

# Weaving the Cosmos: WASM-Powered Interchain Communication for AI Enabled Smart Contracts

Rabimba Karanjai<sup>1</sup>[0000–0002–6705–6506], Lei Xu<sup>2</sup>, and Weidong Shi<sup>1</sup>  
 rkaranjai@uh.edu<sup>1</sup>, lxu12@kent.edu<sup>2</sup>, wshi3@uh.edu<sup>1</sup>

<sup>1</sup> University Of Houston, Houston, USA

<sup>2</sup> Kent State University, USA

**Abstract.** In this era, significant transformations in industries and tool utilization are driven by AI/Large Language Models (LLMs) and advancements in Machine Learning. There’s a growing emphasis on Machine Learning Operations(MLOps) for managing and deploying these AI models. Concurrently, the imperative for richer smart contracts and on-chain computation is escalating. Our paper introduces an innovative framework that integrates blockchain technology, particularly the Cosmos SDK, to facilitate on-chain AI inferences. This system, built on WebAssembly (WASM), enables interchain communication and deployment of WASM modules executing AI inferences across multiple blockchain nodes. We critically assess the framework from feasibility, scalability, and model security, with a special focus on its portability and engine-model agnostic deployment. The capability to support AI on-chain may enhance and expand the scope of smart contracts, and as a result enable new use cases and applications.

**Keywords:** Cosmos · smart contract · AI · LLM · WebGPU · WASM

## 1 Introduction

In today’s technological era, significant strides have been made in the field of artificial intelligence and large language models (LLMs), such as the advancements seen with GPT models [9] and Google’s PaLM2 [7]. These models are now progressively utilized in real-world applications [17,9]. An exciting application of AI and LLMs is their use in generating source code, including smart contracts, directly from natural language descriptions. This application, exemplified by tools like ChatGPT and Github CoPilot [14], translates programmers’ instructions into code across various languages, potentially revolutionizing the programming landscape. Furthermore, LLMs are being applied in forecasting and modeling time series data, outperforming traditional methods by enabling multi-modal predictions that incorporate both time series and other forms of unstructured data. This innovation could greatly enhance financial modeling techniques.

Amidst this backdrop, the increasing use of decentralized ledgers and smart contracts, particularly in the financial sector, is noteworthy. For example, Uniswap’s smart contracts managed transactions worth approximately \$7.17 billion daily in 2021<sup>3</sup>. Given the growing importance of smart contracts in various contexts, such as Confidential Computing [30,18], Decentralized Serverless Functions [19], and Event-based Transactions [21,20], it becomes crucial to investigate whether AI and LLMs can be used effectively for finance modeling and predictions in a distributed and safe way while augmenting a smart contract workflow.

In this research, we explore the feasibility of executing an AI agent or machine learning model through a smart contract on a blockchain. Our focus is on developing a system that integrates with the current smart contract workflows seamlessly, without necessitating disruptive changes to the existing infrastructure. We aim to adhere to established standards while utilizing AI-driven financial models to produce numerical, textual, and multi-modal outputs. The core of our investigation is to determine whether our proposed system can function within well-established frameworks and offer an innovative approach to interacting with AI models, including LLMs, directly on-chain.

In our study, we introduce a framework designed to enable decentralized on-device AI model agnostic inferences for smart contracts. These smart contracts act as triggers for invoking the AI inference engine via a WebAssembly runtime [16]. Our research primarily seeks to address the following questions:

**RQ1: Can smart contracts utilize AI and LLM inferences with reasonable performance?** We try to see if using existing smart contract platforms, we can build a workflow that is AI model and LLM agnostic.

**RQ2: Can we get inference locally in a secure way?** We try to run models locally to see if the smart contracts can execute and invoke the models directly on the device and get the inferences back.

**RQ3: Can we utilize AI accelerators such as GPU for fast inference while executing a smart contract?** Most of the AI and LLM models are compute-intensive and though we can run some of them in CPU, running them using GPU is desirable for reasonable inference throughput. We evaluate if our framework is capable of doing this.

**RQ4: What are the security implications?** We explore the security implications of such a framework its benefits over traditional methods and pitfalls. Here we show how using our framework inherently mitigates some covert channel security attacks, which are possible for native inferences in the same machine.

To address the above, we propose WASM-Powered Interchain Communication for AI Enabled Smart Contracts or WICAS in short. Which uses [24] for the blockchain environment to execute a smart contract and uses our reference implementation of AI inferences with open-source models.

In summary, the contributions of this work are listed below:

<sup>3</sup> <https://decrypt.co/63280/uniswap-trading-volume-exploded-7-billion-heres-why>

- We propose WICAS which provides a way to get AI and LLM invocation and response based on a smart contract execution, a new capability for smart contracts that may enable and power emerging application scenarios in the blockchain and Web3 domain.
- We demonstrate the viability of AI and LLM inferences based on smart contracts on local models.
- We explore the potential security benefits and considerations for the proposed system.
- We demonstrate that our framework is an AI model and underlying inference engine agnostics and can be adapted for other future systems.

