# Developing (experimental) Rust SDK and a Beam engine for IoT devices

Sho Nakatani
TOYOTA Motor Corporation
@laysakura (GitHub)

BEAM SUMMIT
NYC 2023
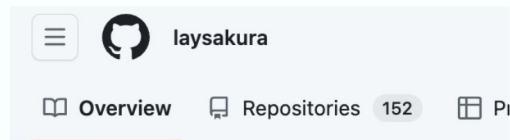
# Goal

- About Beam Rust SDK: Make it the 5th Beam SDK
  - Sharing the <u>motivation</u> behind its development
  - Presenting the <u>current status</u> of the project
  - Encouraging collaboration and <u>gathering contributors</u>
- About SpringQL:
  - Providing a <u>brief overview of SpringQL</u>, a stream processor specifically designed for IoT devices

## About Me

- Research and Development in <u>stream processing</u> for cloud and <u>IoT devices</u>
  - Implementing SpringQL in Rust <u>(GitHub repo)</u>
- Recognizing Beam as <u>standard stream processing model</u> for the next 10 years
  - Desire to <u>support the Beam model for SpringQL</u>
- Active involvement in the development of Beam Rust SDK since February 2023

laysakura

Overview | Repositories 152 | P

**Sho Nakatani**
laysakura

A low-level system developer / backend engineer in Tokyo.

- Rust SDK Development **(17 minutes)**
    - Motivation
    - Design
    - Rust-specific challenges
    - History and future prospects
- Introduction to SpringQL & Integration with Beam **(3 minutes)**

# Rust SDK: <u>Motivation</u>

- For <u>Pipeline Construction</u> (or <u>Programming</u>)
  - Leveraging Rust's <u>statically-typed</u> nature and <u>generics</u>
  - Meeting the demand from Rustaceans for a dedicated Beam Rust SDK
- For <u>Worker</u>
  - <u>Memory safety</u>
  - Performant
    - Comparing to Go: More lightweight runtimes (e.g. no garbage collection)
    - (My interest) High performance single-node SPEs with Beam model?
      - Relevant Research: Scabbard, SABER/LightSaber, StreamBox
      - <u>"Do We Need Distributed Stream Processing?"</u> (blog post)
        - *"a <u>single multicore server</u> can provide <u>better throughput</u> than a multi-node cluster for many streaming applications"*
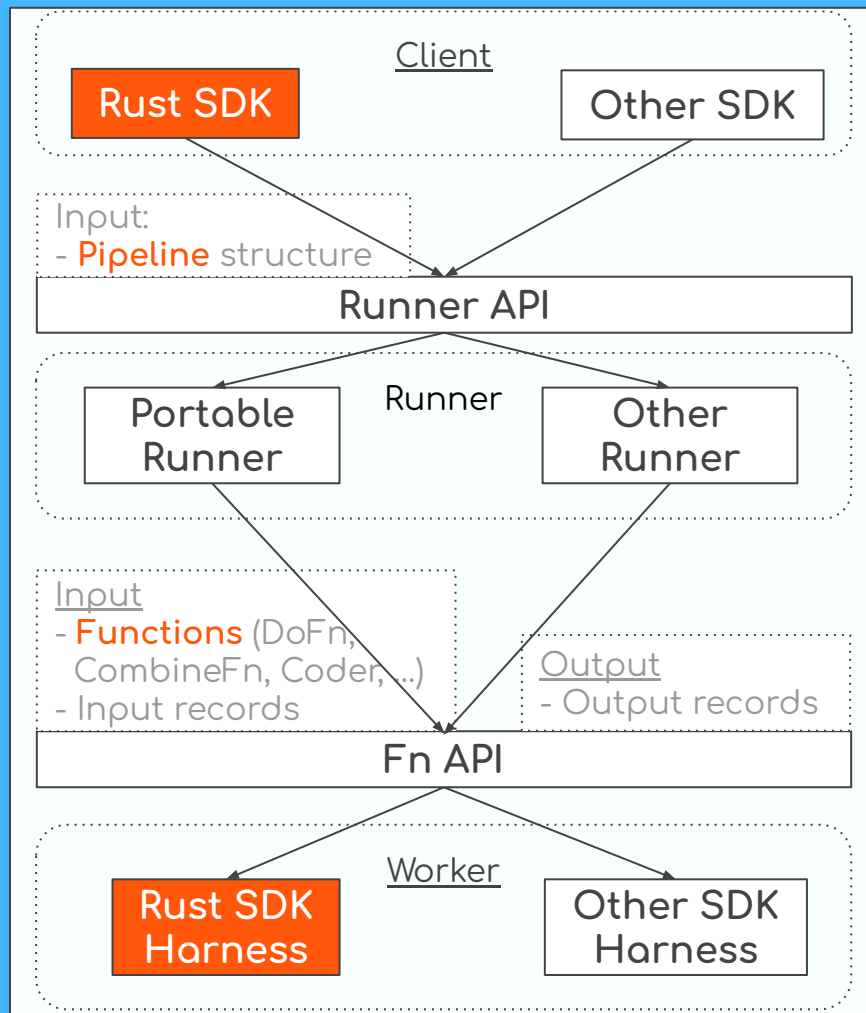
| | Spark | Flink | SABER | Handwritten C++ |
|---|---|---|---|---|
| **Throughput** (million tuples/sec) | 2 | 4.8 | 11.8 | 23 |

