# Proposal: Becoming a Core Maintainer of Dora

Mahmoud Harmouch

January 08, 2025

## Contents

## 1 Proposal: Becoming a Core Maintainer of Dora

### 1.1 Abstract

This proposal aims to outline a plan for improving the **Dora RS Framework**, identifying critical areas for enhancement, and suggesting actionable changes to address current limitations. Dora is a powerful framework for building and managing distributed dataflows such as AI-based robotic applications dataflows, etc. However, as with any growing project, there is plenty room for innovation and maintenance.

Key highlights of this proposal include introducing new backend support for the nodes communication layer, implementing substantial CLI improvements for enhanced user experience, upgrading logging infrastructure to OpenTelemetry OTLP standards, overhauling project documentation, and optimizing the memory data formats used for inter-nodes communication. Additionally, the proposal addresses the need for robust thread safety across APIs and the incorporation of a Text User Interface to visualize running nodes and dataflows.

The proposed roadmap is embedded in the belief that these improvements will not only modernize the Dora RS ecosystem but also enhance its usability, performance, and community adoption. This paper is an in-depth research of the proposed changes.

## 1.2 Communication Layer Backends

Effective communication between nodes is at the heart of any distributed system, and Dora is no exception. The communication layer ensures seamless data transmission and coordination across nodes, enabling the execution of complex dataflows. While the current implementation supports TCP-based communication, this approach, though functional, has limitations. To build a robust, and scalable system, Dora can support a diverse set of communication backends tailored to different use cases. By doing so, it can handle various deployment scenarios, from lightweight IoT setups to high-performance enterprise systems.

This section explores the enhancements to Dora's communication layer, focusing on the integration of additional protocols and technologies that promise to improve performance, scalability, and developer experience.

### 1.2.1 WebSockets

WebSockets, known for their efficient, full-duplex communication capabilities, haven't yet integrated into Dora. This integration proposal leverages the asynchronous capabilities of the `tokio-tungstenite` library which excels in high-throughput scenarios. By implementing async WebSockets, Dora opens up new possibilities for modern, real-time communication between nodes.

**1.2.1.1 Enhanced Implementation Patterns** To optimize WebSocket performance, this proposal suggests employing connection pooling and load balancing techniques. Each node maintains a pool of WebSocket connections to its peers, reducing the overhead of repeatedly opening and closing connections. An example usage pattern might look like this:

```
Node A → WebSocket Pool → Node B, Node C
Node D → WebSocket Pool → Node B, Node C
```

### 1.2.2 MQTT/DDS

Dora stands to benefit greatly from lightweight and scalable messaging protocols like MQTT and DDS. These protocols, commonly used in IoT and distributed systems, bring reliability and efficiency to the communication layer.

**1.2.2.1 Proposed Implementation** The `zenoh` crate serves as a strong foundation for implementing MQTT and DDS in Dora. `Zenoh` bridges traditional pub/sub models with query capabilities, enabling seamless integration of MQTT and DDS alongside Dora's existing mechanisms.

**1.2.2.2 Use Case**   In a factory scenario: - Sensor nodes publish telemetry data (e.g., temperature, pressure) using zenoh topics. - Processing nodes subscribe to specific topics to perform calculations or trigger actions.

```
Sensor A → Topic: "factory/temperature"
Sensor B → Topic: "factory/pressure"
Node C → Subscribes to "factory/*" (all factory-related topics)
```

### 1.2.3 Redis Pub/Sub

Redis's Pub/Sub mechanism offers a high-performance, low-latency option for inter-node communication in distributed systems. Leveraging Redis as a communication backend in Dora introduces several benefits.

#### 1.2.3.1 Use Cases

- **Event Broadcasting**: When a node completes a computation, it can broadcast the result to other nodes via Redis channels.

- **Distributed Task Queues**: Nodes can subscribe to task queues for load-balanced processing.

  ```
  Publisher → Redis Channel: "dataflow_updates"
  Subscriber (Node A) → Receives "Update: Node B completed Task X"
  Subscriber (Node C) → Receives "Update: Node B completed Task X"
  ```

### 1.2.4 gRPC

The inclusion of **gRPC** in Dora brings structured, high-performance RPC capabilities. Using the `tonic` library, gRPC enables seamless communication between nodes, even across different programming languages.

#### 1.2.4.1 Benefits of gRPC

1. **Cross-Language Compatibility**: With support for multiple languages, gRPC allows us developers to integrate Dora with diverse ecosystems.

2. **Protocol Buffers**: gRPC leverages protobufs for defining and serializing structured data.

3. **Streaming Capabilities**: With support for bidirectional streaming, gRPC allows nodes to exchange continuous data streams in real time.

4. **Performance**: gRPC is highly optimized for speed, with built-in multiplexing over HTTP/2.

### 1.2.4.2 Example

```
// schema.proto

service Dataflow {
  rpc SendTask(Task) returns (Status);
  rpc StreamUpdates(stream UpdateRequest) returns (stream Update);
}
```

This schema allows for straightforward RPC methods and streaming updates between nodes.

### 1.2.5 QUIC

QUIC, a modern transport protocol built on UDP, offers exceptional performance and reliability. Integrating QUIC into Dora via the `quinn crate` introduces several advantages:

1. **Reduced Latency**: QUIC's built-in multiplexing minimizes round trips, reducing latency even in congested networks.

2. **Connection Resilience**: QUIC maintains session continuity despite network interruptions, ideal for mobile and unstable environments.

3. **Enhanced Security**: With TLS baked into the protocol, QUIC ensures secure, encrypted communication without additional overhead.

**1.2.5.1 Use Case**   In edge computing, nodes operating in unstable network conditions (e.g., rural IoT devices) can leverage QUIC to maintain reliable connections with minimal downtime.