To the best of our knowledge, we are the first to propose and evaluate a system to (1) run AI and machine learning inference based on smart contract execution on a device and (2) show that it is possible to build such a system in heterogeneous systems with varying CPU/GPU and scalability.

## 2 Background and Motivation

The programming models for smart contracts, such as Solidity, face significant limitations in implementing the complex mathematical models that are commonplace in traditional quantitative finance. These limitations are primarily due to constraints that prevent the use of high-performance computing required for sophisticated quantitative models within smart contracts, thereby impeding innovation in Decentralized Finance (DeFi) and other emerging blockchain applications [31]. Moreover, recent significant advancements in artificial intelligence present new opportunities for blockchain and the decentralized web. The traditional finance sector and the DeFi space are beginning to explore the transformative potential of AI in enhancing financial decision-making, dynamic pricing, prediction, and portfolio optimization [11]. Among various AI technologies, large language models (LLMs) and generative AI are viewed as particularly promising for their potential applications in the finance industry [23]. Generative AI consists of a wide range of technologies that can synthesize new data instances based on statistical patterns. Generative AI models targeting finance applications have the potential to revolutionize both traditional financial modeling and decentralized finance by their capabilities to evaluate massive amounts of numerical and textual data, generate time series, and predict financial performance. Besides data-driven modeling in quantitative finance, large language models present promising opportunities for general-purpose on-chain decision-making and governance where natural language content in conjunction with code can be applied for specifying rules and contract logic.

In Cosmos-based blockchains, smart contracts can be executed as WebAssembly (WASM) code, providing a robust environment for implementing smart contracts beyond traditional approaches. Although Cosmos can facilitate Ethereum

Virtual Machine (EVM) compatibility through Ethermint, allowing for the porting of existing Solidity contracts, this work concentrates on the built-in support for WASM smart contracts within the Cosmos ecosystem. The Cosmos SDK supports the execution of WASM smart contracts through a dedicated module, x/WASM. Thanks to the modular design of the Cosmos SDK, the WASM module can seamlessly interact with other modules, such as the staking and bank modules. Notably, the EVM support within Cosmos is also modular, implemented in a similar fashion to the WASM module. This modular architecture enhances flexibility and integration possibilities within the Cosmos ecosystem.

## 2.1 Transformer Models and Its Applications

The GPT series of language models, which utilize the Transformer architecture [34], operate in an auto-regressive manner [9]. The process starts with a sequence of tokens known as a prompt. The model attentively analyzes this initial input, progressing iteratively. In each iteration, it evaluates the likelihood of possible subsequent tokens. Through intensive training, the model learns the complex process of prediction and selection, thereby refining its ability to generate coherent and contextually appropriate text.

The procedure for processing and sampling to generate a single output token in a language model is known as an iteration. After being trained on a substantial dataset, models like GPT are capable of performing language tasks with notable proficiency. For instance, when presented with the prompt "knowledge is," the model is more likely to predict "power" rather than "apple" as the next word. Following this iteration, the generated token is appended to the initial prompt and re-fed into the model for the generation of the next token. This cycle repeats until an End of Sequence token is emitted, indicating the end of the text, or until a predetermined maximum output length is reached. This inference method, characterized by its iterative nature, contrasts sharply with architectures like ResNet, which typically exhibit a fixed and predictable execution time [15]. While each individual iteration of the language model may be predictable, the total number of iterations—and thus the total inference time—can vary.

Additionally, it is important to mention that large language models (LLMs) are not solely confined to language tasks; they can also be easily configured to model and predict time series data.

## 2.2 Inference Engine

Inference serving systems like TensorFlow Serving [28] and Triton <sup>4</sup> act as middleware for Deep Neural Networks (DNNs), handling resource allocation, task execution, and result delivery. They maximize hardware utilization, especially on GPUs, by batching multiple tasks. This batching improves overall efficiency but increases memory usage compared to single-task processing. This becomes

<sup>4</sup> <https://github.com/triton-inference-server/server>

especially critical for Large Language Models (LLMs) with their heavy memory footprint, limiting their batch size during inferences.

The surge in popularity of GPT models has prompted optimizations in these serving systems to accommodate GPT’s unique design and iterative generation process. At the heart of GPT is the Transformer architecture, featuring the distinctive Masked Self-Attention module. This module generates three key values (query, key, and value) for each input token, allowing each token to “see” how relevant other tokens are, regardless of their position. Masked Self-Attention ensures causality by hiding future tokens from each token during prediction, enabling GPT to focus on the next token in the sequence. The attention mechanism essentially grants each token a global view of the entire input sentence, considering the importance of every other token in its prediction.

Within the GPT inference loop, the attention mechanism relies heavily on key-value pairs from preceding tokens. Traditionally, these are recalculated for each step, incurring significant computational overhead. To address this inefficiency, Fairseq [29] proposes a caching scheme, pre-computing and storing these elements for subsequent use. This two-stage process involves an initial analysis phase where caches are built for each GPT layer based on the prompt. During decoding, only the newly generated token’s key, value, and query are computed, with the cache progressively updated. This optimized scheme significantly reduces per-iteration computational workload compared to the initial computations. Similar caching optimizations are also employed by libraries such as HuggingFace[35] and FasterTransformer[12].

