



CS6650 – Vancouver campus

# Evaluating gRPC performance in real time applications

Comparing gRPC against REST APIs in remote real time string manipulation  
tasks.

Spring 2022

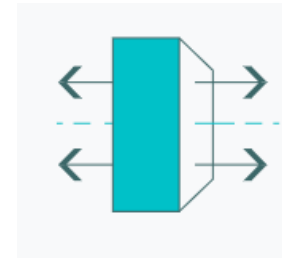
Francisco Proboste

# Hypothesis and research question



The objective of this experiment is to understand to what extent gRPC is a high performant protocol to connect distributed services in real time applications.

1. How does gRPC compares to Rest API in real time string manipulation tasks?
2. How is the distribution of individual-words tasks latency?
3. How useful is bi-directional streaming in this case?



# Tests

A simple remote task requesting invert words will be used to compare the performance of the two mechanisms : API Rest v/s gRPC.

“Speed” → “Deesp”

## Experiment 1:

One by one word inversion.

- 300 words aprox.
- No pause between words.

## Experiment 2:

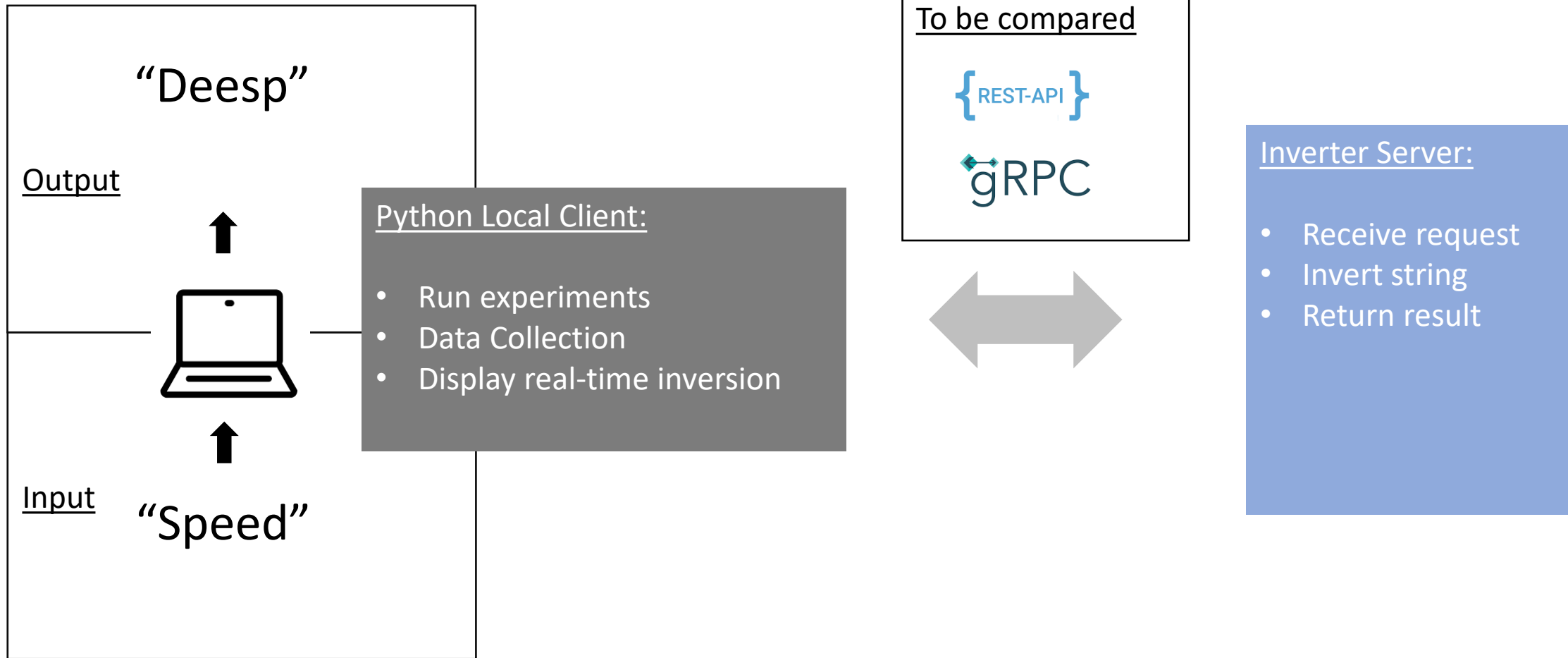
Large text inversion.