### 1.2.6 Potential Impact

By integrating a variety of communication backends, Dora RS can support a broader range of use cases and deployment environments. Each protocol offers unique strengths, allowing developers to choose the best tool for their specific needs. This flexibility, combined with enhanced performance and reliability, makes Dora RS a cutting-edge framework for distributed systems, capable of handling the demands of modern applications at any scale.

### 1.3 Enhancing the CLI

### 1.3.1 Current State and Challenges

The existing CLI for Dora provides the essential functionality to run, start, and stop dataflows but lacks sophistication in terms of usability, feedback mechanisms, and developer-focused features. While functional, the CLI falls short of modern expectations for interactivity, clarity, and actionable feedback, especially when dealing with complex workflows. For example, error messages are often cryptic, leaving users confused about the root cause of issues. A user encountering a failed

dataflow initialization might struggle to identify whether the issue lies in a configuration file, a missing dependency, or a misbehaving node. This lack of clarity creates a steep learning curve for new users and frustrates experienced developers working under tight deadlines.

Moreover, the CLI does not provide real-time visibility into what is happening during the execution of tasks. Users starting a dataflow receive no visual indication of progress or detailed feedback on individual stages, such as configuration loading, node compilation, or connection establishment. This lack of transparency can make the process feel opaque and discouraging, especially when something goes wrong. For example, without timing metrics or progress indicators, it is impossible to determine whether the CLI has stalled or if a particular task is just slow.

One of the most noticeable usability issues is the inability to gracefully terminate the daemon when no coordinators are available. Currently, if a coordinator cannot be found, the CLI enters an infinite retry loop, forcing users to manually kill the process:

```
 error: failed to connect to dora-coordinator. Retrying in 1s..
2025-01-07T20:34:06.298518Z  WARN dora_daemon::coordinator: Could not connect to: 127.0.0.1
2025-01-07T20:34:07.301254Z  WARN dora_daemon::coordinator: Could not connect to: 127.0.0.1
2025-01-07T20:34:08.303026Z  WARN dora_daemon::coordinator: Could not connect to: 127.0.0.1
```

This not only wastes time but also creates unnecessary friction in an otherwise efficient workflow. A developer trying to debug a distributed system might waste valuable time troubleshooting this issue when a simple informative message and a clean shutdown would have been sufficient.

Finally, the CLI does not include advanced features such as log filtering, customizable output formats, or support for configuration profiles, which are increasingly expected in developer tools. This limits its effectiveness in large-scale production environments, where debugging, monitoring, and automation are critical.

### 1.3.2 Proposed Improvements

To address the identified issues, a comprehensive enhancement plan for the CLI is proposed. The goal is to make the CLI intuitive, feature-rich, and capable of meeting the demands of both beginners and advanced users.

#### 1.3.2.1 Enhanced Feedback Mechanisms

1. **Real-Time Progress Indicators**: One of the most significant upgrades to the CLI would be the introduction of real-time progress indicators. These can include animated progress bars, percentage completion markers, and dynamic updates for each stage of execution. For example:

```
dora start dataflow.yaml


Loading configuration... [ done in 0.2s]


Building nodes:
   Node A               100% [ done in 1.2s]
```

```
   Node B                    100% [ done in 1.5s]
   Node C                     30% [ building, 5.0s elapsed]
   Node D                      0%  [ Error, 0.0s elapsed]


Establishing connections:
   Node A   Node B           100% [ done in 0.3s]


Starting dataflow:
   Node A [ running]
   Node B [ running]
   Node C [ running]


Dataflow running successfully!
Tip: Use`dora logs` for real-time monitoring.
```

Each progress indicator is tied to a specific task, with animations providing visual feedback on completion. Additionally, color coding (green for success, yellow for warnings, red for errors) makes it easy to understand the state of the process at a glance. Moreover, icons such as □, □, and □ clearly indicate success, in-progress, and failed states, respectively.

2. **Timing Metrics**: The CLI should also display detailed timing metrics for each step. For instance, the time taken to load configurations, compile nodes, and establish connections should be shown individually. This not only helps developers identify bottlenecks but also creates a sense of transparency and trust in the system.

3. **Color-Coded Output**: Adding color coding to CLI output improves readability and allows users to quickly distinguish between different types of messages. For example:

   - **Green**: Success messages (e.g., "Node A built successfully").
   - **Yellow**: Warnings (e.g., "Node B has high memory usage").
   - **Red**: Errors (e.g., "Node C failed to start").
   - **Blue**: Informational messages (e.g., "Starting dataflow...").

**1.3.2.2 Graceful Shutdowns**  The CLI should handle termination signals (e.g., SIGINT for Ctrl+C) gracefully. When a daemon cannot find a coordinator, instead of entering an infinite retry loop, the CLI should:

1. Display a warning message with detailed context.
2. Retry a configurable number of times (e.g., 2 attempts) before exiting.
3. Offer the user an option to override the default retry limit via a flag.

```
dora daemon
Starting daemon...
No coordinators found. Retrying in 1 seconds... [Attempt 1/2]
No coordinators found. Retrying in 1 seconds... [Attempt 2/2]


Maximum retries reached. Exiting.


Suggested Fix:
   Make sure a coordinator is running and accessible at the specified endpoint.
```

**1.3.2.3 Expanded Command Set**  The CLI should support a set of commands to address common user needs more effectively. Key additions include:

1. **Log Streaming with Filters**: Enable users to view logs in real time with powerful filtering options. For example:

```
dora logs --filter error --tail 10

[2025-01-07T20:34:06.298518Z] [ERROR] [Node B] Connection timeout.
[2025-01-07T20:34:07.298518Z] [ERROR] [Node C] Failed to parse input data.
[2025-01-07T20:34:08.298518Z] [ERROR] [Node A] Dependency not found: `tokio`.
```