Orca[36], on the other hand, introduces a distinct scheduling approach known as iteration-level scheduling. This contrasts with the typical batch processing paradigm by handling just one iteration per batch. While this enables dynamic job entry and completion, it presents challenges in terms of GPU memory limitations and the stringent latency requirements of interactive applications. Therefore, while Orca offers flexibility, its effective utilization requires meticulous resource management.

### 2.3 Blockchain

Blockchain has emerged as a pivotal decentralized framework for data management and storage, serving as the foundation for various distributed applications, including digital currencies[27], decentralized finance, and networks[26]. It functions as a distributed ledger system, enabling transaction recording and verification without relying on central authority or third-party validation. This technology fosters a peer-to-peer network[38] where participants collaboratively maintain a transparent and immutable public ledger[10]. Blockchain’s structure integrates a physical network that supports communication, computing, and data storage. This infrastructure underpins features like the blockchain consensus mechanism, forming a dual-layer system comprising the physical network and the blockchain. Transactions in blockchain encapsulate client information, recorded in blocks that collectively form a chain, delineating the logical sequence of these transactions. The system’s security and efficacy hinge on the consensus

mechanism and smart contracts, automating transaction execution and maintaining system integrity. These attributes position blockchain as a critical technology for enhancing AI’s reliability and trustworthiness.

### 3 Architecture

For our efforts, we have chosen the Cosmos Network[24] as our network of choice for smart contracts. Within the Cosmos network, smart contracts are like programmable gears driving communication and transactions between different blockchains. These contracts act like autonomous robots, following pre-defined instructions without anyone needing to constantly pull the levers. They’re crucial for the whole system to work smoothly.

To build and run these robot gears, Cosmos has two handy tools: CosmWASM<sup>5</sup> and Ethermint<sup>6</sup>. CosmWASM is like a custom Lego set for developers, letting them build all sorts of different contracts from scratch. Ethermint, on the other hand, acts as an adapter, allowing developers to use code already built for the Ethereum blockchain with Cosmos.

In this paper we use CosmWASM, the native Cosmos smart contract support, as our infrastructure for executing the smart contracts.

#### 3.1 CosmWASM

CosmWASM is a specialized platform for building smart contracts within the Cosmos blockchain ecosystem. It uses WebAssembly (WASM) as its engine, allowing developers to write performant and secure contracts. A strength of this is language portability. This "language agnostic" design gives the flexibility and choice. Plus, it runs the contracts in a safe virtual machine sandbox, so everything scales nicely if the project explodes in popularity. That makes CosmWASM ideal for crafting smart contracts that are both speedy and able to handle a ton of activity.

Ethermint acts as a bridge between the Cosmos ecosystem and Ethereum’s smart contract playground. It’s like a translator that lets existing Ethereum contracts seamlessly join the Cosmos party without needing a complete makeover. This fusion leverages Cosmos’s speedy transaction engine, while keeping things modular like the Cosmos SDK, giving developers plenty of building blocks to play with.

The utilization of these smart contract platforms is broad. In decentralized finance (DeFi), both CosmWASM and Ethermint prove invaluable. Ethermint’s Ethereum compatibility simplifies the transfer of existing Ethereum DeFi initiatives to Cosmos.

For applications involving non-fungible tokens and tokenization, CosmWASM’s robust and secure framework is well-suited. Additionally, in areas like supply

<sup>5</sup> <https://github.com/CosmWasm/wasmd>

<sup>6</sup> <https://docs.ethermint.zone>

chain and logistics management, the automation features of smart contracts become beneficial, with both CosmWASM and Ethermint providing appropriate tools for these sectors.

### 3.2 WebAssembly (WASM)

Introduced by the W3C in 2015, WebAssembly (WASM) is a project aimed at creating a standardized, high-performance, and secure machine-independent bytecode. It achieves safety through distinct, isolated memory regions: the stack, global variables, and a linear memory area, accessed via type-safe instructions for assured memory safety during native code compilation. WASM's architecture, incorporating capability-based security, ensures both performance and security, managing operating system resources like networking and threading through stringent security protocols.

In edge computing, WASM's architecture is crucial for fast, secure processing. By eliminating risky features while supporting C/C++ compatibility, WASM addresses numerous technological challenges. A prime example is the automotive industry's supply chain fragility, where increasing demands for features clash with the impracticality of adding more microprocessor-based control units. WASM enables automotive manufacturers to share physical hardware, reducing microprocessor demand and manufacturing costs by simplifying hardware requirements. This software-centric approach allows manufacturers to focus on advances in automation, infotainment, and safety, without supply chain concerns.

Edge computing, particularly applications requiring autonomous computing and distributed collaboration at network edge, benefits from WASM. These applications require low latency, a challenge for cloud computing. Edge computing, by processing data locally, overcomes issues of network latency and connectivity, ensuring timely data availability.