- 300 words aprox.
- Test behavior with a larger payload.

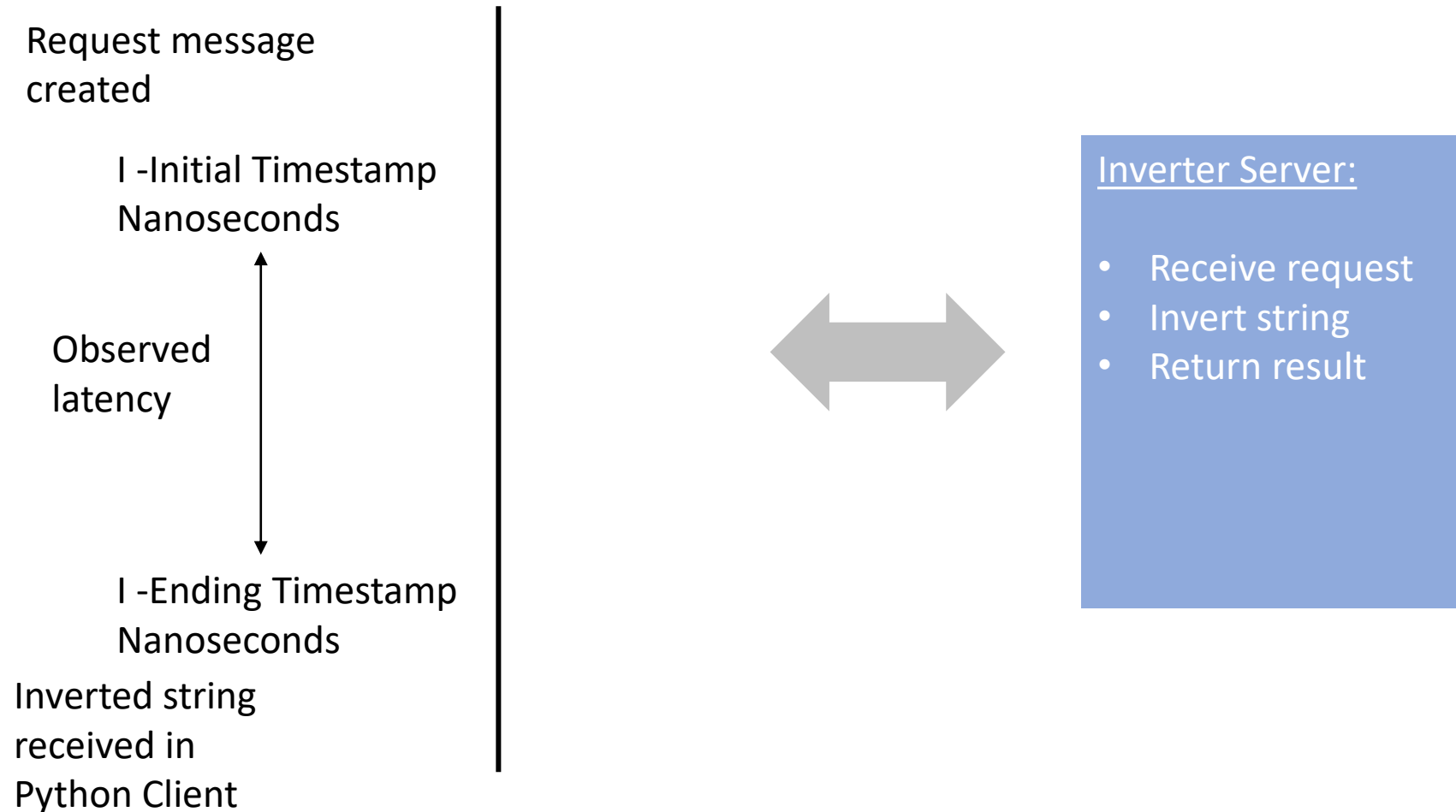
## Experiment 3:

Experiment 1 in a local deployment to understand the differences without transport.