2. **Status Summary**: Introduce a `status` command to provide a high-level overview of all nodes and dataflows. For example:

```
dora status

+----------+---------+------------+---------+---------+
| Node     | Status  | Input Rate | Latency | Errors  |
+----------+---------+------------+---------+---------+
| Node A   | Running | 10k msgs/s | 3ms     | 0       |
| Node B   | Running | 8k msgs/s  | 5ms     | 2       |
| Node C   | Paused  | 0 msgs/s   | -       | 0       |
+----------+---------+------------+---------+---------+
```

3. **Interactive `init` Command**: Guide users through creating a new dataflow configuration interactively (wizard-like experience):

```
dora init

Welcome to Dora! Lets set up your dataflow.

  Name your dataflow: dataflow-name
  Select a runtime environment: local
  Add a node (leave empty to finish): node-a
  Add another node (leave empty to finish): node-b
  Configure logging level [INFO]:
  Save configuration as `dataflow.yaml`? [y/N]: y

Configuration saved! Use `dora start dataflow.yaml` to begin.
```

**1.3.2.4 Improved Error Reporting**  Replace vague error messages with actionable suggestions. For example:

```
dora start dataflow.yaml
Loading configuration... [ done in 0.2s]
Building nodes:
  Node A                100% [ done in 1.2s]
  Node B                100% [ done in 1.5s]
  Node C                  0%  [ Error, 5.0s elapsed]
```

```
Build Failed:
    Node C [failed: Missing dependency `tokio`]

Suggested Fix:
  1. Navigate to the Node C directory.
  2. Install the missing dependency: `cargo add tokio`.
  3. Retry: `dora build`.

Error details logged at: ./out/build.txt
```

**1.3.2.5 CLI Configuration Profiles**   Introduce profiles to simplify repetitive tasks. Example:

```
dora run --profile=staging
```

```
Loaded profile: staging
  Starting dataflow: MyDataflow
```

**1.3.2.6 Extending Output Formats**   Support additional output formats such as JSON and YAML for integration with external tools:

```
dora status --output json
{
  "dataflows": 1,
  "nodes": [
    {"name": "Node A", "status": "Active", "latency": "3ms"},
    {"name": "Node B", "status": "Paused", "latency": null}
  ]
}
```

### 1.3.3 Potential Impact

Enhancing Dora's CLI with these features will transform it into a powerful tool that supports both beginners and advanced users. By improving error feedback, real-time monitoring, and configurability, developers can focus on building and optimizing dataflows rather than troubleshooting opaque processes. These changes will elevate Dora's usability, making it a strong contender in real-time data processing workflows.

## 1.4 Improve Logging

### 1.4.1 Current State and Challenges

Logging is a fundamental standard of modern software systems, serving as a primary tool for developers to gain real-time visibility into system behavior. It provides invaluable insights for debugging, monitoring, and optimizing applications. However, the current logging configuration in Dora presents several challenges that block its effectiveness and limit its potential for advanced observability.

One of the key limitations is **Dora's reliance on Jaeger for distributed tracing**. While Jaeger has historically been a popular choice, it has deprecated its native clients in favor of the **OpenTelemetry OTLP protocol**. This shift towards OpenTelemetry reflects the broader industry trend toward a standardized, vendor-neutral approach to telemetry, covering logs, traces, and metrics in a unified workflow.

Another significant issue is the **lack of advanced log filtering and querying capabilities in the CLI**. The current CLI only supports basic log outputs, making it difficult for users to search logs based on specific criteria, such as date ranges, severity levels, or node-specific events. This lack of sophistication in log interaction limits developers' ability to diagnose issues, particularly when working with large-scale, distributed dataflows.

Finally, **Dora's documentation lacks comprehensive examples of logging integrations** with popular observability tools. While it supports basic logging functionality, there are few real-world examples demonstrating how to set up Dora with centralized log aggregation systems like Logstash or metrics visualization tools like Grafana. This lack of practical guidance can be a barrier for developers looking to implement a robust observability stack with Dora.

### 1.4.2 Proposed Improvements

To address these challenges, we propose several key improvements to Dora's logging infrastructure, focusing on **adopting modern protocols, providing advanced CLI features, and standardizing logging practices**.

#### 1.4.2.1 Adopt OpenTelemetry OTLP
The first and most critical improvement is transitioning Dora's logging and tracing system from **Jaeger's legacy protocol to OpenTelemetry OTLP**.

#### 1.4.2.2 Why OpenTelemetry OTLP

1. **Interoperability**: OTLP is supported by a wide range of observability tools, including **Datadog**, **Grafana**, **Splunk**, and **New Relic**. By adopting OTLP, Dora users can seamlessly integrate their logs and traces with these platforms, enabling powerful visualization and analysis capabilities.

2. **Future-Proofing**: As the industry standard for telemetry, OpenTelemetry ensures long-term compatibility with evolving observability ecosystems. Adopting OTLP now will reduce the need for future migrations and ensure that Dora remains relevant in the long term.

3. **Unified Telemetry**: Unlike Jaeger, which primarily focuses on traces, OpenTelemetry supports **logs, traces, and metrics in a single workflow**. This simplifies the debugging process by allowing developers to correlate logs and traces within the same context.

### 1.4.3 Provide Comprehensive Examples

To encourage adoption and ease integration, Dora's documentation should include **real-world logging examples** that demonstrate how to set up Dora with popular observability tools. These examples should cover various use cases, such as **centralized log aggregation**, **metrics visualization**, and **searchable log analysis**.

