

# One Billion Rows Challenge in Elixir

From 12 minutes to 25 seconds

Raj Rajhans

Code Beam Europe 2024

# Who am I?

Software Engineer at Invideo AI

Elixir, Rust & Javascript

Organizing RustMumbai



*@\_rajrajhans*

*rajrajhans.com*

# THE ONE BILLION ROWS CHALLENGE



...  
Tokyo;16.8  
Cape Town;22.3  
Stockholm;7.1  
Marrakech;29.6  
**Berlin;9.7**  
Auckland;17.5  
Bangkok;32.7  
Reykjavik;3.4  
Vienna;14.6  
Moscow;1.9  
**Berlin;6.3**  
Hanoi;30.2  
Mexico City;20.8  
Cairo;33.5  
Oslo;5.7  
Kyoto;18.1  
Miami;27.6  
**Berlin;11.2**  
Prague;15.9  
Istanbul;23.4  
Edinburgh;10.8  
Seattle;13.7  
Dubai;36.8  
Queenstown;9.9  
...

...  
Tokyo;16.8  
Cape Town;22.3  
Stockholm;7.1  
Marrakech;29.6  
**Berlin;9.7**  
Auckland;17.5  
Bangkok;32.7  
Reykjavik;3.4  
Vienna;14.6  
Moscow;1.9  
**Berlin;6.3**  
Hanoi;30.2  
Mexico City;20.8  
Cairo;33.5  
Oslo;5.7  
Kyoto;18.1  
Miami;27.6  
**Berlin;11.2**  
Prague;15.9  
Istanbul;23.4  
Edinburgh;10.8  
Seattle;13.7  
Dubai;36.8  
Queenstown;9.9  
...

City;Min;Mean;Max



Auckland;17.5;17.5;17.5  
Bangkok;32.7;32.7;32.7  
Berlin;6.3;9.1;11.2  
Cairo;33.5;33.5;33.5  
Cape Town;22.3;22.3;22.3  
Dubai;36.8;36.8;36.8  
Edinburgh;10.8;10.8;10.8  
Hanoi;30.2;30.2;30.2  
Istanbul;23.4;23.4;23.4  
Kyoto;18.1;18.1;18.1  
Marrakech;29.6;29.6;29.6  
Mexico City;20.8;20.8;20.8  
Miami;27.6;27.6;27.6  
Moscow;1.9;1.9;1.9  
Oslo;5.7;5.7;5.7  
Prague;15.9;15.9;15.9  
Queenstown;9.9;9.9;9.9  
Reykjavik;3.4;3.4;3.4  
Seattle;13.7;13.7;13.7  
Stockholm;7.1;7.1;7.1  
Tokyo;16.8;16.8;16.8  
Vienna;14.6;14.6;14.6

...	...	...	...	...	...	...
Nairobi;28.3	Madrid;24.6	Sydney;23.7	Tokyo;16.8	Rome;23.5	Madrid;22.7	<b>Berlin;9.1</b>
Tokyo;15.7	Hanoi;29.8	Helsinki;4.2	Cape Town;22.3	Lima;19.8	Bangkok;33.8	Barcelona;23.5
<b>Berlin;7.9</b>	Queenstown;11.3	Havana;29.5	Dubai;38.9	Oslo;7.2	Reykjavik;4.1	Tokyo;17.8
Lisbon;19.4	Cairo;33.5	Johannesburg;20.1	Stockholm;7.1	Tokyo;18.4	Mumbai;31.5	Dubai;38.6
Cairo;32.1	<b>Berlin;8.7</b>	Shanghai;18.6	Marrakech;29.6	Delhi;33.6	Cairo;34.2	Oslo;6.2
Vancouver;12.6	Kuala Lumpur;31.2	Vancouver;13.9	<b>Berlin;9.7</b>	<b>Berlin;9.1</b>	Vancouver;14.3	Cape Town;20.9
Mumbai;30.5	Stockholm;6.9	<b>Berlin;7.4</b>	Auckland;17.5	Miami;28.7	Helsinki;6.9	Bangkok;32.7
Stockholm;5.2	Brisbane;25.1	Cairo;32.8	Bangkok;32.7	Cairo;35.2	Rio de Janeiro;28.6	Reykjavik;3.8
Rio de Janeiro;27.8	Nairobi;22.7	Oslo;6.5	Reykjavik;3.4	Bern;12.5	<b>Berlin;8.5</b>	Mexico City;21.3
Moscow;-2.1	Mexico City;19.4	Rio de Janeiro;27.3	Vienna;14.6	Seoul;20.9	Dubai;39.1	Sydney;25.6
<b>Berlin;11.3</b>	Lisbon;20.8	Mumbai;31.6	Moscow;1.9	Riga;8.3	Prague;16.2	Marrakech;29.4
Cape Town;22.6	Anchorage;-1.5	Edinburgh;11.2	<b>Berlin;6.3</b>	Dubai;40.1	Sydney;24.8	Vancouver;13.5
Bangkok;33.9	Singapore;28.9	Bangkok;33.9	Hanoi;30.2	Paris;16.7	Wellington;15.7	<b>Berlin;7.6</b>
Dublin;13.0	Warsaw;13.6	<b>Berlin;10.8</b>	Mexico City;20.8	<b>Berlin;6.8</b>	Marrakech;30.3	Helsinki;5.7
Sydney;24.2	Marrakech;27.3	Reykjavik;2.7	Cairo;33.5	Accra;29.4	Anchorage;-2.1	Istanbul;25.1
Toronto;8.7	<b>Berlin;11.9</b>	Prague;16.3	Oslo;5.7	Quito;15.3	Seoul;19.5	Prague;15.4
Dubai;38.4	Seattle;14.2	San Francisco;17.8	Singapore;28.4	Sofia;14.8	<b>Berlin;10.9</b>	Singapore;30.2
Amsterdam;9.6	Athens;26.5	Wellington;15.4	Lisbon;19.3	Lagos;31.2	Istanbul;25.1	Edinburgh;11.8
Singapore;29.8	Reykjavik;3.8	Marrakech;28.2	Vancouver;12.5	<b>Berlin;11.3</b>	Singapore;29.7	Cairo;33.9
<b>Berlin;6.5</b>	Vienna;15.4	Taipei;25.9	Kyoto;18.1	Porto;19.5	Lisbon;20.6	Lisbon;19.7
Seoul;18.2	Cape Town;23.1	<b>Berlin;5.1</b>	Miami;27.6	Doha;36.8	<b>Berlin;7.3</b>	Kyoto;18.6
Vienna;14.9	Moscow;2.6	Montreal;8.3	<b>Berlin;11.2</b>	Hanoi;30.7	Toronto;11.8	Buenos Aires;24.8
Buenos Aires;26.3	Dubai;38.7	Rome;22.6	Prague;15.9	Minsk;9.6	Athens;26.9	Toronto;12.1
Helsinki;3.8	Kyoto;18.9	Seoul;19.7	Istanbul;23.4	Tunis;25.4	Oslo;5.4	Istanbul;26.5
San Francisco;17.5	<b>Berlin;5.3</b>	Lisbon;21.4	Edinburgh;10.8	Lyon;15.9	San Francisco;17.9	Auckland;16.9
Oslo;6.1	Amsterdam;12.1	Anchorage;-0.8	Seattle;13.7	Surat;32.1	Kyoto;18.3	Moscow;2.4
Mexico City;21.7	Buenos Aires;26.8	Singapore;30.1	Dubai;36.8	Boston;13.7	Buenos Aires;26.2	Vancouver;14.3
Athens;23.5	Toronto;9.5	Kyoto;17.5	Queenstown;9.9	Perth;22.8	Hanoi;30.8	<b>Berlin;10.3</b>
...	...	...	...	...	...	...

# Version 1: Simple & Idiomatic Elixir

- Read the entire file into memory.
- Split it by newlines.
- Create a weather station to values map.
- Aggregate the result.

```
{:ok, content} = File.read(file_path)
acc =
  content
  ▷ String.split("\n")
  ▷ Enum.map(&String.split(&1, ";"))
  ▷ Enum.reject(fn value → value ▷ Enum.at(0) == "" end)
  ▷ Enum.reduce(%{}, fn [key, value], acc →
    {val, _} = Float.parse(value)
    Map.update(acc, key, [val], fn v → [val | v] end)
  end)
result =
  acc
  ▷ Enum.map(fn {key, values} →
    min = Enum.min(values) ▷ round_to_single_decimal()
    max = Enum.max(values) ▷ round_to_single_decimal()
    mean =
      (Enum.sum(values) / length(values))
      ▷ round_to_single_decimal()

    {key, %{min: min, max: max, mean: mean}}
  end)
```

# 1BRC in Elixir: Version 1



**10 Million Rows**

**11.5 sec**



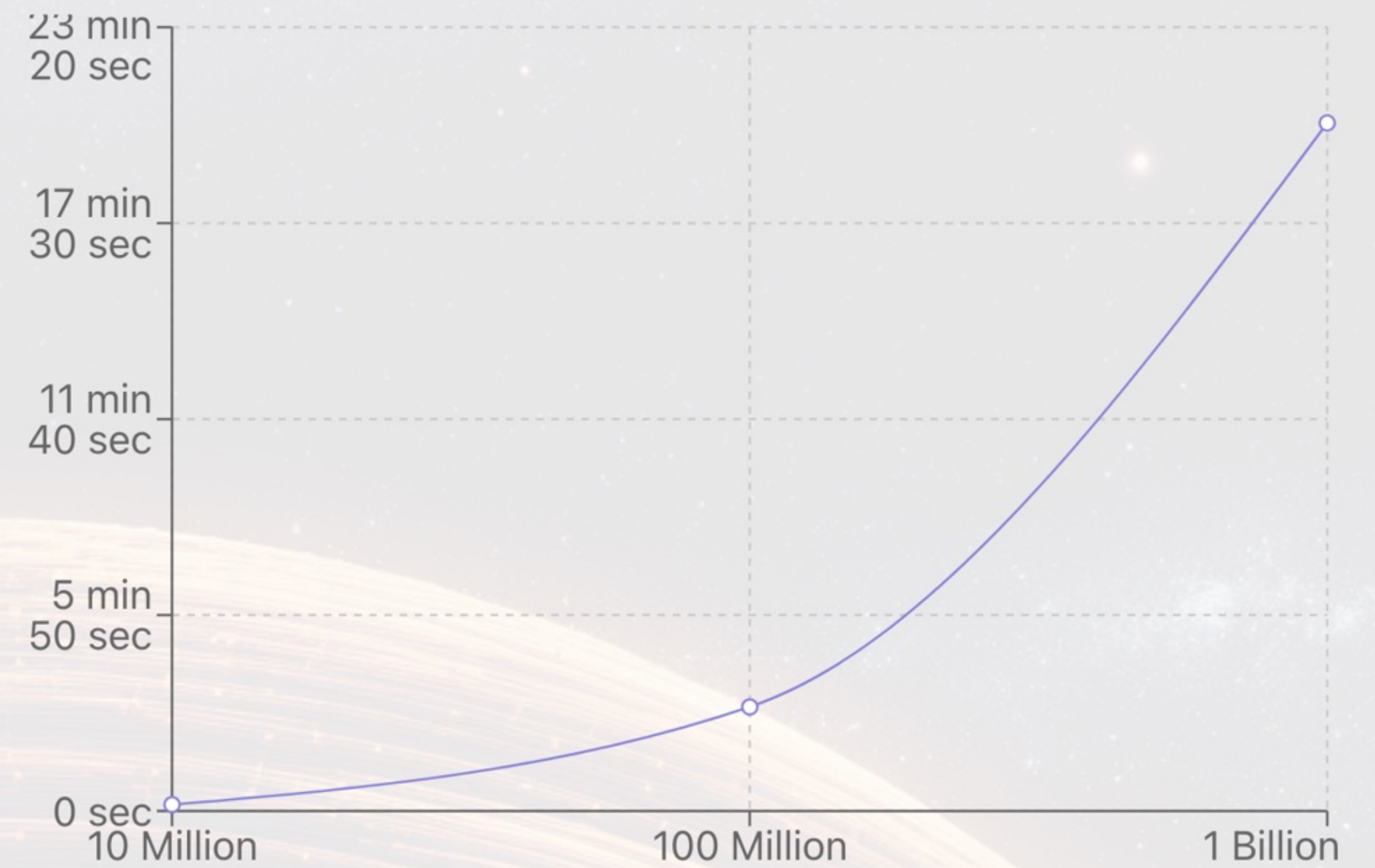
**100 Million Rows**

**3 min 6 sec**



**1 Billion Rows**

**X**





## Version 2: Quick wins: Streaming, Incremental Processing

- `File.stream!/2` for efficient file reading
- Incremental processing: updating min, max, sum and count on-the-fly.
- Single pass over final map to calculate mean temperatures.
- Reduced memory usage by not storing all temperature values.

## Version 2: Quick wins: Streaming, Incremental Processing

- File.stream!/2 for efficient file reading
- Incremental processing: updating min, max, sum and count on-the-fly.
- Single pass over final map to calculate mean temperatures.
- Reduced memory usage by not storing all temperature values.

```
result =
  File.stream!(file_path)
  ▷ Stream.map(&String.split(&1, ";"))
  ▷ Enum.reduce(%{}, fn [key, value], acc →
    {val, _} = Float.parse(value)

    default = %{
      min: val,
      max: val,
      sum: val,
      count: 1
    }

    Map.update(acc, key, default, fn record →
      min = if val < record.min, do: val, else: record.min
      max = if val > record.max, do: val, else: record.max
      sum = record.sum + val
      count = record.count + 1

      %{
        min: min,
        max: max,
        sum: sum,
        count: count
      }
    end)
  end)
end)
```

## Version 2: Quick wins: Streaming, Incremental Processing

- File.stream!/2 for efficient file reading
- Incremental processing: updating min, max, sum and count on-the-fly.
- Single pass over final map to calculate mean temperatures.
- Reduced memory usage by not storing all temperature values.

```
result =
  result
  ▷ Enum.map(fn {key,
    %{min: min, max: max, sum: sum, count: count}} →
      mean = (sum / count) ▷ round_to_single_decimal()

      {key, %{min: min, max: max, mean: mean}}
  end)
```

```
result_txt =
  result
  ▷ Enum.sort_by(fn {key, _} → key end)
  ▷ Enum.reduce("", fn {key,
    %{min: min, max: max, mean: mean}}, acc →
    acc ◇ "#{key};#{min};#{mean};#{max}\n"
  end)
```

# 1BRC in Elixir: Version 2

10 Million Rows



7.8 sec

↘ 32.2%

100 Million Rows



1 min 24 sec

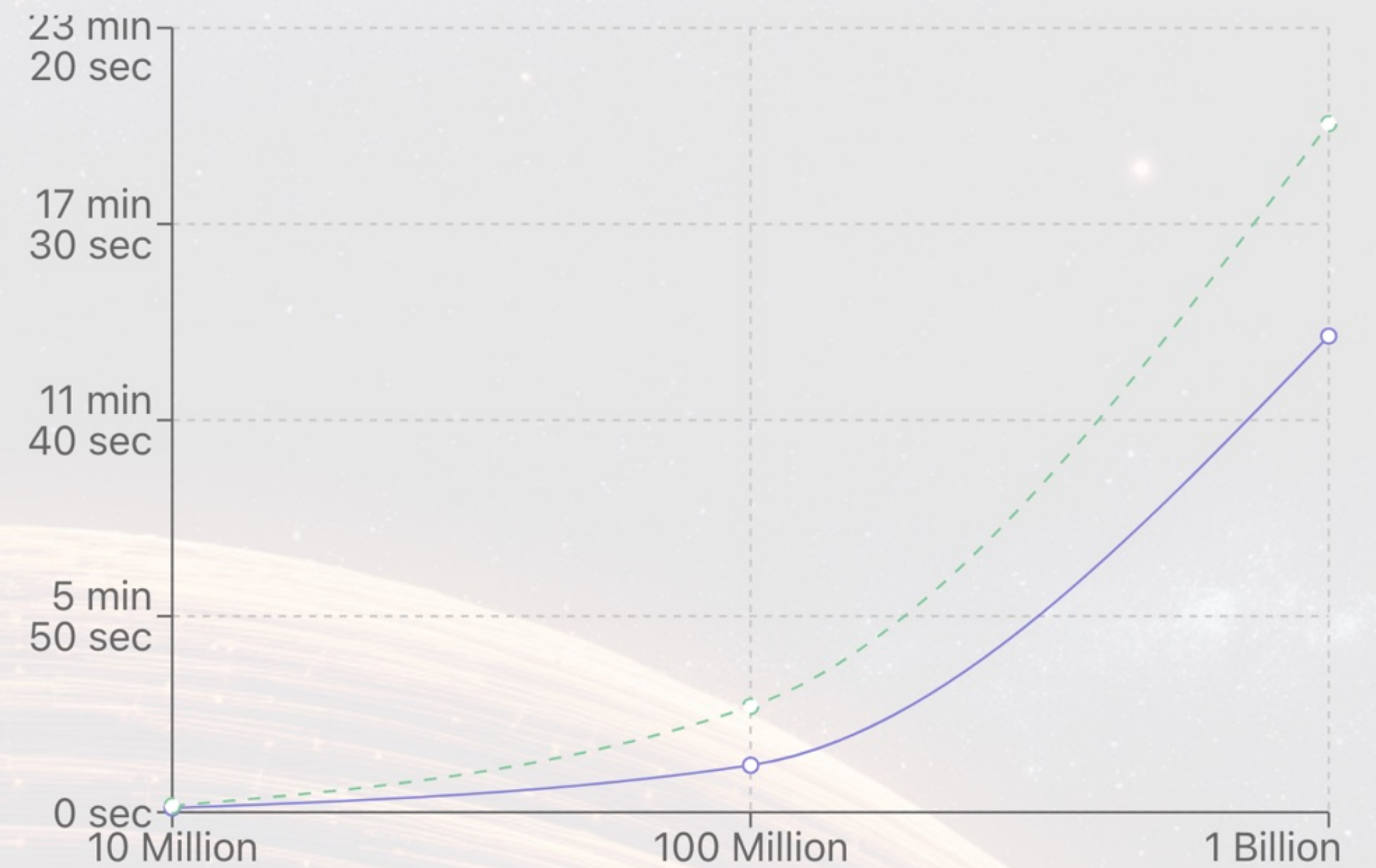
↘ 55.6%

1 Billion Rows



14 min 10 sec

↘ 30.8%



## Version 2: Quick wins: Streaming, Incremental Processing

- File.stream!/2 for efficient file reading
- Incremental processing: updating min, max, sum and count on-the-fly.
- Reduced memory usage by not storing all temperature values.
- Single pass over final map to calculate mean temperatures.

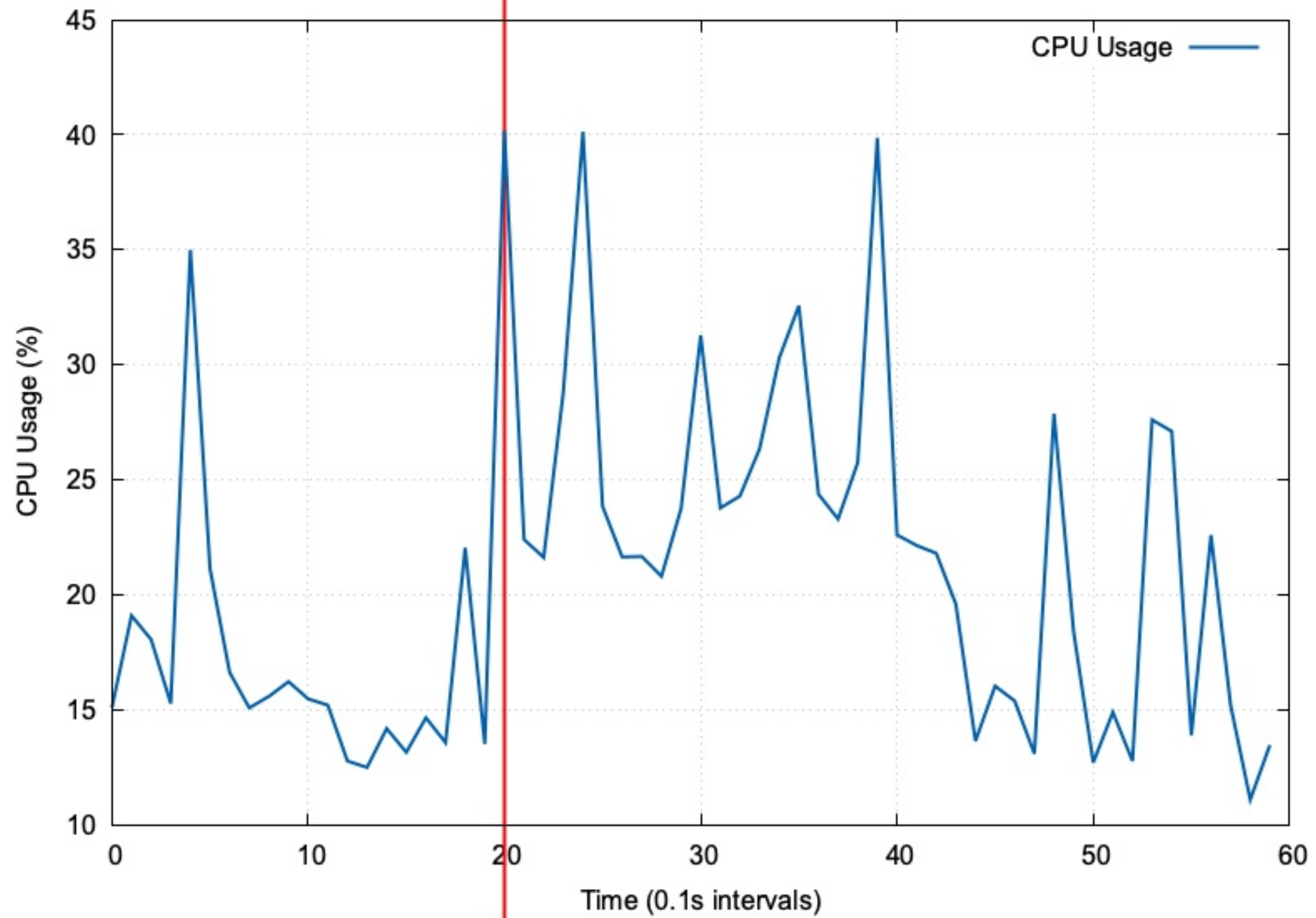
```
result =
  File.stream!(file_path)
  ▷ Stream.map(&String.split(&1, ";"))
  ▷ Enum.reduce(%{}, fn [key, value], acc →
    {val, _} = Float.parse(value)

    default = %{
      min: val,
      max: val,
      sum: val,
      count: 1
    }

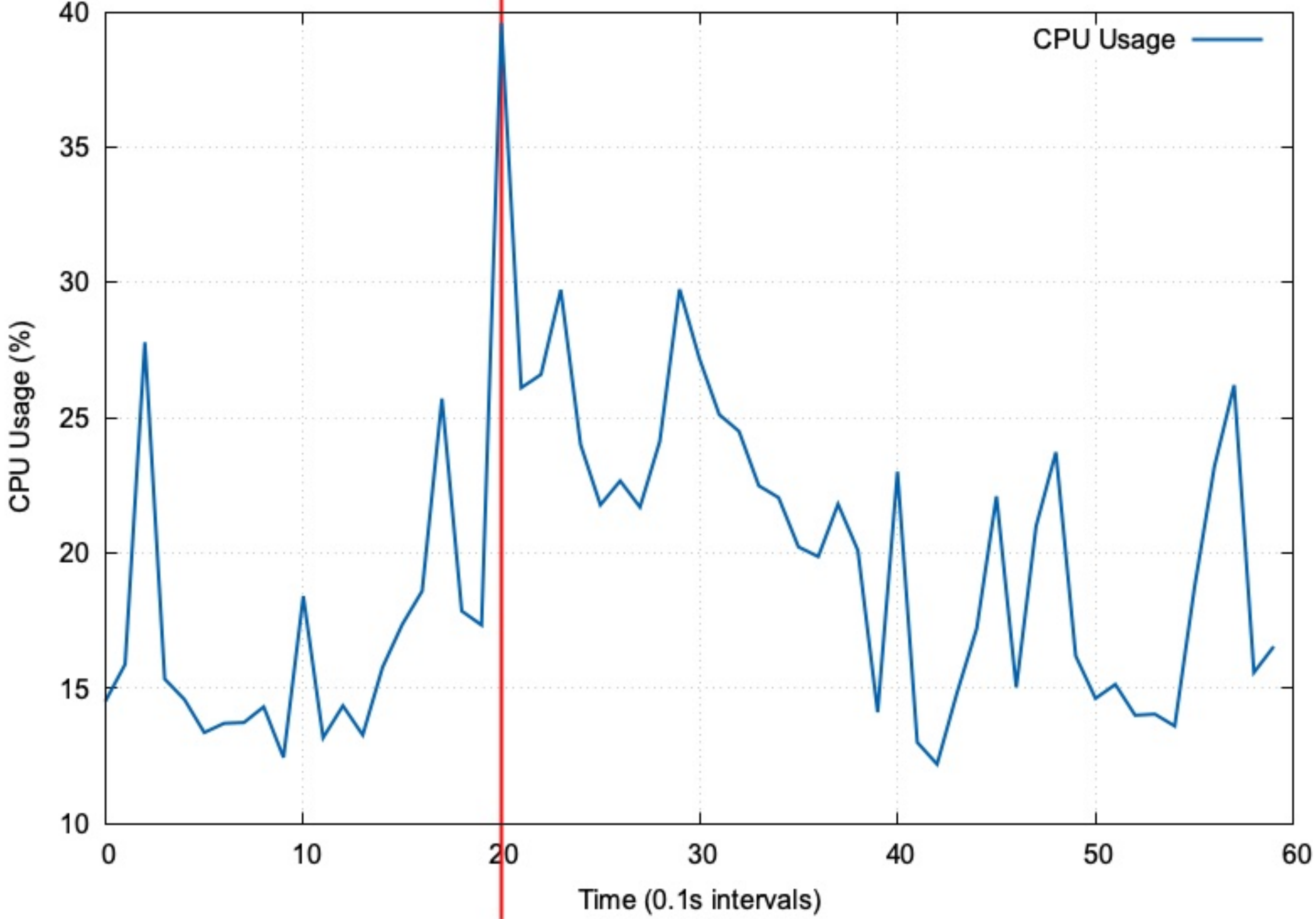
    Map.update(acc, key, default, fn record →
      min = if val < record.min, do: val, else: record.min
      max = if val > record.max, do: val, else: record.max
      sum = record.sum + val
      count = record.count + 1

      %{
        min: min,
        max: max,
        sum: sum,
        count: count
      }
    end)
  end)
```

Overall CPU Usage Over Time (V1, 100M measurements)

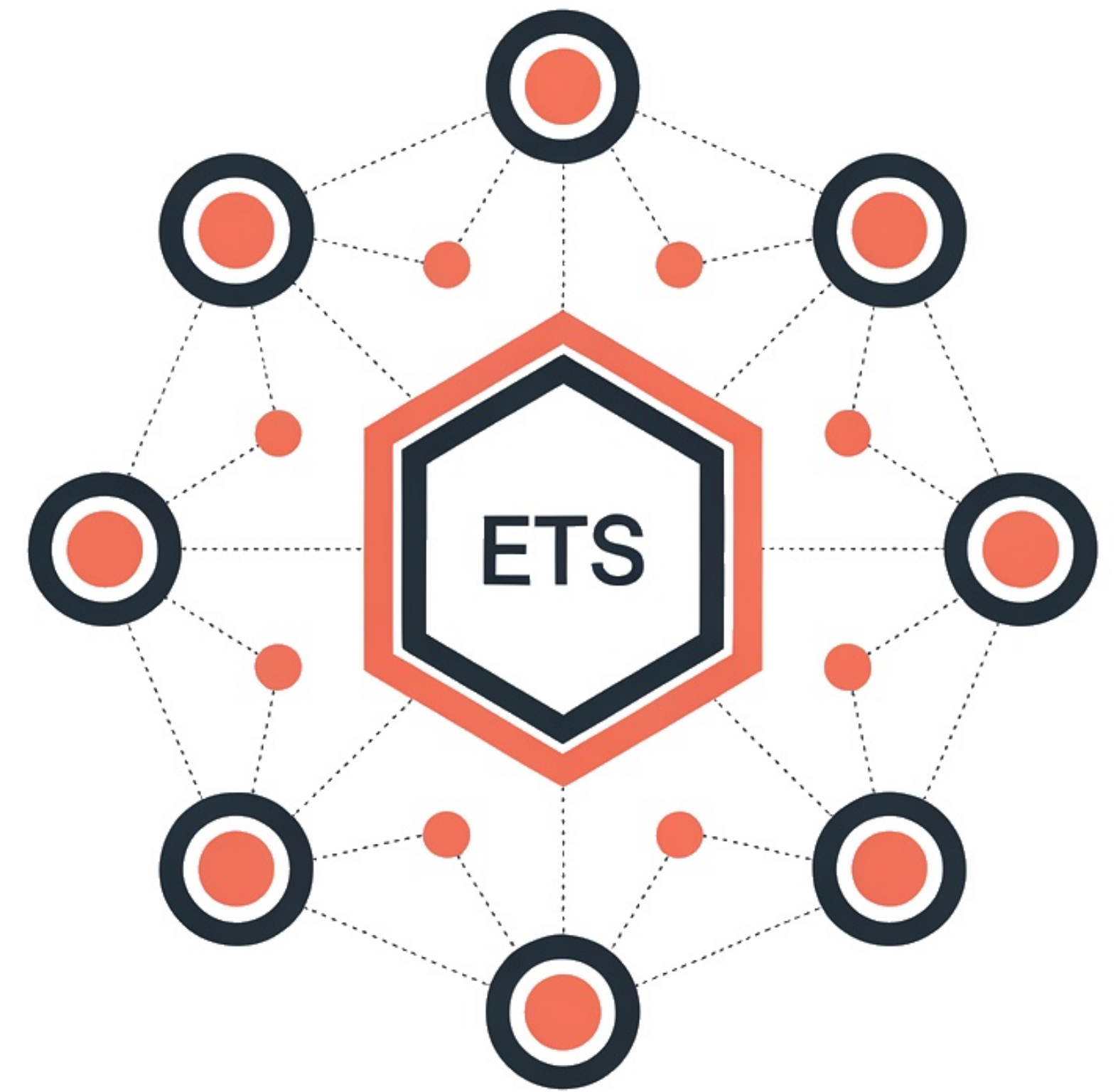


Overall CPU Usage Over Time (V2, 100M measurements)