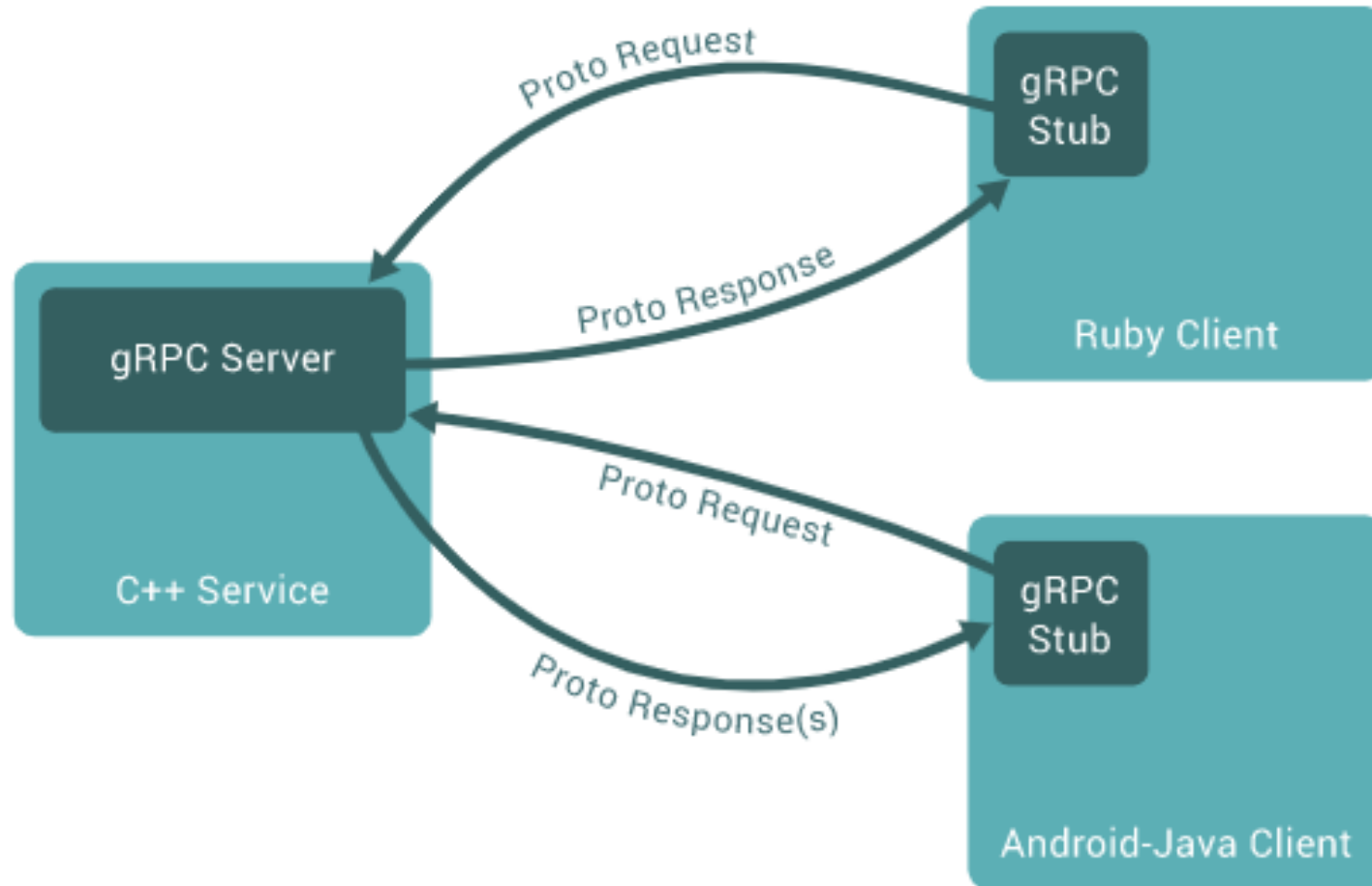
# Architecture Diagram



# For each string to be inverted:



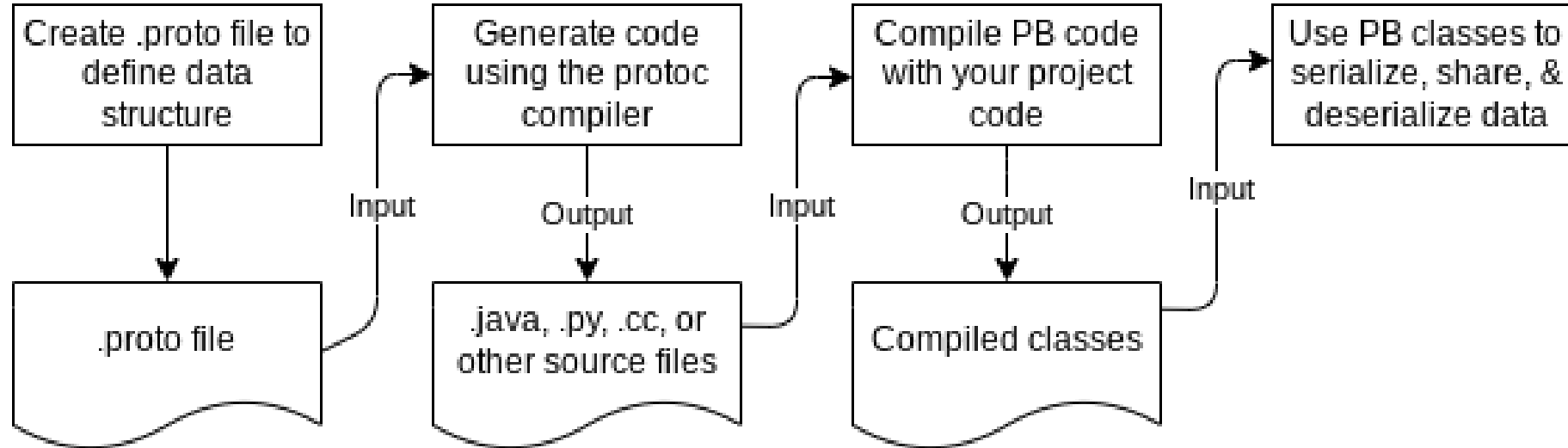
# gRPC implementation



Source: <https://grpc.io/docs/what-is-grpc/introduction/>

# Leverages use of protocol buffers

Multi-language, bidirectional data serialization mechanism.



Source: <https://developers.google.com/protocol-buffers/docs/overview>

# textSender.proto

```
1  syntax = "proto3";  
2  
3  message Word {  
4      string value = 1;  
5  }  
6  
7  service wordInverter {  
8      rpc Invert(Word) returns (Word) {}  
9  }
```