#### 1.4.3.1 Example Integrations

1. **Logstash** for centralized log aggregation.
2. **Prometheus + Grafana** for metrics visualization and alerting.
3. **Elasticsearch + Kibana** for advanced log search and analysis.
4. **Hybrid Systems**: Examples for combining Dora with `axum`-based microservices, providing end-to-end traceability.

Basic example setup for docker compose **Elasticsearch + Kibana** would look like:

```yaml
# docker-compose.yml

services:
  dora:
    image: dora-rs/dora:latest
    environment:
      - OTEL_EXPORTER_OTLP_METRICS_ENDPOINT=http://localhost:4318/v1/metrics
  elasticsearch:
    image: docker.elastic.co/elasticsearch/elasticsearch:latest
    ports:
      - "9200:9200"
  kibana:
    image: docker.elastic.co/kibana/kibana:latest
    ports:
      - "5601:5601"
    depends_on:
      - elasticsearch
```

### 1.4.4 Real-Time Log Insights

One of the most impactful improvements to the logging system would be **real-time log visualization**. By integrating real-time log dashboards, Dora can offer developers immediate feedback on system behavior, enabling faster diagnosis and resolution of issues.

#### 1.4.4.1 Proposed Features

1. **Live Search**: Allow users to search for specific keywords in logs as they are streamed.
2. **Log Grouping**: Automatically group logs by node, severity, or timestamp to make them easier to analyze.

3. **Export Options**: Enable users to export logs in **JSON**, **CSV**, or **YAML** formats for further analysis.

### 1.4.5 Standardize Logging Practices

Consistency in logging is crucial for effective observability. Dora should enforce **standard logging practices** across all modules and nodes to ensure that logs are easy to read, parse, and correlate.

#### 1.4.5.1 Standardized Logging Formats

1. **Timestamps**: Use **ISO 8601 format** for all log timestamps to ensure consistency and compatibility with most log management systems.
2. **Structured Logging**: Use **key-value pairs** for log metadata, making logs machine-readable and easier to analyze.
3. **Trace Context**: Attach **trace IDs** to every log entry to correlate logs with distributed traces.

Example standardized log entry:

```
{
  "timestamp": "2025-01-07T20:34:07.298518Z",
  "level": "ERROR",
  "node": "Node B",
  "message": "Connection timeout",
  "trace_id": "abc"
}
```

### 1.4.6 Potential Impact

By addressing the current logging challenges, Dora will significantly improve its observability capabilities, making it more suitable for **mission-critical deployments**. Enhanced logging features will:

1. Simplify **debugging and troubleshooting**.
2. Provide **real-time visibility** into system behavior.
3. Ensure **long-term sustainability** by adopting industry standards.
4. Enable **seamless integration** with popular observability tools.

These improvements will position Dora as a **developer-friendly framework** that prioritizes observability and modern telemetry practices, ultimately increasing its adoption in production environments.

## 1.5 Proper Documentation

### 1.5.1 Current State and Challenges

Documentation is the cornerstone of any open-source project. It plays a crucial role in enabling us developers to understand, use, and contribute to a framework effectively. Despite Dora's potential as a powerful framework for building distributed dataflows, its current documentation is incomplete and lacks the depth and structure needed to support a wide range of users, from beginners to advanced developers.

One of the most pressing issues is **incomplete coverage**. Most core concepts, such as multi-node workflows, error handling strategies, and performance optimizations, are either touched on briefly or completely missing from the current documentation. Without detailed explanations of these advanced topics, users are left to experiment through trial and error, which can be time-consuming and frustrating, especially when dealing with complex distributed systems.

In addition to coverage gaps, **the lack of real-world examples** significantly restrict the practical usability of the documentation. Many developers prefer learning by example, yet Dora's documentation provides few realistic use cases. For instance, there are no comprehensive examples of deploying Dora in a production environment, handling large-scale data processing tasks, or integrating with third-party services like databases or cloud platforms. This lack of practical context makes it difficult for users to visualize how Dora can be applied to solve real-world problems.

Another critical challenge is the **absence of beginner-friendly tutorials**. Dora's documentation assumes a certain level of familiarity with distributed systems and dataflow concepts, which can make it intimidating for new users. Without step-by-step guides, new users struggle to grasp the basics, from setting up their environment to creating their first dataflow. This steep learning curve discourages adoption and limits Dora's reach to a broader developer community.

Lastly, the **API documentation is underdeveloped**. While the APIs are functional, the documentation often lacks detailed explanations of parameters, return values, and potential use cases. In some cases, developers have to rely on source code or external discussions to fully understand how to use certain APIs. This inconsistency in API documentation creates confusion and increases the likelihood of mistakes, particularly for developers who are new to the framework.

### 1.5.2 Proposed Improvements

Addressing these documentation gaps is essential to improving Dora's usability and adoption. Below are detailed recommendations for enhancing core documentation, creating comprehensive tutorials, and improving API references.

**1.5.2.1 Improve Core Documentation**  The core documentation should be improved to cover both **basic** and **advanced** concepts comprehensively. This improvement should follow a structured approach, gradually introducing users to more complex topics as they progress.

1. **Getting Started Guide**: The first thing new users need is a straightforward, step-by-step guide to set up Dora, create a basic dataflow, and deploy it. This guide should assume no prior knowledge of distributed systems and focus on simplicity and clarity.

   - **Setup**: Explain how to install Dora and its dependencies on different platforms (Windows, macOS, Linux).
   - **Creating a Basic Dataflow**: Walk users through creating a simple dataflow with two nodes exchanging messages.
   - **Deployment**: Show how to run the dataflow locally and in a Dockerized environment.