**Table 1**: Single CPU core throughput for Yahoo Streaming Benchmark

# Rust SDK: [Design](#)

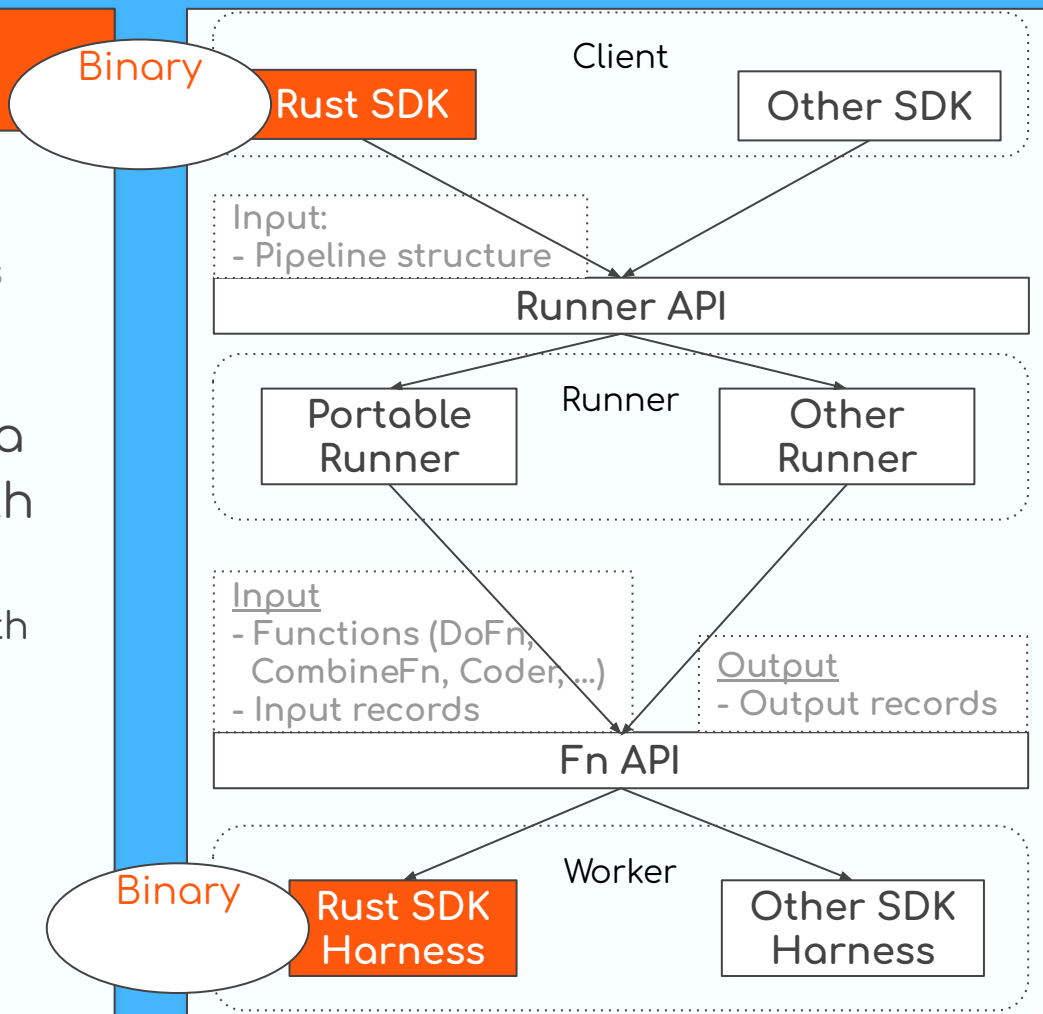# Where Rust SDK Works

- Rust SDK works in:
  - <u>Client</u> to construct pipelines
  - <u>Workers</u> to execute Rust-specific functions
- An application is built as a binary statically linked with the Rust SDK
  - Binaries are deployed to both Client and Workers
  - Different binaries are built from the same app (source)

# Where Rust SDK Works

- Rust SDK works in:
  - <u>Client</u> to construct pipelines
  - <u>Workers</u> to execute Rust-specific functions
- An application is built as a binary statically linked with the Rust SDK
  - Binaries are deployed to both Client and Workers
  - Different binaries are built from the same app (source)

**Binary**

Client

**Rust SDK**   Other SDK

Input:
- Pipeline structure

Runner API

Runner

Portable Runner   Other Runner

Input
- Functions (DoFn, CombineFn, Coder, ...)
- Input records

<u>Output</u>
- Output records

Fn API

Worker

**Binary**

**Rust SDK Harness**   Other SDK Harness

# Where Rust SDK Works

- Rust SDK works in:
  - Client to construct pipelines
  - Workers to execute Rust-specific functions
- An application is built into binary statically linked with Rust SDK
  - Binaries are deployed to both Client and Workers
  - Different binaries are built from the same app (source)

**Binary** aarch64 / macOS

**Client**

Rust SDK

Other SDK

Input:
- Pipeline structure

Runner API

**Runner**

Portable Runner

Other Runner

Input
- Functions (DoFn, CombineFn, Coder, ...)
- Input records

Output
- Output records

Fn API

**Binary** x86-64 / Linux

**Worker**

Rust SDK Harness

Other SDK Harness

# Design Concepts

- Mainly influenced by TypeScript (features) and Go (compilation & deployment)
- Statically-typed pipeline construction
- Removal of Pipeline APIs (explained later)
- Asynchronous execution of workers

Note: The design concepts may require further synchronization with other contributors.

- Mainly influenced by <u>TypeScript (features)</u> and <u>Go (compilation & deployment)</u>
- <u>Statically-typed pipeline</u> construction
- <u>Removal of Pipeline APIs</u>
- <u>Asynchronous</u> execution of <u>workers</u>

Note: The design concepts may require further synchronization with other contributors.

## Show the concepts via a word-count pipeline

## Word-count pipeline...

```rust
fn word_count(lines: PCollection<String>) -> PCollection<KV<String, i32>> {
    lines
        .apply(ParDo::from_map(
            // convert lines to lowercase
            |line| line.to_lowercase(),
        ))
        .apply(ParDo::from_flat_map(|line| {
            // split a line into words
            line.split_whitespace()
        }))
        .apply(ParDo::from_map(|word| {
            // count each word
            KV::new(word, 1)
        }))
        .apply(GroupByKey::default())
        .apply(Combine::per_key(|values| values.count()))
}
```

## and its usage from DirectRunner

```rust
#[tokio::test]
async fn main() {
    DirectRunner::new()
        .run(|root| {
            let lines = root.apply(Create::new(vec![
                "And God said, Let there be light: and there was light",
            ]));
            let result = word_count(lines);

            result.apply(AssertEqualUnordered::new(&[
                KV::new("and".to_string(), 2),
                KV::new("god".to_string(), 1),
                KV::new("said".to_string(), 1),
                KV::new("let".to_string(), 1),
                KV::new("there".to_string(), 2),
                KV::new("be".to_string(), 1),
                KV::new("light".to_string(), 2),
                KV::new("was".to_string(), 1),
            ]))
        })
        .await;
}
```

# Word-count pipeline…

```rust
fn word_count(lines: PCollection<String>) -> PCollection<KV<String, i32>> {
    lines
        .apply(ParDo::from_map(
            // convert lines to lowercase
            |line| line.to_lowercase(),
        ))
        .apply(ParDo::from_flat_map(|line| {
            // split a line into words
            line.split_whitespace()
        }))
        .apply(ParDo::from_map(|word| {
            // count each word
            KV::new(word, 1)
        }))
        .apply(GroupByKey::default())
        .apply(Combine::per_key(|values| values.count()))
}
```