### 3.3 WASMEdge

WASMEdge facilitates the deployment of WebAssembly[3] (WASM) in edge computing environments. It allows for the integration of serverless functions, which are WASM executable, into various software platforms. Notably, WASMEdge is versatile in its application, being usable at the edge of cloud networks, serving as an API endpoint in a Function as a Service (FaaS) model, operating within Node command line interfaces, and even within embedded systems [3].

Artificial Intelligence (AI) advancements, Natural Language Processing (NLP), and machine learning (ML) are important fields of research and use. Before we get into ML implementation with WASMEdge, let's go over some basic knowledge.

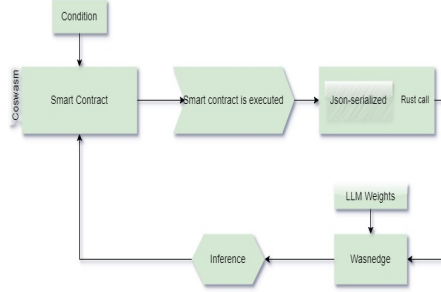
TensorFlow Lite runs on WASMEdge and is optimized for embedded systems. It makes on-device machine learning easier. By removing server communication, this method protects data privacy and helps to prevent network latency and connectivity problems.

### 3.4 WebAssembly Binary

The interface between WebAssembly (WASM) and the operating system is facilitated by the WebAssembly System Interface (WASI) [13]. The interface is standardized and resembles POSIX, with an emphasis on cross-platform compatibility and ease of use. Because of this, it can be used in Trusted Execution contexts (TEEs), edge devices, and Internet of Things contexts. WASM applications may access file systems, networks, and other OS services thanks to WASI, which has 46 functionalities. POSIX calls can be effectively converted into WASI calls by languages like C and Rust.

Additionally, WASI uses a capability-based security paradigm, which establishes a sandbox environment by requiring WASM runtime to approve resource access. By creating an intermediary layer and restricting applications to particular file system components, it improves security.

For our framework we use WASMEdge since it is the fastest WASM virtual machine available at the moment [37,3]. We explain the high level architecture of WICAS as shown in Fig. 1.



**Fig. 1.** Overview of WICAS.

For the contract, we have taken a simple name service sample code from Cosmos SDK and have modified it to pass a prespecified input dynamically created using the name lookup as an input prompt to the AI or LLM inference engine. Here for our AI or LLM inference engine we have chosen WASMEdge. Once the smart contract has been executed successfully, the lookup value is passed as a rust system call to WASMEdge.

## 4 Implementation Highlights

For WICAS we decided to use the fairly large Llama2[33] family of models. Before we can pass the input to the model, the CosmWasm modules has to compile and run the contract.

The compilation process is same for both the smart contract and the Rust inference engine we wrote for loading the model.



## 4.1 Compiling the Contract

The contract written in Rust has to be compiled to WASM code to be executed. The easiest way for us to do that was to use "cargo". Cargo is Rust's build system and package manager which allows us to easily get packages created in the Rust ecosystem. We use the following to optimize and compile it to a WASM runtime

**Listing 1.1.** WasmCosm Contract Compile

```
RUSTFLAGS= -C link-arg=-s cargo wasm
```

Once we have the compiled WASM file, this was deployed in the blockchain.

**Listing 1.2.** WasmCosm Contract Store

```
wasmd tx wasm store nameservice.wasm --from <your-key>
--chain-id <chain-id> --gas auto
```

## 4.2 Llama 2 Inference

For our inference engine we chose to use WasmEdge as to use WASM runtime. However to do that we needed a way for us to first have a program to run LLM inferences (for proof-of-concept and feasibility evaluation). For the inference engine we take inspiration from the excellent llama.cpp<sup>7</sup> project and create a simple Rust program for running the inferences. Here we chose to use Rust as a programming language to remove friction between our smart contract execution on CosmosWasm and here, since our Rust code can directly be integrated into the smart contract and does not need any external message passing scheme to work.

## 4.3 AI Inferences across Nodes

While our solution works completely fine running in a single node, to be useful in real-world blockchain setting we also explore how it performs if we have multiple Cosmos validator nodes running inferences (use the same AI model). One of the challenging use cases of such a setting will be multiple node running an inference with the same machine learning model and an identical input.

To achieve consensus, it requires AI inference reproducibility. Since LLMs are not-deterministic, we need a way to support determinism for this use case. Though inference reproducibility is a research topic of its own, existing results in the literature are enough for the purpose of our research.

Running advanced ML and LLM inference across multiple blockchain nodes could introduce variability in inference results due to several factors. These factors include using deterministic AI code, floating number resolution, node hardware capabilities, and the stochastic nature of some LLMs if they utilize mechanisms like dropout during inference for regularization purposes. It is crucial to

<sup>7</sup> <https://github.com/ggerganov/llama.cpp>

understand that blockchain architectures are inherently designed to ensure consensus despite potential discrepancies across nodes. Typically, mechanisms like Proof of Work (PoW) or Proof of Stake (PoS) are employed to achieve agreement among nodes regarding the validity of transactions, which could analogously be applied to agreeing on LLM inference results.