2. **Intermediate Tutorials**: Once users are comfortable with the basics, the documentation should introduce more complex topics, such as:

   - Multi-node workflows: Explain how to manage multiple interconnected nodes.
   - Multi-coordinator setups: Guide users through deploying dataflows across multiple coordinators for fault tolerance.
   - Integration with third-party services: Provide examples of integrating with databases, cloud storage, and message brokers.

3. **Advanced Tutorials**: For power users, the documentation should include advanced tutorials covering:

   - Distributed deployments: Show how to deploy Dora across multiple servers or cloud environments.
   - Performance tuning: Explain strategies for optimizing node performance and minimizing latency.
   - Writing custom operators and plugins: Guide users through extending Dora's functionality by creating their own nodes and operators.

In addition to tutorials, the documentation should include detailed conceptual explanations, accompanied by diagrams and code snippets where applicable.

1. **Dataflow Architecture**: Provide a clear, visual representation of Dora's architecture, illustrating how nodes interact, how data is transferred, and how workflows are executed. Diagrams can be created using **Mermaid.js**, currently supported in Dora, or similar tools like **diagrams** python package to help users visualize complex interactions.

2. **Error Handling**: Explain best practices for detecting, logging, and recovering from errors in distributed dataflows. Include examples of implementing retry mechanisms, handling node failures, and using structured logging to diagnose issues effectively.

3. **Extensibility**: Provide a comprehensive guide on how to write custom nodes, operators, and plugins. Include a sample project showcasing a custom node and explain how to package and distribute it for reuse.

**1.5.2.2 Improve API Documentation**   A well-structured API documentation is crucial for helping developers use Dora efficiently. The following are some key improvements that can make the API reference more user-friendly:

1. **Detailed Descriptions**: Each method, function, or struct should include a clear explanation of its purpose and usage.

2. **Parameter and Return Value Documentation**: Describe what each parameter represents and what the function or method returns. This helps developers avoid confusion and misuse.

3. **Code Examples**: Provide at least one usage example for each API. Include practical examples that demonstrate common use cases.

4. **Common Pitfalls**: Highlight potential mistakes or "gotchas" that developers might encounter when using specific APIs, along with solutions.

Example (Proposed API Doc Entry):

```
/// Initializes a `DoraNode` using a given `NodeId`.
///
/// This function establishes a connection to the Dora daemon, requests the node's configur
/// and initializes the node based on the received configuration.
///
/// # Parameters
/// - `node_id` - The unique identifier of the node to be initialized (`NodeId`).
///
/// # Returns
/// A result containing:
/// - `Self` - An instance of the initialized `DoraNode`.
/// - `EventStream` - A stream of events from the node.
///
/// If an error occurs during the process, an `eyre::Error` is returned.
///
/// # Behavior
/// - Connects to the Dora daemon using a TCP channel.
/// - Sends a request to fetch the node's configuration.
/// - Handles the daemon's response to either initialize the node or propagate an error.
///
/// # Errors
/// - Returns an error if the connection to the daemon fails.
/// - Returns an error if the daemon responds with an invalid configuration or an unexpecte
///
/// # Example
/// ```no_run
/// use dora_node_api::DoraNode;
/// use dora_node_api::dora_core::config::NodeId;
///
/// fn main() -> eyre::Result<()> {
///     // Initialize the node with a given NodeId
///     let (mut node, mut events) = DoraNode::init_from_node_id(NodeId::from("plot".to_str
///         .expect("Could not initialize node with ID 'plot'");
///
///     while let Some(event) = events.recv() {
///         match event {
///             Event::Input {
///                 id,
```

```
///                 metadata,
///                 data,
///             } => match id.as_str() {
///                 "speech" => {
///                     let message: &str = (&data).try_into()?;
///                     println!("I heard: {message} from {id}");
///                 }
///                 other => println!("Received input `{other}`"),
///             },
///             _ => {}
///         }
///     }
///
///     Ok(())
/// }
/// ```
///
/// # Implementation Details
/// - The function establishes a TCP connection to the daemon at `127.0.0.1:DORA_DAEMON_LOC
/// - It sends a `DaemonRequest::NodeConfig` to the daemon to fetch the node's configuratio
/// - The response is matched to either initialize the node or return an error based on the
///
/// # See Also
/// - [`DoraNode::init`](#method.init) for initializing a node directly with a `NodeConfig`
pub fn init_from_node_id(node_id: NodeId) -> eyre::Result<(Self, EventStream)> {
}
```

**1.5.2.3 Add Real-World Examples**   Real-world examples can bridge the gap between theory and practice, helping users understand how Dora can be applied to solve real problems.

1. **Handling Large-Scale Data Processing**: Provide an example of a dataflow that processes millions of data points in parallel.

2. **Integration with Cloud Services**: Show how to integrate Dora with AWS S3 for storage or Google Pub/Sub for messaging.

3. **Using External Libraries**: Demonstrate how to extend Dora by integrating external Rust crates for specialized functionality.

**1.5.3 Potential Impact**

Improved documentation will have a **significant impact on Dora's adoption and usability**. By providing comprehensive tutorials, detailed conceptual guides, and well-structured API references, Dora will become much more accessible to new users. Developers will spend less time troubleshooting and more time building, which can accelerate the adoption of Dora in production environments.

Moreover, a well-documented framework attracts contributors from the community, promoting an ecosystem of shared knowledge and innovation. By reducing the learning curve and addressing common pain points, Dora can position itself as

a leading framework for distributed dataflows in Rust, appealing to both individual developers and enterprise users.

## 1.6 Exploring Enhanced In-Memory Data Formats

### 1.6.1 Current State and Challenges

Dora's data processing efficiency heavily relies on its **in-memory data format**. Since dataflows are at the core of the framework, the way data is stored and exchanged between nodes directly impacts performance. In scenarios demanding **high throughput**, **low latency**, or **real-time analytics**, the choice of data format can be a decisive factor in ensuring seamless operations.