# and its usage from DirectRunner

<u>Statically-typed</u> (w/ automatic type-inference)

line: **String**

line.split_whitespace(): **Vec<String>**
→ flat-mapped into **String**

word: **String**

(output PCollection): **KV<String, i32>**

(output PCollection): **KV<String, Vec<i32>>**

(output PCollection): **KV<String, i32>**

# Word-count pipeline...

```rust
fn word_count(lines: PCollection<String>) -> PCollection<KV<String, i32>> {
    lines
        .apply(ParDo::from_map(
            // convert lines to lowercase
            |line| line.to_lowercase(),
        ))
        .apply(ParDo::from_flat_map(|line| {
            // split a line into words
            line.split_whitespace()
        }))
        .apply(ParDo::from_map(|word| {
            // count each word
            KV::new(word, 1)
        }))
        .apply(GroupByKey::default())
        .apply(Combine::per_key(|values| values.count()))
}
```

# and its usage from DirectRunner

Statically-typed (w/ generics)

fn from_map<F, In, Out>(func: F) -> **ParDo**
  where
    F: Fn(&In) -> Out,
    In: ElemType, Out: ElemType

fn from_flat_map<F, In, Out>(func: F) -> **ParDo**
  where
    F: Fn(&In) -> Vec<Out>,
    In: ElemType, Out: ElemType

fn per_key<F, In, Out>(func: F) -> **Combine**
  where
    F: Fn(&In) -> Vec<Out>,
    In: ElemType, Out: ElemType

# Word-count pipeline...

## Runner.run() instead of Pipeline.run()

- Same API as TypeScript SDK.

  Runner.run() introduce pipeline root (PValue)

- Proposed in a [design doc](#).

**Simplifying Apache Beam**
*or* Pipelines Considered Harmful

https://s.apache.org/no-beam-pipeline

Robert Bradshaw (robertwb@google.com)

# and its usage from DirectRunner

```rust
#[tokio::test]
async fn main() {
    DirectRunner::new()
        .run(|root| {
            let lines = root.apply(Create::new(vec![
                "And God said, Let there be light: and there was light",
            ]));
            let result = word_count(lines);

            result.apply(AssertEqualUnordered::new(&[
                KV::new("and".to_string(), 2),
                KV::new("god".to_string(), 1),
                KV::new("said".to_string(), 1),
                KV::new("let".to_string(), 1),
                KV::new("there".to_string(), 2),
                KV::new("be".to_string(), 1),
                KV::new("light".to_string(), 2),
                KV::new("was".to_string(), 1),
            ]))
        })
        .await;
}
```

# Rust SDK:
# Rust-specific Challenges

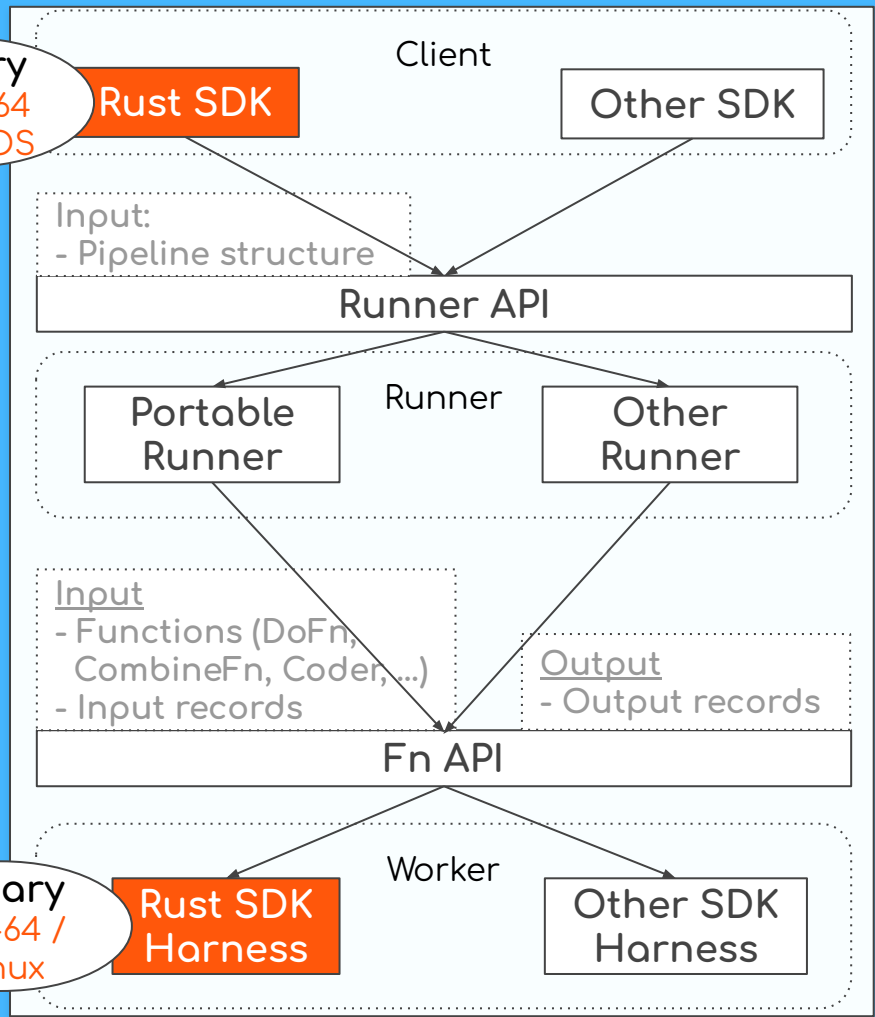# How to share functions?
(between client and worker)

- <u>Functions</u> (and <u>closures</u>)
  - User-defined ParDo, CombineFn, Coder, ...
- Both binaries contain the same functions, but how does a <u>worker</u> determine <u>which functions to execute</u>?

