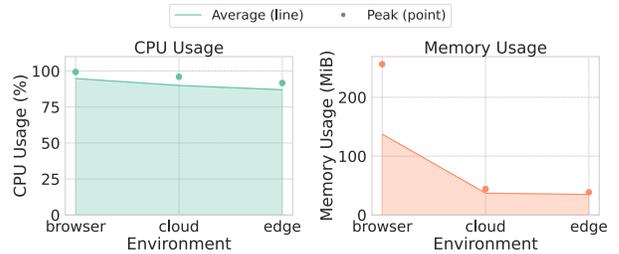


Figure 7: Resource usage.

Resource usage. We sample the CPU (average/peak) and memory (average/peak RSS) of the executor process, with a sampling interval of 20 ms on edge/cloud (readings from /proc) and 20 ms in the browser. In the browser, CPU is obtained by recording a Performance session from in-browser DevTools and aggregating the profile samples; the same tools on the Web Worker side estimate memory. As shown in Figure 7, CPU usage is high in all environments (indicating computational bottlenecks), with average values slightly lower on the edge than in the cloud and browser. Memory usage differs more significantly: the browser has higher average/peak RSS (JS engine stack, Worker, GC management), while the cloud and edge remain more compact thanks to the native Wasm runtime and the absence of browser components.

6 Conclusion

We presented a comparative analysis of the execution of the same serverless workflow in browsers, edge, and cloud, keeping Wasm and ABI artifacts unchanged. The results indicate that warm and AOT significantly reduce startup times and variability. The browser is competitive on small loads and on steps dominated by orchestration/aggregation, while edge/cloud in AOT prevail with increasing payloads when the regime becomes compute/memory-bound. Makespan and throughput consistently reflect these trends, and CPU/RSS profiles explain the differences. This evidence should be interpreted within the scope of the setup (runtimes, hardware, browser, workload) and does not claim to be generalizable beyond the experiment.

Acknowledgment

This work received partial funding from the Italian Ministry of University and Research (MUR) “Research projects of National Interest (PRIN-PNRR)” call through the project “Cloud Continuum aimed at On-Demand Services in Smart Sustainable Environments” (CUP: J53D23015080001- IC: P2022YNBHP), and “SEcurity and RIghts in the Cyberspace (SERICS)” project (PE00000014), under the MUR National Recovery and Resilience Plan funded by the European Union - NextGenerationEU (CUP: D43C22003050001), the European Union’s Horizon Europe research and innovation program, grant agreements 101093202 (Graph-Massivizer), and the Austrian Research Promotion Agency (FFG, grant agreement 909989, AIM AT Stiftungsprofessur für Edge AI)