# How the server is instantiated

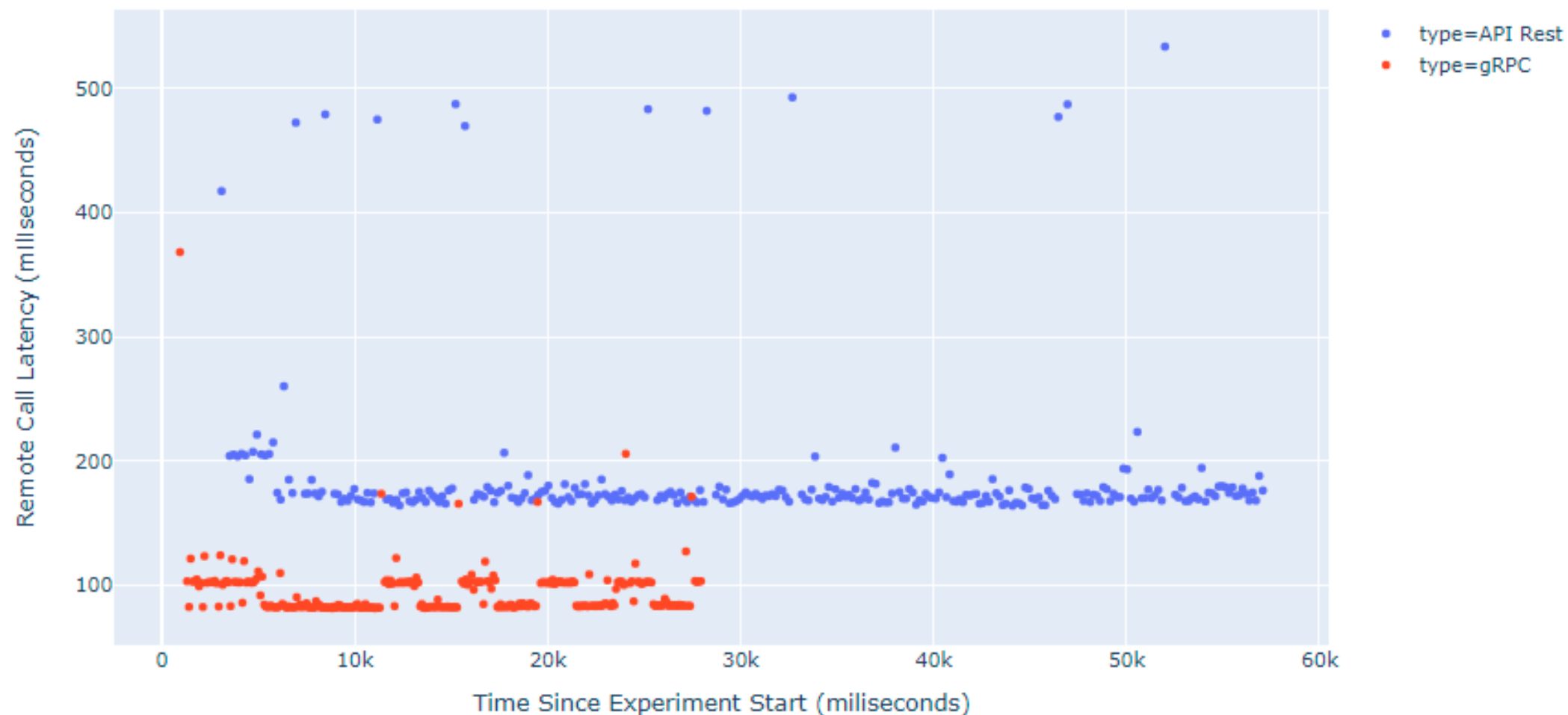
```
40 def serve():  
41     server = grpc.server(futures.ThreadPoolExecutor(max_workers=10))  
42     textSender_pb2_grpc.add_wordInverterServicer_to_server(WordInverterServicer(), server)  
43     server.add_insecure_port('[::]:50051')  
44     server.start()  
45     server.wait_for_termination()
```

# How the client is instantiated

```
6      # Import classes
7      import textSender_pb2
8      import textSender_pb2_grpc
9
10     # Open a GRPC channel
11     channel = grpc.insecure_channel('54.209.72.104:50051')
12
13     # Create the stub from the compiled pb2 file
14     wordStub = textSender_pb2_grpc.wordInverterStub(channel)
15
16     # Create a valid request message (see proto file)
17     word = textSender_pb2.Word(value="hola")
18
19     # Make the call
20     response = wordStub.Invert(word)
```

# Experiment 1

One by one inversion of 300 words aprox.



# Experiment 1

One by one inversion of 300 words aprox.

		API Rest	gRPC
Average	[ms]	187.2	93.7
99 percentile	[ms]	497.2	171.4
SD	[ms]	62.17	22.5
Var. Coef.	[%]	33	24

# Experiment 2

Large payload inversion: 300 words in one call



33

A dagger of the mind, a false creation  
Proceeding from the heat-oppressed brain?  
I see thee yet, in form as palpable  
As this which now I draw. *He draws his dagger.*  
Thou marshal'st me the way that I was going,  
And such an instrument I was to use.  
Mine eyes are made the fools o' th' other senses  
Or else worth all the rest. I see thee still,  
And, on thy blade and dudgeon, gouts of blood,  
Which was not so before. There's no such thing.  
It is the bloody business which informs  
Thus to mine eyes. Now o'er the one-half world  
Nature seems dead, and wicked dreams abuse  
The curtained sleep. Witchcraft celebrates  
Pale Hecate's off'rings, and withered murder,  
Alarumed by his sentinel, the wolf,  
~~Whose howl's his watch, thus with his stealthy pace~~

One Remote Call  
To invert the whole text.