**Binary** aarch64 / macOS

**Client**

Rust SDK

Other SDK

Input:
- Pipeline structure

Runner API

**Runner**

Portable Runner

Other Runner

<u>Input</u>
- **Functions (DoFn, CombineFn, Coder, ...)**
- Input records

<u>Output</u>
- Output records

Fn API

**Binary** x86-64 / Linux

**Worker**

Rust SDK Harness

Other SDK Harness

# How to share functions?
(between client and worker)

- <u>Functions</u> (and <u>closures</u>)
  - User-defined ParDo, CombineFn, Coder, …
- Both binaries contain the same functions, but how does a <u>worker</u> determine <u>which functions to execute</u>?
  - From Fn API, worker receives:

```
message FunctionSpec {

    // (Required) A URN that describes the accompanying payload.
    // For any URN that is not recognized (by whomever is inspecting
    // it) the parameter payload should be treated as opaque and
    // passed as-is.
    string urn = 1;

    // (Optional) The data specifying any parameters to the URN. If
    // the URN does not require any arguments, this may be omitted.
    bytes payload = 3;

}
```
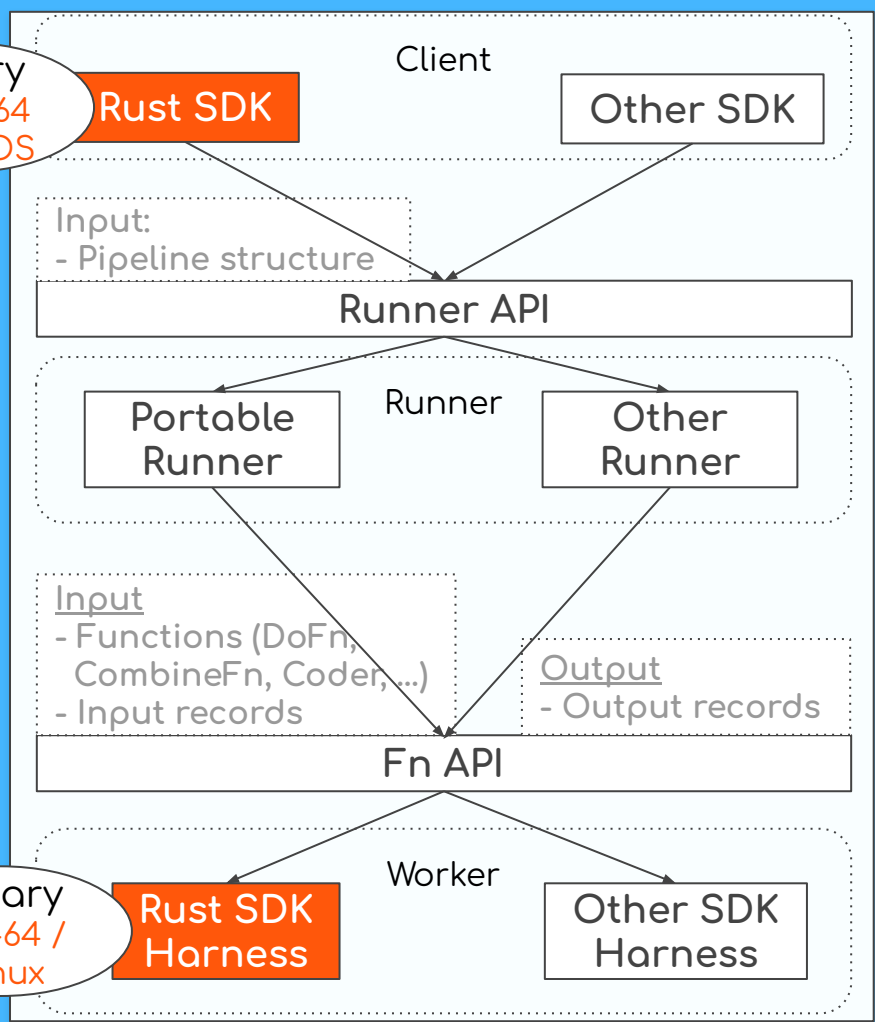
**Binary** aarch64 / macOS

Client

**Rust SDK**     Other SDK

Input:
- Pipeline structure

Runner API

Runner

Portable Runner     Other Runner

Input
- Functions (DoFn, CombineFn, Coder, …)
- Input records

Output
- Output records

Fn API

**Binary** x86-64 / Linux

Worker

**Rust SDK Harness**     Other SDK Harness

# How to share functions?
(between client and worker)

- <u>How does a worker decide which function to execute?</u>
- Deserialize function body from payload?
  - **Cannot serialize functions** in Rust (especially for generic ones).
    - See [discussion in a design doc](#) for detail

```
message FunctionSpec {

  // (Required) A URN that describes the accompanying payload.
  // For any URN that is not recognized (by whomever is inspecting
  // it) the parameter payload should be treated as opaque and
  // passed as-is.
  string urn = 1;

  // (Optional) The data specifying any parameters to the URN. If
  // the URN does not require any arguments, this may be omitted.
  bytes payload = 3;

}
```

**Binary** aarch64 / macOS

**Rust SDK**

Client

**Other SDK**

Input:
- Pipeline structure

**Runner API**

**Portable Runner**   Runner   **Other Runner**

Input
- Functions (DoFn, CombineFn, Coder, ...)
- Input records

Output
- Output records

**Fn API**

**Binary** x86-64 / Linux

Worker

**Rust SDK Harness**

**Other SDK Harness**

# How to share functions?
(between client and worker)

- **How does a worker decide which function to execute?**
- Function symbols in URN?
  - **No reflection** in Rust (cannot call function from its symbol)
  - Closures are unnamed
  - Different from Go SDK

```
message FunctionSpec {

  // (Required) A URN that describes the accompanying payload.
  // For any URN that is not recognized (by whomever is inspecting
  // it) the parameter payload should be treated as opaque and
  // passed as-is.
  string urn = 1;

  // (Optional) The data specifying any parameters to the URN. If
  // the URN does not require any arguments, this may be omitted.
  bytes payload = 3;
}
```

**Binary** aarch64 / macOS

**Rust SDK**

Client

Other SDK

Input:
- Pipeline structure

Runner API

Portable Runner

Runner

Other Runner

**Input**
- Functions (DoFn, CombineFn, Coder, ...)
- Input records

**Output**
- Output records

Fn API

**Binary** x86-64 / Linux

**Rust SDK Harness**

Worker

Other SDK Harness

# How to share functions?
(between client and worker)

- **How does a worker decide which function to execute?**
- Regis
  "URN
  - In
    m
    - *Note that function pointers differ in Client and Worker (different binary)*
  - Requires macro and further implementation efforts, but seems not a bad idea

We are currently working on the development of safe serialization for functions.

**Binary**
aarch64 / macOS

Rust SDK

Client

Other SDK

Input:
- Pipeline structure

Input
- Functions (DoFn, CombineFn, Coder, ...)
- Input records

Output
- Output records

Fn API

**Binary**
x86-64 / Linux

Rust SDK Harness

Worker

Other SDK Harness

# Rust SDK:
# Development history and future

# Why History?

- While I currently serve as the repository owner of the experimental Beam Rust SDK, I am not the project's original contributor.
- It is important for me to <u>acknowledge and honor the contributions of past and current individuals</u> involved in the project.