Based on experiments, the best practice to achieve inference reproducibility across nodes, is to use deterministic implementation of machine learning algorithm, manage random seeds with on-chain support like using on-chain random beacon (a well studied problem in the recent years), ran AI models using 64 bit floating numbers. It is important to highlight that inference reproducibility has been successfully demonstrated in one application use case where LLMs are used as a covert communication channel between a sender and a receiver. In [8], the authors demonstrate how LLMs can be configured to produce consistent LLM outputs across devices. Our case directly benefits from the findings and best practice identified in the LLM covert communication research for achieving reproducible inferences across validator nodes.

If the LLM inferences do vary across nodes, the blockchain protocol would need to specify a method for selecting the 'correct' or 'agreed-upon' inference. This is based on a majority rule where the most frequently generated result is chosen, or more sophisticated approaches like weighted scoring based on node reputability or stake can be employed.

## 5 Portability

One of the key benefits to our proposed solution is model portability and inference engine (as well as model) agnostic nature. Even though we had chosen WasmEdge to showcase our proposed framework, this can easily be replaced by any other code that can communicate directly with a Web Assembly runtime. The primary reason we chose WasmEdge or any Web Assembly runtime because of their ability to communicate with the WebGPU [22] API. We have also tested our implementation using the well documented WGPU [6] API.

### 5.1 Inference Engine Portability

Rust is used to construct the demo inference code, which is compiled to WebAssembly. This core Rust code is only forty lines long and is surprisingly small. It can process inputs from the user, log the flow of the conversation, modify the text to make it compatible with the llama2 input template, and do inference tasks using the WASI NN API. This simplified method demonstrates how well the code manages these operations and how efficient it is.

**Listing 1.3.** Rust Inference

```

1 fn main() {
2     let arguments: Vec<String> = env::args().collect();
3     let model_identifier: &str = &arguments[1];
4
5     let neural_network =
6         wasi_nn::GraphBuilder::new(wasi_nn::GraphEncoding::Ggml, wasi_nn::ExecutionTarget::AUTO)

```

```

7      .assemble_from_cache(model_identifier)
8      .expect("Successful-build");
9      let mut execution_context = neural_network.create_execution_context().expect("Context-initialization");
10
11      .....
12      .....
13      .....
14
15      // Input processing.
16      let input_tensor_data = combined_prompt.as_bytes().to_vec();
17      execution_context
18      .assign_input(0, wasi_nn::TensorType::U8, &[1], &input_tensor_data)
19      .expect("Input-set");
20
21      // Inference execution.
22      execution_context.execute().expect("Successful-computation");
23
24      // Output retrieval.
25      let mut result_buffer = vec![0u8; 1000];
26      let result_size = execution_context.obtain_output(0, &mut result_buffer).expect("Output-retrieval");
27      let response = String::from_utf8_lossy(&result_buffer[..result_size]).to_string();
28      println!("Response:\n{}", response.trim());
29
30      historical_prompt.push_str(&format!("{}", combined_prompt, response.trim()));
31  }
32 }

```

This can be compiled using

#### Listing 1.4. Rust Compile

```

curl --proto '=https' --tlsv1.2 -sSf https://sh.rustup.rs
| sh
rustup target add wasm32-wasi

```

Our handwritten Rust code is also definitely not a requirement. To showcase portability we had directly taken the llama2.c [2] engine and compiled it into WASM file.

#### Listing 1.5. llama2.c wasm

```

$ run.c -D_WASLEMULATEDMMAN -lwasi-emulated-mmman -
D.WASLEMULATED_PROCESS_CLOCKS -lwasi-emulated-process
-clocks -o run.wasm

```

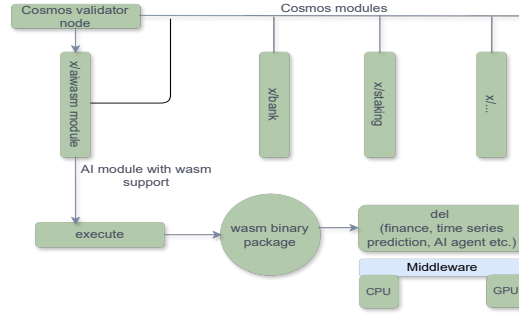
## 5.2 Model Portability

Even though we have primarily targeted the Llama 2 family of models. We are not bounded by them. Since we base our implementation on the llama.cpp, the model formats have to be ggml. For which we utilize the GGML plugin for WasmEdge [4]. We also have tested the framework with LLAVA [25] which is a multi-modal model. However our inference code is not equipped to handle non-text inputs, hence we only tested the text input and output of the model. However, a base64 encoder and decoder is certainly possible to pass multi-modal (image, numbers, text) input to the model to get its feedback.

WICAS works as shown in Fig 2.

The validator node running in Cosmwasm generates the input message to be passed to the x/aiwasm runtime (which in this case is WasmEdge) which runs the inferences.