Currently, Dora leverages **Apache Arrow** as its primary in-memory data format. Arrow's columnar data structure offers a modern, efficient way to handle large datasets in memory, minimizing the need for repeated serialization and deserialization. While this format has numerous advantages, particularly for analytics and batch processing, it also presents limitations in **resource-constrained environments**.

Some of the key challenges with relying solely on a single format include:

1. **One-Size-Fits-All Limitation:** While Apache Arrow excels at processing large data batches and analytical workloads, it may not be the best fit for every use case. Tasks requiring **real-time processing**, **low-latency communication**, or operating in **memory-constrained environments** might benefit from alternative formats that prioritize speed, memory efficiency, or schema flexibility.

2. **Lack of Configurability:** Currently, developers using Dora have limited flexibility to choose or configure data formats based on their specific workloads. The inability to switch between formats or tailor data structures to fit performance needs makes Dora less adaptable to diverse domains, such as **IoT**, **finance**, or **machine learning pipelines**.

Addressing these challenges by introducing **support for multiple in-memory data formats** would make Dora more versatile, allowing us developers to optimize out workflows for both **general-purpose** and **performance-critical** applications.

### 1.6.2 Proposed Enhancements

To improve performance and adaptability, we propose extending Dora's support for **multiple in-memory data formats** beyond Apache Arrow. Additionally, **user-configurable options** should be introduced to allow developers to select and fine-tune the data format that best fits their application's requirements.

The proposed enhancements include:

1. **Support for Alternative In-Memory Data Formats**

2. **User Configuration for Flexible Data Handling**

3. **Performance Benchmarks for Informed Decision-Making**

4. **Comprehensive Documentation and Use Cases**

### 1.6.3 Expanding In-Memory Data Formats

Dora can greatly benefit from supporting a range of data formats, each optimized for specific scenarios. The following are some alternative formats and their potential use cases:

**1.6.3.1 Cap'n Proto** Cap'n Proto is a **high-performance serialization format** that eliminates the need for serialization and deserialization. It allows data to be accessed directly from its serialized form, reducing latency and memory usage.

**Key Benefits:**

- **Near-Zero Serialization Overhead:** Data can be read directly without additional parsing.
- **Fast Communication:** Ideal for distributed systems with frequent, real-time node communication.

**Ideal Use Cases:**

- Low-latency systems such as **HFT platforms**.
- Real-time messaging applications.
- **Gaming servers** or **IoT devices**.

**1.6.3.2 FlatBuffers** FlatBuffers is a compact, **zero-copy** serialization format designed by Google. It is optimized for **resource-constrained environments** and provides direct access to serialized data.

**Key Benefits:**

- **Low Memory Overhead:** Data can be accessed directly without unpacking.
- **Efficient for Embedded Systems:** Ideal for scenarios where memory and CPU resources are limited.

**Ideal Use Cases:**

- Mobile and embedded systems.
- IoT devices with limited hardware resources.
- AR/VR applications.

**1.6.3.3 Apache Avro**   Apache Avro is a **schema-based** format widely used in **big data** systems. It ensures **data consistency** and provides efficient serialization for dynamic, evolving data structures.

**Key Benefits:**

- **Schema Evolution:** Supports backward and forward compatibility.
- **Compact Serialization:** Efficient for storing large datasets in distributed environments.

**Ideal Use Cases:**

- **Event-driven architectures** like Kafka.
- **Big data pipelines** using tools like Hadoop and Spark.
- Systems requiring **dynamic schema updates**.

### 1.6.4  User Configuration Options

To make the adoption of different data formats seamless, Dora should allow developers to configure their preferred data format at both the **global level** and the **node level**.

### 1.6.4.1  Configuration Example

```
default_data_format: arrow

nodes:
  - id: node-a
    data_format: capnp

  - id: node-b
    data_format: flatbuffers

  - id: node-c
    data_format: arrow
```

- **Global Configuration:** The `default_data_format` specifies the format used across all nodes unless overridden.
- **Node-Level Configuration:** Each node can override the global setting with a custom data format.

This flexibility enables developers to optimize performance for different nodes based on their specific roles within the dataflow.

### 1.6.5  Performance Benchmarks

Providing **performance benchmarks** will help developers make informed decisions when selecting data formats. These benchmarks should compare serialization and deserialization times, memory usage, and suitability for different workloads.

| Data Format | Serialization Time | Deserialization Time | Memory Usage | Best Use Case |
|---|---|---|---|---|
| Apache Arrow | Low | Low | Moderate | Real-time analytics |
| Cap'n Proto | Very Low | Very Low | Low | Low-latency systems |
| FlatBuffers | Low | Very Low | Very Low | IoT/mobile systems |
| Apache Avro | Moderate | Moderate | Moderate | Big data pipelines |

### 1.6.6 Comprehensive Documentation and Use Cases

Dora's documentation should provide **detailed guides** for each supported format, including:

- **Benefits and trade-offs of each format.**
- **Step-by-step instructions for configuration.**
- **Example scenarios demonstrating optimal use cases.**

**Example:**

```
### Choosing the Right Data Format

| Use Case                     | Recommended Format |
|------------------------------|--------------------|
| Real-Time Analytics          | Apache Arrow       |
| Low-Latency Applications     | Cap'n Proto        |
| Resource-Constrained Devices | FlatBuffers        |
| Big Data Processing          | Apache Avro        |
```

### 1.6.7 Potential Impact

Introducing **multiple in-memory data formats** will make Dora a more flexible and powerful framework for developers across diverse industries. By reducing serialization overhead and optimizing memory usage, these enhancements will improve Dora's suitability for:

- **High-throughput systems** (e.g., HFT trading platforms).
- **Low-latency applications** (e.g., gaming, IoT).
- **Big data pipelines** (e.g., event-driven architectures).