# Version 3: Introducing Concurrency

- Concurrency using `Task.async_stream/3`
- Using ETS table for storing intermediate results





```
defp process_row([key, value], ets_table) do
  {val, _} = Float.parse(value)
  existing_record = :ets.lookup(ets_table, key)
  new_record =
    case existing_record do
      [] →
        %{
          min: val,
          max: val,
          sum: val,
          count: 1
        }
      [{^key, record}] →
        min = if val < record.min, do: val, else: record.min
        max = if val > record.max, do: val, else: record.max
        sum = record.sum + val
        count = record.count + 1
        %{
          min: min,
          max: max,
          sum: sum,
          count: count
        }
    end
  :ets.insert(ets_table, {key, new_record})
end
```

```
lib/measurements_processor.ex  +9 -1
```

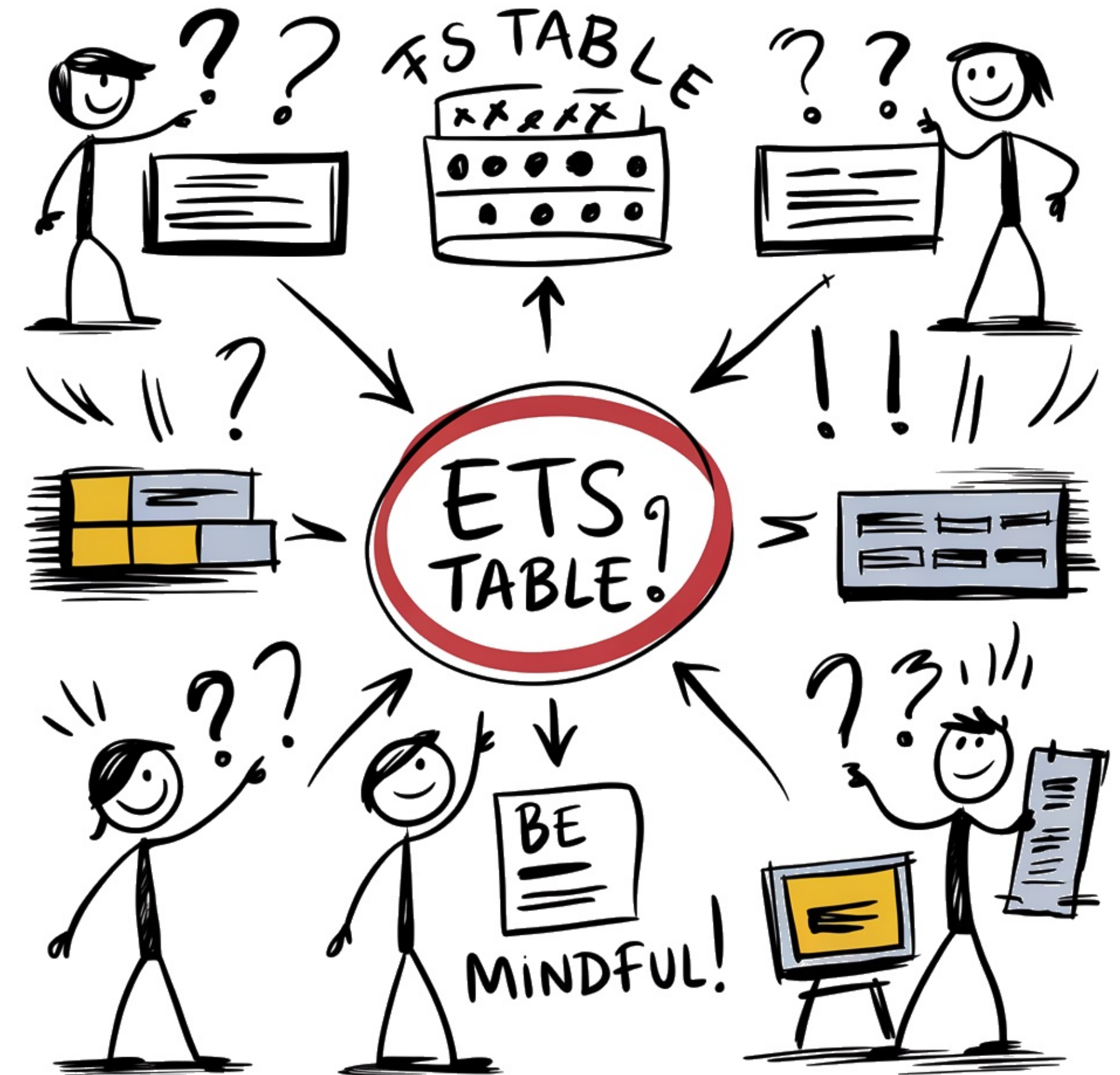
```
@@ -28,7 +28,15 @@ defmodule OneBRC.MeasurementsProcessor do
  28     fs
  29     |> Stream.map(&String.split(&1, ";"))
  30     |> Stream.reject(fn value -> value |> Enum.at(0) == "" end)
  31 -   |> Stream.map(fn val -> process_row(val, ets_table) end)
  32     |> Stream.run()
  33
  34     t2 = System.monotonic_time(:millisecond)
  28     fs
  29     |> Stream.map(&String.split(&1, ";"))
  30     |> Stream.reject(fn value -> value |> Enum.at(0) == "" end)
  31 +   |> Stream.chunk_every(10_000)
  32 +   |> Task.async_stream(
  33 +     fn val ->
  34 +       Enum.map(val, fn row -> process_row(row, ets_table)
  35 +         end)
  36 +     end,
  37 +     max_concurrency: System.schedulers_online() * 5,
  38 +     ordered: false,
  39 +     timeout: :infinity
  40     )
  41     |> Stream.run()
  42     t2 = System.monotonic_time(:millisecond)
```

Replacing `Stream.map` with `Task.async_stream` to add concurrency 🤪



# Version 3: Introducing Concurrency

- Using ETS table for storing intermediate results
- Lesson learned: be mindful of shared state and potential race conditions with concurrent processes.



```
lib/measurements_processor.ex  +18 -8
```

56	end	54	end
57		55	
58	- defp process_row([key, value], ets_table) do	56	+ defp parse_row(row) do
59	-   {val, _} = Float.parse(value)	57	+   item = row  > String.split(";")
		58	+
		59	+   case item do
		60	+     ["" ] ->
		61	+     nil
		62	+
		63	+     [key, value] ->
		64	+       {parsed_value, _} = Float.parse(value)
		65	+       [key, parsed_value]
		66	+     end
		67	+   end
		68	+
		69	+ defp process_row([key, val], ets_table) do
60	existing_record = :ets.lookup(ets_table, key)	70	existing_record = :ets.lookup(ets_table, key)
61		71	
62	new_record =	72	new_record =

Separating parsing from the actual calculation



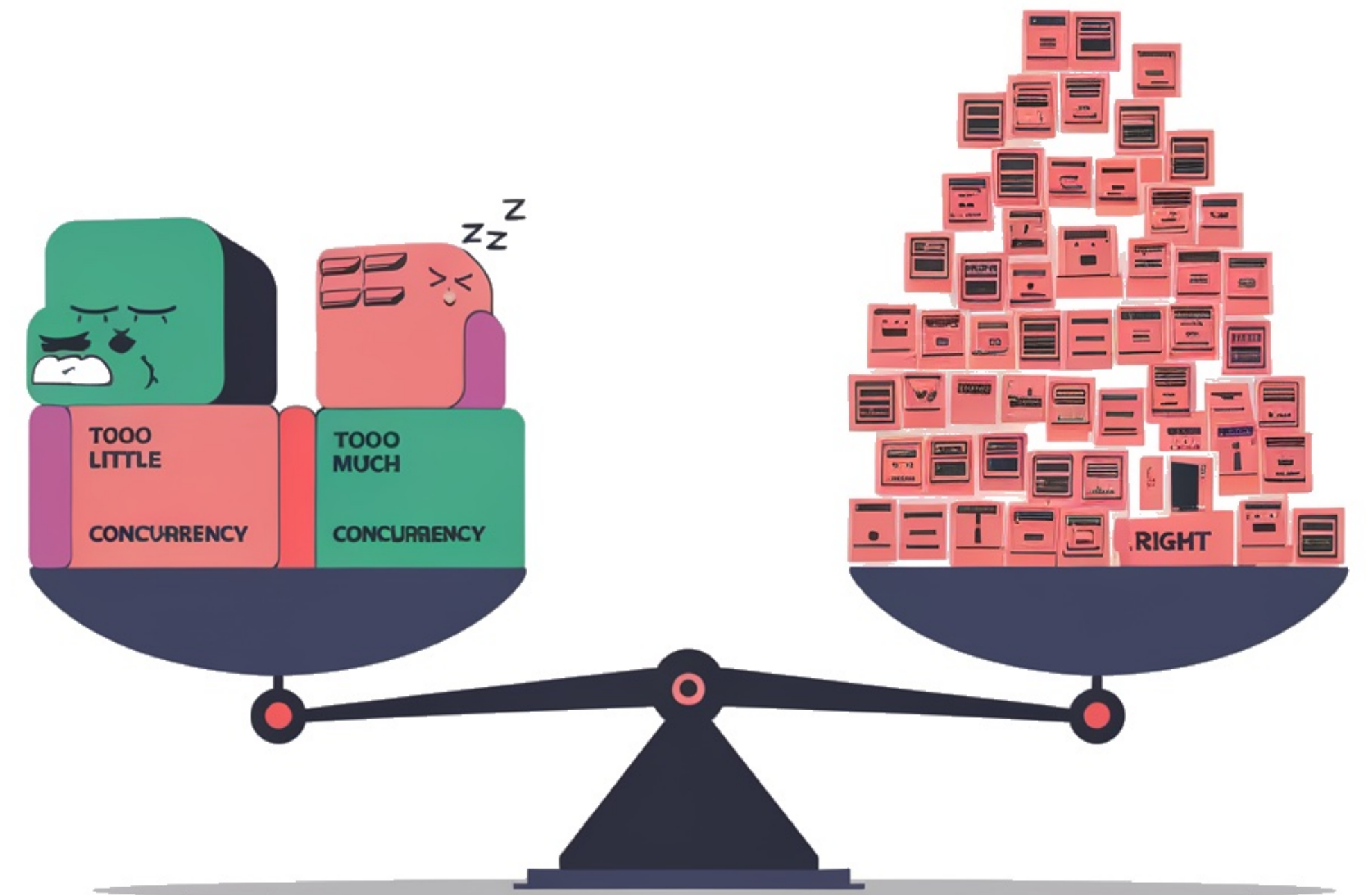
lib/measurements_processor.ex		
32	> Task.async_stream( 33 - fn val -> 34 - Enum.map(val, fn row -> process_row(row, ets_table) end) 35 - end, 36 max_concurrency: System.schedulers_online() * 5, 37 ordered: false, 38 timeout: :infinity 39 ) 40  > Stream.run() 41 42 t2 = System.monotonic_time(:millisecond) @@ -55,8 +53,20 @@ defmodule OneBRC.MeasurementsProcessor do 55 result 56 end	+18 -8
30	> Task.async_stream( 31 + fn val -> Enum.map(val, &parse_row/1)  > Enum.reject(&is_nil/1) end, 32 max_concurrency: System.schedulers_online() * 5, 33 ordered: false, 34 timeout: :infinity 35 ) 36 +  > Stream.flat_map(&elem(&1, 1)) 37 +  > Stream.map(&process_row(&1, ets_table)) 38  > Stream.run() 39 40 t2 = System.monotonic_time(:millisecond) 53 result 54 end	

*Parse rows concurrently, then process synchronously*



# Version 3: Introducing Concurrency

- Using ETS table for storing intermediate results
- Lesson 1: be mindful of shared state and potential race conditions with concurrent processes.
- Lesson 2: Too little concurrency underutilizes resources, while too much leads to overhead



```
lib/one_brc/measurements_processor.ex  +1 -1  ...
@@ -28,7 +28,7 @@ defmodule OneBRC.MeasurementsProcessor do
28     ets_table = :ets.new(:station_stats, [:set, :public])
29
30     fs
31 -   |> Stream.chunk_every(2000)
31 +   |> Stream.chunk_every(10000)
32   |> Task.async_stream(
33     fn val -> Enum.map(val, &parse_row/1) end,
34     max_concurrency: System.schedulers_online(),
```

*Process spawning overhead matters!*



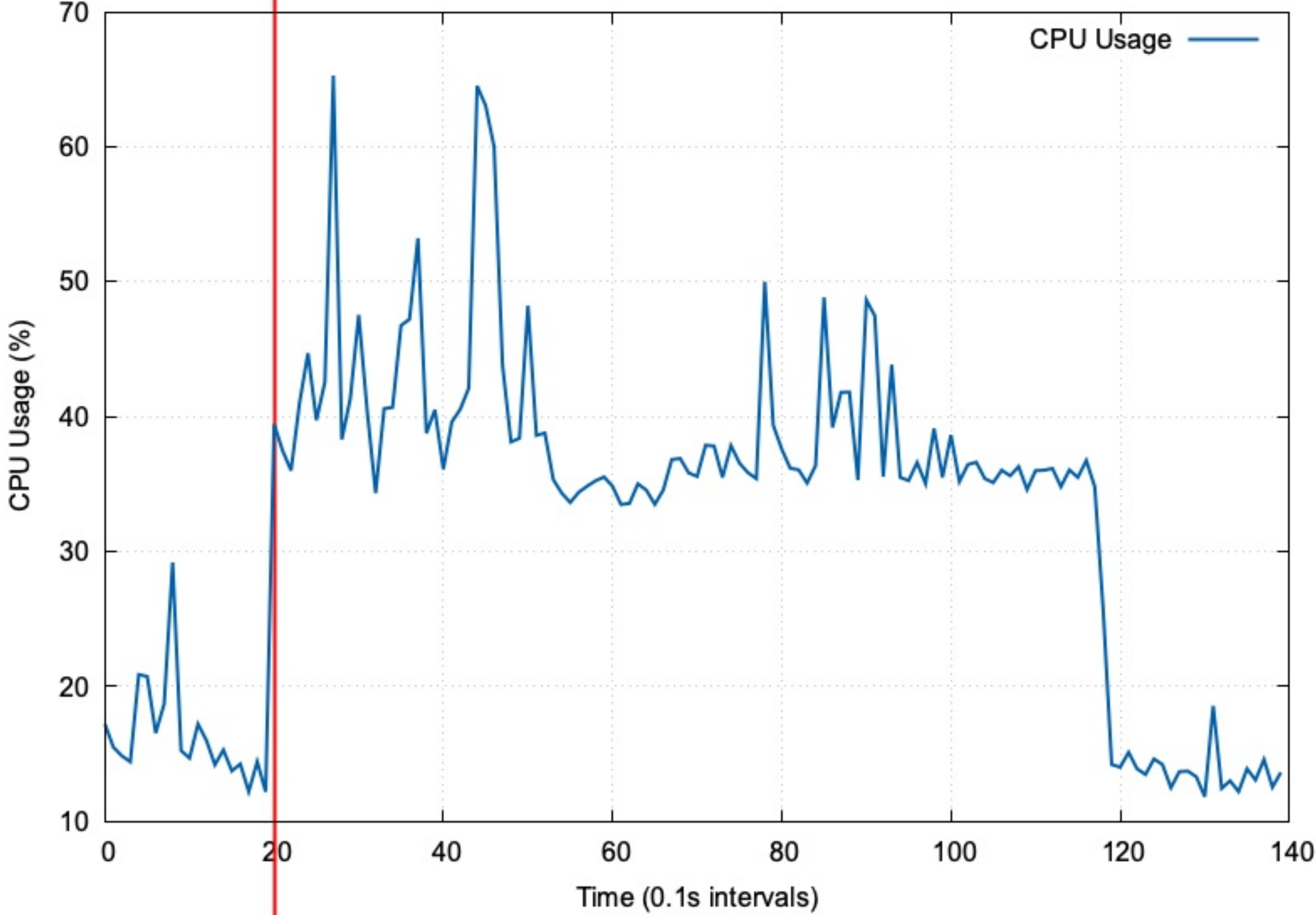
lib/measurements_processor.ex		+9 -3	
45	- mean = (sum / count)  > round_to_single_decimal()	45	+ mean = (sum / (count * 10.0))  > round_to_single_decimal()
46		46	
47	- {key, %{min: min, max: max, mean: mean}}	47	+ {key, %{min: min / 10.0, max: max / 10.0, mean: mean}}
48	end)	48	end)
49		49	
50	Logger.info("Processing data 1: #{t2 - t1} ms")	50	Logger.info("Processing data 1: #{t2 - t1} ms")
@@ -60,7 +60,13 @@ defmodule OneBRC.MeasurementsProcessor do			
60		60	
61	row ->	61	row ->
62	[key, value] = String.split(row, ";")	62	[key, value] = String.split(row, ";")
63	- parsed_value = String.to_float(value  > String.trim_trailing())	63	+  64 + parsed_value = 65 + value  > String.trim_trailing() 66 + 67 + [a, b] = parsed_value  > String.split(".") 68 + parsed_value = (a <> b)  > String.to_integer() 69 +
64	[key, parsed_value]	70	[key, parsed_value]
65	end	71	end
--	.	--	.

Getting around float parsing



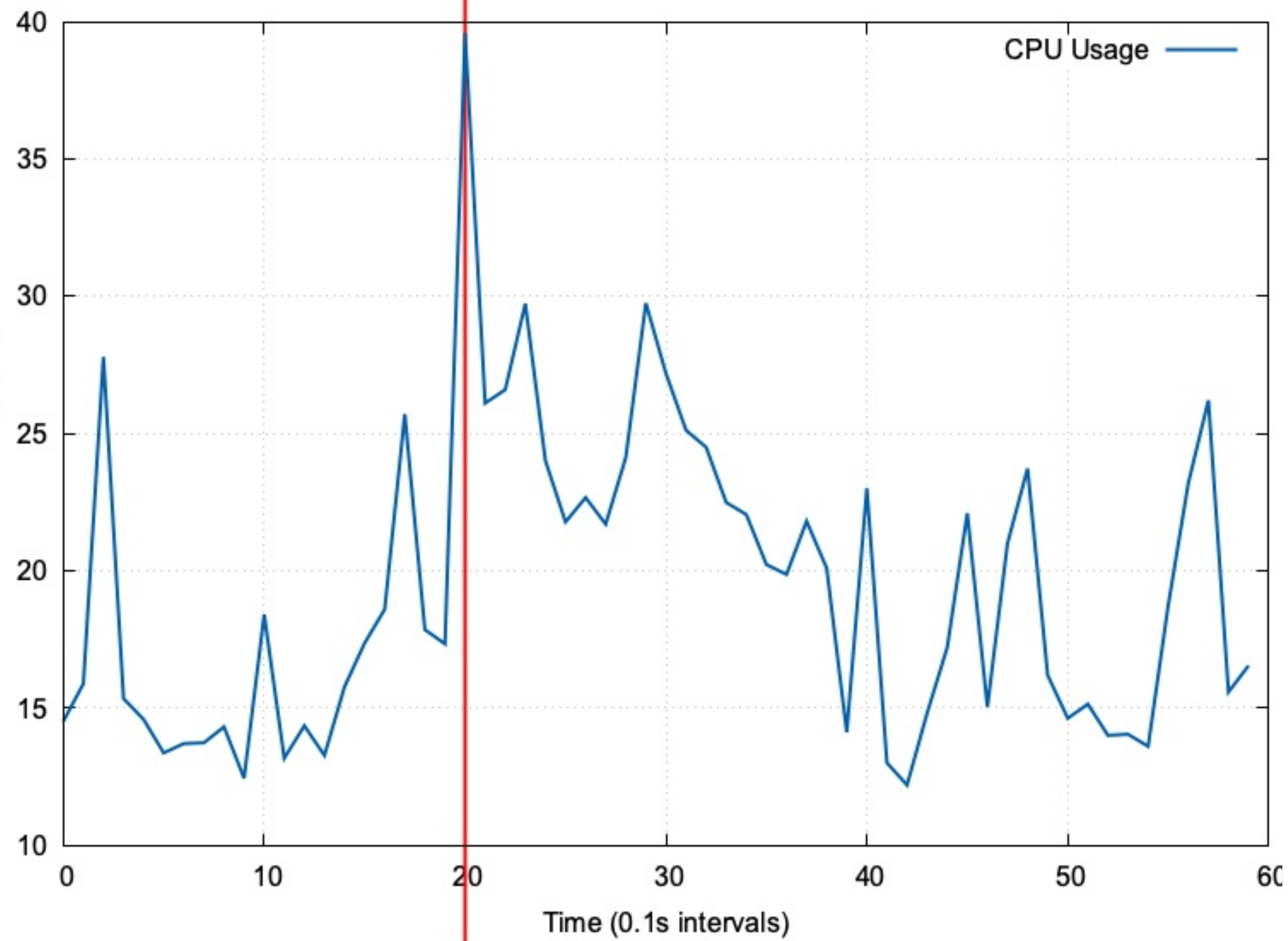


Overall CPU Usage Over Time (V3, 100M measurements)

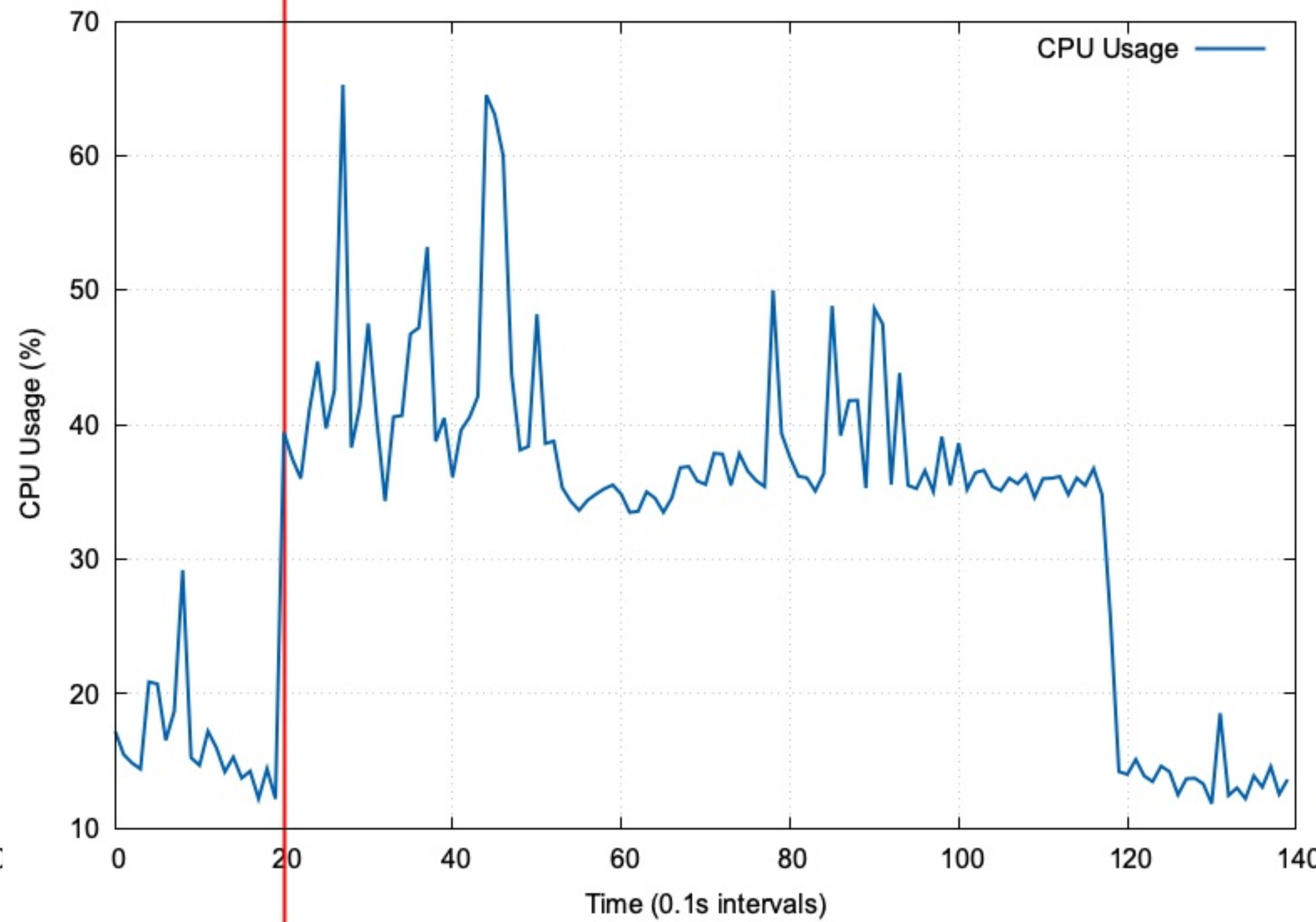


*CPU usage of version 3*

Overall CPU Usage Over Time (V2, 100M measurements)



Overall CPU Usage Over Time (V3, 100M measurements)