WasmEdge can utilize both GPU or CPU based on what is available in the host system. It uses WebGPU API to see if there is a compatible GPU available,



**Fig. 2.** Working of WICAS.

and runs the inference if it is available. If not then tries to run the inference in CPU. WICAS is compatible with other Cosmos modules since it does not modify anything within the CosmWasm, but only individual contract.

## 6 Discussion

We have run the LLM/AI execution and Cosmwasm in a machine equipped with Intel I9 Processor and NVIDIA 4080 Mobile graphics processor with 16GB of RAM. We intentionally chose a moderately configured device to note any potential limitation of the system. In our testing both of the models ran in reasonable amount of time. The llama2 model was able to generate a 35 token per second and llava 30 token per second.

### 6.1 CosmWasm and WASM Contract

We had modified the name space example contract from the Cosmos example repository in Github to incorporate our Rust code. And we were able to run the compiled WASM smart contract as described in Section 4.1.

**RQ1** *Can smart contracts utilize LLM and AI inferences with reasonable performance?*

We were able to write a proof of concept code and call a LLM/AI inference engine we to produce a model output in a moderately configured commodity hardware. Within our limited scope smart contracts, or more specifically smart contract executed using CosmWasm is able to generate LLM inferences using our framework. Note that LLMs are more demanding in computing resources than other AI models.

### 6.2 Security

Since the input is executed directly from the contract and the contract is compiled into a WASM file before running. The data should be safe as long as the

host machine is safe. Also since the inference is being run in either that machine or in another clustered machine, the models and inferences on device and protected as long as the system integrity is protected.

**RQ2** *Can we get inference in local in secure way?*

The inference is done in local since the compiled models are also in local. It assumes that all the validator nodes have access to the same models.

### 6.3 Model Performance

Since we use WasmEdge for our model evaluation, we are able to utilize WebGPU API completely. That gives us access to GPU in the host system if its available, or switch to a CPU runtime based on the inference engine.

**RQ3** *Can we utilize GPU for fast inferences while executing a smart contract?*

WebGPU allows us to access and utilize the GPU runtime to run our models in the GPU in an efficient way.

### 6.4 Security Consideration

Even though the model is running locally, as the model is loaded into memory, attacks are possible. We specifically look into one specific family of attack that makes most of the local inference system today vulnerable if using native solution.

The VU#446598<sup>8</sup> dubbed as *LeftoverLocals* affects all Apple, AMD, Qualcomm family of GPUs.

In the context of GPU computing, it has been identified that one GPU kernel is capable of accessing memory values from another kernel, even when these kernels are segregated across different applications, processes, or user environments. The memory area where this interaction occurs is known as local memory. This local memory functions akin to a software-managed cache, comparable to the L1 cache found in CPUs. The size of this local memory varies between GPUs, ranging from tens of kilobytes to several megabytes. Trail of Bits demonstrated this vulnerability's presence<sup>9</sup> through various programming interfaces, including Metal, Vulkan, and OpenCL, across diverse combinations of operating systems and drivers. This finding highlights a significant security concern in GPU architecture.

However their findings report that even in vulnerable machines, the models when accessed through WebGPU did not dump any secrets apart from zeros [1].

<sup>8</sup> <https://kb.cert.org/vuls/id/446598>

<sup>9</sup> <https://blog.trailofbits.com/2024/01/16/leftoverlocals-listening-to-llm-responses-through-leaked-gpu-local-memory/>

This makes WICAS resilient to these type of memory dump attacks. However, these kind of attacks show us more work is needed to continually protect the system from future threat attacks, even though the present one did not pose a threat.

In case model confidentiality is a concern, AI models can be executed by enclaves and TEEs (trusted execution environments). TEEs such as Intel SGX are capable of running inferences for machine learning models like DNNs. The new GPU TEEs may soon support running LLM inferences in concealed environment.

**RQ4** *What are the security implications?*

Attacks like LeftoverLocals are mitigated in this framework. but future work is necessary to understand the security and privacy needs of on-chain AI inferences and assess the attack surface.

## 6.5 Potential Applications

The capability of running AI inference engine as part of contract execution could open new opportunities for application areas such as DeFi, decentralized insurance, DAO governance, voting, prediction market, AI agent based decision making, to name just few.

**Open Research Avenues AI based DeFi models:** With the kind of support described in this work, machine learning based finance models can be incorporated on-chain as part of a DeFi contract. This will certainly expand the scope of DeFi to a new level where DeFi contracts can directly apply machine learning on-chain in decision making such as pricing, liquidity optimization, risk reduction.

**On-chain governance based on natural language texts:** With LLM support, the concept "code is law" can be expanded to use both code and nature language to regulate the behavior of smart contracts. Complex business logic could be described in natural language and enforced as smart contracts.

**On-chain AI agents:** With integration of AI into smart contracts, smart contracts can act as AI agents. This may redefine to our current understanding of smart contracts and significantly increase the scope of smart contracts in future.