References

- [1] Zhiyang Chen, Yun Ma, Haiyang Shen, and Mugeng Liu. 2025. WeInfer: Unleashing the Power of WebGPU on LLM Inference in Web Browsers. In *Proceedings of the ACM on Web Conference 2025* (Sydney NSW, Australia) (WWW ’25). Association for Computing Machinery, New York, NY, USA, 4264–4273. doi:10.1145/3696410.3714553
- [2] Joao De Macedo, Rui Abreu, Rui Pereira, and João Saraiva. 2021. On the Runtime and Energy Performance of WebAssembly: Is WebAssembly Superior to Javascript Yet?. In *2021 36th IEEE/ACM International Conference on Automated Software Engineering Workshops (ASEW)*. IEEE, 255–262.
- [3] Giuseppe De Palma, Saverio Giallorenzo, Jacopo Mauro, Matteo Trentin, and Gianluigi Zavattaro. 2024. WebAssembly at the Edge: Benchmarking a Serverless Platform for Private Edge Cloud Systems. *IEEE Internet Computing* 28, 6 (2024), 37–44. doi:10.1109/MIC.2024.3513035
- [4] Abdallah Elshamy, Ahmed Alquraan, and Samer Al-Kiswani. 2023. A Study of Orchestration Approaches for Scientific Workflows in Serverless Computing. In *Proceedings of the 1st Workshop on Serverless Systems, Applications and Methodologies* (Rome, Italy) (SESAME ’23). Association for Computing Machinery, New York, NY, USA, 34–40. doi:10.1145/3592533.3592809
- [5] Reza Farahani, Dragi Kimovski, Sashko Ristov, Alexandru Iosup, and Radu Prodan. 2023. Towards Sustainable Serverless Processing of Massive Graphs on the Computing Continuum. In *Companion of the 2023 ACM/SPEC International Conference on Performance Engineering*. 221–226.
- [6] Reza Farahani, Frank Loh, Dumitru Roman, and Radu Prodan. 2024. Serverless Workflow Management on the Computing Continuum: A Mini-Survey. In *Companion of the 15th ACM/SPEC International Conference on Performance Engineering*. 146–150.
- [7] Reza Farahani, Narges Mehran, Sashko Ristov, and Radu Prodan. 2024. Heftless: A Bi-objective Serverless Workflow Batch Orchestration on the Computing Continuum. In *2024 IEEE International Conference on Cluster Computing (CLUSTER)*. IEEE, 286–296.
- [8] Reza Farahani and Radu Prodan. 2025. EnergyLess: An Energy-Aware Serverless Workflow Batch Orchestration on the Computing Continuum. In *2025 IEEE 18th International Conference on Cloud Computing (CLOUD)*. IEEE, 243–254.
- [9] Muhammed Golec, Guneet Kaur Walia, Mohit Kumar, Felix Cuadrado, Sukhpal Singh Gill, and Steve Uhlig. 2024. Cold Start Latency in Serverless Computing: A Systematic Review, Taxonomy, and Future Directions. *Comput. Surveys* 57, 3 (2024), 1–36.
- [10] Muhammad Hamza, Muhammad Azeem Akbar, and Rafael Capilla. 2024. Understanding Cost Dynamics of Serverless Computing: An Empirical Study. In *Software Business*. Sami Hyrynsalmi, Jürgen Münch, Kari Smolander, and Jorge Melegati (Eds.). Springer Nature Switzerland, Cham, 456–470.
- [11] Sangeeta Kakati and Mats Brorsson. 2024. A Cross-Architecture Evaluation of WebAssembly in the Cloud-Edge Continuum. In *2024 IEEE 24th International Symposium on Cluster, Cloud and Internet Computing (CCGrid)*. 337–346. doi:10.1109/CCGrid59990.2024.00046
- [12] Mugeng Liu, Haiyang Shen, Yixuan Zhang, Hong Mei, and Yun Ma. 2025. WebAssembly for Container Runtime: Are We There Yet? *ACM Trans. Softw. Eng. Methodol.* 34, 6, Article 174 (July 2025), 22 pages. doi:10.1145/3712197
- [13] Samaneh Hajy Mahdizadeh and Saeid Abrishami. 2024. An Assignment Mechanism for Workflow Scheduling in Function as a Service Edge Environment. *Future Generation Computer Systems* 157 (2024), 543–557. doi:10.1016/j.future.2024.04.003
- [14] Pankaj Mendki. 2020. Evaluating WebAssembly Enabled Serverless Approach for Edge Computing. In *2020 IEEE Cloud Summit*. IEEE, 161–166.
- [15] Biju R Mohan et al. 2022. Comparative Analysis of JavaScript And WebAssembly in the Browser Environment. In *2022 IEEE 10th Region 10 Humanitarian Technology Conference (R10-HTC)*. IEEE, 232–237.
- [16] Jeffrey M Perkel. 2024. No Installation Required: How WebAssembly Is Changing Scientific Computing. *Nature* 627, 8003 (2024), 455–456.
- [17] Steven Pham, Kaue Oliveira, and Chung-Hornng Lung. 2023. WebAssembly Modules as Alternative to Docker Containers in IoT Application Development. In *2023 IEEE 3rd International Conference on Electronic Communications, Internet of Things and Big Data (ICEIB)*. IEEE, 519–524.
- [18] Sashko Ristov, Reza Farahani, and Radu Prodan. 2023. Large-scale Graph Processing and Simulation with Serverless Workflows in Federated FaaS. In *Companion of the 2023 ACM/SPEC International Conference on Performance Engineering*. 227–231.
- [19] Charlie F Ruan, Yucheng Qin, Xun Zhou, Ruihang Lai, Hongyi Jin, Yixin Dong, Bohan Hou, Meng-Shiun Yu, Yiyang Zhai, Sudeep Agarwal, et al. 2024. WebLLM: A High-Performance In-Browser LLM Inference Engine. *arXiv preprint arXiv:2412.15803* (2024).
- [20] Yixuan Zhang, Mugeng Liu, Haoyu Wang, Yun Ma, Gang Huang, and Xuanzhe Liu. 2025. Research on WebAssembly Runtimes: A Survey. *ACM Trans. Softw. Eng. Methodol.* (Jan. 2025). doi:10.1145/3714465
- [21] Mircea Țălu. 2025. A Comparative Study of WebAssembly Runtimes: Performance Metrics, Integration Challenges, Application Domains, and Security Features. *Archives of Advanced Engineering Science* (Apr. 2025), 1–13. doi:10.47852/bonvi ewAAES52024965