CPU usage of version 2 vs version 3

# 1BRC in Elixir: Version 3

10 Million Rows



4.7 sec

↘ 39.7%

100 Million Rows



46 sec

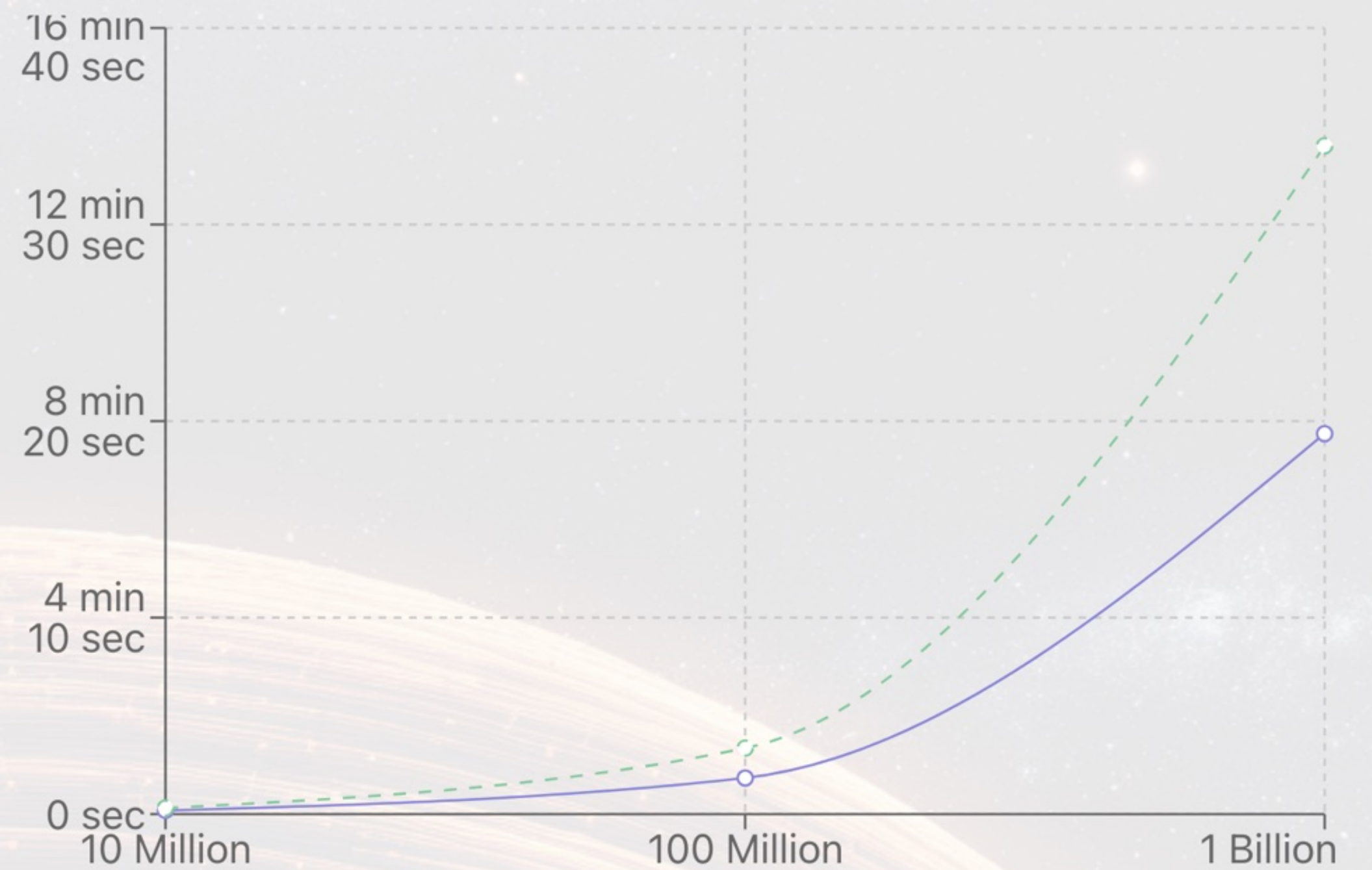
↘ 45.2%

1 Billion Rows



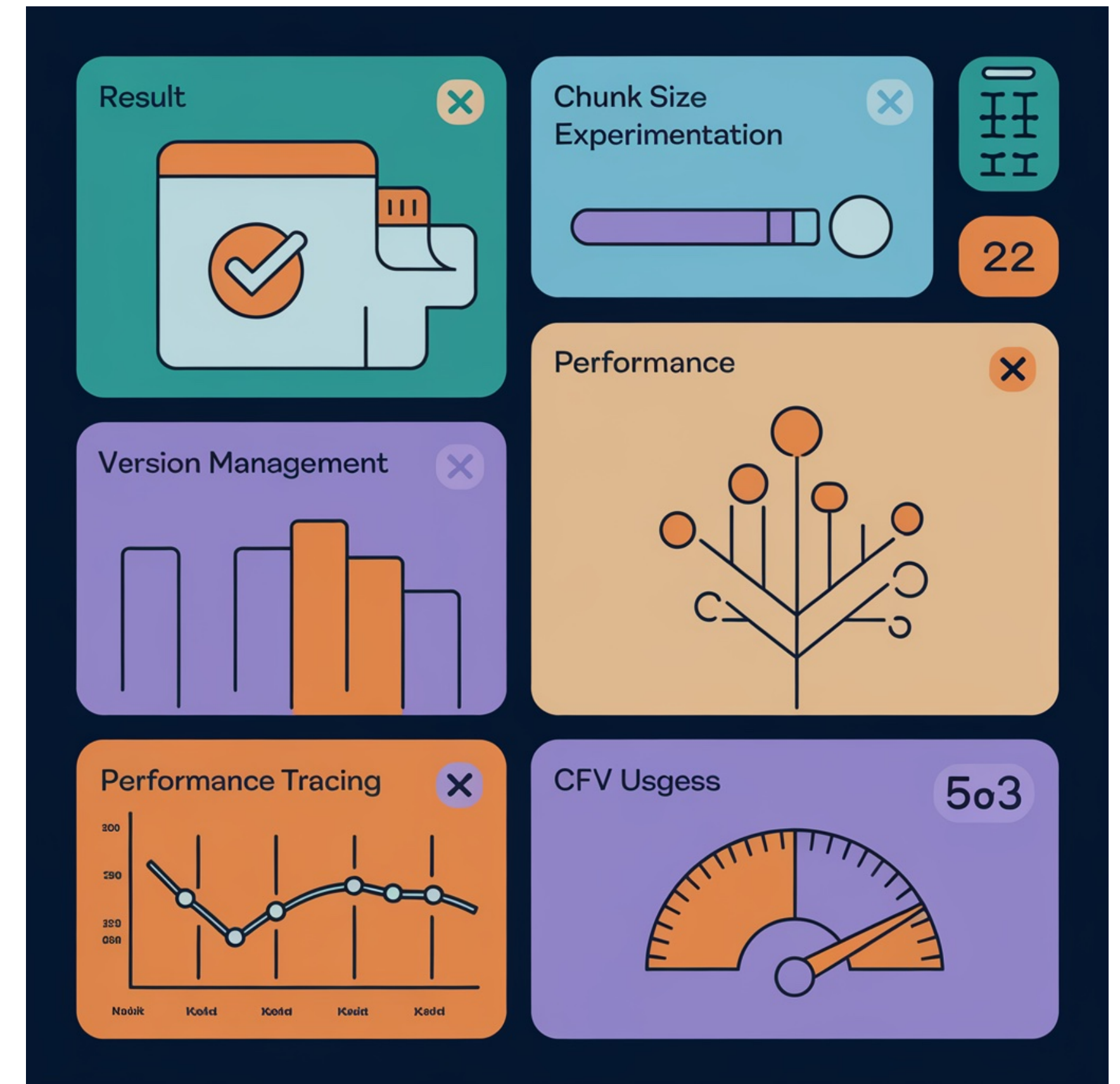
8 min 4 sec

↘ 43.1%



# Need for better instrumentation

- Key needs:
  - Result correctness verification
  - Chunk size experimentation
  - Version management, performance comparison
  - Performance tracing
  - CPU usage measurement



***Demo: maintaining & executing different versions, measuring cpu usage***

# Profiling Elixir Code

- Pinpointing time-consuming code sections to detect perf bottlenecks.
- Frequency and duration of function calls, where they spend their time.
- A flame graph like visualisation would be nice.
- Measure the impact of optimisations



# Profiling Elixir Code

- Pinpointing time-consuming code sections to detect perf bottlenecks.
- Frequency and duration of function calls, where they spend their time.
- A flame graph like visualisation would be nice.
- Measure the impact of optimisations
- cprof, prof, eflambe, benchee.



# 1. cprof

- Counts function calls, not execution time.
- Provides a quick overview of most called functions.
- Lightweight, minimal impact on code performance.





```
> run process_measurements.profile.cprof --count=1000 --version=3
```

```
Warmup ...
```

```
...
Total CNT
61928
:lists 6054 ←
  :lists.keyfind/3 4013
  ...
Stream 5032 ←
  Stream.do_resource/5 1002
  Stream.do_element_resource/6 1000
  ...
String 5008 ←
  String.split/3 2000
  String.split/2 2000
  String.trim_trailing/1 1000
  String.do_replace/4 2
  ...
String.Break 5000 ←
  String.Break.trim_trailing/2 2000
  ...
OneBRC.MeasurementsProcessor.Version3 4503 ←
  OneBRC.MeasurementsProcessor.Version3.process_row/2 1000
  OneBRC.MeasurementsProcessor.Version3.parse_row/1 1000
  ...
```

```
Profile done over 21915 matching functions
```

# 1. cprof

- Counts function calls, not execution time.
- Provides a quick overview of most called functions.
- Lightweight, minimal impact on code performance.

# 2. eprof

- Counts function calls and execution time per function.
- Provides a clearer picture of time consuming functions
- Output shows percentage function calls, percentage ti per function and microseconds / call



```
> run process_measurements.profile.eprof --count=1000 --version=3
```

```
Warmup ...
```

```
...
```


```
Profile results of #PID<0.228.0>
```

#	CALLS	% TIME	μS/CALL
Total	37216	100.	3009 0.08
...			
IO.each_binstream/2	1001	2.82	85 0.08
Stream.do_resource/5	1002	2.82	85 0.08
anonymous fn/4 in Stream.chunk_while_fun/2	1000	3.09	93 0.09
String.Chars.Float.to_string/1	1126	3.59	108 0.10
Enum."-map/2-lists^map/1-1-"/2	1188	3.82	115 0.10
OneBRC.MeasurementsProcessor.Version3.process_row/2	1000	3.99	120 0.12
:ets.lookup/2	1018	4.05	122 0.12
:prim_file.read_line_1/4	1002	4.12	124 0.12
:ets.insert/2	1000	5.38	162 0.16
:prim_file.read_line/1	1001	5.95	179 0.18

*eprof output*

# 3. eFlambe

README Apache-2.0 license



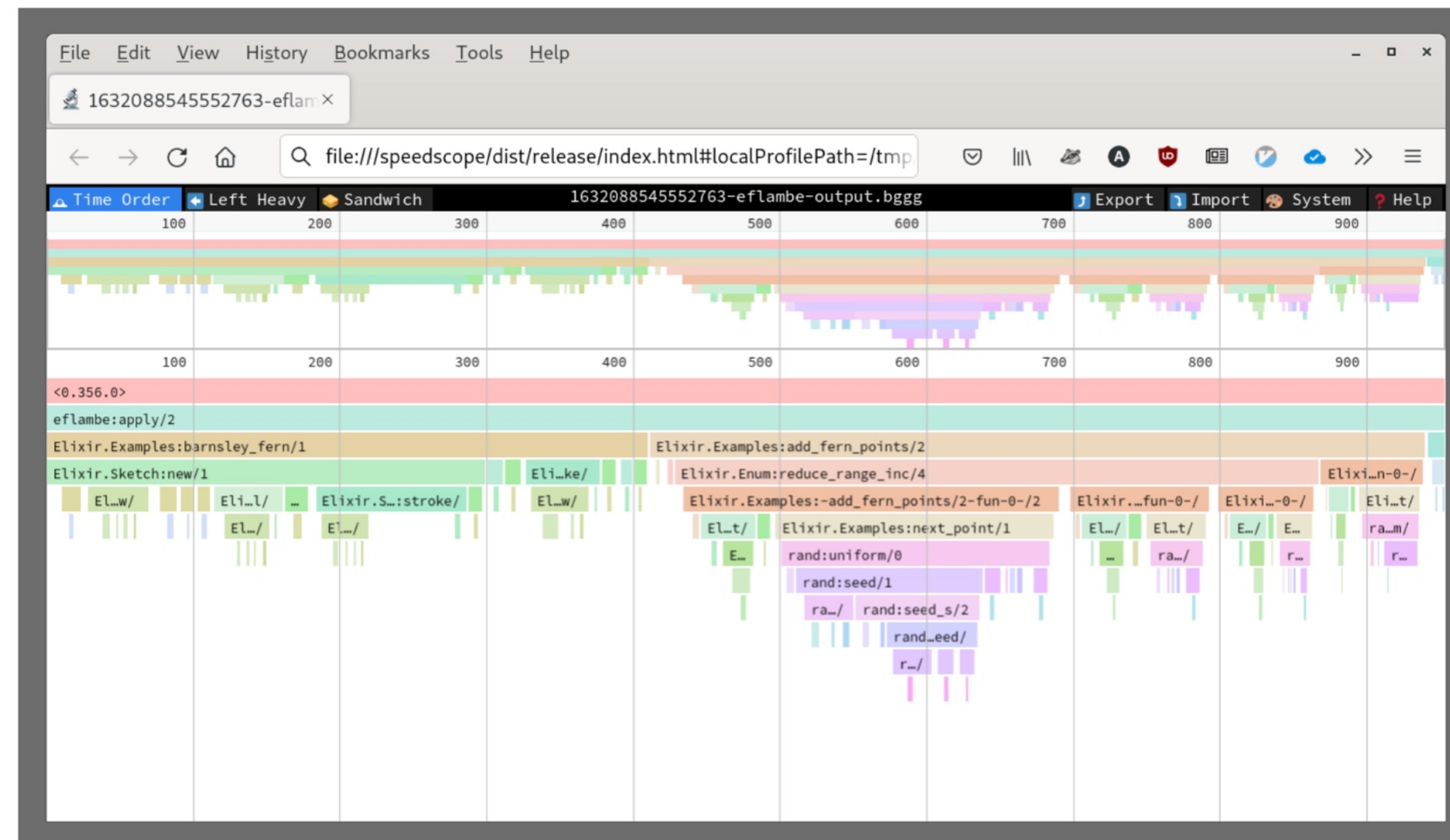
## eFlambe!

A tool for rapid profiling of Erlang and Elixir applications

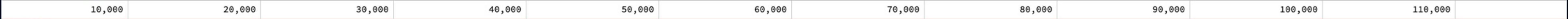
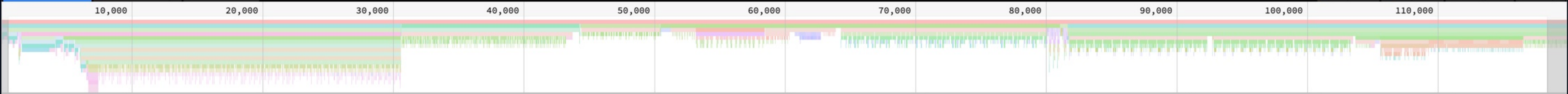
CI failing hex v0.3.1 documentation

- Designed for rapid generation of flamegraphs - no need to edit source code or recompile!
- Generates output in several common formats - perf and Brendan Gregg's
- Interface inspired by recon\_trace to automatically stop flamegraph generation
- Can automatically open data in your flamegraph viewer of choice

### Output Viewed in Speedscope



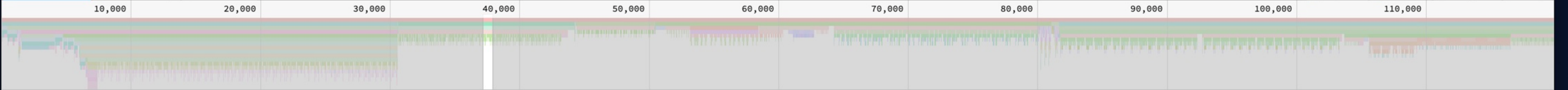
[Watch a quick demo here](#)



eflambe:apply/2  
Elixir.OneBRC.MeasurementsProcessor:process/2  
timer:tc/2  
Elixir.OneBRC.MeasurementsProcessor.Version5:process/1  
Elixir.Stream:run/1  
Elixir.Task.Supervised:stream\_reduce/7  
Elixir.Enumerable.Stream:do\_each/4

Elixir.Task.Supervised:reply/4  
Elixir...-1-/2  
Elixir.OneBRC.MeasurementsProcessor.Version5:process/1  
Elixir.Task.Su...ed:invoke\_mfa/2  
Elixir...-0-/3  
E... Elixi...map/2 Elix...y/3  
Elixir.Enum:-reduc...ists^foldl/2-0-/3  
Elixir.OneBRC.MeasurementsProcessor:verify\_result/1  
Elixir.Enum:-ma...sts^map/1-1-/2  
maps:...d\_1/4 E...  
Elixir.OneBrc.ResultVerification:verify\_result/2  
Elixir.Enum:...^map/1-1-/2 Elixir.Enum:...^map/1-1-/2 E... Elixir.OneBr...mpare\_data/2 E...  
Eli...t/2  
Eli...e/3

24,524 (20%) Elixir.Task.Supervised:stream\_reduce/7

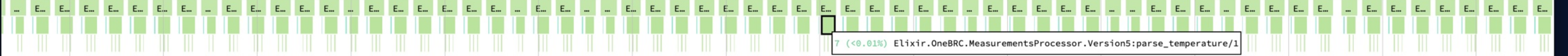


<0.98.0>

Elixir.Task.Supervised:reply/4

Elixir.Task.Supervised:invoke\_mfa/2

Elixir.Enum:-map/2-lists^map/1-1-/2



7 (<0.01%) Elixir.OneBRC.MeasurementsProcessor.Version5:parse\_temperature/1

# 4. Benchee

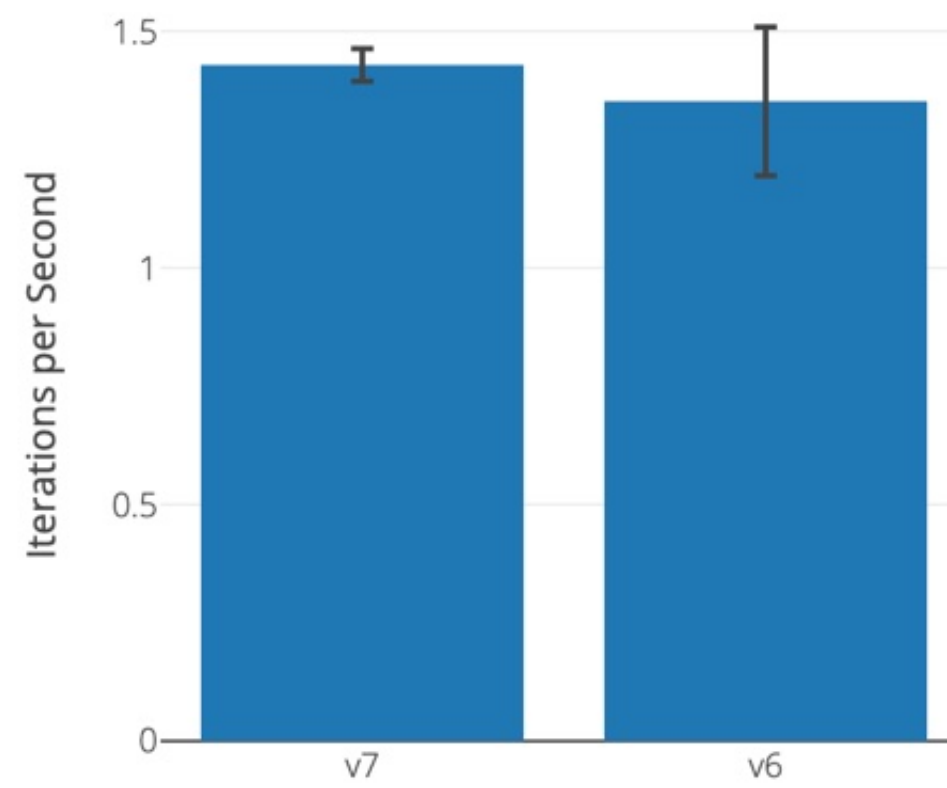
- Library for easy benchmarking in Elixir.
- Compares the performance of different pieces of code.
- Runs each of your functions for a given amount of time after an initial warmup.
- Measures execution time, memory consumption.
- Provides a plethora of statistics - avg, ips, median, 99p.
- Has different visualisation plugins for creating reports - HTML, Markdown, JSON.



## Run Time Comparison ?

Name	Iterations per Second	Average	Deviation	Median	Mode	Minimum	Maximum	Sample size
v7	1.43	699.57 ms	±2.40%	707.20 ms	none	674.78 ms	714.78 ms	8
v6	1.35	739.26 ms	±11.60%	780.37 ms	none	577.56 ms	799.37 ms	7

Average Iterations per Second



Run Time Boxplot

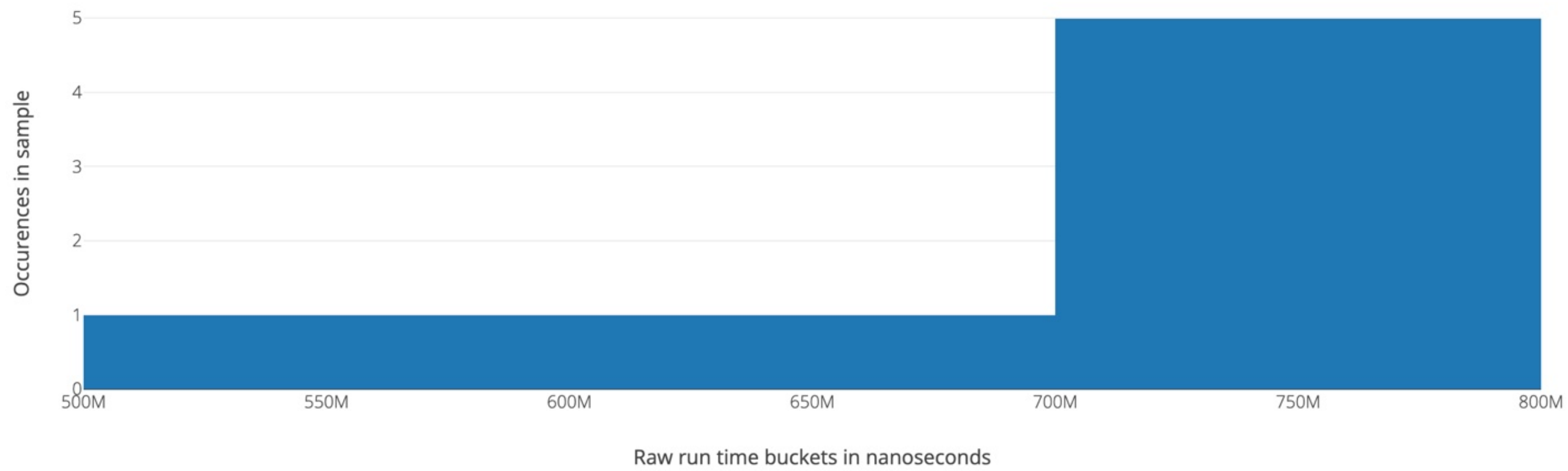




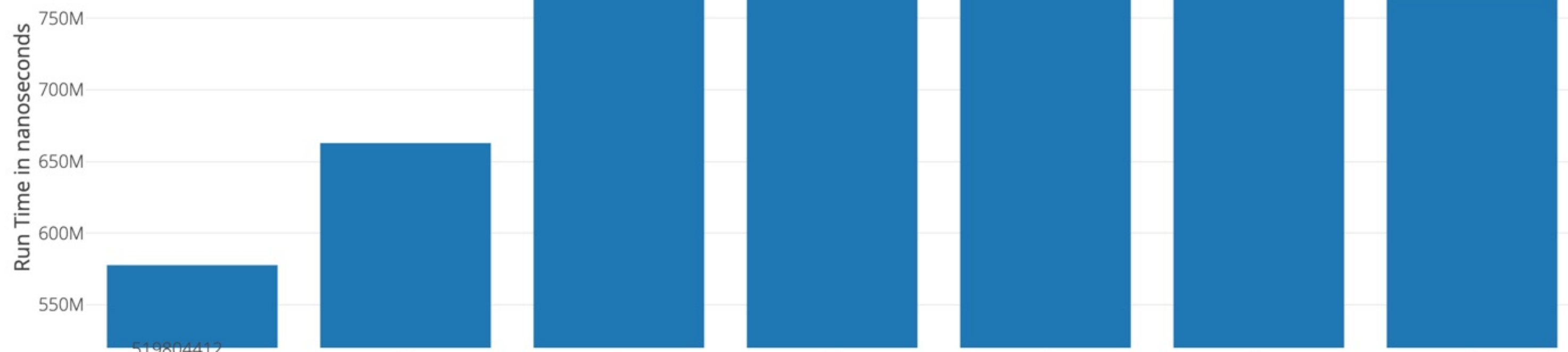
Name	Iterations per Second	Average	Deviation	Median	Mode	Minimum	Maximum	Sample size
v6	1.35	739.26 ms	±11.60%	780.37 ms	none	577.56 ms	799.37 ms	7



v6 Run Times Histogram

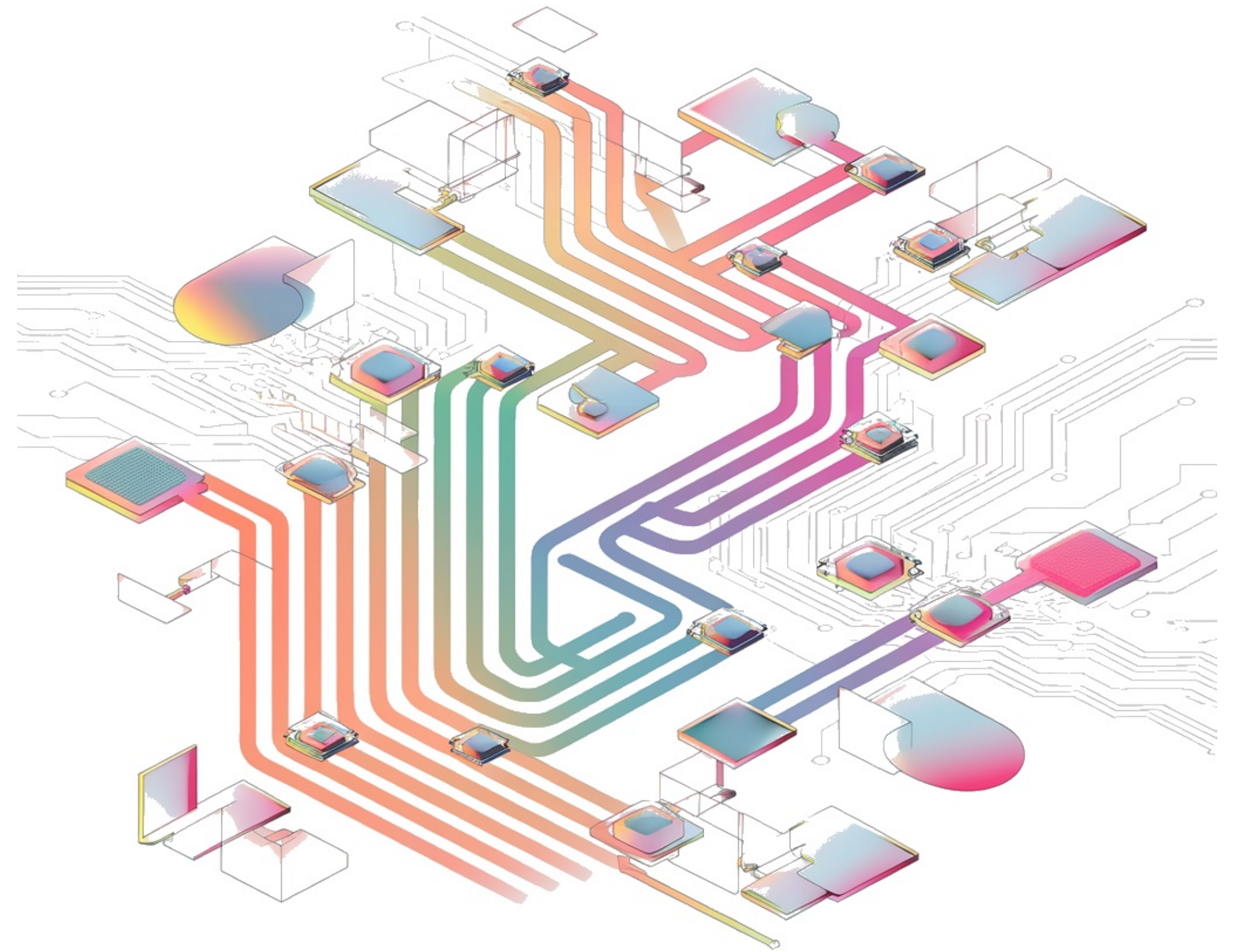


v6 Raw Run Times



## Version 4: Making processing rows concurrent

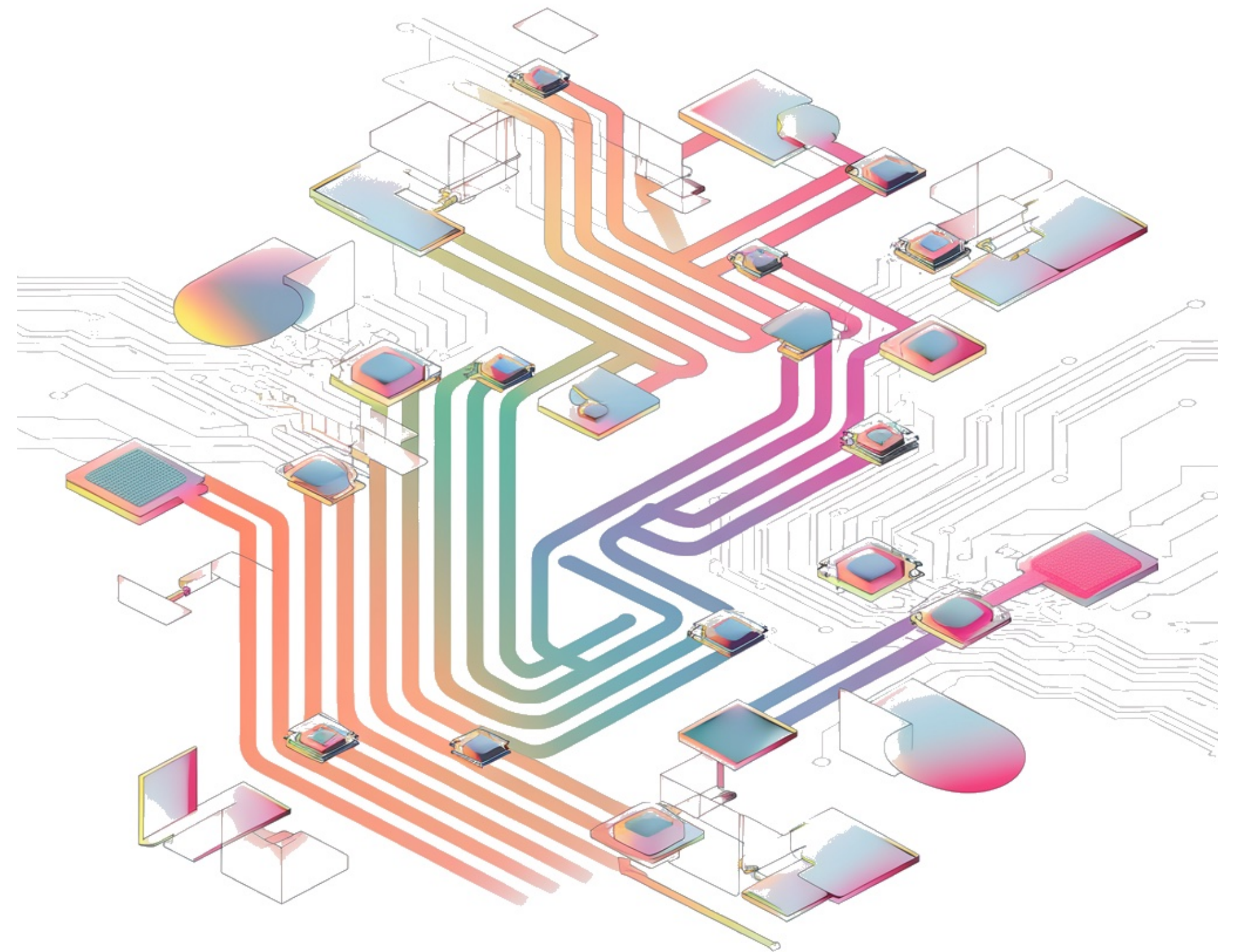
- Concurrency currently limited to the parsing step.
- Processing step was sequential due to shared state.



```
defp process_row([key, value], ets_table) do
  {val, _} = Float.parse(value)
  existing_record = :ets.lookup(ets_table, key)
  new_record =
    case existing_record do
      [] →
        %{
          min: val,
          max: val,
          sum: val,
          count: 1
        }
      [{^key, record}] →
        min = if val < record.min, do: val, else: record.min
        max = if val > record.max, do: val, else: record.max
        sum = record.sum + val
        count = record.count + 1
        %{
          min: min,
          max: max,
          sum: sum,
          count: count
        }
    end
  :ets.insert(ets_table, {key, new_record})
end
```