## 6.6 Remarks

Depending on the types of blockchain, for instance, permissioned or permissionless, running AI inferences may require gas fee. It is outside the scope of this work to specify a detailed gas model for AI inferences. It is plausible to compute gas fee based on model size and the number of LLM tokens (similar to how fee is calculated by OpenAI for using ChatGPT). In case of Cosmos module, aiwasm module can charge gas fee using a specific fee model implemented by the module. A simple approach is to apply a fee model similar to EVM pre-compiles.

Another remark is that although the current work focuses on enabling AI for WASM based smart contracts, it is plausible to apply the same concept to EVM. This could be achieved using dedicated pre-compiles. Last but not the least, it is not assumed that AI inferences must be executed on the main chain. The framework could be applied to support AI modeling in side-chain, application specific chain (app chain), or layer 2 network.

## 7 Future Work

### 7.1 WebGPU Native

Having received universal support from top browser developers, WebGPU<sup>10</sup> is an impressive step forward. There have been some work which tries to run AI inferences directly using WebGPU like tokenhawk<sup>11</sup> which tries to implement models by hand, which is not scalable. The other more mature methods all use WASM as a backend like mlc-llm [32] and WebONNX [5].

### 7.2 Usage of WebAssembly (WASM)

The large language model (LLM) and AI inference engine in our architecture run on both CPUs and GPUs and is written in Rust. This dual compatibility addresses common limitations in web-based computation by facilitating a balance between memory consumption, data transfer volume, and CPU/GPU utilization. However we still need to load the model every time. That is not suitable for large models and specially for execution on-chain and from a smart contract.

In our future work, we want to explore these two venues, of (a) coming up with a WebGPU framework that can run inferences in native speed in a programming language agnostic way, communicated just by message passing. Similar to what WebGPU allows us but not limited to Rust. We also want to explore if (b) we can develop a better way to *stream* model weights instead of loading the model completely. That will allow us to have more control for performance optimization.

## 8 Conclusion

In conclusion, our research addresses vital questions regarding the integration of AI and Large Language Models (LLMs) with smart contracts. We demonstrated through a proof of concept that smart contracts, particularly those executed using CosmWasm, can effectively generate AI/LLM inferences. The use of WebGPU technology allows for efficient GPU-accelerated inferences within a smart contract environment. Our findings affirm the feasibility, performance efficiency, and security of utilizing AI models in blockchain environments. This opens up new avenues for blockchain applications, leveraging the power of AI and LLMs while maintaining security standards.

<sup>10</sup> <https://developer.chrome.com/blog/webgpu-io2023>

<sup>11</sup> <https://github.com/kayvr/token-hawk>

## References

1. Leftoverlocalsrelease: The public release of leftoverlocals code (2024), <https://github.com/trailofbits/LeftoverLocalsRelease>
2. llama2.c: Inference llama 2 in one file of pure c (2024), <https://github.com/karpathy/llama2.c>
3. Wasmedge (2024), <https://webassembly.org/>
4. Wasmedge-wasinn (2024), <https://github.com/second-state/WasmEdge-WASINN-examples/tree/master/wasmedge-ggml-llama-interactive>
5. webonnx (2024), <https://github.com/webonnx/wonnx>
6. wgpu: Cross-platform, safe, pure-rust graphics api. (2024), <https://github.com/gfx-rs/wgpu>
7. Anil, R., Dai, A.M., Firat, O., Johnson, M., Lepikhin, D., Passos, A., Shakeri, S., Taropa, E., Bailey, P., Chen, Z., et al.: Palm 2 technical report. arXiv preprint arXiv:2305.10403 (2023)
8. Bauer, L.: Leveraging generative models for covert messaging: Challenges and tradeoffs for “dead-drop” deployments. In: Proceedings of the Thirteenth ACM Conference on Data and Application Security and Privacy (2024)
9. Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Nee-lakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models are few-shot learners. *Advances in neural information processing systems* **33**, 1877–1901 (2020)
10. Cao, B., Wang, Z., Zhang, L., Feng, D., Peng, M., Zhang, L., Han, Z.: Blockchain systems, technologies and applications: A methodology perspective. *IEEE Communications Surveys & Tutorials* (2022)
11. Cao, L.: Ai in finance (2020)
12. Chelba, C., Chen, M., Bapna, A., Shazeer, N.: Faster transformer decoding: N-gram masked self-attention. arXiv preprint arXiv:2001.04589 (2020)
13. Clark, L.: Standardizing wasi: A system interface to run webassembly outside the web. Mozilla Hacks—the Web developer blog (2019)
14. Finnie-Ansley, J., Denny, P., Becker, B.A., Luxton-Reilly, A., Prather, J.: The robots are coming: Exploring the implications of openai codex on introductory programming. In: Proceedings of the 24th Australasian Computing Education Conference. p. 10–19. ACE ’22, Association for Computing Machinery, New York, NY, USA (2022). <https://doi.org/10.1145/3511861.3511863>, <https://doi.org/10.1145/3511861.3511863>
15. Gujarati, A., Karimi, R., Alzayat, S., Hao, W., Kaufmann, A., Vigfusson, Y., Mace, J.: Serving {DNNs} like clockwork: Performance predictability from the bottom up. In: 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20). pp. 443–462 (2020)
16. Haas, A., Rossberg, A., Schuff, D.L., Titzer, B.L., Holman, M., Gohman, D., Wagner, L., Zakai, A., Bastien, J.: Bringing the web up to speed with webassembly. In: Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation. pp. 185–200 (2017)
17. Jang, J., Kim, S., Ye, S., Kim, D., Logeswaran, L., Lee, M., Lee, K., Seo, M.: Exploring the benefits of training expert language models over instruction tuning (2023). <https://doi.org/10.48550/ARXIV.2302.03202>, <https://arxiv.org/abs/2302.03202>
18. Karanjai, R., Gao, Z., Chen, L., Fan, X., Suh, T., Shi, W., Xu, L.: DhTEE: Decentralized infrastructure for heterogeneous tees. In: 2023 IEEE International Conference on Blockchain and Cryptocurrency (ICBC). pp. 1–3 (2023). <https://doi.org/10.1109/ICBC56567.2023.10174906>