I apologize if I have unintentionally omitted mentioning any specific contributor names.

# Started from a JIRA Ticket

- The Rust SDK issue was created in July 2021 on JIRA
- There was a recommendation to learn from the TypeScript SDK
- An initial concept of pipeline construction was shared in a Gist
- Contributor
  - jayendra13
- Advisers
  - kennknowles
  - robertwb
  - lostluck



Beam / BEAM-12658
## Rust SDK

Start Progress   Resolve issue   Need more information                    Export ▾

▾ Details
| | | | |
|---|---|---|---|
| Type: | ⬆ Wish | Status: | OPEN |
| Priority: | ⩔ P3 | Resolution: | Unresolved |
| Affects Version/s: | None | Fix Version/s: | None |
| Component/s: | sdk-ideas | | |
| Labels: | None | | |
| Language: | Rust | | |

▾ People
Assignee:
? Unassigned

Reporter:

Votes:
1 Vote for this issue

Watchers:
8 Start watching this issue

▾ Description
It would be great to have Rust SDK in order to create very high-performant yet safe pipelines.

GitHub Gist    Search...    All gists   Back to GitHub

jayendra13 / **rust_beam.rs**                    Subscribe   ☆ Star  0   Fork  0
Created last year · Report abuse

‹› Code    Revisions 1                    Embed ▾  <script src="https://g   Download ZIP

‹› **rust_beam.rs**                                                    Raw

```
 1  #[derive(Debug)]
 2  struct Pipeline {
 3      context: String,
 4  }
 5
 6  impl Pipeline {
 7    fn new(context: String) -> Pipeline {
 8      Pipeline {
 9        context: context,
10      }
11    }
12
13    fn apply_transform(&self, transform: &impl PTransform, input: &impl PValue) -> PCollection {
14      transform.expand(input)
15    }
16  }
17
18  trait PTransform {
19    fn expand(&self, input: &impl PValue) -> PCollection;
20  }
21
```

# Issue Migrated to GitHub

- The issue on GitHub is still active to this day
- Experimental implementation repos:
  - kennknowles/beam [old]
  - ↓ (merged into)
  - nivaldoh/beam [old]
  - ↓ (forked to)
  - laysakura/beam [current]
- Organizer: brucearctor

## [Old repo] kennknowles/beam

- Project initiation: January 2023
- The Google Cloud Dataflow team started a Rust SDK development
- Later merged into nivaldoh/beam repository
- Contributors
  - antonbobkov
  - robertwb
  - JayDosunmu
  - y1chi

robertwb commented on Jan 7 · edited ▾          Contributor   ···

I *just* saw this, there's actually an effort to build a Rust SDK this week from the Dataflow team. What we have is at https://github.com/kennknowles/beam/tree/rust/sdks/rust ; it would be great to combine efforts. Though that one looks much further along.

robertwb commented on Jan 7          Contributor   ···

IMHO, @nivaldoh's repo is further along, and better structured, so I think it makes sense to start there. In the next day or two we'll probably be pushing willy-nilly to the one at kennknowles, in the spirit of the hackathon to explore ideas, but next week I suggest we start creating pull requests to https://github.com/nivaldoh/beam/tree/rust_sdk to carry anything over that has value (and isn't already in the latter) and continue there.

# [Old repo] nivaldoh/beam

- Project initiation: November 2022
- Added:
  - Codes for pipeline construction (partial)
  - Worker codes (partial)
- Development activities ceased since February 2023
- Contributors
  - nivaldoh
  - sjvanrossum
  - laysakura (me)
  - Miuler



nivaldoh commented on Nov 2, 2022 · edited ▾                    Contributor  ...

Hi, I would like to express interest in working on the Rust SDK. I'll create an incubator fork soon.

☺  👍 4   🎉 3

nivaldoh commented on Nov 2, 2022                              Contributor  ...

.take-issue

☺  👍 2

🔘 github-actions bot assigned nivaldoh on Nov 2, 2022

nivaldoh commented on Nov 12, 2022                            Contributor  ...

Work is underway here. Progress may be slow, and early code will look quite rough. I'll be really happy to receive any feedback or collaboration opportunities.

☺  👍 6   ❤️ 2

# [Current repo] laysakura/beam

- Project initiation: February 2023
- Forked from nivaldoh/beam
- Added:
  - Coder serialization (partial)
  - More worker codes (partial)
  - General function serialization (doing)
  - The Beam Programming Guide for Rust (doing)
- Contributors
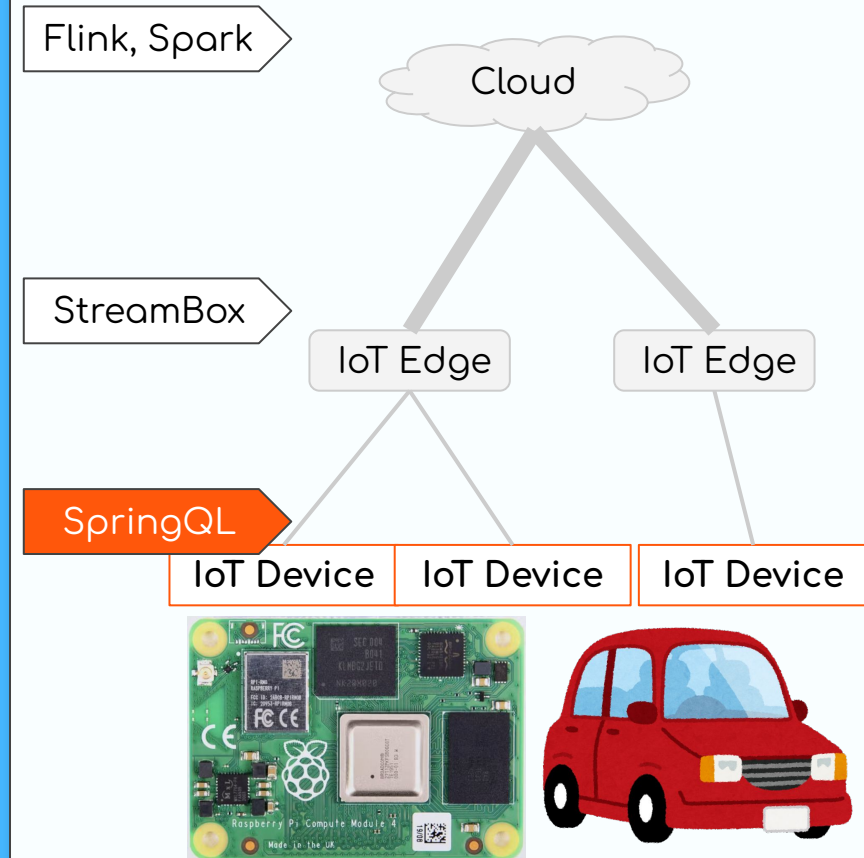  - dahlbaek
  - sjvanrossum
  - Kelvinyu1117
  - laysakura (me)