## Version 4: Making processing rows concurrent

- Concurrency currently limited to the parsing step.
- Processing step was sequential due to shared state.
- Proposed solution: Divide data into chunks.
- Each chunk produces intermediate results independently.
- A single pass at the end combines intermediate results for the final output.



```
74 - defp process_row([key, val], ets_table) do
75 -   existing_record = :ets.lookup(ets_table, key)
76
77   new_record =
78     case existing_record do
79 -     [] ->
80       %{
81         min: val,
82         max: val,
83         mean: val,
84         count: 1
85       }
86
87 -     [{^key, record}] ->
88 -     min = if val < record.min, do: val, else: record.min
89 -     max = if val > record.max, do: val, else: record.max
90 -     mean = (record.mean * record.count + val) / (record.count + 1)
91
92     %{
93       min: min,
94       max: max,
95       mean: mean,
96 -     count: record.count + 1
97     }
98   end
99
100 - :ets.insert(ets_table, {key, new_record})
```

```
135 + defp process_row([key, val], acc) do
136 +   existing_record = Map.get(acc, key, nil)
137
138   new_record =
139     case existing_record do
140 +     nil ->
141       %{
142         min: val,
143         max: val,
144         mean: val,
145         count: 1
146       }
147
148 +     %{count: count, min: min, max: max, mean: mean} ->
149 +     min = if val < min, do: val, else: min
150 +     max = if val > max, do: val, else: max
151 +     mean = (mean * count + val) / (count + 1)
152 +
153     %{
154       min: min,
155       max: max,
156       mean: mean,
158 +     count: count + 1
159     }
160   end
161
162 +   Map.put(acc, key, new_record)
```

```
30     |> Task.async_stream(  
31       fn val -> Enum.map(val, &parse_row/1) end,  
32 -     max_concurrency: System.schedulers_online() * 5,
```

```
33     ordered: false,  
34     timeout: :infinity  
35   )  
36 -   |> Stream.flat_map(&elem(&1, 1))  
37 -   |> Stream.map(&process_row(&1, ets_table))  
38   |> Stream.run()  
39
```

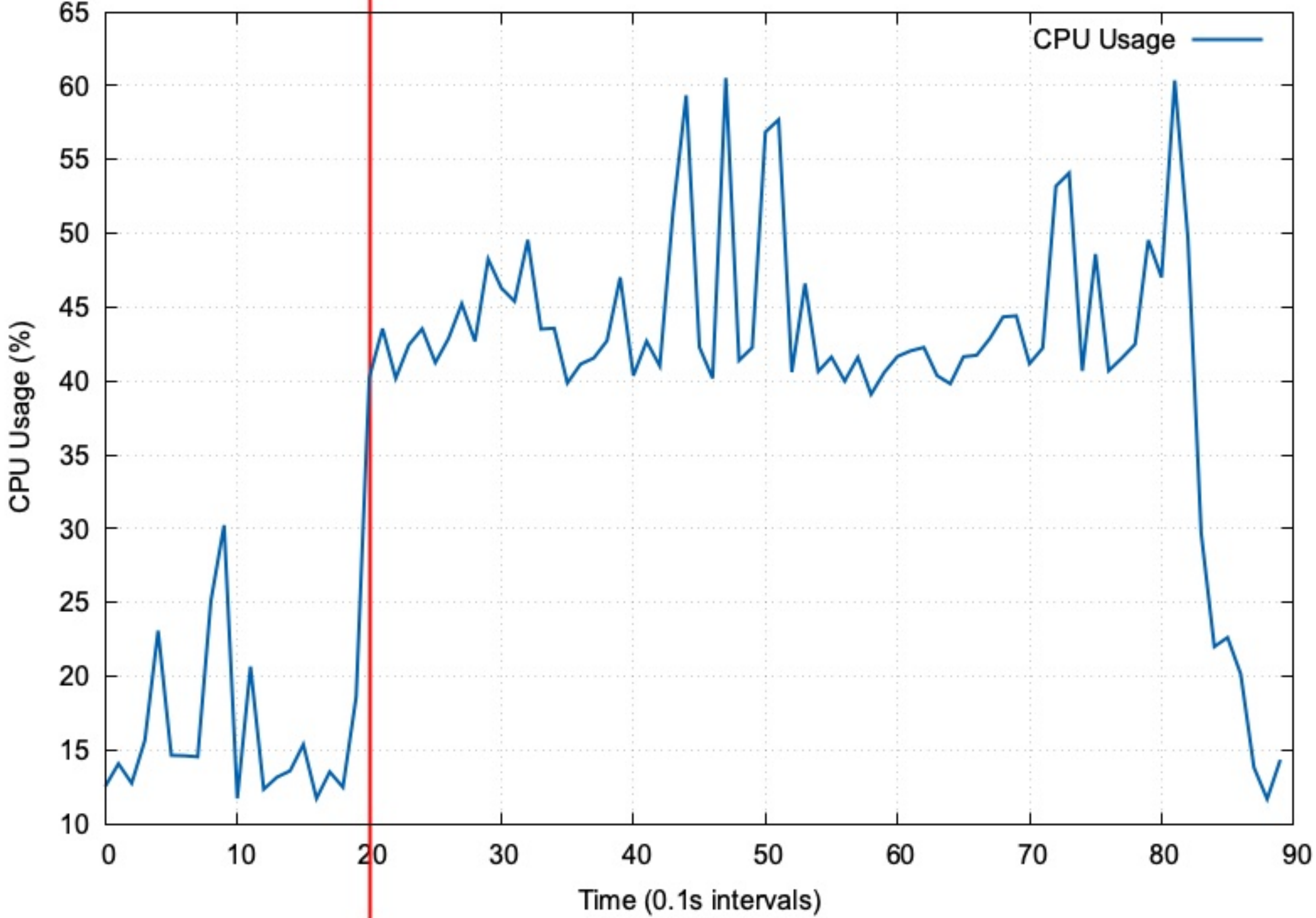
```
+ 30     |> Task.async_stream(  
31       fn val -> Enum.map(val, &parse_row/1) end,  
32 +     max_concurrency: System.schedulers_online(),  
33 +     ordered: false,  
34 +     timeout: :infinity  
35 +   )  
36 +   |> Stream.with_index()  
37 +   |> Task.async_stream(  
38 +     fn {:ok, parsed_rows}, row_index ->  
39 +       interim_records =  
40 +         Enum.reduce(parsed_rows, %{}, fn row, acc ->  
41 +           process_row(row, acc)  
42 +         end)  
43 +  
44 +         :ets.insert(ets_table, {row_index, interim_records})  
45 +     end,  
46 +     max_concurrency: System.schedulers_online(),  
47     ordered: false,  
48     timeout: :infinity  
49   )  
50   |> Stream.run()  
51
```



```
initialize final accumulator
for each entry in ETS table do
  intermediate accumulator = entry.value
  for each key (city) in intermediate accumulator do
    if key exists in final accumulator then
      merge records:
      update min to the lesser of the two min values.
      update max to the greater of the two max values.
      calculate the new sum by adding both sums.
      calculate the new count by adding both counts.
    else
      add record to final accumulator
    end
  end
end
end
```



Overall CPU Usage Over Time (V4, 100M measurements)



CPU usage of version 4



# 1BRC in Elixir: Version 4



10 Million Rows

3.3 sec

↘ 29.8%



100 Million Rows

33 sec

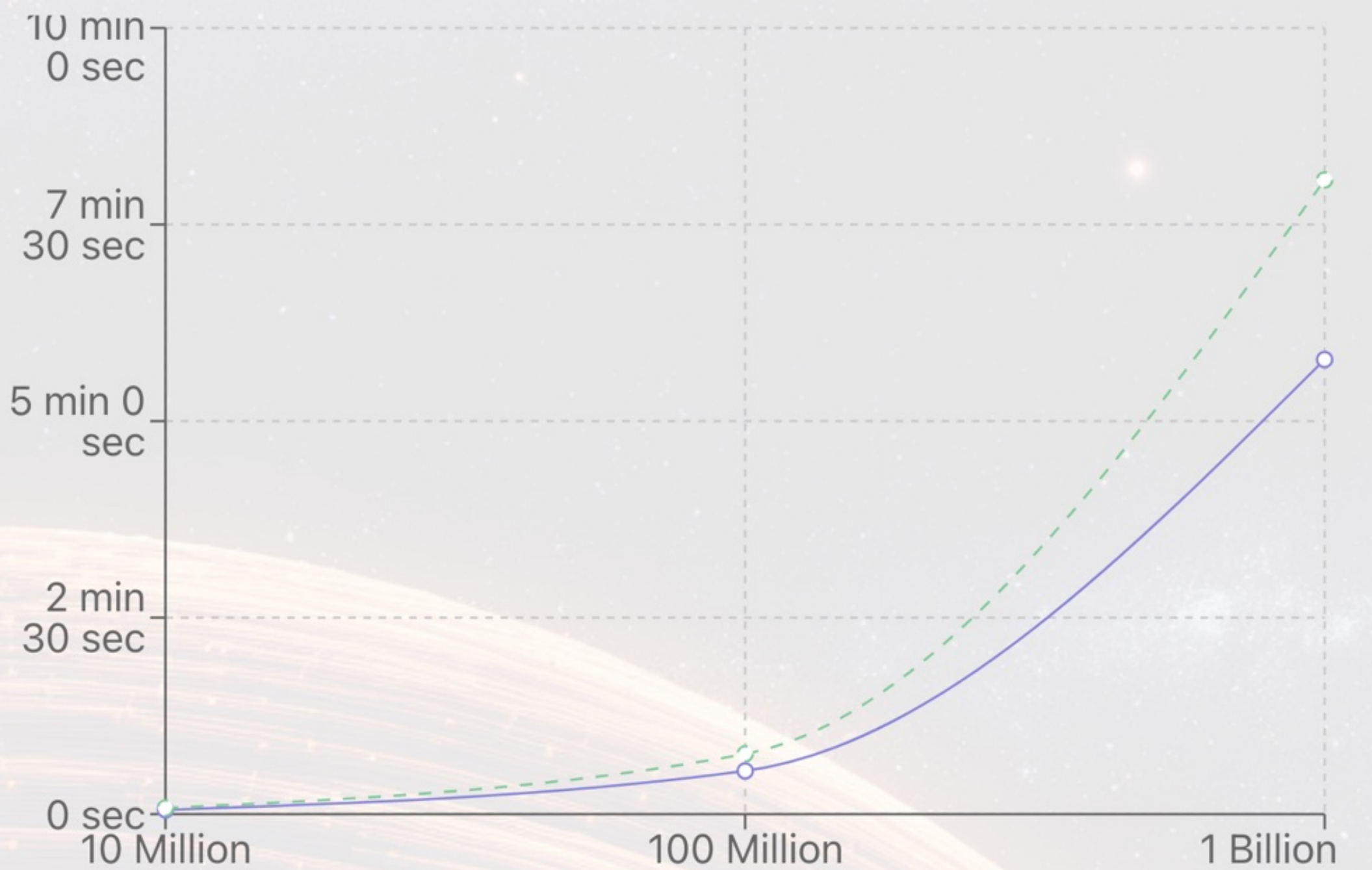
↘ 28.3%

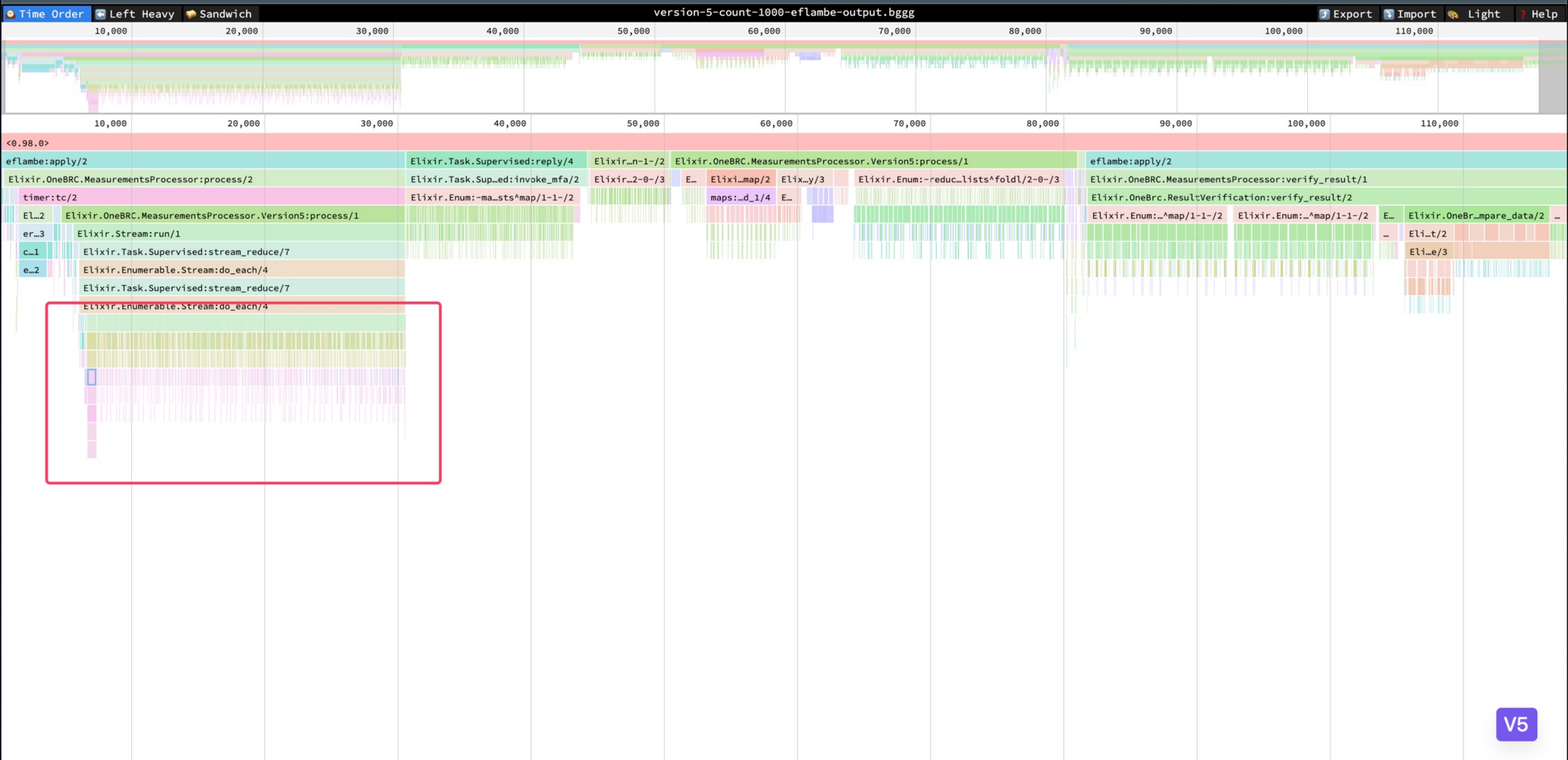


1 Billion Rows

5 min 47 sec

↘ 28.3%

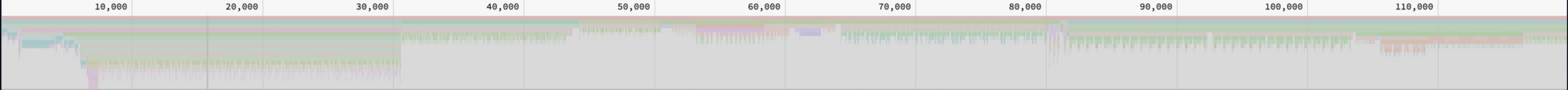




V5

This Instance		All Instances	
Total	Self	Total	Self
10	4	10,837	4,023
<0.01%	<0.01%	9.0%	3.4%

- prim\_file:read\_line/1
- > Elixir.IO:binread/2
- > Elixir.IO:each\_binstream/2
- > Elixir.Stream:do\_resource/5
- > Elixir.Enumerable.Stream:do\_each/4
- > Elixir.Task.Supervised:stream\_reduce/7
- > Elixir.Enumerable.Stream:do\_each/4
- > Elixir.Task.Supervised:stream\_reduce/7
- > Elixir.Stream:run/1
- Elixir.OneBRC.MeasurementsProcessor.Version5:process/1
- > timer:tc/2
- > Elixir.OneBRC.MeasurementsProcessor:process/2



12 15,814 15,816 15,818 15,820 15,822 15,824 15,826 15,828 15,830 15,832 15,834 15,836 15,838 15,840 15,842 15,844 15,846 15,848

<0.98.0>

eflambe:apply/2

Elixir.OneBRC.MeasurementsProcessor:process/2

timer:tc/2

Elixir.OneBRC.MeasurementsProcessor.Version5:process/1

Elixir.Stream:run/1

Elixir.Task.Supervised:stream\_reduce/7

Elixir.Enumerable.Stream:do\_each/4

Elixir.Task.Supervised:stream\_reduce/7

Elixir.Enumerable.Stream:do\_each/4

Elixir.Stream:do\_resource/5

Elixir.Stream:do\_element\_resource/6

Elixir.Stream:do\_resource/5

Elixir.IO:each\_binstream/2

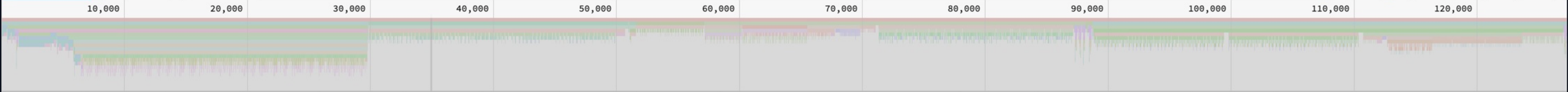
Elixir.IO:binread/2

13 (0.01%) Elixir.IO:binread/2

prim\_file:read\_line/1

prim\_file:read\_line\_1/4

prim\_file\_line\_1/4



34,982 34,984 34,986 34,988 34,990 34,992 34,994 34,996 34,998 35,000 35,002 35,004 35,006 35,008 35,010 35,012 35,014 35,016 35,018

<0.98.0>

Elixir.Task.Supervised:reply/4 127,324 (100%) <0.98.0>

Elixir.Task.Supervised:invoke\_mfa/2

Elixir.Enum:-map/2-lists^map/1-1-/2

Elixir.OneBRC.MeasurementsProcessor.Version4:parse\_row/1 Elixir.OneBRC.MeasurementsProcessor.Version4:parse\_row/1 Elixir.OneBRC.MeasurementsProcessor.Version4:parse\_row/1

Elixir.String.Break:trim\_trailing/2 bi.../2 erlang:bin...integer/1 bi.../2 EL.../1 EL.../1 Elixir.String.Break:trim\_trailing/2 bi.../2 erlang:bin...integer/1 bi.../2 EL.../1 EL.../1

EL.../1 EL.../1 er.../2 EL.../1 EL.../1 er.../2 er.../2

# Version 5: Some micro-optimisations to parsing

- String.split/3, String.trim\_trailing/1 are being called a lot of times.
- Pattern matching to the rescue!

This Instance		All Instances		
Total	Self	Total	Self	
52	7	18,098	5,044	Elixir.OneBRC.MeasurementsProcessor.Version4:parse_row/1
0.04%	<0.01%	14%	4.0%	> Elixir.Enum:-map/2-lists^map/1-1-/2
				> Elixir.Task.Supervised:invoke_mfa/2
				> Elixir.Task.Supervised:reply/4
				> <0.98.0>

```
# sample row: `Mumbai;30.5\n`
```

```
# BEFORE:
```

```
[key, t_value] = :binary.split(row, ";")
```

```
[a, b] = t_value ▷ String.trim_trailing() ▷ :binary.split(".")
```

```
parsed_temp = (a ◊ b) ▷ String.to_integer()
```



main ▾

elixir / lib / elixir / lib / string.ex

↑ Top

Code

Blame

3177 lines (2388 loc) · 94.1 KB



Raw



2949

2950

`String.to_integer("invalid data")`

2951

`** (ArgumentError) argument error`

2952

2953

`"""`

2954

`@spec to_integer(String.t()) :: integer`

...

2955

`def to_integer(string) when is_binary(string) do`

2956

 `:erlang.binary_to_integer(string)`

2957

`end`

2958

2959

`@doc """`

2960

`Returns an integer whose text representation is `string` in base `base`.`

2961

2962

`Inlined by the compiler.`

2963

2964

`## Examples`

2965

2966

`iex> String.to_integer("3FF", 16)`

2967

`1023`

2968

2969

`"""`

2970

`@spec to_integer(String.t(), 2..36) :: integer`

2971

`def to_integer(string, base) when is_binary(string) and is_integer(base) do`

2972

 `:erlang.binary_to_integer(string, base)`

2973

`end`

2974





Code

Blame

12723 lines (10733 loc) · 445 KB



Raw



```
991 as list_to_integer/1 .
992
993 Failure: `badarg` if `Binary` contains a bad representation of an integer.
994 """
995 -doc(#{since => <<"OTP R16B">>}).
996 -doc #{ group => terms }.
997 -spec binary_to_integer(Binary) -> integer() when
998     Binary :: binary().
999 binary_to_integer(Binary) ->
1000     case erts_internal:binary_to_integer(Binary, 10) of
1001         N when erlang:is_integer(N) ->
1002             N;
1003         big ->
1004             case big_binary_to_int(Binary, 10) of
1005                 N when erlang:is_integer(N) ->
1006                     N;
1007                 Reason ->
1008                     error_with_info(Reason, [Binary])
1009             end;
1010         badarg ->
1011             badarg_with_info([Binary])
1012     end.
1013
1014 %% binary_to_integer/2
1015 -doc """
1016 Returns an integer whose text representation in base `Base` is `Binary`.
1017
1018 For example:
```







Code

Blame

3244 lines (2895 loc) · 70.4 KB



Raw



```
3118     }
3119
3120     BIF_RETYPE erts_internal_binary_to_integer_2(BIF_ALIST_2)
3121     {
3122         const byte *temp_alloc = NULL, *bytes;
3123         Uint size;
3124         Uint base;
3125         Eterm res;
3126
3127         if (!is_small(BIF_ARG_2)) {
3128             BIF_RET(am_badarg);
3129         }
3130
3131         base = (Uint)signed_val(BIF_ARG_2);
3132
3133         if (base < 2 || base > 36) {
3134             BIF_RET(am_badarg);
3135         }
3136
3137         bytes = erts_get_aligned_binary_bytes(BIF_ARG_1, &size, &temp_alloc);
3138         if (bytes == NULL) {
3139             BIF_RET(am_badarg);
3140         }
3141
3142         res = chars_to_integer(bytes, size, base);
3143         erts_free_aligned_binary_bytes(temp_alloc);
3144         BIF_RET(res);
3145     }
```





Code

Blame

3244 lines (2895 loc) · 70.4 KB



Raw

3051 `static Eterm chars_to_integer(const byte *bytes, Uint size, const Uint base)`

```
3052 {
3053     Sint i = 0;
3054     int neg = 0;
3055
3056     if (size == 0) {
3057         return am_badarg;
3058     }
3059
3060     if (bytes[0] == '-') {
3061         neg = 1;
3062         bytes++;
3063         size--;
3064     } else if (bytes[0] == '+') {
3065         bytes++;
3066         size--;
3067     }
3068
3069     if (size == 0) {
3070         return am_badarg;
3071     }
3072
3073     /* Trim leading zeroes */
3074     while (*bytes == '0') {
3075         bytes++;
3076         size--;
3077         if (size == 0) {
3078             /* All zero! */
```



Code

Blame

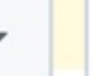
3244 lines (2895 loc) · 70.4 KB



Raw



```
3091     * case.
3092     */
3093     while (size--) {
3094         Uint digit = *bytes++ - '0';
3095         if (digit >= base) {
3096             return am_badarg;
3097         }
3098         i = i * base + digit;
3099     }
3100 } else {
3101     while (size) {
3102         byte b = *bytes++;
3103         size--;
3104
3105         if (c2int_is_invalid_char(b, base)) {
3106             return am_badarg;
3107         }
3108
3109         i = i * base + c2int_digit_from_base(b);
3110     }
3111 }
3112
3113 if (neg) {
3114     i = -i;
3115 }
3116 ASSERT(IS_SSMALL(i));
3117 return make_small(i);
3118 }
```





Code

Blame

3244 lines (2895 loc) · 70.4 KB



```
2956     return (ch >= '0' && ch < ('0' + base));
2957 else
2958     return (ch >= '0' && ch <= '9')
2959         || (ch >= 'A' && ch < ('A' + base - 10))
2960         || (ch >= 'a' && ch < ('a' + base - 10));
2961 }
2962
2963 static ERTS_INLINE int c2int_is_invalid_char(byte ch, int base) {
2964     return !c2int_is_valid_char(ch, base);
2965 }
2966
2967 static ERTS_INLINE byte c2int_digit_from_base(byte ch) {
2968     return ch <= '9' ? ch - '0'
2969         : (10 + (ch <= 'Z' ? ch - 'A' : ch - 'a'));
2970 }
2971
2972 /*
2973  * How many bits are needed to store 1 digit of given base in binary
2974  * Wo.Alpha formula: Table [log2[n], {n,2,36}]
2975  */
2976 static const double lg2_lookup[36-1] = {
2977     1.0, 1.58496, 2.0, 2.32193, 2.58496, 2.80735, 3.0, 3.16993, 3.32193,
2978     3.45943, 3.58496, 3.70044, 3.80735, 3.90689, 4.0, 4.08746, 4.16993, 4.24793,
2979     4.32193, 4.39232, 4.45943, 4.52356, 4.58496, 4.64386, 4.70044, 4.75489,
2980     4.80735, 4.85798, 4.90689, 4.9542, 5.0, 5.04439, 5.08746, 5.12928, 5.16993
2981 };
2982
2983 /*
```



iex(1)> ?A  
65

iex(2)> ?B  
66

iex(3)> ?0  
48

iex(4)> ?1  
49

iex(5)> ?2  
50

iex(6)> ?3  
51

```
iex(1)> ?5 - ?0  
5
```

- One digit before the decimal ( '4' )
- Decimal point ( '.' )
- One digit after the decimal ( '5' )

4.5

- One digit before the decimal ( '4' )
- Decimal point ( '.' )
- One digit after the decimal ( '5' )
- <<d1, ?., d2, \_::binary>>
  - *Also eliminates the need for String.trim\_trailing*

4.5

```
defp parse_temperature(<<d1, ?., d2, _::binary>>) do
  char_to_num(d1) * 10 + char_to_num(d2)
end
```



```
# sample row: `Mumbai;30.5\n`
```

```
# AFTER:
```

```
parsed_temp = t_value ▷ parse_temperature()
```

```
...
```

```
# ex: -4.5
```

```
defp parse_temperature(<<?-, d1, ?., d2, _::binary>>) do
```

```
  -(char_to_num(d1) * 10 + char_to_num(d2))
```

```
end
```

```
# ex: 4.5
```

```
defp parse_temperature(<<d1, ?., d2, _::binary>>) do
```

```
  char_to_num(d1) * 10 + char_to_num(d2)
```

```
end
```

```
# ex: -45.3
```

```
defp parse_temperature(<<?-, d1, d2, ?., d3, _::binary>>) do
```

```
  -(char_to_num(d1) * 100 + char_to_num(d2) * 10 + char_to_num(d3))
```

```
end
```

```
# ex: 45.3
```

```
defp parse_temperature(<<d1, d2, ?., d3, _::binary>>) do
```

```
  char_to_num(d1) * 100 + char_to_num(d2) * 10 + char_to_num(d3)
```

```
end
```

```
defp char_to_num(char) do
```

```
  char - ?0
```