This evolution will position Dora as a **next-generation dataflow framework**, capable of adapting to the growing demands of **real-time**, **resource-efficient**, and **high-performance** applications across various domains.

## 1.7 TUI

### 1.7.1 Overview

A Text User Interface for Dora would provide users with a convenient way to interact with and monitor dataflows in real time without requiring a graphical interface. The TUI could be particularly useful for developers working in terminal-based environments or managing workflows on remote servers.

### 1.7.2 Current State and Challenges

Dora currently lacks any form of interactive interface for monitoring or managing dataflows. Users must rely on logs or external tools, which: 1. Limits real-time visibility into active workflows. 1. Makes debugging and performance monitoring cumbersome. 1. Provides no intuitive way to interact with the system in real time.

### 1.7.3 Proposed Features

#### 1.7.3.1 Core Functionalities

1. **Real-Time Dataflow Monitoring**:
   - Display a list of active dataflows, including their status (e.g., running, paused, or completed).
   - Show key metrics such as throughput, latency, and error rates for each node.

2. **Node-Level Details**:
   - Allow users to drill down into individual nodes to view:
     - Input/output data rates.
     - Active connections.
     - Memory usage and processing time.

3. **Interaction Capabilities**:
   - Pause, resume, or restart specific nodes or entire dataflows.
   - Dynamically update configurations for nodes (e.g., change data formats, input/output channels).

4. **Error Logs**:
   - Display recent errors and warnings, grouped by severity and source.

5. **Customizable Views**:
   - Let users filter or reorder displayed information (e.g., focus on nodes with the highest error rates or latency).

### 1.7.3.2 Example TUI Layout

```
+-------------------------------------------------+
|        Dora TUI - Real-Time Dataflow Manager    |
+-------------------------------------------------+
| Dataflows:                                      |
|  1. OrderProcessing (Running)                   |
|  2. AnalyticsPipeline (Paused)                  |
|  3. LogAggregator (Completed)                   |
+-------------------------------------------------+
| Selected Dataflow: OrderProcessing              |
| Status: Running | Nodes: 5                      |
+--------------------+----------------------------+
| Node Name  | Status | Input Rate | Latency      |
| -----------|--------|------------|------------- |
| Ingest     | Active | 10k msgs/s | 3ms          |
| Processor  | Active | 8k msgs/s  | 5ms          |
| Aggregator | Active | 8k msgs/s  | 4ms          |
| ErrorLogger | Idle  | -          | -            |
| Output     | Active | 8k msgs/s  | 2ms          |
+-------------------------------------------------+
| Logs:                                           |
| [WARN] Processor: High latency detected (5ms).  |
| [ERROR] Ingest: Failed to process batch #1234.  |
+-------------------------------------------------+
| Actions: (p)ause (r)esume (l)ogs (q)uit         |
+-------------------------------------------------+
```

### 1.7.4 Implementation

1. **Framework Selection**:
   - Use **TUI libraries** such as Ratatui to create the interface.
   - Libraries like `tokio` can provide asynchronous updates for real-time metrics.

2. **Integration with Dora**:
   - Expose an API endpoint in Dora for real-time metric streaming.
   - Use WebSockets or similar protocols to push updates to the TUI.

3. **Configuration**:
   - Allow users to customize the TUI (e.g., themes, displayed metrics) through a configuration file.

### 1.7.5 Potential Impact

A TUI will: - Improve user experience, especially for developers who manage workflows in terminal environments. - Enhance real-time debugging and monitoring capabilities. - Reduce reliance on external tools, making Dora more self-sufficient.

### 1.8 Use Safe Rust

To maximize the **safety, reliability, and concurrency** of applications built with Dora, all APIs should be rewritten or verified to ensure they adhere to **100% safe Rust** practices. By eliminating all `unsafe` code, Dora will align with Rust's core philosophy: providing **memory safety, data race prevention**, and **thread-safe abstractions** without sacrificing performance. This will enhance **developer confidence** and make the framework more appealing for **mission-critical systems**, particularly in industries like **finance**, **aerospace**, **IoT**, and **machine learning**.

The key to achieving this goal is leveraging Rust's **ownership model**, **type system**, and **concurrency guarantees** while enforcing safe coding practices through the use of the `#![forbid(unsafe_code)` directive. This ensures that **all code** within Dora, including core APIs and extensions, is memory-safe and free from undefined behavior.

#### 1.8.1 Current State and Challenges

While Rust offers robust mechanisms for ensuring **memory safety** and **concurrency guarantees**, Dora currently has APIs that rely on **unsafe Rust** to achieve certain low-level optimizations. This presents several challenges:

**1.8.1.1 Potential Memory Safety Risks**  APIs that use `unsafe` code may bypass Rust's built-in safety checks, introducing risks such as:

- **Dangling pointers**: References to deallocated memory.
- **Use-after-free errors**: Accessing memory that has already been freed.
- **Data races**: Concurrent access to shared data without proper synchronization.

These risks are particularly concerning in **multithreaded environments**, where concurrent workflows are common. Even a single unsafe operation can lead to **undefined behavior**, making applications unpredictable and difficult to debug.

**1.8.1.2 Developer Friction**  The presence of `unsafe` code increases **developer friction**, as users must carefully audit their own code to ensure it interacts safely with Dora's APIs. This goes against Rust's promise of **providing memory safety by default**, which is one of the language's primary selling points.

#### 1.8.2 Proposed Enhancements

To address these challenges, we propose a comprehensive strategy to make **all Dora APIs 100% Rust safe**, ensuring that **unsafe code is completely removed** from the framework.