- Technically challenging implementations
  - Serialization/deserialization of functions (including closures), led by sjvanrossum
- Align design considerations for non-trivial features
  - Registration of user-defined objects (possibly through init function w/ macros)
  - Coders (custom coders, row coders, etc.)
  - Artifact staging service
- Completion of the Programming Guide and working examples
- **Call for more contributors!**
  - Will create good-first issues in laysakura/beam

# SpringQL: Introduction and integration with Beam

# SpringQL's Target

- ## Stream Processing Engine for <u>IoT devices</u>
  - Targeting middle-to-high end devices
    - Raspberry Pi
    - Connected vehicles
- ## Support semi-realtime stream processing
  - Input:
    - Sensor data
    - UI
  - Output:
    - Device actuation
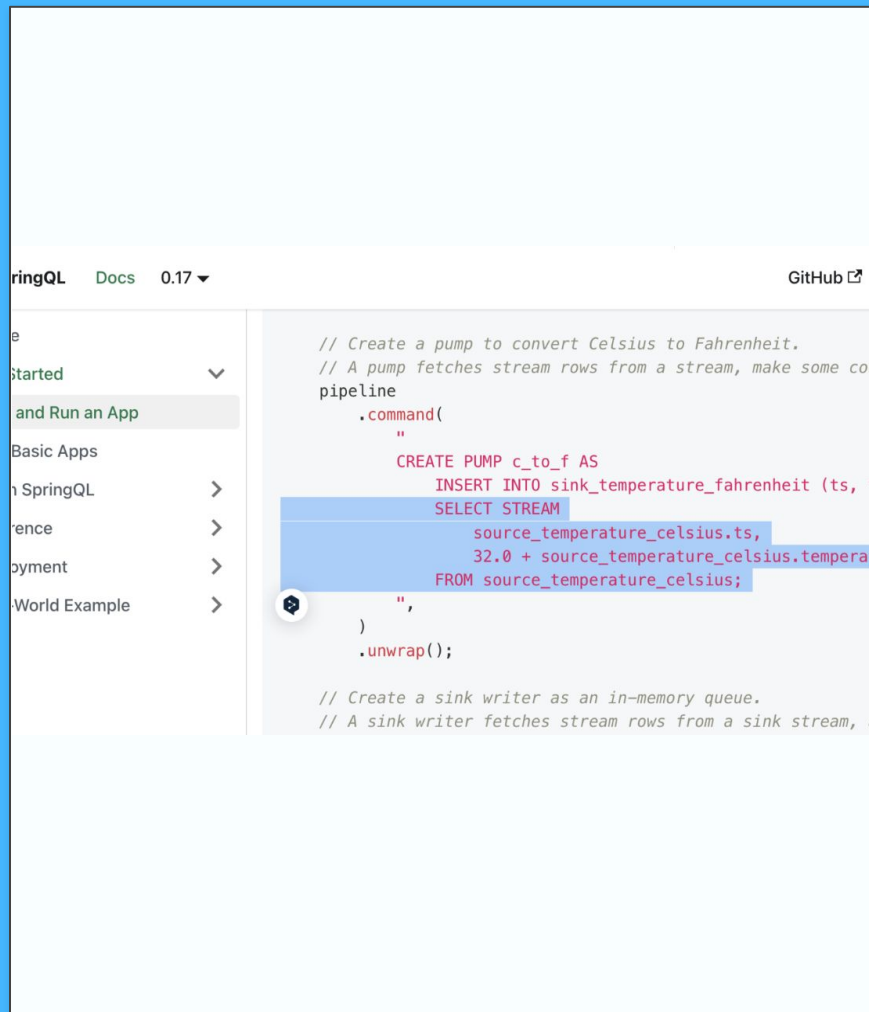    - Aggregated data (sent to edge/cloud)
    - UI (display, sound, …)

Flink, Spark

Cloud

StreamBox

IoT Edge          IoT Edge

SpringQL

IoT Device | IoT Device | IoT Device

https://www.raspberrypi.com/products/compute-module-4/?variant=raspberry-pi-cm4001000