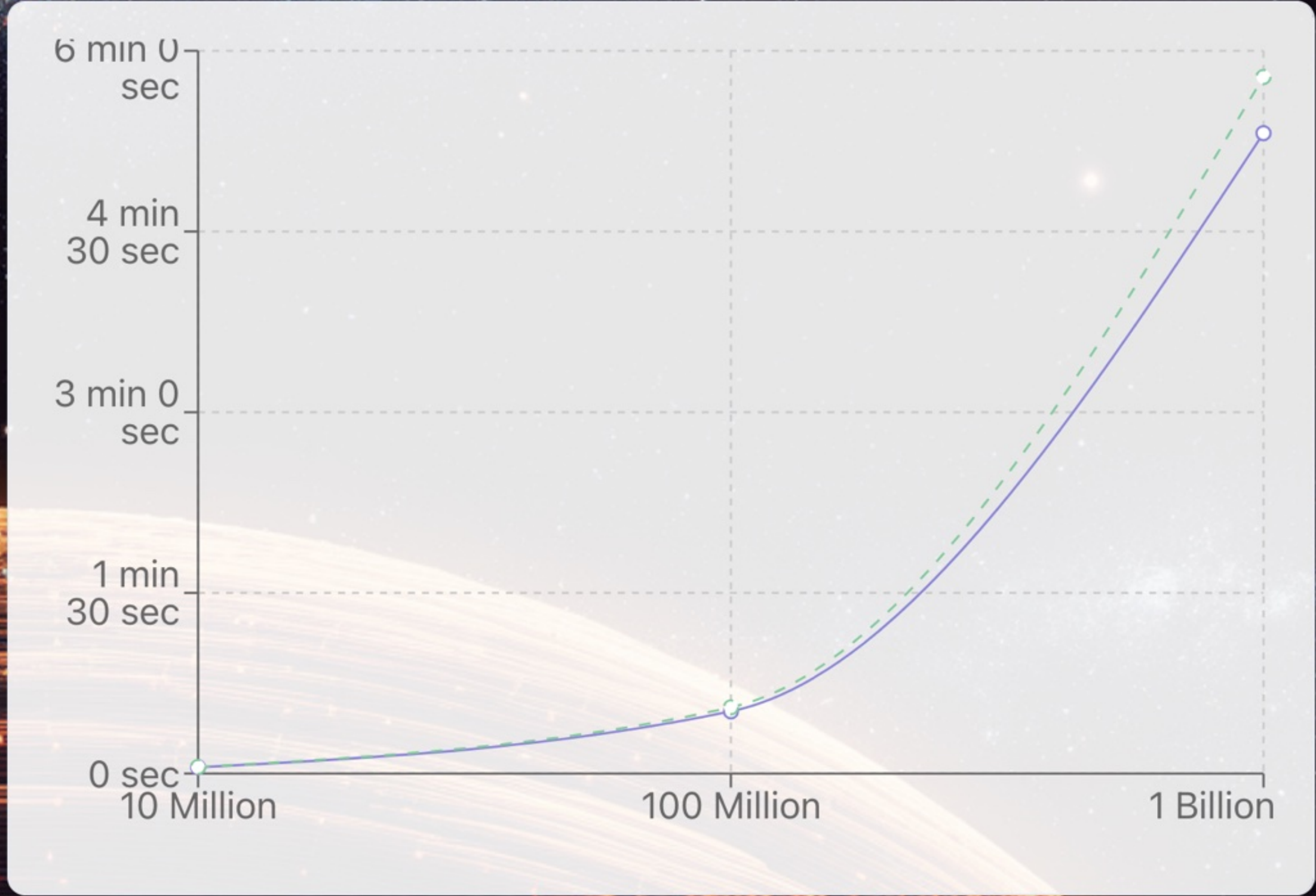
```
end
```

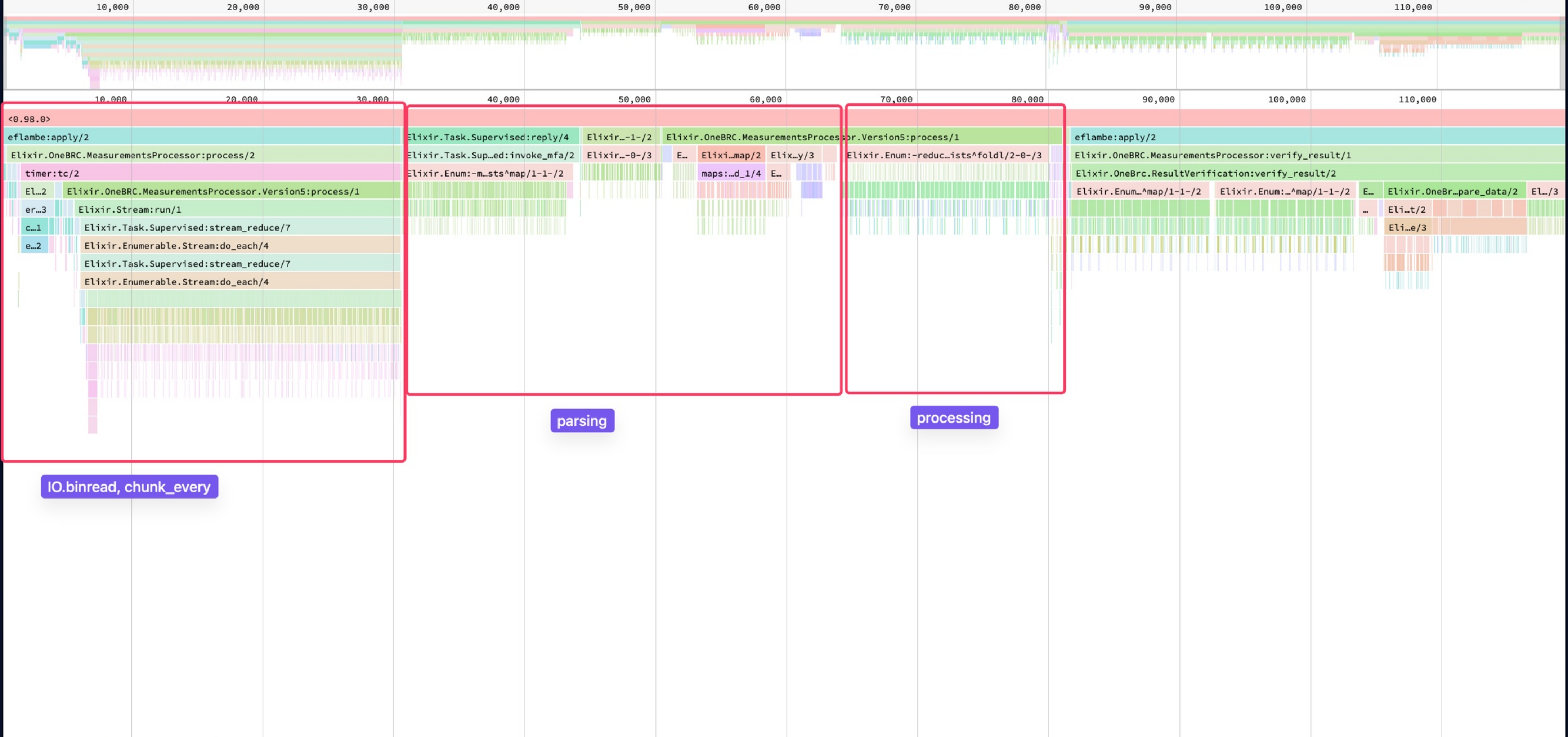
# 1BRC in Elixir: Version 5

**10 Million Rows**  
🕒 **3.1 sec**  
↘ **6.1%**

**100 Million Rows**  
🕒 **31 sec**  
↘ **6.1%**

**1 Billion Rows**  
🕒 **5 min 19 sec**  
↘ **8.1%**





IO.binread, chunk\_every

parsing

processing

This Instance		All Instances	
Total	Self	Total	Self
1	1	1,005	1,005
<0.01%	<0.01%	0.84%	0.84%

- > file:read\_line/1
- > Elixir.IO:binread/2
- > Elixir.IO:each\_binstream/2
- > Elixir.Stream:do\_resource/5
- > Elixir.Enumerable.Stream:do\_each/4
- > Elixir.Task.Supervised:stream\_reduce/7
- > Elixir.Enumerable.Stream:do\_each/4
- > Elixir.Task.Supervised:stream\_reduce/7
- > Elixir.Stream:run/1
- > Elixir.OneBRC.MeasurementsProcessor.Version5:process/1
- > timer:tc/2
- > Elixir.OneBRC.MeasurementsProcessor:process/2

```
stream!(path, line_or_bytes, modes)
```

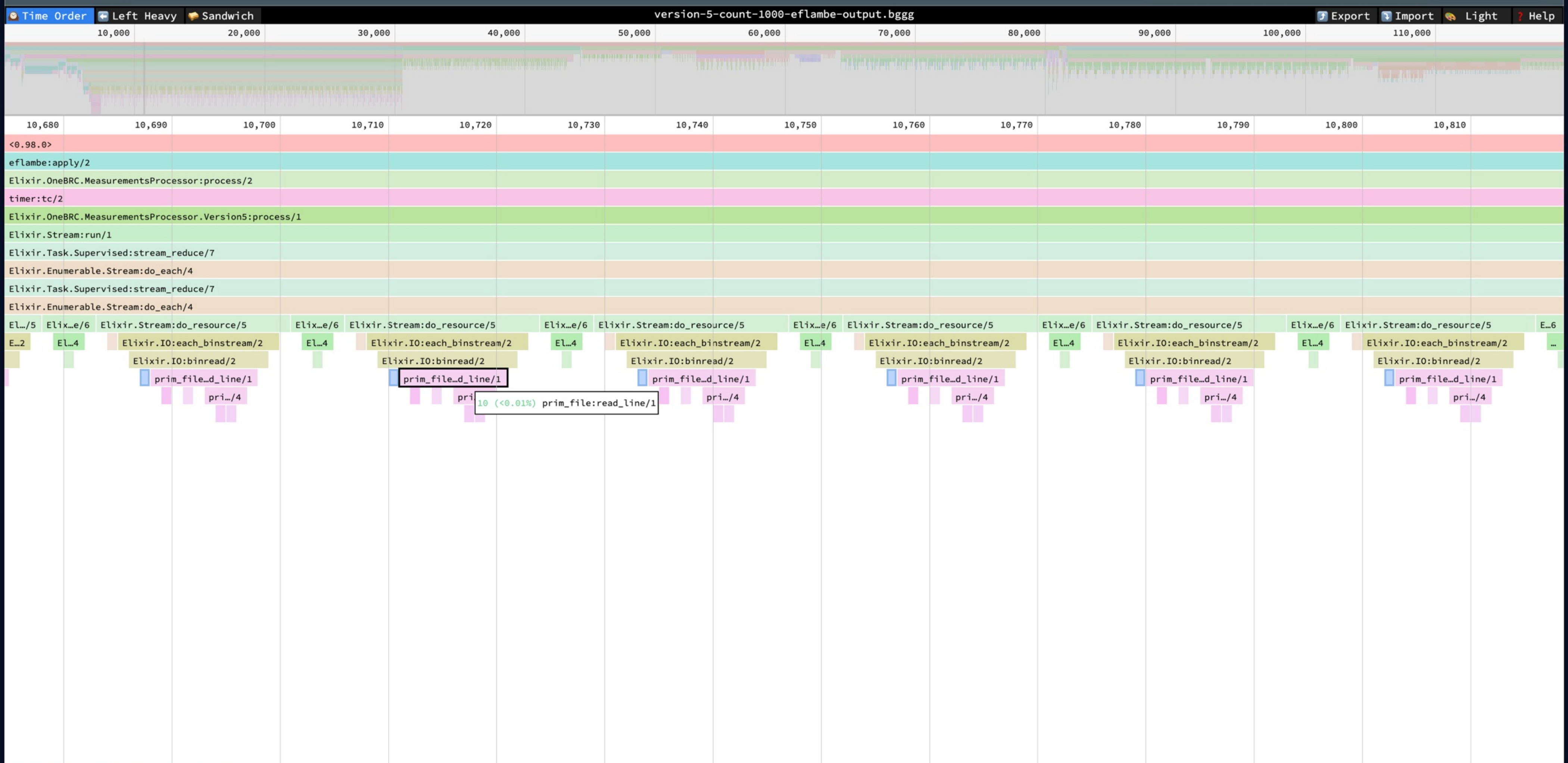
```
</>
```

```
@spec stream!(Path.t(), :line | pos_integer(), [stream_mode()]) :: File.Stream.t()
```

Returns a `File.Stream` for the given `path` with the given `modes`.

The stream implements both `Enumerable` and `Collectable` protocols, which means it can be used both for read and write.

The `line_or_bytes` argument configures how the file is read when streaming, by `:line` (default) or by a given number of bytes. When using the `:line` option, CRLF line breaks (`"\r\n"`) are normalized to LF (`"\n"`).



This Instance		All Instances		
Total	Self	Total	Self	
1	1	1,005	1,005	file:read_line/1
<0.01%	<0.01%	0.84%	0.84%	> Elixir.IO:binread/2
				> Elixir.IO:each_binstream/2
				> Elixir.Stream:do_resource/5
				> Elixir.Enumerable.Stream:do_each/4
				> Elixir.Task.Supervised:stream_reduce/7
				> Elixir.Enumerable.Stream:do_each/4
				> Elixir.Task.Supervised:stream_reduce/7
				> Elixir.Stream:run/1
				> Elixir.OneBRC.MeasurementsProcessor.Version5:process/1
				> timer:tc/2
				> Elixir.OneBRC.MeasurementsProcessor:process/2

...

Tokyo;16.8\nCape Town;22.3\nStockholm;7.1\nMarrakech;29.6\nBerlin;9.7\nAuckland;17.5\nBangkok;32.7\nReykjavik;3.4\nVienna;14.6\nMoscow;1.9\nBerlin;6.3\nHanoi;30.2\nCairo;33.5\nOslo;5.7\nKyoto;18.1\nMiami;27.6\nBerlin;11.2\nPrague;15.9\nIstanbul;23.4\nEdinburgh;10.8\nSeattle;13.7\nDubai;36.8\nQueenstown;9.9

...

...

Tokyo;16.8\n	Cape Town;22.3\n	Stockholm;7.1\n
Marrakech;29.6\n	Berlin;9.7\n	Auckland;17.5\n
Bangkok;32.7\n	Reykjavik;3.4\n	Vienna;14.6\n

Moscow;1.9\nBerlin;6.3\nHanoi;30.2\n

Cairo;33.5\nOslo;5.7\nKyoto;18.1\n

Miami;27.6\nBerlin;11.2\nPrague;15.9\n

Istanbul;23.4\nEdinburgh;10.8\nSeattle;13.7\n

Dubai;36.8\nQueenstown;9.9\n

...

...

Tokyo;16.8	Cape Town;22.3	Stockholm;7.1
Marrakech;29.6	Berlin;9.7	Auckland;17.5
Bangkok;32.7	Reykjavik;3.4	Vienna;14.6

Moscow;1.9

Cairo;33.5

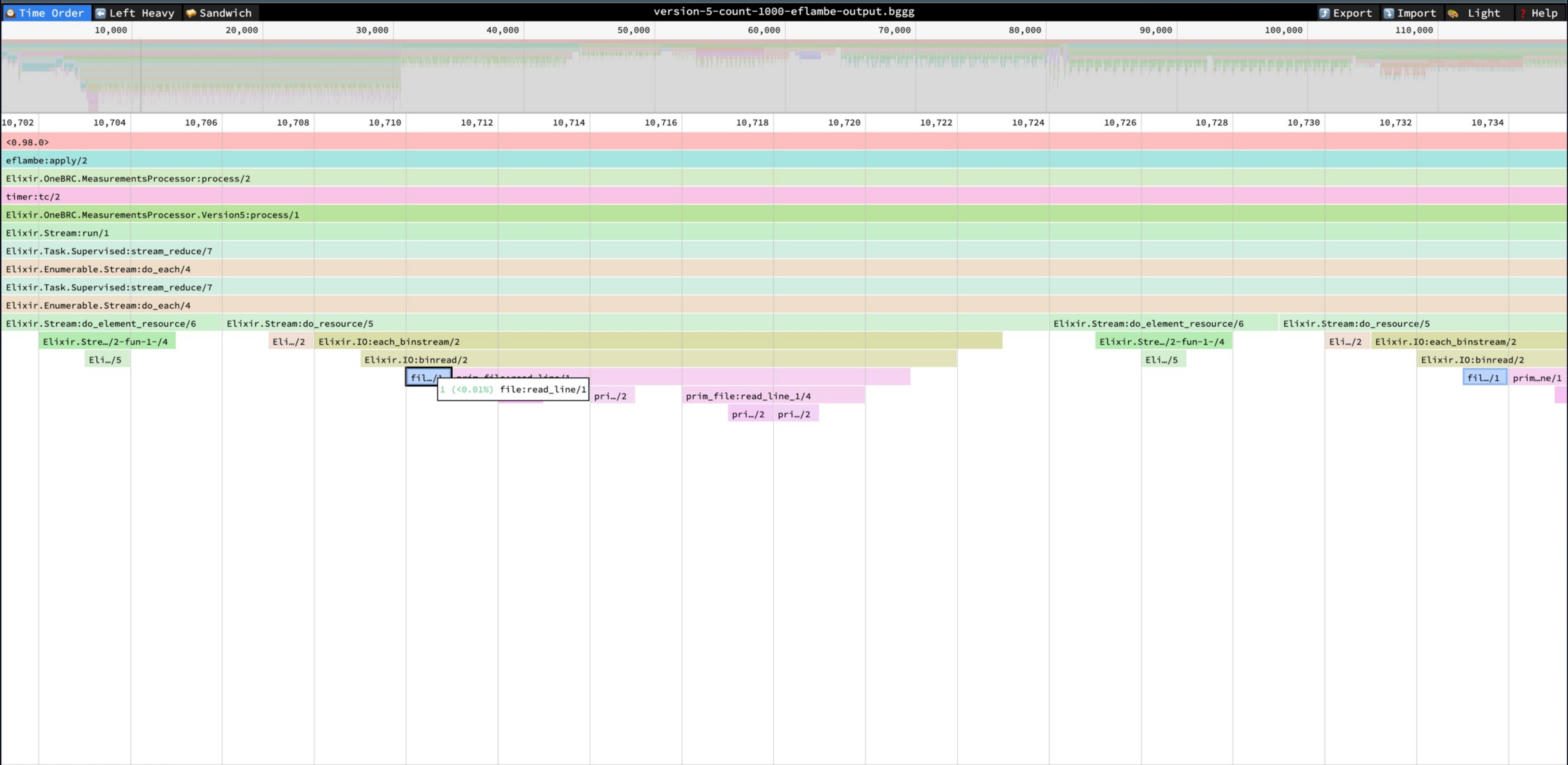
Miami;27.6

Istanbul;23.4

Dubai;36.8

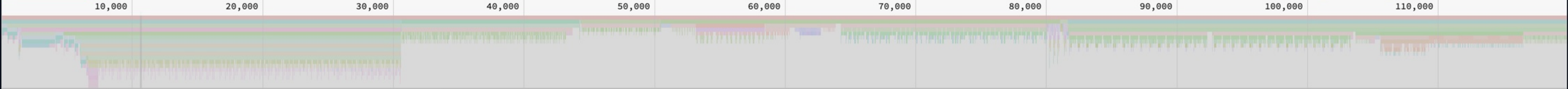
...





This Instance		All Instances	
Total	Self	Total	Self
1	1	1,005	1,005
<0.01%	<0.01%	0.84%	0.84%

- file:read\_line/1
- > Elixir.IO:binread/2
- > Elixir.IO:each\_binstream/2
- > Elixir.Stream:do\_resource/5
- > Elixir.Enumerable.Stream:do\_each/4
- > Elixir.Task.Supervised:stream\_reduce/7
- > Elixir.Enumerable.Stream:do\_each/4
- > Elixir.Task.Supervised:stream\_reduce/7
- > Elixir.Stream:run/1
- > Elixir.OneBRC.MeasurementsProcessor.Version5:process/1



10,700	10,705	10,710	10,715	10,720	10,725	10,730	10,735	10,740
<code>&lt;0.98.0&gt;</code>								
<code>eflambe:apply/2</code>								
<code>Elixir.OneBRC.MeasurementsProcessor:process/2</code>								
<code>timer:tc/2</code>								
<code>Elixir.OneBRC.MeasurementsProcessor.Version5:process/1</code>								
<code>Elixir.Stream:run/1</code>								
<code>Elixir.Task.Supervised:stream_reduce/7</code>								
<code>Elixir.Enumerable.Stream:do_each/4</code>								
<code>Elixir.Task.Supervised:stream_reduce/7</code>								
<code>Elixir.Enumerable.Stream:do_each/4</code>								
<code>Elixir.Stream:do_resource/5</code>			<code>Elixir.Stream:do_resource/6</code>			<code>Elixir.Stream:do_resource/5</code>		
<code>Elixir.IO:binread/2</code>		<code>Elixir.Stream:do_resource/6</code>		<code>Elixir.IO:binread/2</code>		<code>Elixir.Stream:do_resource/5</code>		<code>Elixir.IO:binread/2</code>
<code>Elixir.IO:binread/2</code>		<code>Elixir.S.un-1-/4</code>		<code>Elixir.S.un-1-/4</code>		<code>Elixir.IO:binread/2</code>		<code>Elixir.IO:binread/2</code>
<code>prim_/1</code>		<code>EL5</code>		<code>EL5</code>		<code>EL5</code>		<code>EL5</code>
<code>fi...1 prim_file:read_line/1</code>		<code>EL2 Elixir.IO:each_binstream/2</code>		<code>Elixir.S.un-1-/4</code>		<code>EL2 Elixir.IO:each_binstream/2</code>		<code>EL2 Elixir.IO:each_binstream/2</code>
<code>pr...10 (&lt;0.01%) prim_file:read_line/1_line_1/4</code>		<code>Elixir.IO:binread/2</code>		<code>Elixir.S.un-1-/4</code>		<code>Elixir.IO:binread/2</code>		<code>Elixir.IO:binread/2</code>
<code>pr...2 pr...2</code>		<code>EL5</code>		<code>EL5</code>		<code>EL5</code>		<code>EL5</code>
<code>fi...1 prim_file:read_line/1</code>		<code>EL2 Elixir.IO:each_binstream/2</code>		<code>Elixir.S.un-1-/4</code>		<code>EL2 Elixir.IO:each_binstream/2</code>		<code>EL2 Elixir.IO:each_binstream/2</code>
<code>pr...1 pr...2 prim_file:r_d_line_1/4</code>		<code>Elixir.IO:binread/2</code>		<code>Elixir.S.un-1-/4</code>		<code>Elixir.IO:binread/2</code>		<code>Elixir.IO:binread/2</code>
<code>pr...2 pr...2</code>		<code>EL5</code>		<code>EL5</code>		<code>EL5</code>		<code>EL5</code>

This Instance		All Instances	
Total	Self	Total	Self
1	1	1,005	1,005
<0.01%	<0.01%	0.84%	0.84%

- file:read\_line/1
- > Elixir.IO:binread/2
- > Elixir.IO:each\_binstream/2
- > Elixir.Stream:do\_resource/5
- > Elixir.Enumerable.Stream:do\_each/4
- > Elixir.Task.Supervised:stream\_reduce/7
- > Elixir.Enumerable.Stream:do\_each/4
- > Elixir.Task.Supervised:stream\_reduce/7
- > Elixir.Stream:run/1
- > Elixir.OneBRC.MeasurementsProcessor.Version5:process/1



nato

Aug/2023

I swapped out some `file:read_file` and `file:write_files` today, for some `prim_file` ones, and I noticed an astounding speedup. Other than the the `don't use this module` standard reply, can I get some context on this module, and why it's there (along with the other erts prim modules)?

Appreciate the info!!

4  

629 views 18 likes 2 links 7 users     



mikpe

Aug/2023

`prim_file` is a building block for the so-called "file server". The proper question to ask is why does Erlang need the file server. (I don't *know* the answer to that, and I prefer not to speculate.)

1  

Aug 2023

1 of 13  
Aug 2023

Aug/2023

# Prim\_file speedup



...adding need the file server (and then the answer to that, and I prefer not to speculate)

1

Aug 2023

2 of 13  
Aug 2023



raimo Erlang Core Team

Aug/2023

The purpose of this indirection is to be able to run on a diskless machine, so a node can be configured to run with a remote file server. Therefore `file` calls a file server process for all its operations. The file server uses `prim_file` as backend to its local file system.

Applications then does not need to know on which machine the file system is, as long as they are using `file` to access it.

There is a `raw` file mode option that can be used for most file operations, e.g. for `file:write_file/3`. Try it and see if the performance is close enough to `prim_file:write_file/3`.

Unfortunately; `file:read_file/1` does not today have an arity that takes a `Mode` argument...