19. Karanjai, R., Xu, L., Diallo, N., Chen, L., Shi, W.: Defaas: Decentralized function-as-a-service for emerging dapps and web3. In: 2023 IEEE International Conference on Blockchain and Cryptocurrency (ICBC). pp. 1–3 (2023). <https://doi.org/10.1109/ICBC56567.2023.10174945>
20. Karanjai, R., Xu, L., Gao, Z., Chen, L., Kaleem, M., Shi, W.: On conditional cryptocurrency with privacy. In: 2021 IEEE International Conference on Blockchain and Cryptocurrency (ICBC). pp. 1–3 (2021). <https://doi.org/10.1109/ICBC51069.2021.9461133>
21. Karanjai, R., Xu, L., Gao, Z., Chen, L., Kaleem, M., Shi, W.: Privacy preserving event based transaction system in a decentralized environment. In: Proceedings of the 22nd International Middleware Conference. p. 286–297. Middleware '21, Association for Computing Machinery, New York, NY, USA (2021). <https://doi.org/10.1145/3464298.3493401>, <https://doi.org/10.1145/3464298.3493401>
22. Kenwright, B.: Introduction to the webgpu api. In: ACM SIGGRAPH 2022 Courses, pp. 1–184 (2022)
23. Krause, D.: Large language models and generative ai in finance: An analysis of chatgpt, bard, and bing ai (2023)
24. Kwon, J., Buchman, E.: Cosmos whitepaper. A Netw. Distrib. Ledgers **27** (2019)
25. Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual instruction tuning. In: NeurIPS (2023)
26. Luo, H., Wu, Y., Sun, G., Yu, H., Xu, S., Guizani, M.: Escm: An efficient and secure communication mechanism for uav networks. arXiv preprint arXiv:2304.13244 (2023)
27. Nakamoto, S.: Bitcoin: A peer-to-peer electronic cash system. Decentralized business review (2008)
28. Olston, C., Fiedel, N., Gorovoy, K., Harmsen, J., Lao, L., Li, F., Rajashekhar, V., Ramesh, S., Soyke, J.: Tensorflow-serving: Flexible, high-performance ml serving. arXiv preprint arXiv:1712.06139 (2017)
29. Ott, M., Edunov, S., Baevski, A., Fan, A., Gross, S., Ng, N., Grangier, D., Auli, M.: fairseq: A fast, extensible toolkit for sequence modeling. arXiv preprint arXiv:1904.01038 (2019)
30. Rabimba, K., Xu, L., Chen, L., Zhang, F., Gao, Z., Shi, W.: Lessons learned from blockchain applications of trusted execution environments and implications for future research. In: Workshop on Hardware and Architectural Support for Security and Privacy. HASP '21, Association for Computing Machinery, New York, NY, USA (2022). <https://doi.org/10.1145/3505253.3505259>, <https://doi.org/10.1145/3505253.3505259>
31. Shah, K., Lathiya, D., Lukhi, N., Parmar, K., Sanghvi, H.: A systematic review of decentralized finance protocols. International Journal of Intelligent Networks (2023)
32. team, M.: MLC-LLM (2023), <https://github.com/mlc-ai/mlc-llm>
33. Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al.: Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971 (2023)
34. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. Advances in neural information processing systems **30** (2017)
35. Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., et al.: Huggingface’s transformers: State-of-the-art natural language processing. arXiv preprint arXiv:1910.03771 (2019)

36. Yu, G.I., Jeong, J.S., Kim, G.W., Kim, S., Chun, B.G.: Orca: A distributed serving system for {Transformer-Based} generative models. In: 16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22). pp. 521–538 (2022)
37. Zheng, S., Wang, H., Wu, L., Huang, G., Liu, X.: Vm matters: a comparison of wasm vms and evms in the performance of blockchain smart contracts. arXiv preprint arXiv:2012.01032 (2020)
38. Zyskind, G., Nathan, O., et al.: Decentralizing privacy: Using blockchain to protect personal data. In: 2015 IEEE security and privacy workshops. pp. 180–184. IEEE (2015)