https://www.irasutoya.com/2019/10/blog-post_57.html

# SpringQL's Current Status

- Implemented in Rust ([repo](#))
- Distributed as libraries:
  - Rust (static)
  - C (static / dynamic)
- User interface
  - Client: Rust / C
  - Pipeline construction: SQL-like
  - Operation: Streaming SQL
- Problems
  - Difficulty in constructing DAGs using SQL-like language
  - Limited operations available through streaming SQL
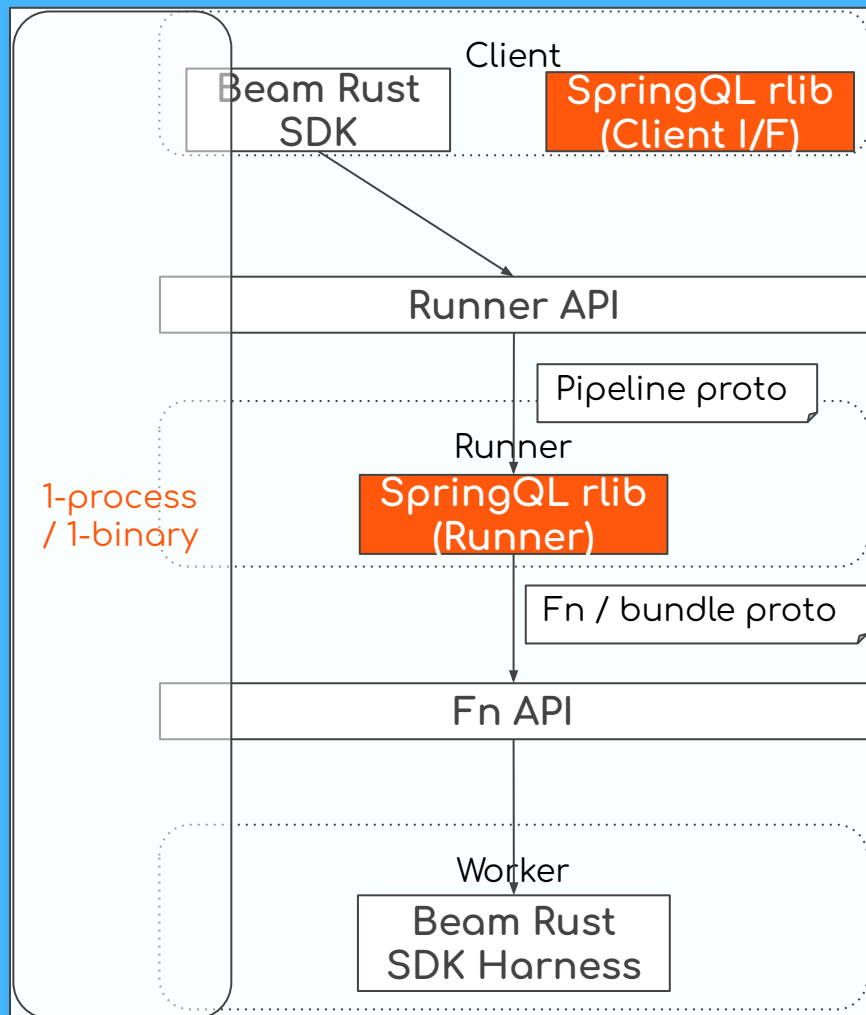
Desire to utilize Beam for U/I



```
ringQL   Docs   0.17 ▾                                    GitHub ⧉

              // Create a pump to convert Celsius to Fahrenheit.
 Started      // A pump fetches stream rows from a stream, make some co
              pipeline
 and Run an App    .command(
Basic Apps            "
                      CREATE PUMP c_to_f AS
 SpringQL         >       INSERT INTO sink_temperature_fahrenheit (ts,
                          SELECT STREAM
rence             >           source_temperature_celsius.ts,
                              32.0 + source_temperature_celsius.tempera
oyment            >       FROM source_temperature_celsius;
                      ",
World Example     >    )
                      .unwrap();

              // Create a sink writer as an in-memory queue.
              // A sink writer fetches stream rows from a sink stream,
```
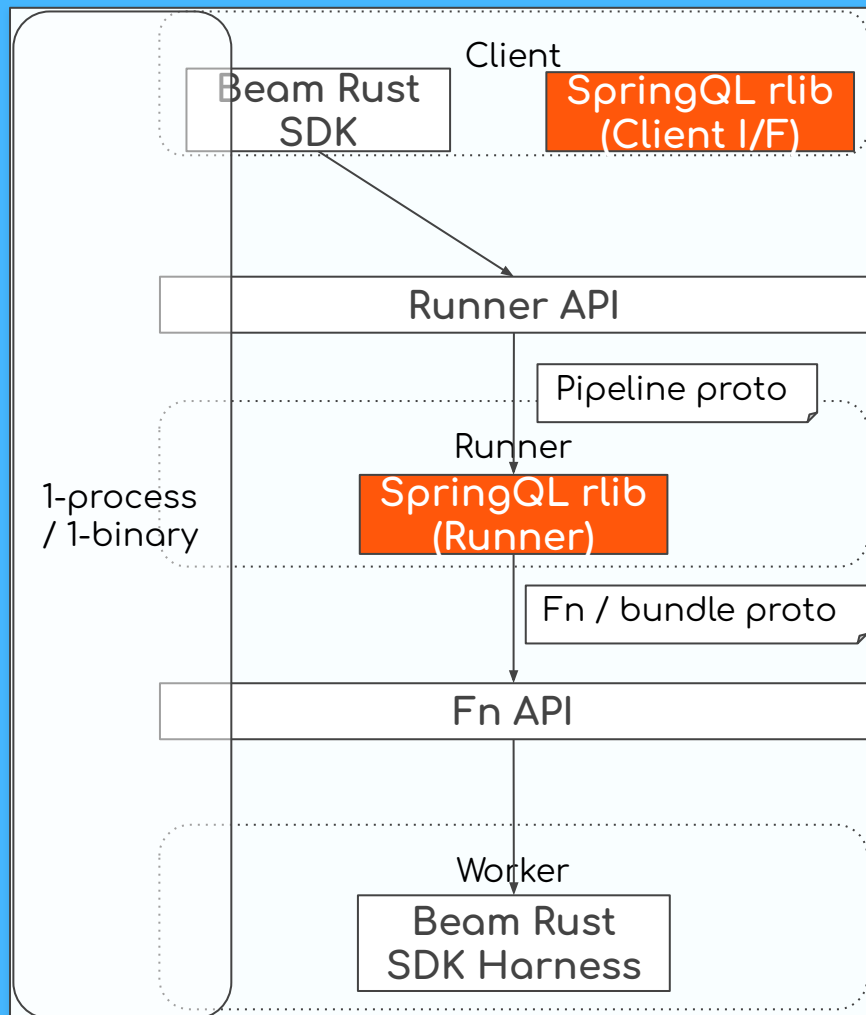
# Initial Idea: Integration with Beam

- App, Beam SDK, and SpringQL library are all within the same process and binary
- SpringQL library serves as:
  - Client interface
  - Dedicated runner
- SpringQL Runner receives pipeline graph via Runner API in protobuf format
- SpringQL runner calls SDK Harness to execute UDFs
  - May use "LOOPBACK" SDK Harness (config doc)



Client
Beam Rust SDK
SpringQL rlib (Client I/F)

Runner API

Pipeline proto

Runner
1-process / 1-binary
SpringQL rlib (Runner)