2 Replies

5

Aug/2023



mikpe

raimo Aug/2023

It's one thing to have an indirection to add flexibility. it's another to indirect through a single (\*)



Code

Blame

893 lines (790 loc) · 27.9 KB



Raw



203 read\_line(Fd) -&gt;



```
204     try
205         #{ handle := FRef,
206           r_ahead_size := RASz,
207           r_buffer := RBuf } = get_fd_data(Fd),
208         SearchResult = prim_buffer:find_byte_index(RBuf, $\n),
209         LineSize = max(?MIN_READLINE_SIZE, RASz),
210         read_line_1(FRef, RBuf, SearchResult, LineSize)
211     catch
212         error:badarg -> {error, badarg}
213     end.
214
215 -spec read_line_1(FRef, RBuf, SearchResult, LineSize) -> Result when
216     FRef :: prim_file_ref(),
217     RBuf :: term(),
218     SearchResult :: not_found | {ok, non_neg_integer()},
219     LineSize :: non_neg_integer(),
220     Result :: eof | {ok, binary()} | {error, Reason :: atom()}.
221 read_line_1(FRef, RBuf, not_found, LineSize) ->
222     case read_nif(FRef, LineSize) of
223     {ok, Data} ->
224         prim_buffer:write(RBuf, [Data]),
225         SearchResult = prim_buffer:find_byte_index(RBuf, $\n),
226         read_line_1(FRef, RBuf, SearchResult, LineSize);
227     eof ->
228         case prim_buffer:size(RBuf) of
229         Size when Size > 0 -> {ok, prim_buffer:read(RBuf, Size)};
230         Size when Size == 0 -> eof
```

# :prim\_file

- **open/2:** `open(Name, Modes)`  
Opens a file and returns a file descriptor.
- **read\_line/1:** `read_line(Fd)`  
Reads a line from a file descriptor until a newline (`\n`) character is encountered.
- **read/2:** `read(Fd, Size)`  
Reads a specified number of bytes from an open file descriptor.

```
{ok, Fd} = prim_file:open("file.txt", [read, binary]),  
% Do something with Fd  
prim_file:close(Fd).
```

# :prim\_file

- **open/2:** `open(Name, Modes)`  
Opens a file and returns a file descriptor.
- **read\_line/1:** `read_line(Fd)`  
Reads a line from a file descriptor until a newline (`\n`) character is encountered.
- **read/2:** `read(Fd, Size)`  
Reads a specified number of bytes from an open file descriptor.

```
{ok, Fd} = prim_file:open("file.txt", [read, binary]),
case prim_file:read_line(Fd) of
  {ok, Line} → io:format("Read line: ~p~n", [Line]);
  eof → io:format("End of file reached~n")
end,
prim_file:close(Fd).
```

# :prim\_file

- **open/2:** `open(Name, Modes)`  
Opens a file and returns a file descriptor.
- **read\_line/1:** `read_line(Fd)`  
Reads a line from a file descriptor until a newline (`\n`) character is encountered.
- **read/2:** `read(Fd, Size)`  
Reads a specified number of bytes from an open file descriptor.

```
{ok, Fd} = prim_file:open("file.txt", [read, binary]),  
{ok, Data} = prim_file:read(Fd, 100),  
prim_file:close(Fd).
```



# :prim\_file

- **open/2:** `open(Name, Modes)`  
Opens a file and returns a file descriptor.
- **read\_line/1:** `read_line(Fd)`  
Reads a line from a file descriptor until a newline (`\n`) character is encountered.
- **read/2:** `read(Fd, Size)`  
Reads a specified number of bytes from an open file descriptor.
- Taking the middle road:
  - use `read/2` to take in a fixed chunk of bytes
  - then use `read_line/1` to finish till the next `\n`

```
{ok, Fd} = prim_file:open("file.txt", [read, binary]),  
{ok, Data} = prim_file:read(Fd, 100),  
prim_file:close(Fd).
```



Code

Blame

194 lines (156 loc) · 5.22 KB



Raw



```
1  defmodule OneBRC.MeasurementsProcessor.Version5 do
13  def process(count) do
16    fs = File.stream!(file_path)
17
18    ets_table = :ets.new(:station_stats, [:set, :public])
19
20    fs
21    |> Stream.chunk_every(10000)
22    |> Task.async_stream(
23      fn val -> Enum.map(val, &parse_row/1) end,
24      max_concurrency: System.schedulers_online(),
25      ordered: false,
26      timeout: :infinity
27    )
28    |> Stream.with_index()
29    |> Task.async_stream(
30      fn {{:ok, parsed_rows}, row_index} ->
31        interim_records =
32          Enum.reduce(parsed_rows, %{}, fn row, acc ->
33            process_row(row, acc)
34          end)
35
36        :ets.insert(ets_table, {row_index, interim_records})
37      end,
38      max_concurrency: System.schedulers_online(),
39      ordered: false,
40      timeout: :infinity
41    )
```





Code

Blame

239 lines (190 loc) · 6.13 KB



Raw



```
107
108 ✓ defp read_and_process(file, ets_table, tasks) do
109     chunk_size = 1024 * 1024 * 1
110
111     data =
112     case :prim_file.read(file, chunk_size) do
113         :eof ->
114             nil
115
116         {:ok, data} ->
117             case :prim_file.read_line(file) do
118                 {:ok, line} ->
119                     <<data::binary, line::binary>>
120
121                 :eof ->
122                     data
123             end
124     end
125
126 ✓ if is_nil(data) do
127     tasks
128 else
129     task = Task.async(fn -> process_chunk(data, ets_table) end)
130
131     read_and_process(file, ets_table, [task | tasks])
132 end
133 end
134
```





Code

Blame

239 lines (190 loc) · 6.13 KB



Raw



```
11     require Logger
12
13     def process(count) do
14         t1 = System.monotonic_time(:millisecond)
15         file_path = measurements_file(count)
16
17         {:ok, file} = :prim_file.open(file_path, [:raw, :binary, :read])
18
19         ets_table = :ets.new(:station_stats, [:duplicate_bag, :public])
20
21         tasks = read_and_process(file, ets_table, [])
22
23         Task.await_many(tasks, :infinity)
24         # old way ->
25         # fs
26         # |> Stream.chunk_every(10000)
27         # |> Task.async_stream(
28         #   fn val -> Enum.map(val, &parse_row/1) end,
29         #   max_concurrency: System.schedulers_online(),
30         #   ordered: false,
31         #   timeout: :infinity
32         # )
33         # |> Stream.with_index()
34         # |> Task.async_stream(
35         #   fn {{:ok, parsed_rows}, row_index} ->
36         #     interim_records =
37         #       Enum.reduce(parsed_rows, %{}, fn row, acc ->
38         #         process_row(row, acc)
```





Code

Blame

239 lines (190 loc) · 6.13 KB



```
13     def process(count) do
51         result =
52             :ets.tab2list(ets_table)
53             |> Enum.reduce([], fn {_, m}, acc -> [acc | Enum.into(m, [])] end)
54             |> List.flatten()
55             |> Enum.reduce(%{}, fn {key, val}, acc ->
56                 existing_record = Map.get(acc, key, nil)
57
58                 new_record =
59                 case existing_record do
60                     nil ->
61                         val
62
63                     %{count: count, min: min, max: max, mean: mean} ->
64                         min = if val.min < min, do: val.min, else: min
65                         max = if val.max > max, do: val.max, else: max
66                         new_c = count + val.count
67
68                         mean = (mean * count + val.mean * val.count) / new_c
69
70                         %{
71                             min: min,
72                             max: max,
73                             mean: mean,
74                             count: new_c
75                         }
76             end
77
```

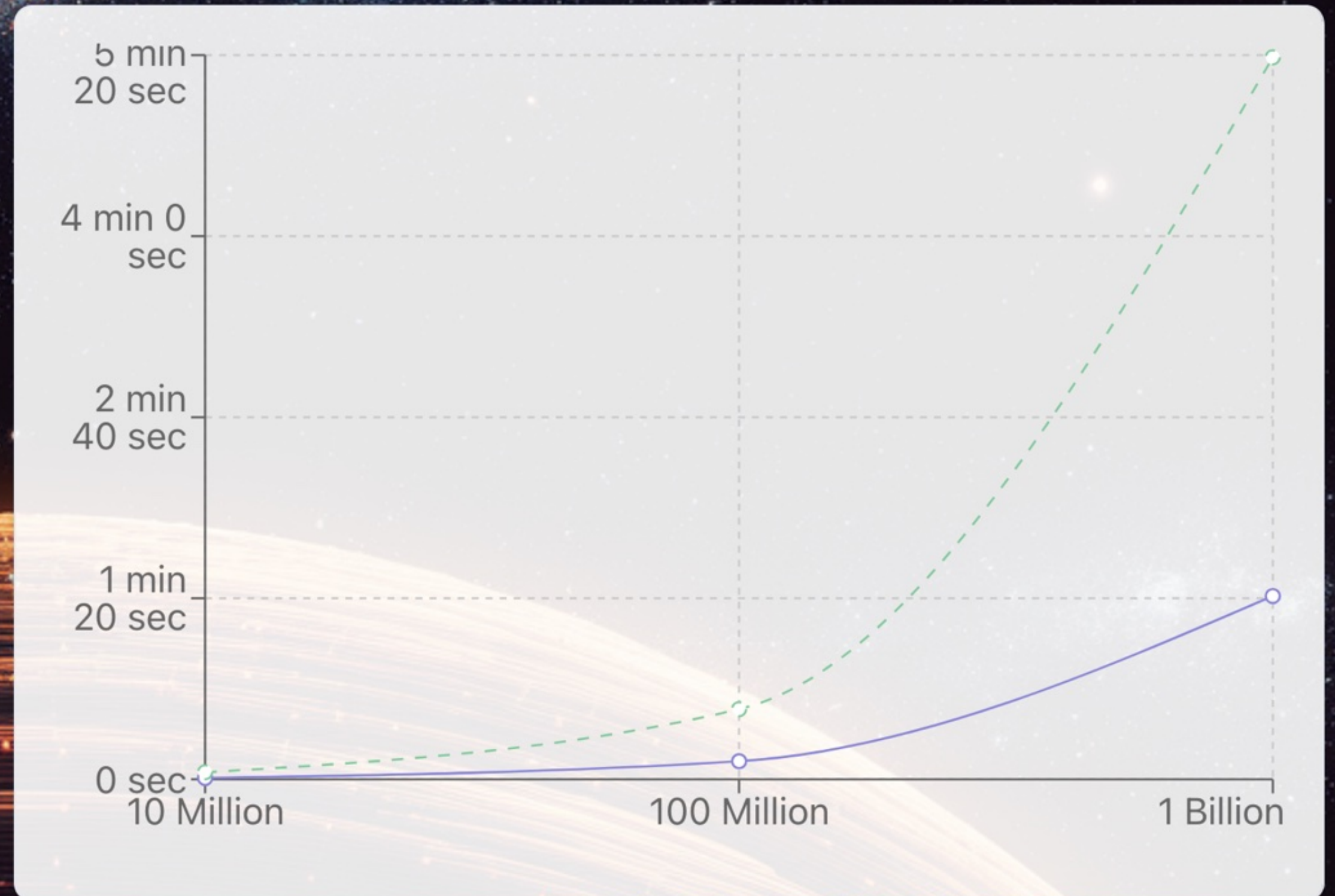


# 1BRC in Elixir: Version 6

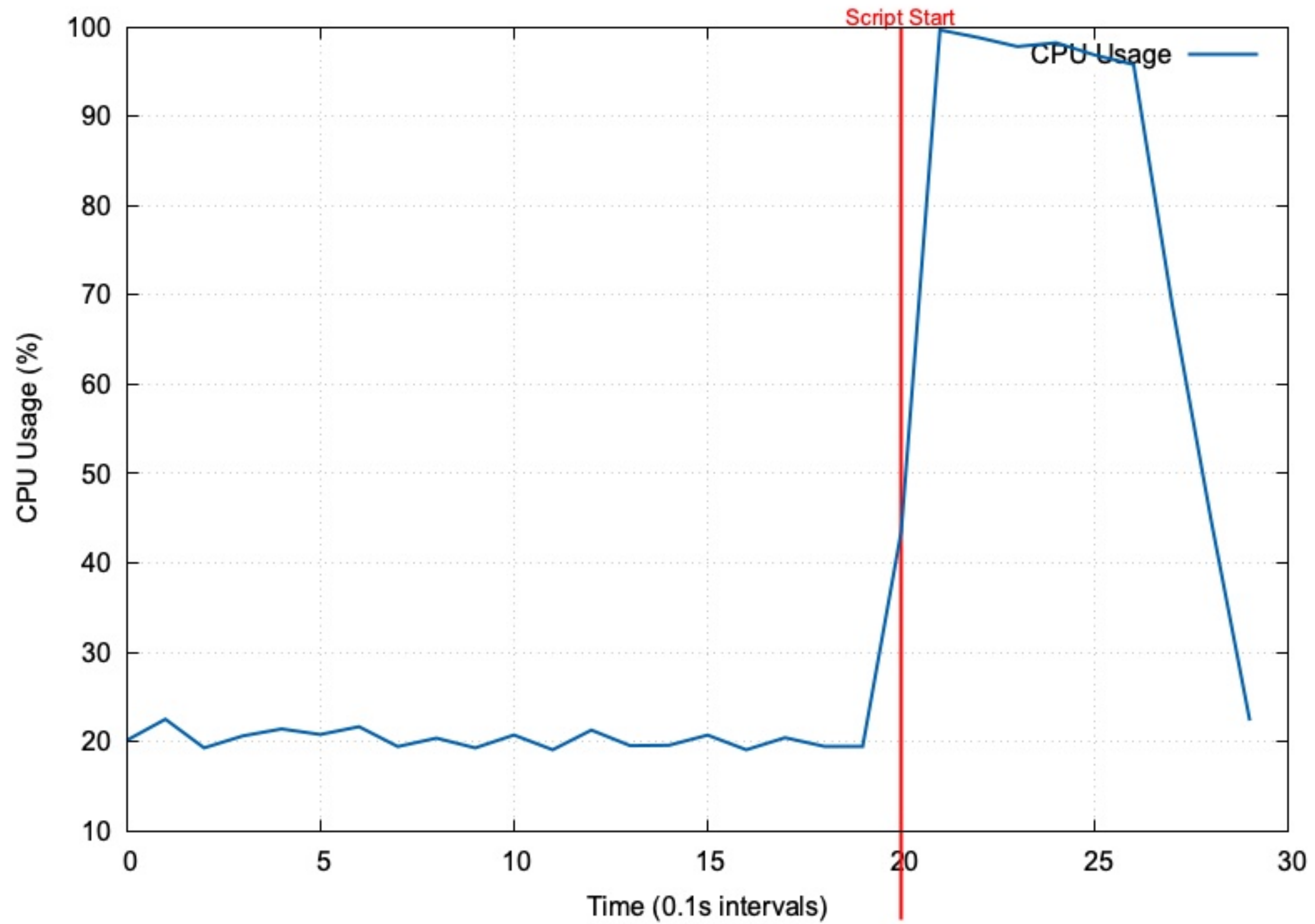
**10 Million Rows**  
🕒 **0.6 sec**  
↘ **80.6%**

**100 Million Rows**  
🕒 **8 sec**  
↘ **74.2%**

**1 Billion Rows**  
🕒 **1 min 21 sec**  
↘ **74.6%**

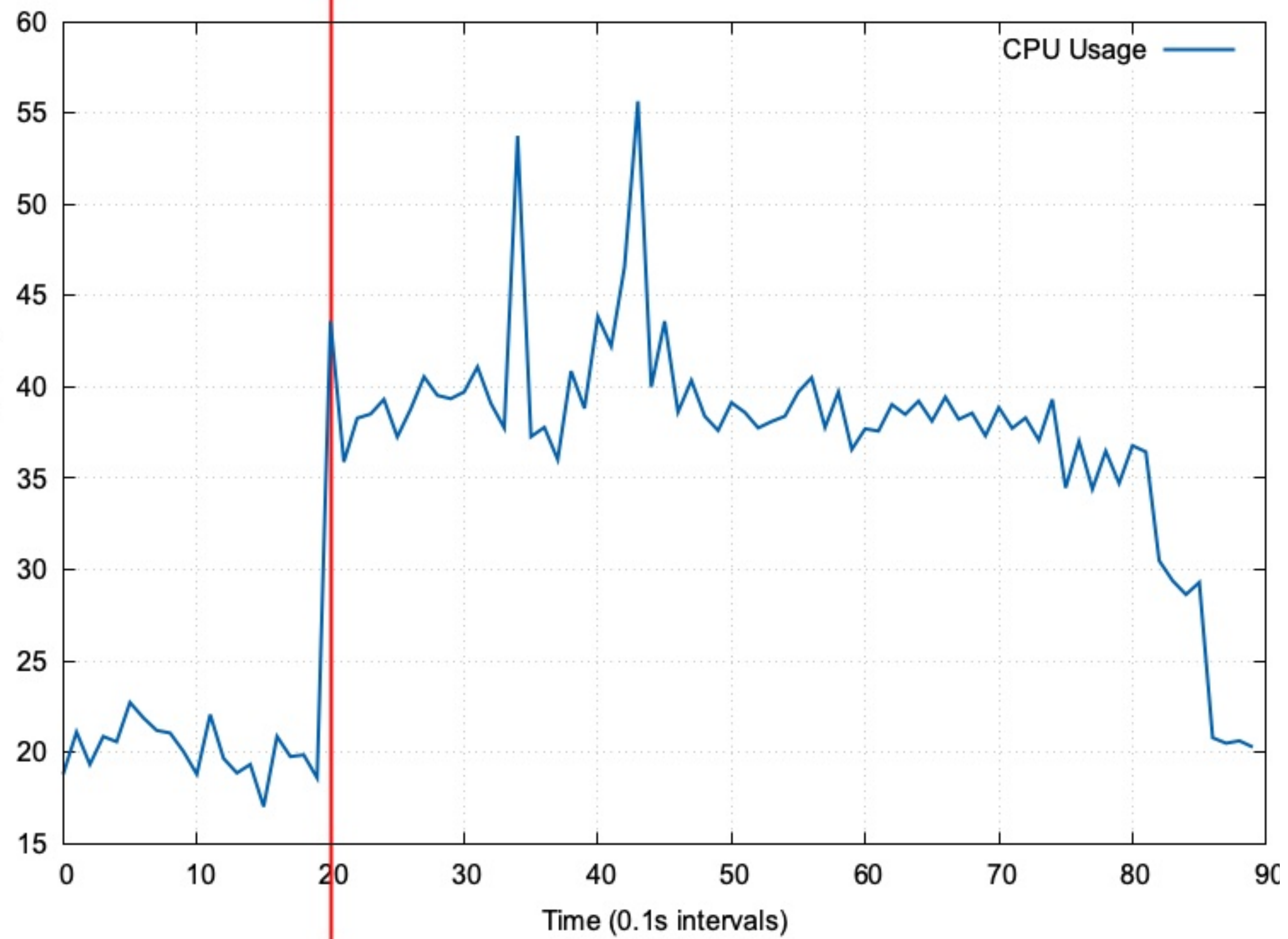


Overall CPU Usage Over Time (V6, 100M measurements)

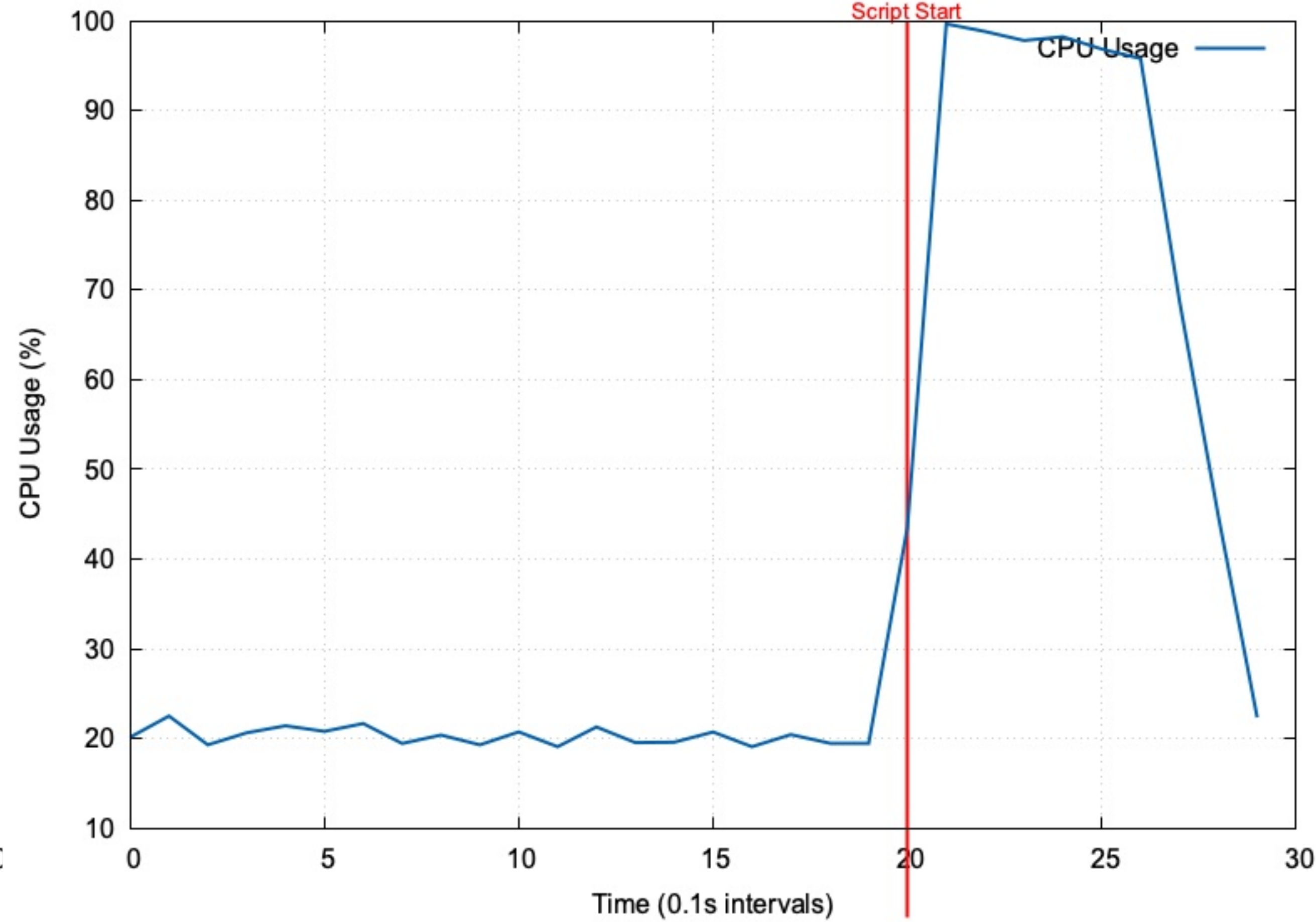


*CPU usage of version 6*

Overall CPU Usage Over Time (V5, 100M measurements)



Overall CPU Usage Over Time (V6, 100M measurements)

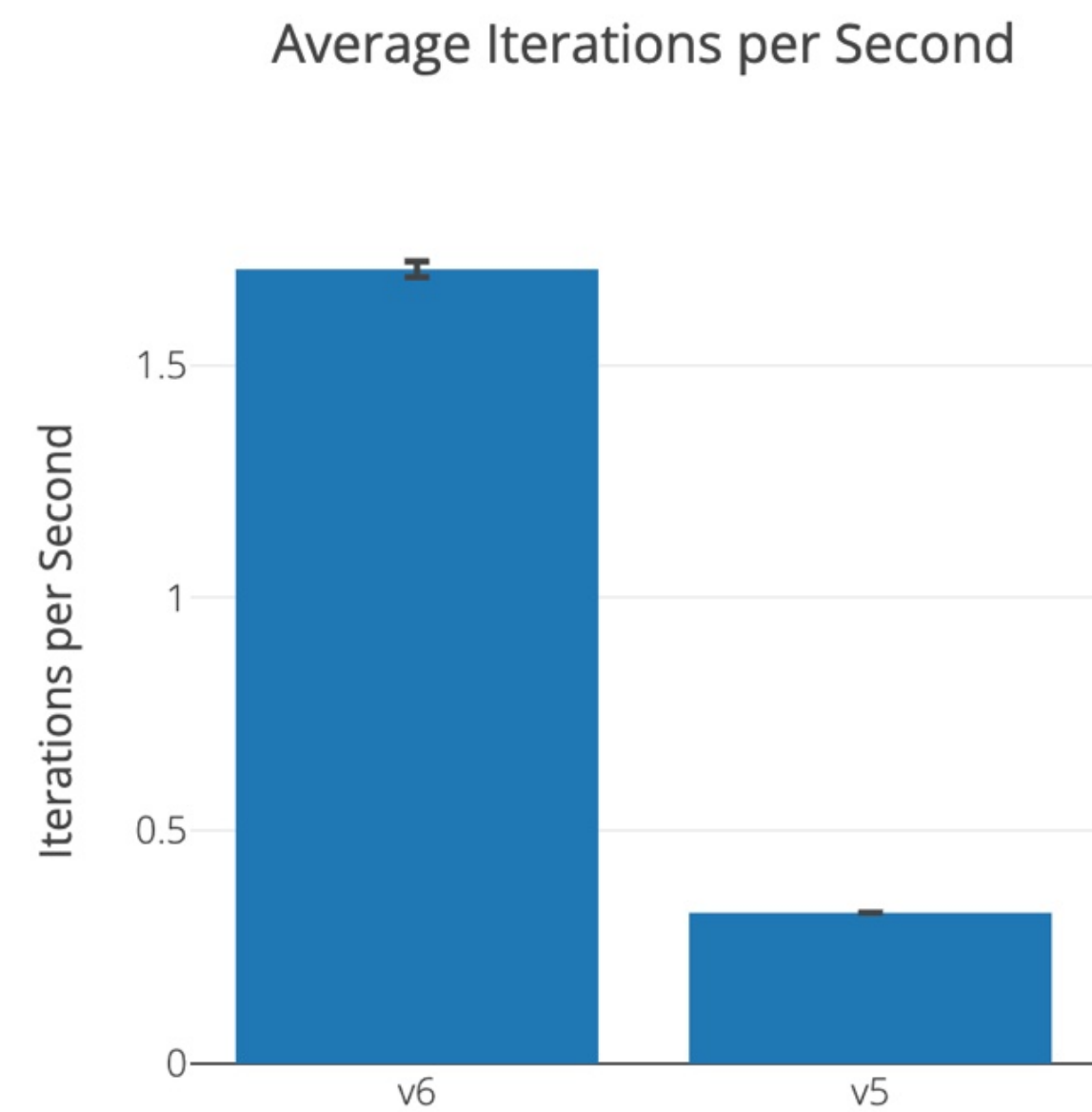


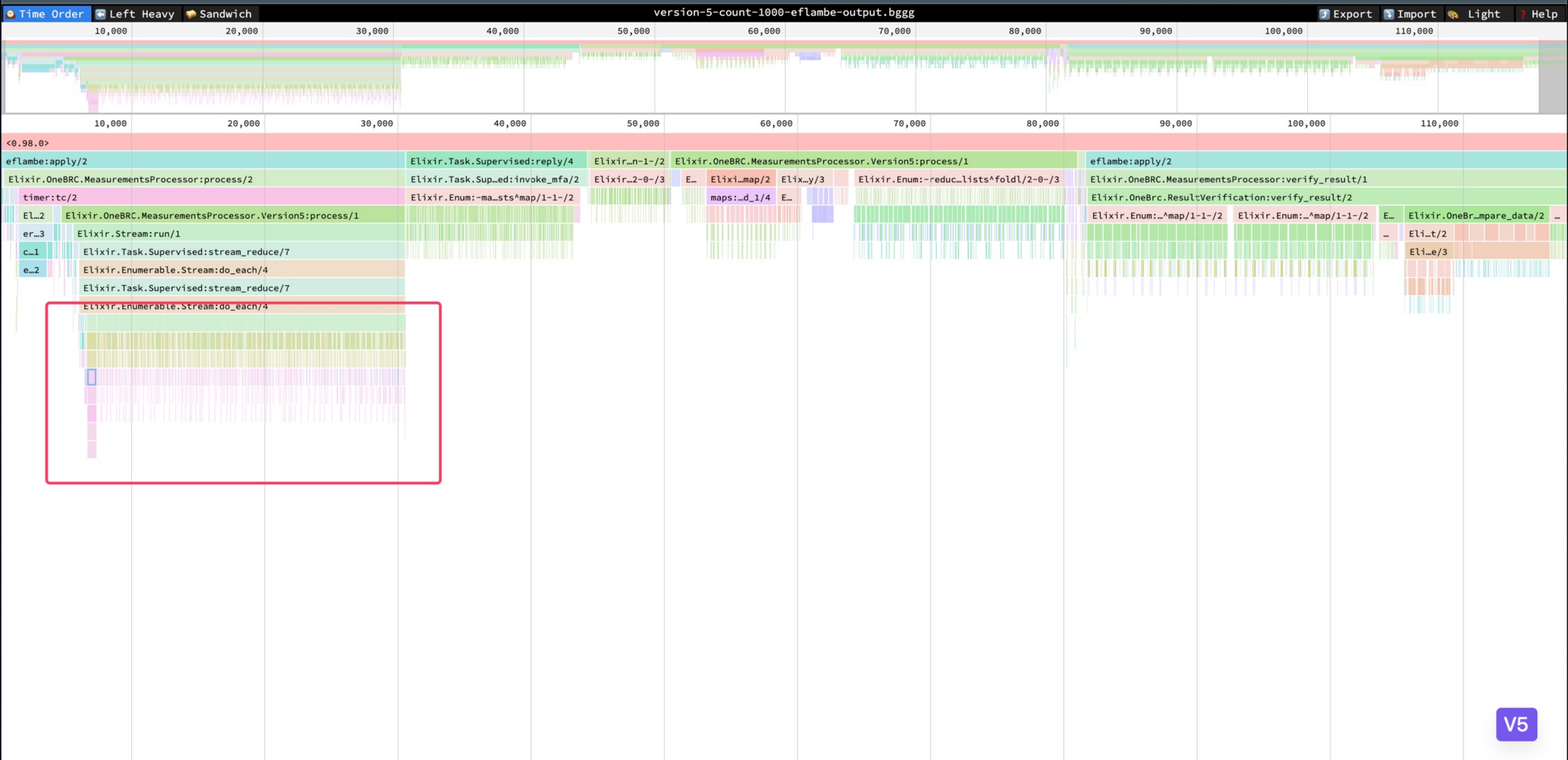
*CPU usage of version 5 vs version 6*



## Run Time Comparison <sup>?</sup>

Name	Iterations per Second	Average	Deviation	Median	Mode	Minimum	Maximum	Sample size
v6	1.71	0.59 s	±1.00%	0.59 s	none	0.58 s	0.60 s	9
v5	0.32	3.08 s	±0.01%	3.08 s	none	3.08 s	3.08 s	2

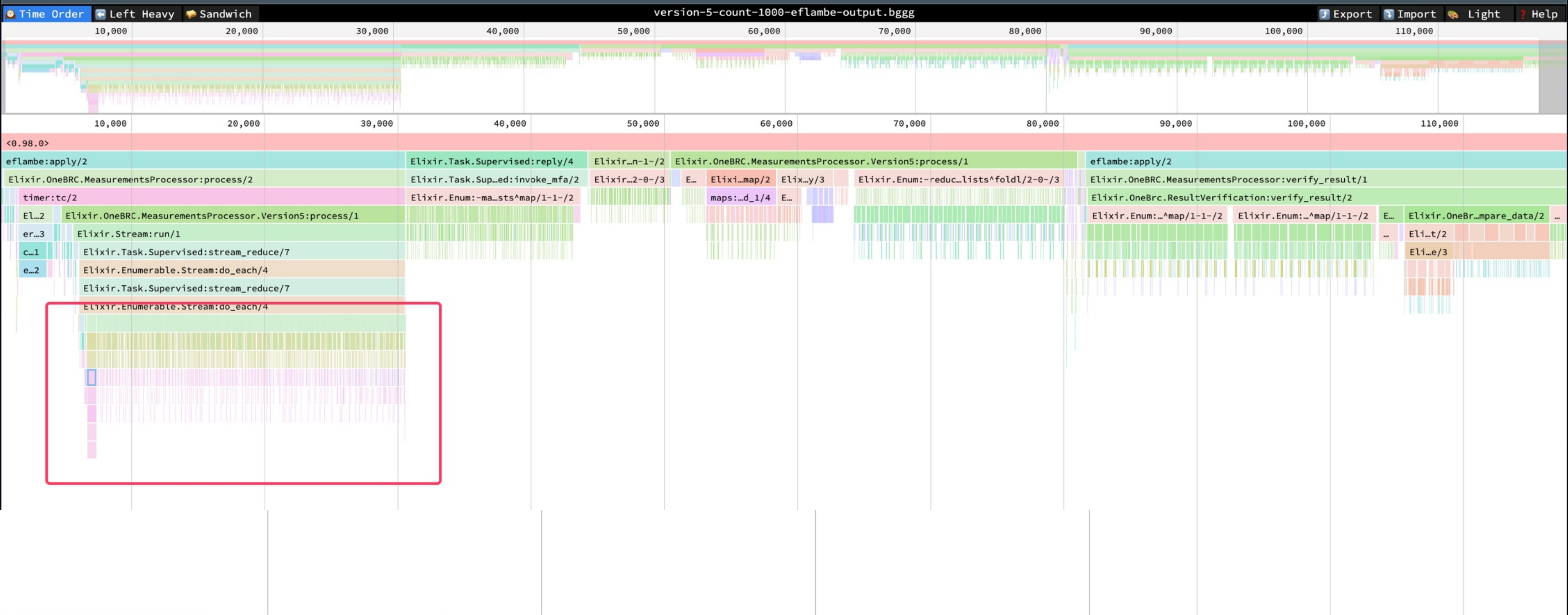




This Instance		All Instances	
Total	Self	Total	Self
10	4	10,837	4,023
<0.01%	<0.01%	9.0%	3.4%