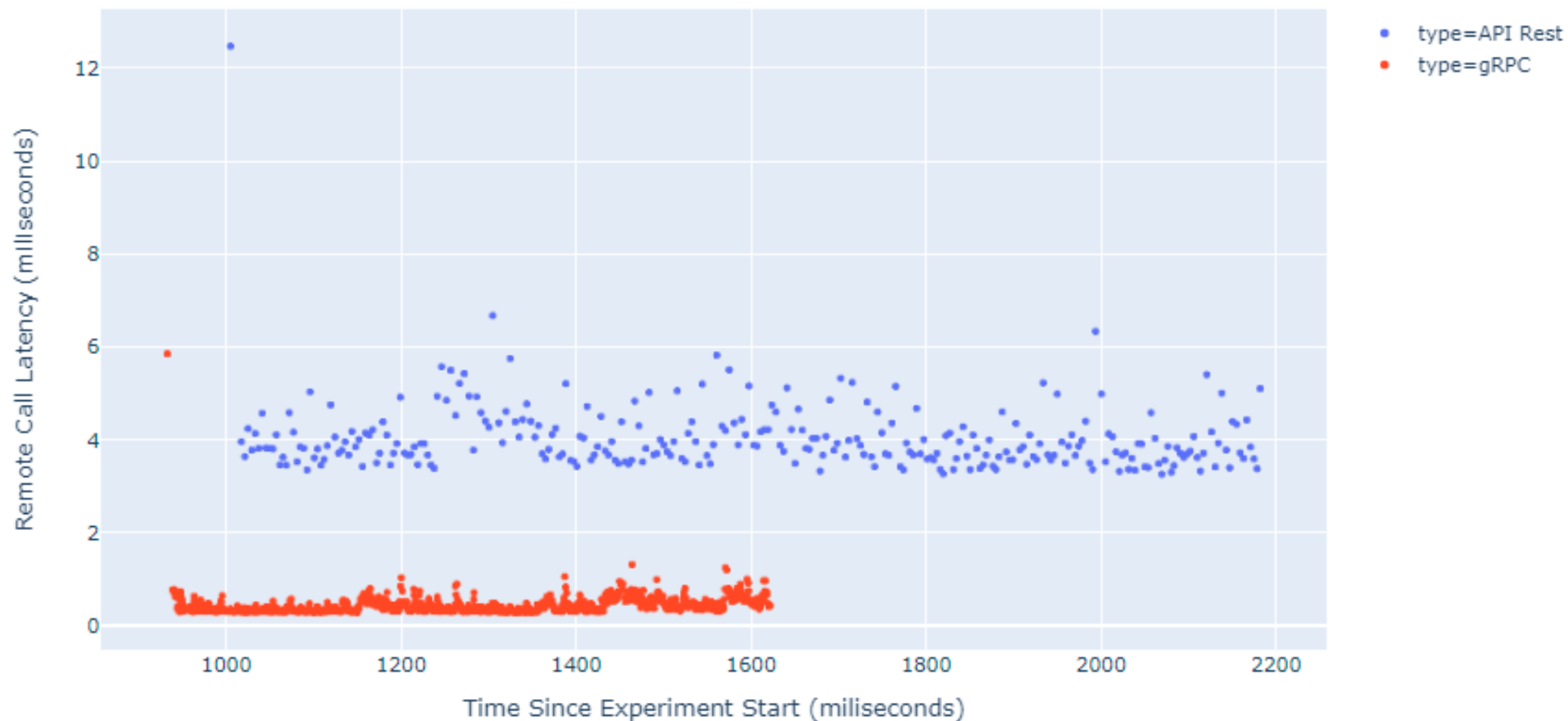
API Rest: 195.5 ms

gRPC: 101.1 ms

# Experiment 3

One by one inversion of 300 words aprox.

No Transport Case



# Experiment 3

One by one inversion of 300 words aprox.

No Transport Case

		API Rest	gRPC
Average	[ms]	4.04	0.41
99 percentile	[ms]	5.89	0.89
SD	[ms]	0.76	0.19
Var. Coef.	[%]	19	46

# Conclusions

1. gRPC is at least two times faster than API Rest in a simple low size load manipulation task on the cloud.
2. gRPC could be up to one order of magnitude faster in short distance applications.
3. Difficult to conclude which shows more variability in speeds.



# Future questions

1. How do they compare when servers are near capacity?
2. How do they compare in term of parallel requests handling?
3. How do they compare in resources requirements? Which is more suitable for IoT communication?

Thanks for your attention.