Fn / bundle proto

Fn API

Worker
Beam Rust SDK Harness

# Initial Idea: Integration with Beam

- App, Beam SDK, and SpringQL library are all within the same process and binary
- SpringQL library serves as:
  - Client interface
  - Dedicated runner
- SpringQL Runner receives pipeline graph via Runner API in protobuf format
- SpringQL runner calls SDK Harness to execute UDFs
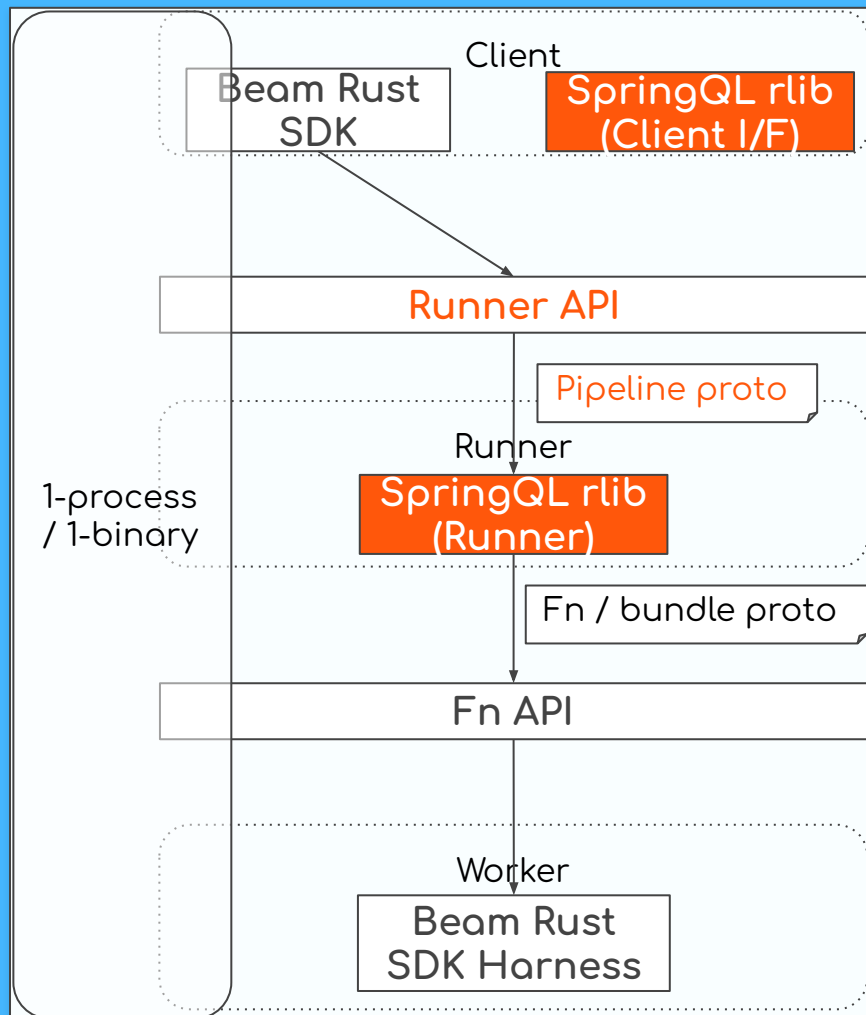  - May use "LOOPBACK" SDK Harness (config doc)

# Initial Idea: Integration with Beam

- App, Beam SDK, and SpringQL library are all within the same process and binary
- SpringQL library serves as:
  - Client interface
  - Dedicated runner
- **SpringQL Runner receives pipeline graph via Runner API in protobuf format**
- SpringQL runner calls SDK Harness to execute UDFs
  - May use "LOOPBACK" SDK Harness (config doc)

Client
Beam Rust SDK
SpringQL rlib (Client I/F)

Runner API

Pipeline proto

1-process / 1-binary

Runner
SpringQL rlib (Runner)

Fn / bundle proto

Fn API

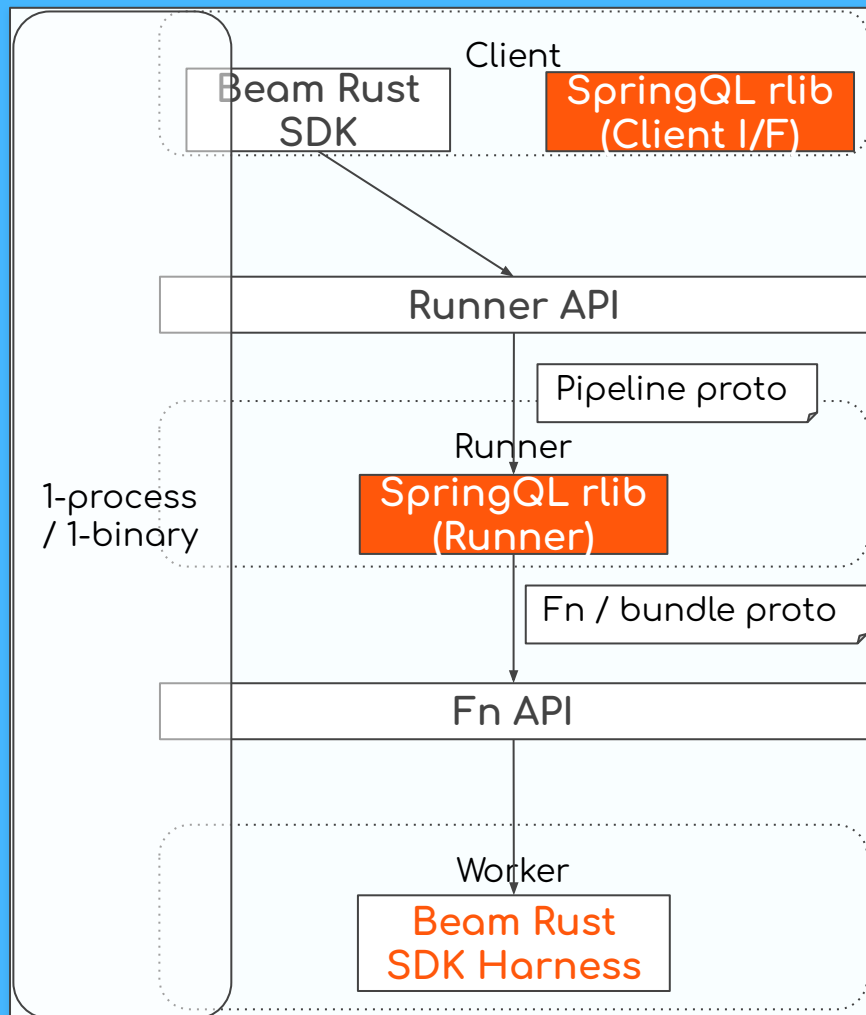Worker
Beam Rust SDK Harness

# Initial Idea: Integration with Beam

- App, Beam SDK, and SpringQL library are all within the same process and binary
- SpringQL library serves as:
  - Client interface
  - Dedicated runner
- SpringQL Runner receives pipeline graph via Runner API in protobuf format
- SpringQL runner calls SDK Harness to execute UDFs
  - May use "LOOPBACK" SDK Harness (config doc)

Client
Beam Rust SDK
SpringQL rlib (Client I/F)

Runner API

Pipeline proto

Runner
1-process / 1-binary
SpringQL rlib (Runner)

Fn / bundle proto

Fn API

Worker
Beam Rust SDK Harness

- ## About Beam Rust SDK
  - Motivation behind its development
  - Current status of the project
  - Call for contributions
- ## About SpringQL
  - SpringQL's target systems and architecture
  - Integration idea with Beam