- prim\_file:read\_line/1
- > Elixir.IO:binread/2
- > Elixir.IO:each\_binstream/2
- > Elixir.Stream:do\_resource/5
- > Elixir.Enumerable.Stream:do\_each/4
- > Elixir.Task.Supervised:stream\_reduce/7
- > Elixir.Enumerable.Stream:do\_each/4
- > Elixir.Task.Supervised:stream\_reduce/7
- > Elixir.Stream:run/1
- Elixir.OneBRC.MeasurementsProcessor.Version5:process/1
- > timer:tc/2
- > Elixir.OneBRC.MeasurementsProcessor:process/2

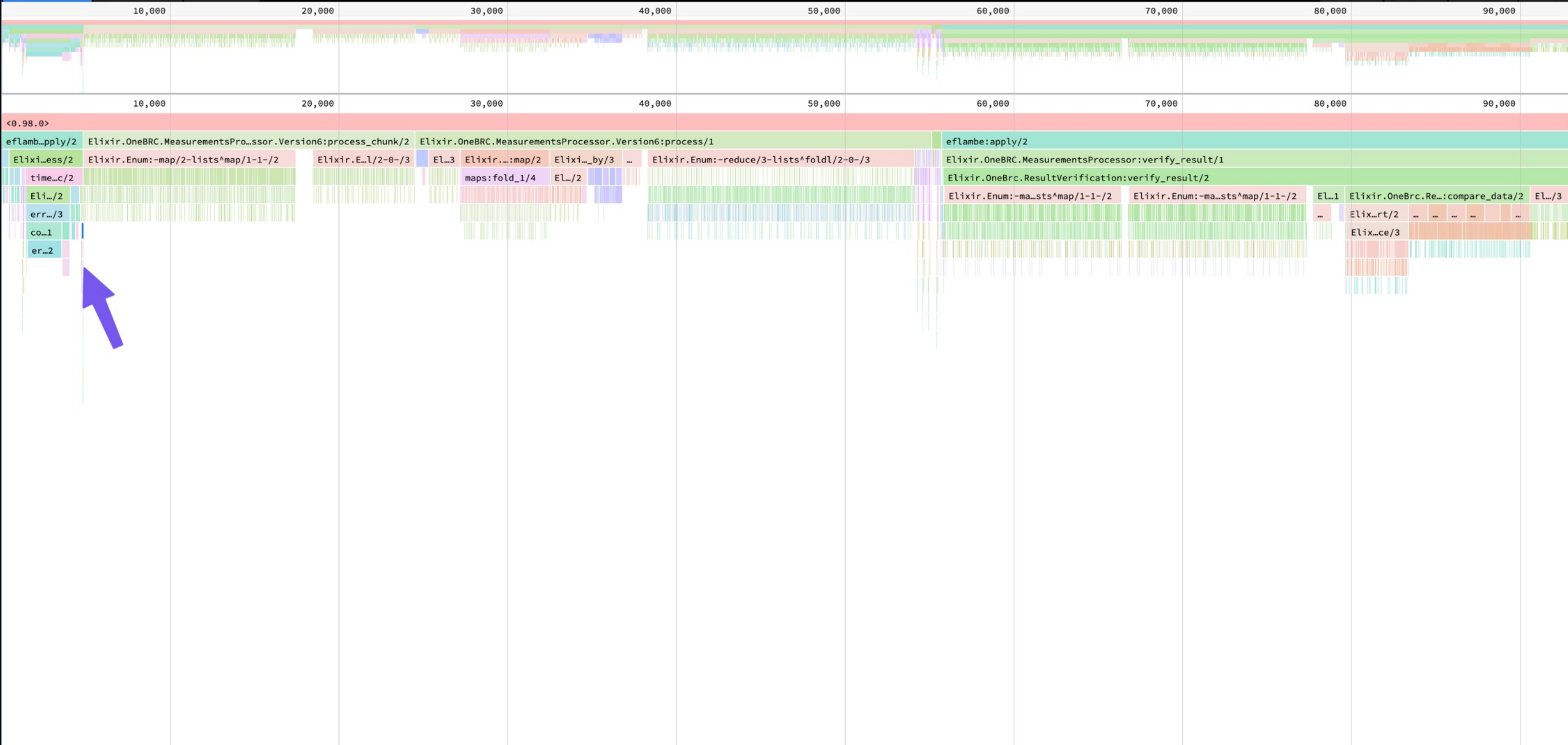
V5



This Instance		All Instances	
Total	Self	Total	Self
10	4	10,837	4,023
<0.01%	<0.01%	9.0%	3.4%

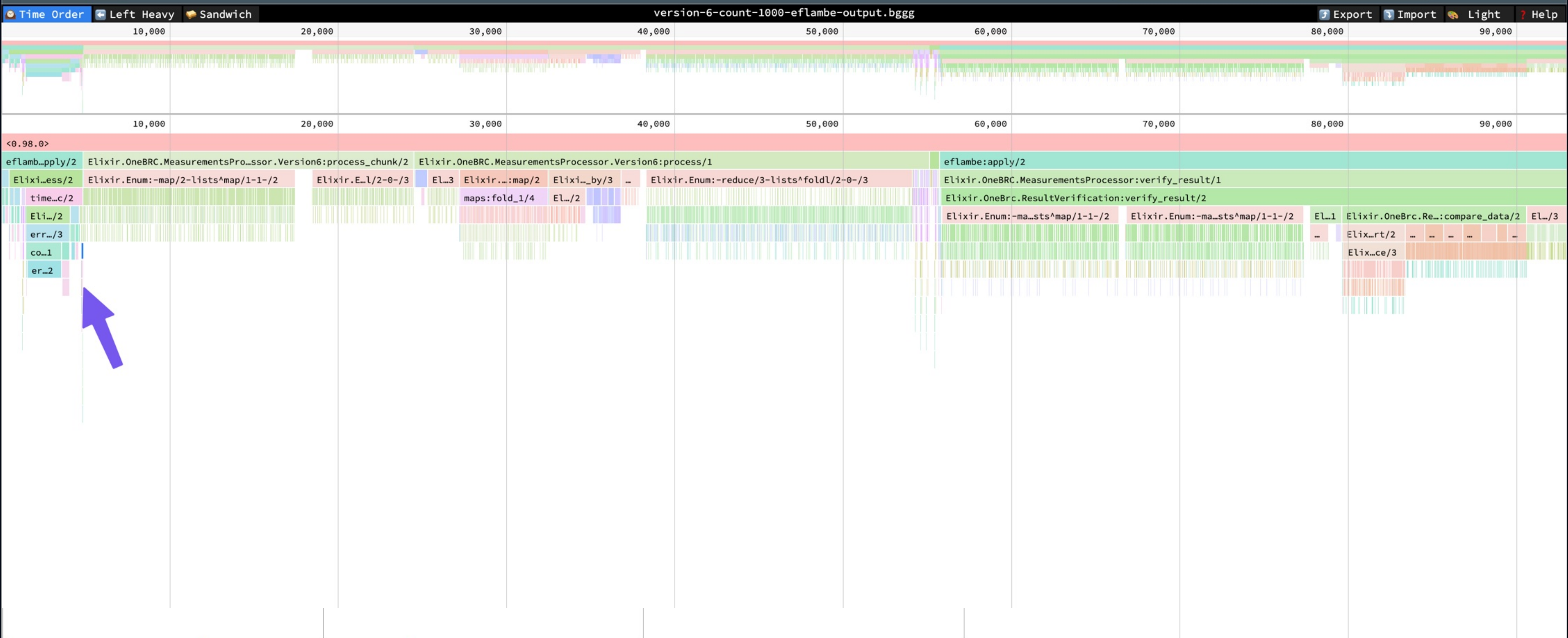
- prim\_file:read\_line/1
- > Elixir.IO:binread/2
- > Elixir.IO:each\_binstream/2
- > Elixir.Stream:do\_resource/5
- > Elixir.Enumerable.Stream:do\_each/4
- > Elixir.Task.Supervised:stream\_reduce/7
- > Elixir.Enumerable.Stream:do\_each/4
- > Elixir.Task.Supervised:stream\_reduce/7
- > Elixir.Stream:run/1
- > Elixir.OneBRC.MeasurementsProcessor.Version5:process/1
- > timer:tc/2
- > Elixir.OneBRC.MeasurementsProcessor:process/2

V5



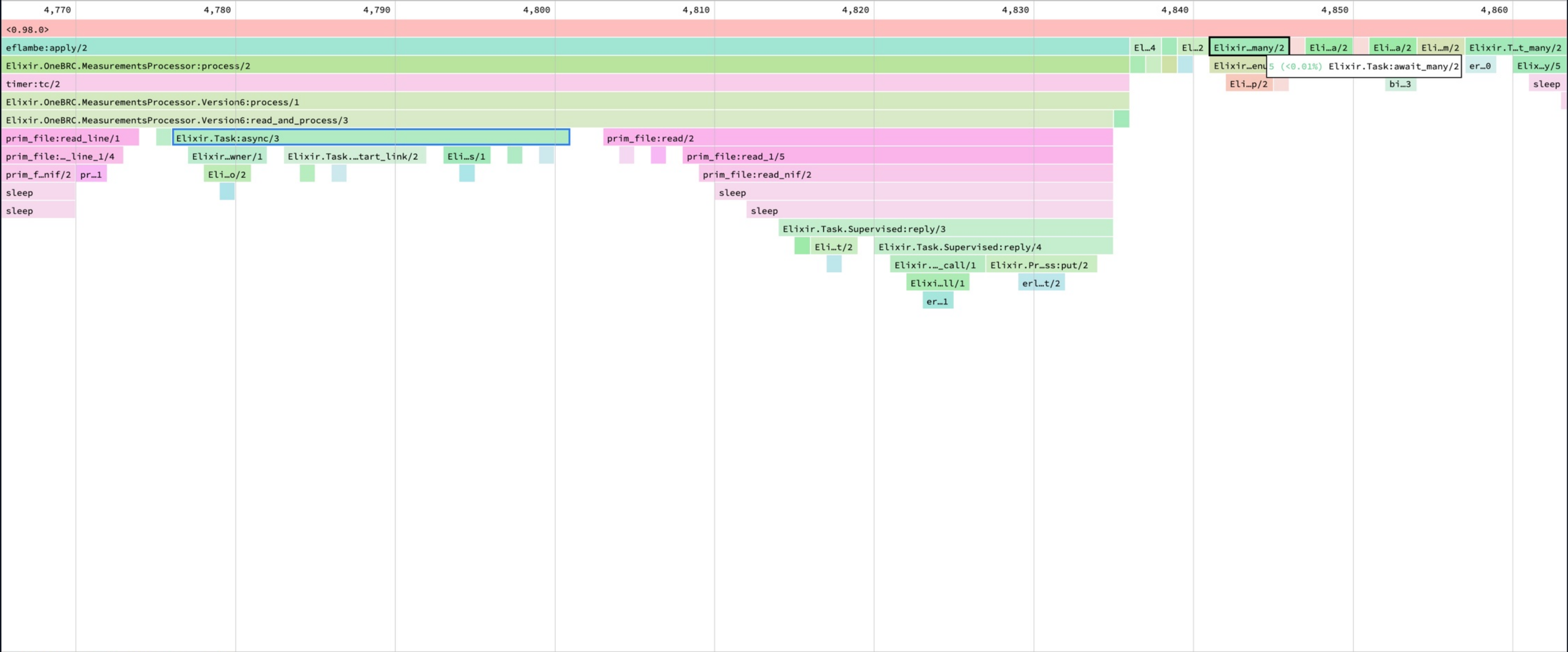
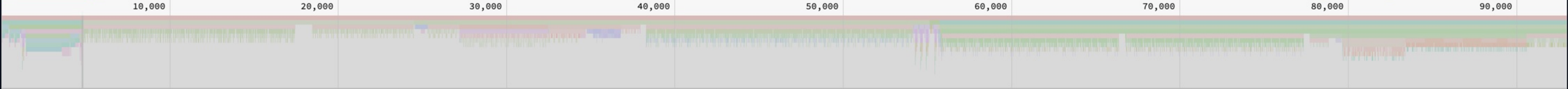
This Instance		All Instances	
Total	Self	Total	Self
40	7	40	7
0.04%	<0.01%	0.04%	<0.01%

- prim\_file:read\_line/1
- > Elixir.OneBRC.MeasurementsProcessor.Version6:read\_and\_process/3
- > Elixir.OneBRC.MeasurementsProcessor.Version6:process/1
- > timer:tc/2
- > Elixir.OneBRC.MeasurementsProcessor:process/2
- > eflambe:apply/2
- > <0.98.0>



This Instance		All Instances	
Total	Self	Total	Self
40	7	40	7
0.04%	<0.01%	0.04%	<0.01%

- prim\_file:read\_line/1
- > Elixir.OneBRC.MeasurementsProcessor.Version6:read\_and\_process/3
- > Elixir.OneBRC.MeasurementsProcessor.Version6:process/1
- > timer:tc/2
- > Elixir.OneBRC.MeasurementsProcessor:process/2
- > eflambe:apply/2
- > <0.98.0>

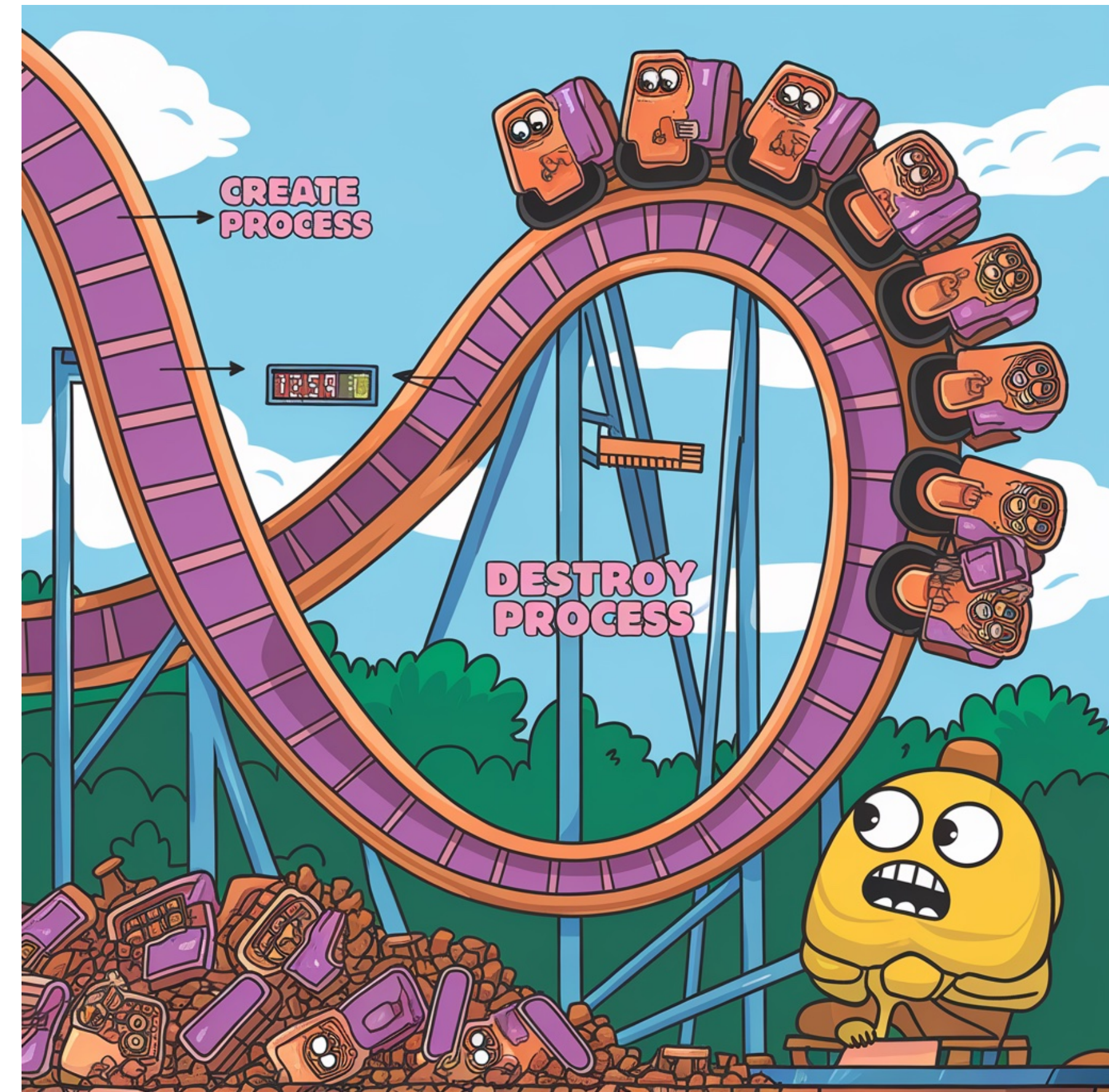


This Instance		All Instances	
Total	Self	Total	Self
25	6	25	6
0.03%	<0.01%	0.03%	<0.01%

- Elixir.Task:async/3
- > Elixir.OneBRC.MeasurementsProcessor.Version6:read\_and\_process/3
- > Elixir.OneBRC.MeasurementsProcessor.Version6:process/1
- > timer:tc/2
- > Elixir.OneBRC.MeasurementsProcessor:process/2
- > eflambe:apply/2
- > <0.98.0>

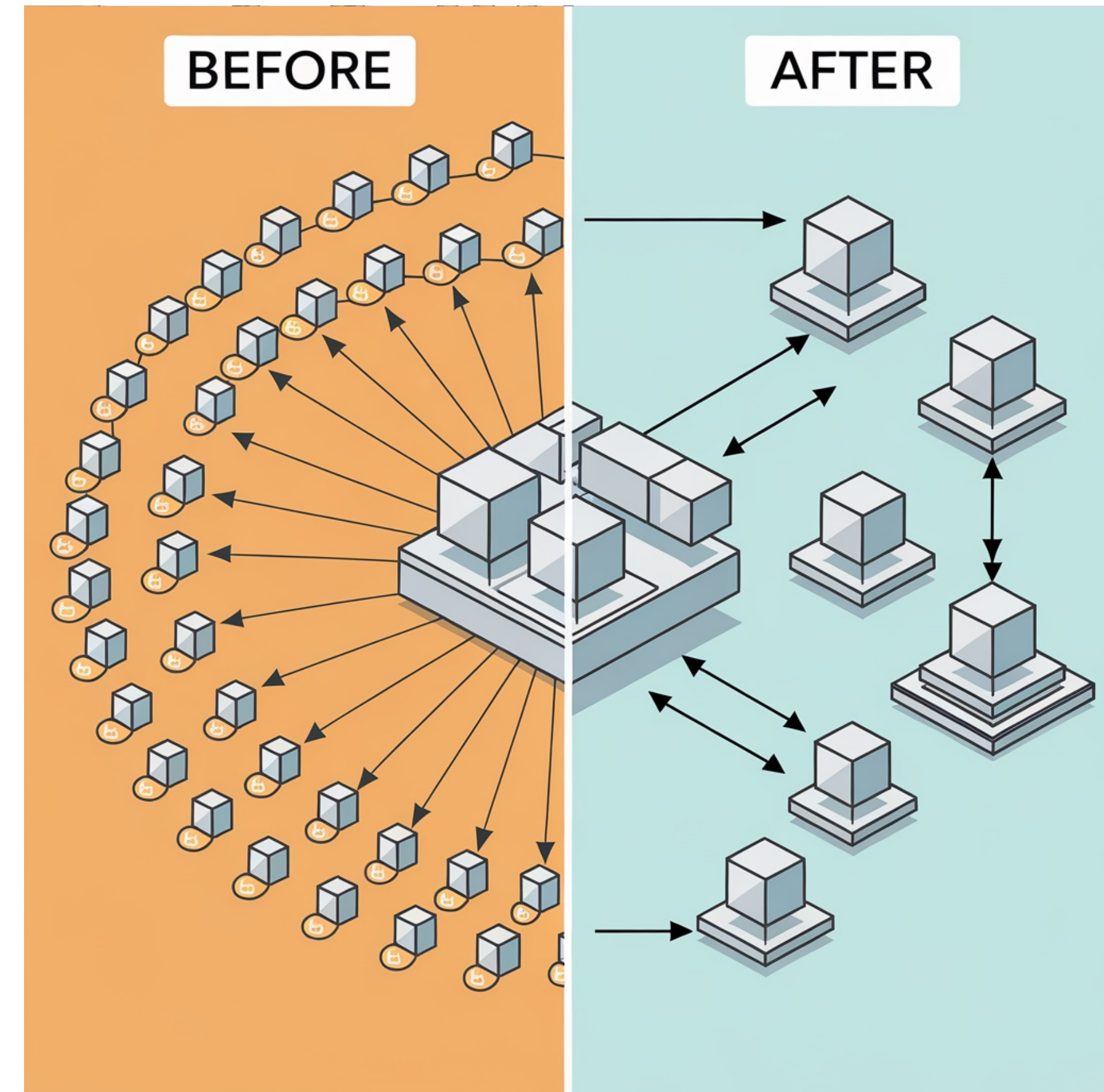
# Optimising Task overhead

- Task.async for each chunk and then Task.await\_many helped 🚀 CPU utilisation
- Task.async creates a process, runs our code, destroys the process.



# Optimising Task overhead

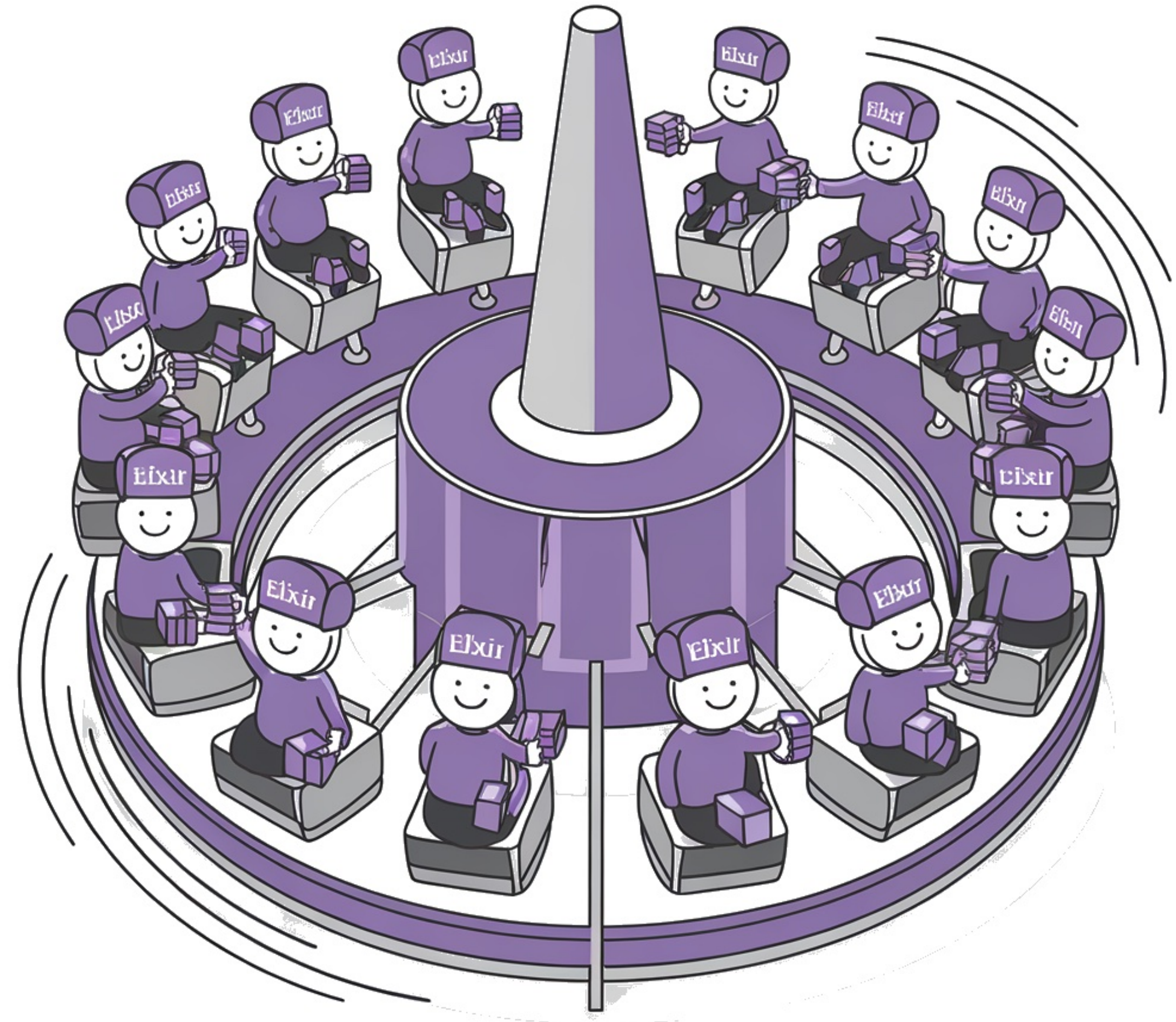
- Task.async for each chunk and then Task.await\_many helped 🚀 CPU utilisation
- Task.async creates a process, runs our code, destroys the process.
- Process creation overhead can be avoided.
- We can use a fixed worker pool for handling jobs.
- Reuse workers, minimise overhead.