**1.8.2.1 Audit and Refactor APIs for Safety** The first step is to **audit all existing APIs** and **remove any use of unsafe Rust**. This involves identifying areas where unsafe code is currently being used and refactoring them to use **safe abstractions** provided by the Rust standard library or third-party crates.

Key actions include:

- **Replace raw pointers** with **smart pointers** like `Rc`, `Arc`, or `Mutex` for safe memory management.
- **Use atomic types** from the `std::sync` module for thread-safe operations without locks.

Refactoring Dora's APIs to comply with these traits will guarantee that **all shared data structures** are safe to use in **concurrent workflows**, without the need for manual synchronization.

### 1.8.3 Promote Immutable Data Sharing

Wherever possible, Dora's APIs should promote the use of **immutable data structures**. Immutable data can be safely shared across threads without requiring locks, reducing **synchronization overhead** and improving performance.

Rust's **ownership and borrowing model** makes it easy to enforce immutability at compile time. By defaulting to **immutable APIs**, developers are encouraged to write **lock-free, thread-safe code**.

### 1.8.4 Forbid Unsafe Code with `#![forbid(unsafe_code)]`

To enforce a **zero-unsafe policy** across Dora, we recommend adding the `#![forbid(unsafe_code)]` directive at the **crate level**. This ensures that no unsafe code can be introduced, either intentionally or accidentally.

### 1.8.5 Update Documentation

Comprehensive documentation is essential to help developers understand the **thread safety guarantees** provided by Dora's APIs. The documentation should clearly indicate:

- **Which APIs are thread-safe.**
- **How shared resources are safely handled.**
- **Best practices for building concurrent applications with Dora.**

### 1.8.6 Potential Impact

**1.8.6.1 Improved Safety and Reliability** By eliminating unsafe code, Dora will offer **strong memory safety guarantees**. This reduces the risk of **critical vulnerabilities** such as **data races**, **null pointer dereferences**, and **buffer overflows**, making Dora more suitable for **production-grade applications**.

**1.8.6.2 Enhanced Performance Through Concurrency**   Ensuring that APIs are **thread-safe** will unlock the potential for **parallel execution**. Developers can safely create **multi-threaded workflows**, improving **throughput** and **reducing latency** in performance-critical applications.

Use cases that will benefit include:

- **Real-time data processing** (e.g., HFT trading systems).
- **IoT pipelines** with millions of concurrent devices.
- **Machine learning workflows** requiring parallel computation.

**1.8.6.3 Increased Developer Confidence**   Removing unsafe code will **reduce cognitive load** for developers using Dora, allowing them to focus on building their applications without worrying about **low-level safety concerns**. This will make Dora more appealing to **Rust developers** who value safety guarantees and are looking for frameworks that align with **Rust's best practices**.

**1.8.6.4 Broader Adoption in Mission-Critical Systems**   By ensuring **100% safe Rust**, Dora will become a more attractive choice for industries that require **high reliability** and **predictable performance**, such as:

- **Finance** (e.g., risk analysis, trading systems).
- **Aerospace and defense** (e.g., control systems, simulations).
- **Healthcare** (e.g., medical devices, data analysis).
- **Automotive** (e.g., autonomous systems).

Making all Dora APIs **Rust safe** is a critical step toward ensuring **memory safety**, **thread safety**, and **concurrency guarantees**. By enforcing **safe Rust practices** and eliminating `unsafe` code, Dora will offer a **robust, reliable, and high-performance framework** for building **modern, concurrent applications** in **mission-critical domains**.

## 1.9 Conclusion

The proposed enhancements to Dora aim to transform it into a **more robust, flexible, and developer-friendly framework** for real-time dataflows, positioning it as a top choice for building **high-performance distributed systems**. By introducing **real-time monitoring through a Terminal User Interface**, developers will gain **deep visibility into running dataflows**, enabling them to identify bottlenecks, debug issues, and optimize performance in real time. This capability will reduce downtime, improve reliability, and make Dora more practical for **production environments** where system responsiveness is critical. Combined with **improved documentation**, including step-by-step guides and best practices for both beginners and advanced users, Dora will significantly lower the barrier to entry for new users and encourage the **growth of an engaged developer community**.

Another key enhancement is the **support for multiple in-memory data formats**, which will make Dora more adaptable to a variety of workloads. While the framework's current reliance on Apache Arrow provides strong foundations for real-time

analytics, introducing alternatives like **Cap'n Proto**, **FlatBuffers**, and **Avro** will allow developers to **tailor dataflows for specific performance needs**. For instance, low-latency applications can benefit from Cap'n Proto's near-zero serialization overhead, while resource-constrained devices may prefer the lightweight structure of FlatBuffers. By allowing users to configure the data format at both the **global** and **node levels**, Dora will offer remarkable flexibility in **balancing memory usage, speed, and scalability**, making it suitable for use cases ranging from **IoT pipelines** to **big data workflows**.

Lastly, ensuring that **all APIs are implemented using 100% safe Rust** will make Dora a **reliable choice for building concurrent, multi-threaded applications**. By forbidding unsafe code and enforcing Rust's **ownership and type safety principles**, developers will be able to build **memory-safe and thread-safe workflows** without worrying about data races or unexpected behavior in distributed systems. Emphasizing the use of `Send` **and** `Sync` **traits** and leveraging **immutable data structures** will further reduce synchronization overhead and improve performance in parallel workloads. This focus on thread safety will **increase developer confidence**, particularly for those building **mission-critical applications** in industries like **finance**, **machine learning**, and **real-time analytics**. By addressing key technical and user experience challenges, Dora will evolve into a **future-proof framework** that meets the needs of **modern, high-throughput systems** across a wide range of industries.