# Reusable Worker pool of processes

- Parent process spawns worker processes
- Each worker process sends give\_work message to it's parent
- Parent process, when received a give\_work message from worker, sends a chunk for the worker to process
- Worker process processes the chunk, writes result to it's dict, and sends give\_work message to parent again
- Repeat
- At the end, parent collects intermediate results from each worker.



```
defmodule Worker do
  def run(parent_pid) do
    send(parent_pid, {:give_work, self()})

    receive do
      {:do_work, chunk} →
        process_chunk(chunk)
        run(parent_pid)

      :result →
        send(parent_pid, {:result, :erlang.get()})
        # die
    end
  end
end

defp process_chunk(bin) do
  :binary.split(bin, "\n", [:global])
  ▷ Enum.map(&parse_row/1)
  ▷ Enum.map(fn row →
    process_row(row)
  end)
end

...
```

```
worker_count = System.schedulers_online() * 2
# boot up workers
parent = self()

wpids =
  Enum.map(1..worker_count, fn _ →
    spawn_link(fn →
      Worker.run(parent)
    end)
  end)

{:ok, file} = :prim_file.open(file_path, [:raw, :binary, :read])

:ok = read_and_process(file)

...
```

```
defp read_and_process(file) do
  chunk_size = 1024 * 1024 * 1
  data =
    case :prim_file.read(file, chunk_size) do
      :eof →
        nil

      {:ok, data} →
        case :prim_file.read_line(file) do
          {:ok, line} →
            <<data::binary, line::binary>>

          :eof →
            data
        end
    end

  end

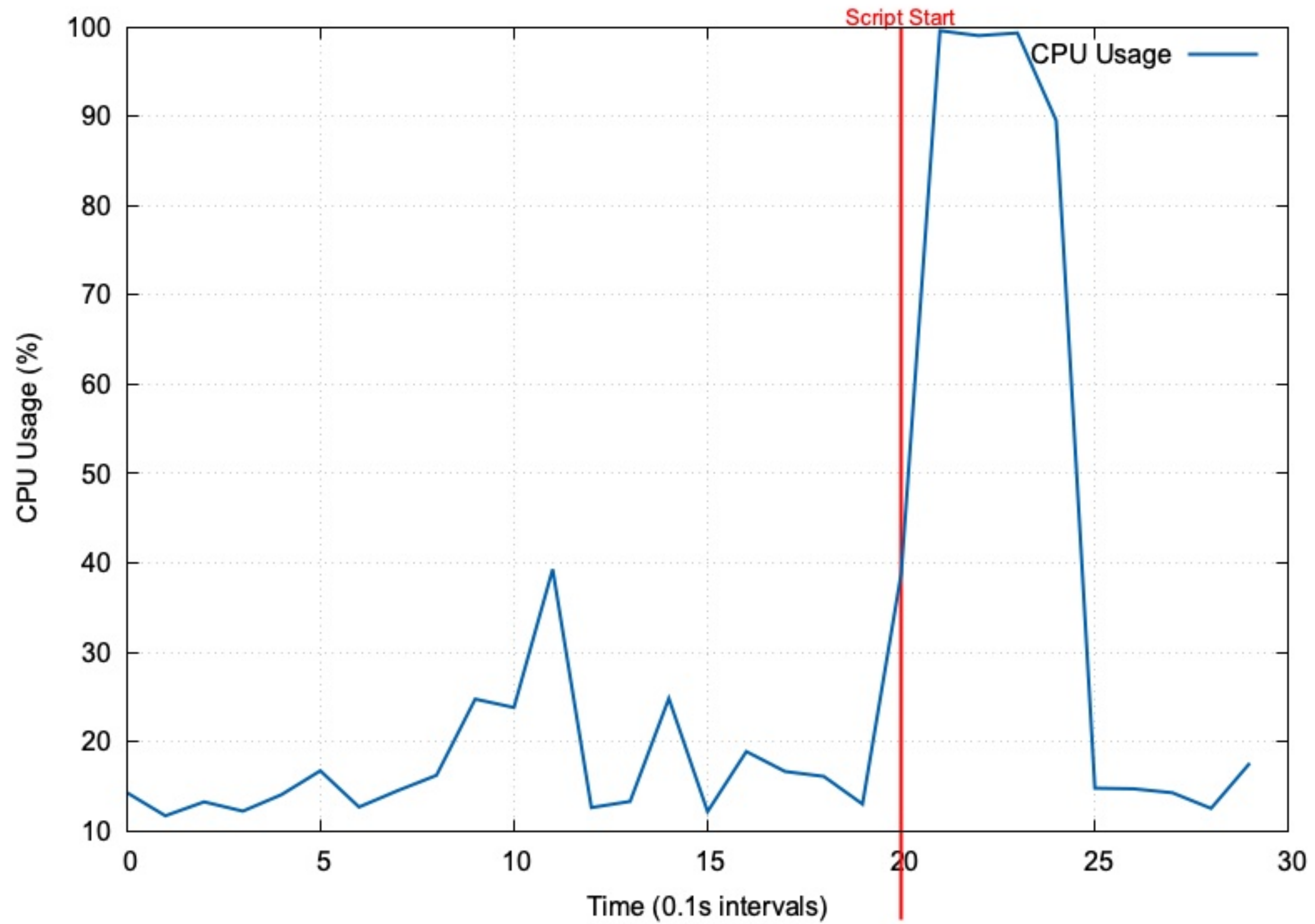
  if !is_nil(data) do
    receive do
      {:give_work, worker_pid} →
        send(worker_pid, {:do_work, data})
    end
    read_and_process(file)
  else
    :ok
  end
end
```

```
# wait for all workers to finish
Enum.map(1..worker_count, fn _ →
  receive do
    {:give_work, _worker_pid} →
      :ok
  end
end)

results =
  wpids
  ▷ Enum.map(fn wpid →
    send(wpid, :result)
  end)
  ▷ Enum.map(fn _ →
    receive do
      {:result, result} →
        result
    end
  end)

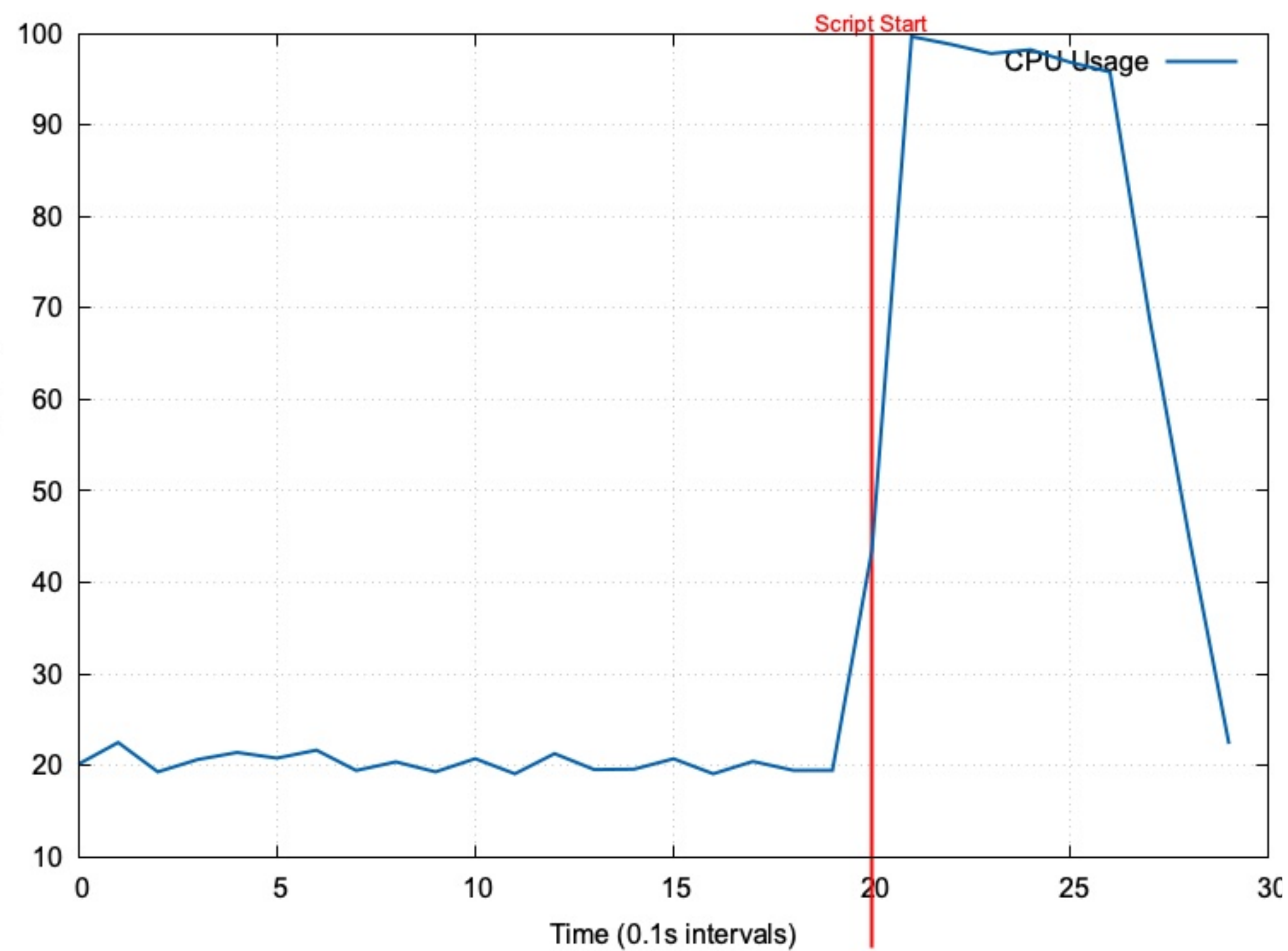
:prim_file.close(file)
```

Overall CPU Usage Over Time (V7, 100M measurements)

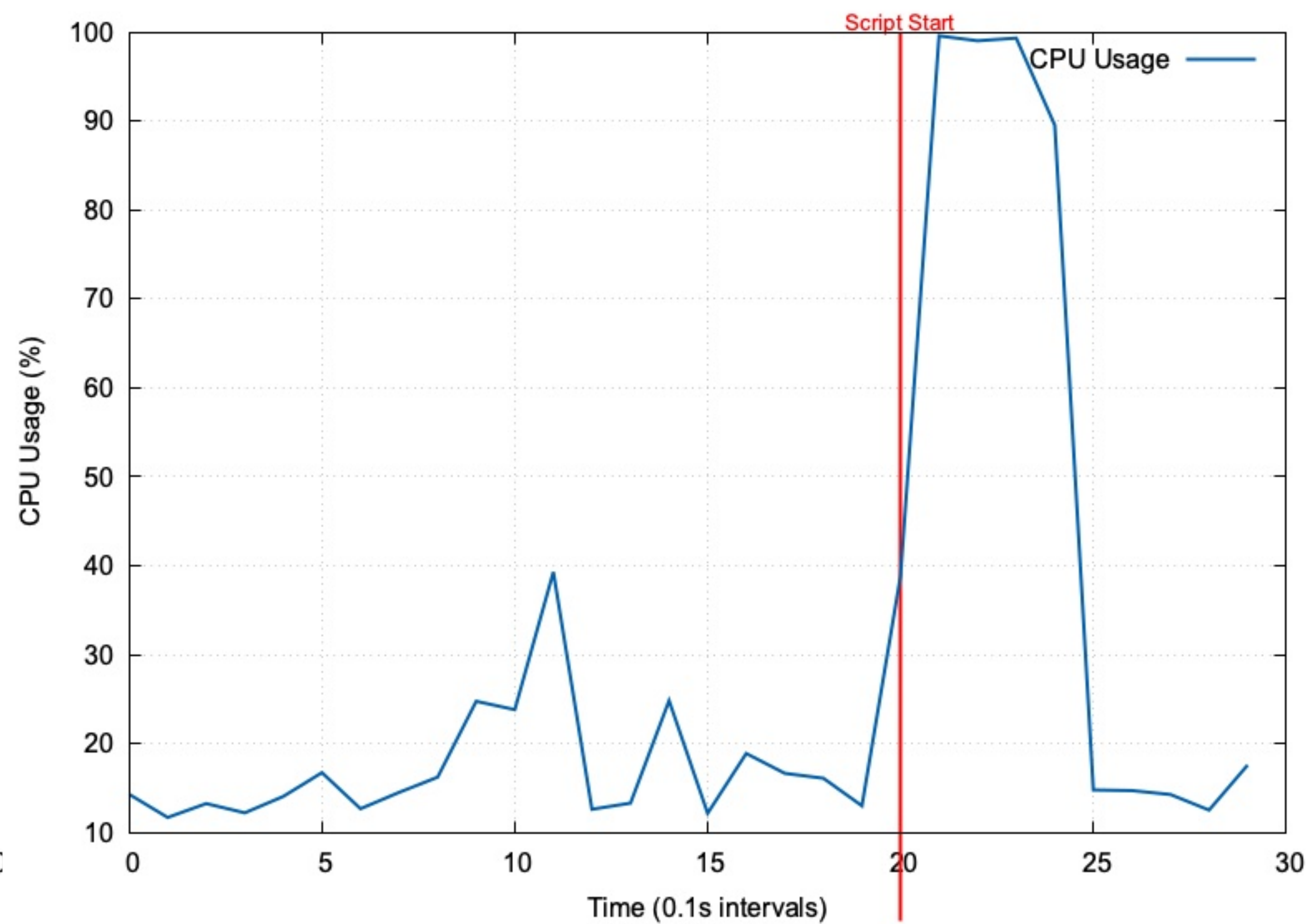


*CPU usage of version 7*

Overall CPU Usage Over Time (V6, 100M measurements)



Overall CPU Usage Over Time (V7, 100M measurements)



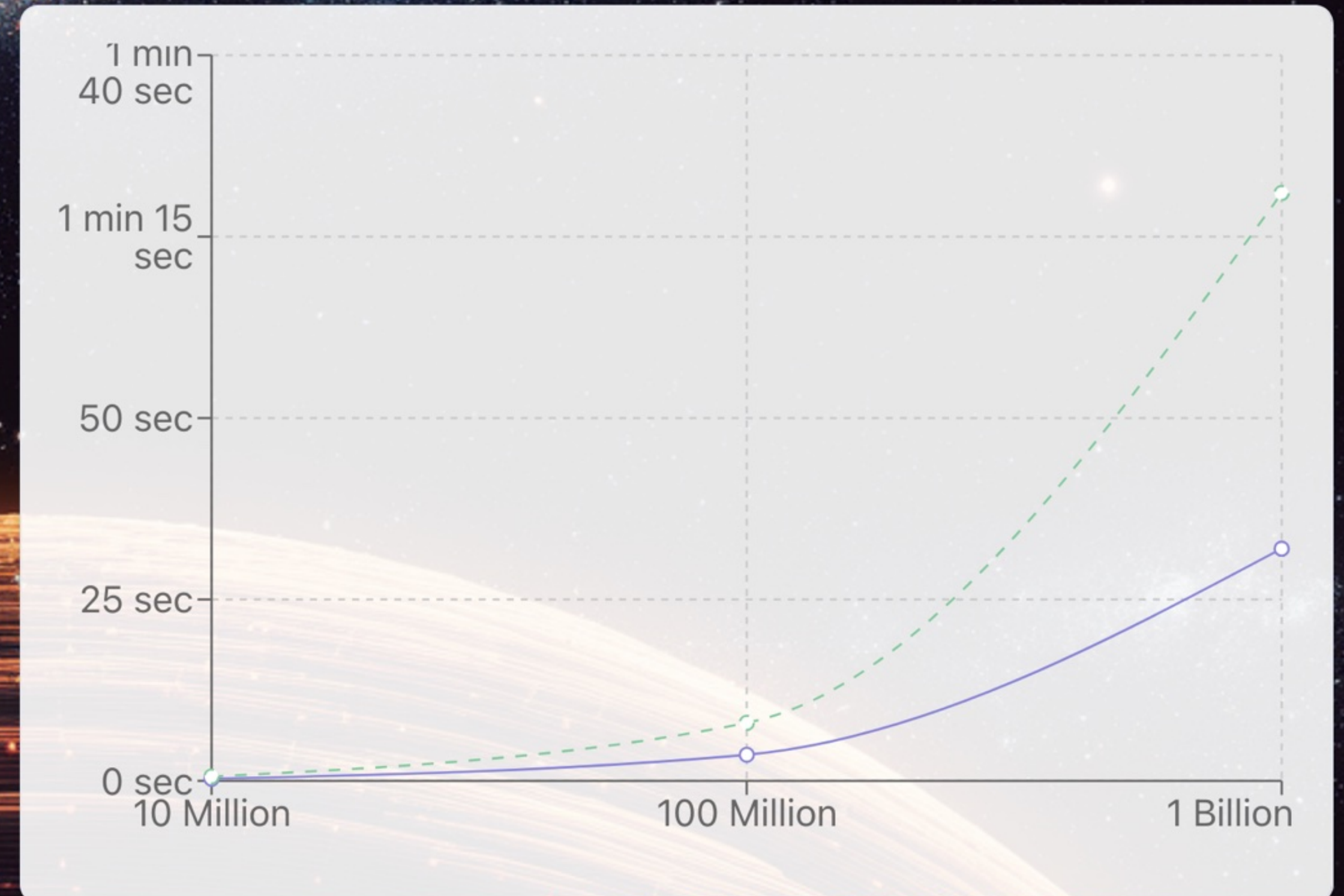
*CPU usage of version 6 vs version 7*

# 1BRC in Elixir: Version 7

**10 Million Rows**  
🕒 **0.3 sec** ↘ 50.0%

**100 Million Rows**  
🕒 **3.6 sec** ↘ 55.0%

**1 Billion Rows**  
🕒 **32 sec** ↘ 60.5%

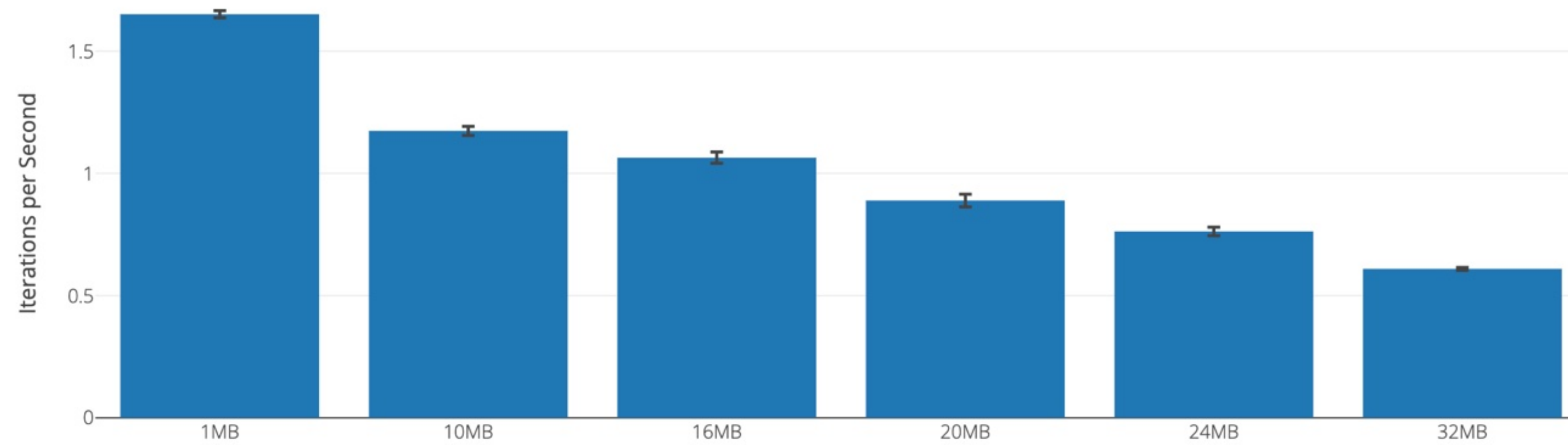




## Run Time Comparison <sup>®</sup>

Name	Iterations per Second	Average	Deviation	Median	Mode	Minimum	Maximum	Sample size
1MB	1.65	0.61 s	±0.88%	0.61 s	none	0.60 s	0.61 s	9
10MB	1.17	0.85 s	±1.58%	0.85 s	none	0.83 s	0.87 s	6
16MB	1.06	0.94 s	±2.16%	0.94 s	none	0.91 s	0.96 s	6
20MB	0.89	1.12 s	±2.91%	1.13 s	none	1.09 s	1.17 s	5
24MB	0.76	1.31 s	±2.24%	1.32 s	none	1.27 s	1.34 s	4
32MB	0.61	1.64 s	±0.98%	1.63 s	none	1.63 s	1.66 s	4

Average Iterations per Second



Run Time Boxplot

*Using Benchee to find the optimal chunk size*

```
defmodule Worker do
  def run(parent_pid) do
    send(parent_pid, {:give_work, self()})

    receive do
      {:do_work, chunk} →
        process_chunk(chunk)
        run(parent_pid)

      :result →
        send(parent_pid, {:result, :erlang.get()})
        # die
    end
  end
end

defp process_chunk(bin) do
  :binary.split(bin, "\n", [:global])
  ▷ Enum.map(&parse_row/1)
  ▷ Enum.map(fn row →
    process_row(row)
  end)
end)
...

```

```

# ex: -4.5
defp parse_temp(<<?-, d1, ?., d2, "\n", rest::binary>>, key) do
  temp = -(char_to_num(d1) * 10 + char_to_num(d2))
  process_row(key, temp)
  process_chunk_lines(rest)
end

# ex: 4.5
defp parse_temp(<<d1, ?., d2, "\n", rest::binary>>, key) do
  temp = char_to_num(d1) * 10 + char_to_num(d2)
  process_row(key, temp)
  process_chunk_lines(rest)
end

# ex: -45.3
defp parse_temp(<<?-, d1, d2, ?., d3, "\n", rest::binary>>, key) do
  temp = -(char_to_num(d1) * 100 + char_to_num(d2) * 10 + char_to_num(d3))
  process_row(key, temp)
  process_chunk_lines(rest)
end

# ex: 45.3
defp parse_temp(<<d1, d2, ?., d3, "\n", rest::binary>>, key) do
  temp = char_to_num(d1) * 100 + char_to_num(d2) * 10 + char_to_num(d3)
  process_row(key, temp)
  process_chunk_lines(rest)
end

```

```
defp process_chunk_lines(<<>>) do
  :ok
end
```

```
defp process_chunk_lines(bin) do
  parse_weather_station(bin, bin, 0)
end
```

```
# ↓ found ";" - split and parse
```

```
defp parse_weather_station(bin, <<";", _rest::binary>>, count) do
  <<key::binary-size(count), ";", temp_bin::binary>> = bin
  parse_temp(temp_bin, key)
end
```

```
# → no ";" yet - move to next char
```

```
defp parse_weather_station(bin, <<_c, rest::binary>>, count) do
  parse_weather_station(bin, rest, count + 1)
end
```

```
iex(1)> h <<>>
```

```
defmacro <<args>>
```

```
...
```

For binaries, the default is the size of the binary.

```
iex> <<name::binary-size(5), " the ", species::binary>> = <<"Frank the Walrus">>
"Frank the Walrus"
iex> {name, species}
{"Frank", "Walrus"}
```

The size can be a variable or any valid guard expression:

```
iex> name_size = 5
iex> <<name::binary-size(^name_size), " the ", species::binary>> = <<"Frank the Walrus">>
iex> {name, species}
{"Frank", "Walrus"}
```

```
...
```

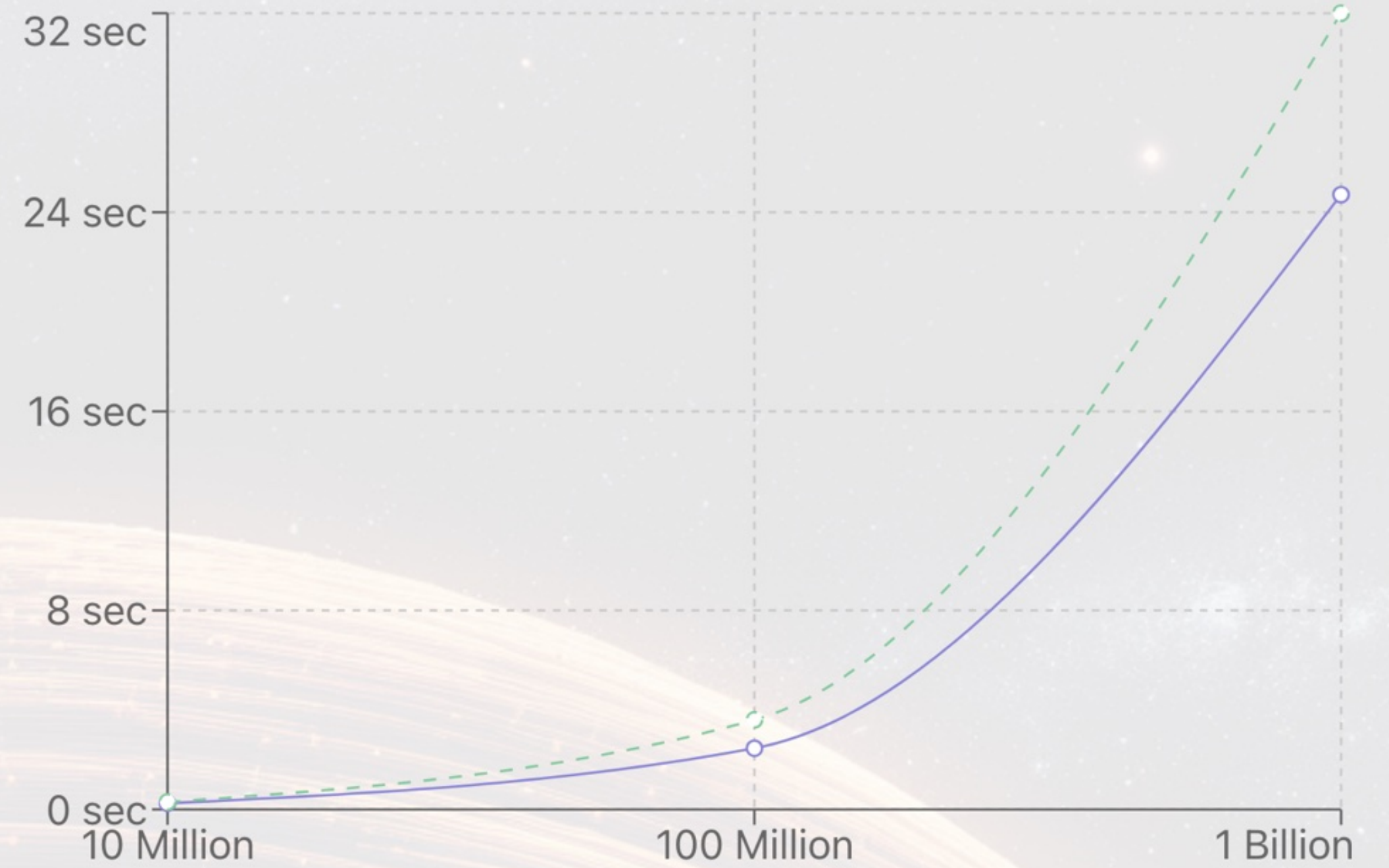
*Output of `h <<>>` in iex*

# 1BRC in Elixir: Version 8

**10 Million Rows**  
🕒 **0.25 sec**  
↗️ **16.7%**

**100 Million Rows**  
🕒 **2.46 sec**  
↗️ **31.7%**

**1 Billion Rows**  
🕒 **24.71 sec**  
↗️ **22.8%**

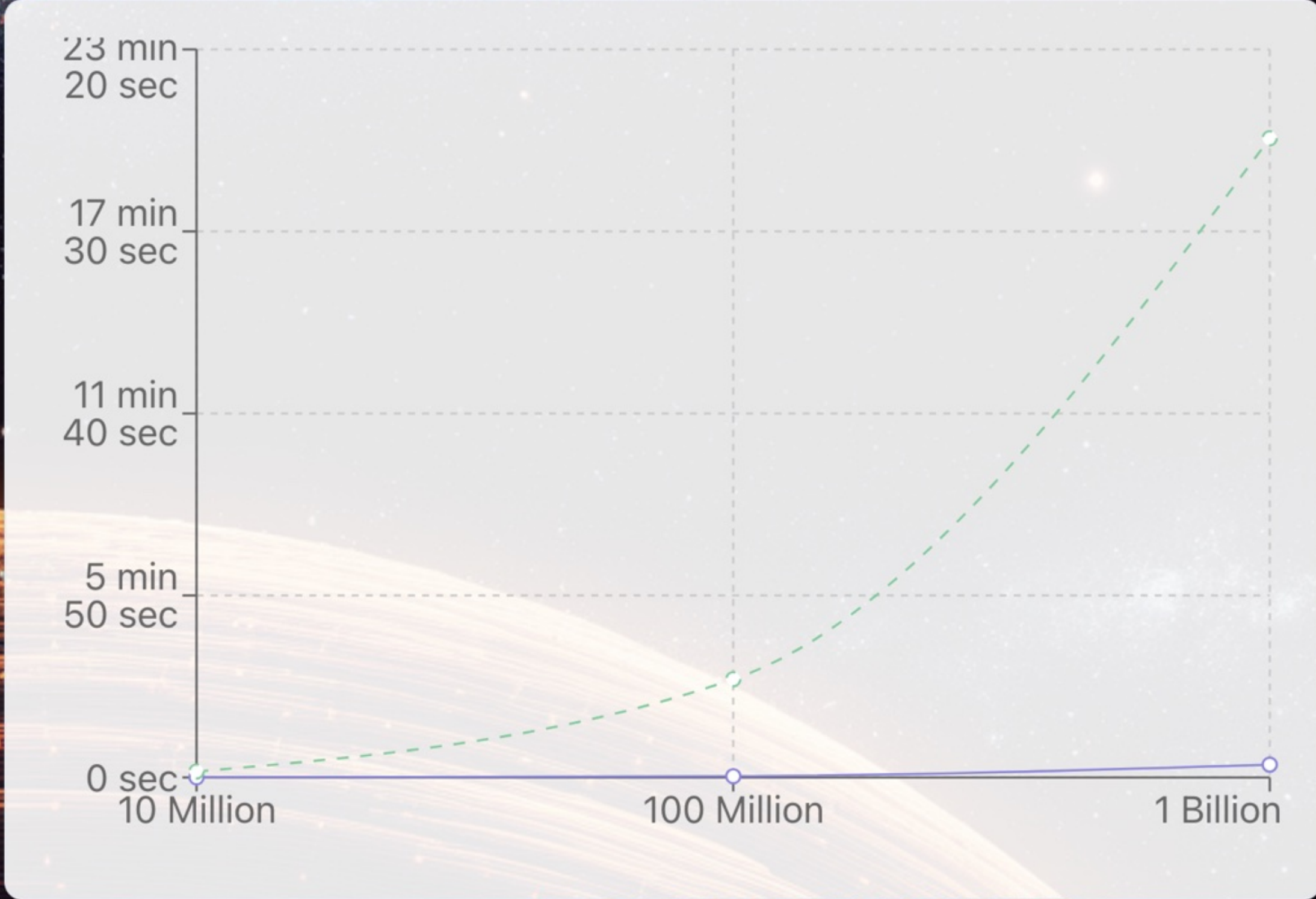


# 1BRC in Elixir: From Version 1 to Version 8

**10 Million Rows**  
🕒 **0.25 sec** ↘ 97.8%

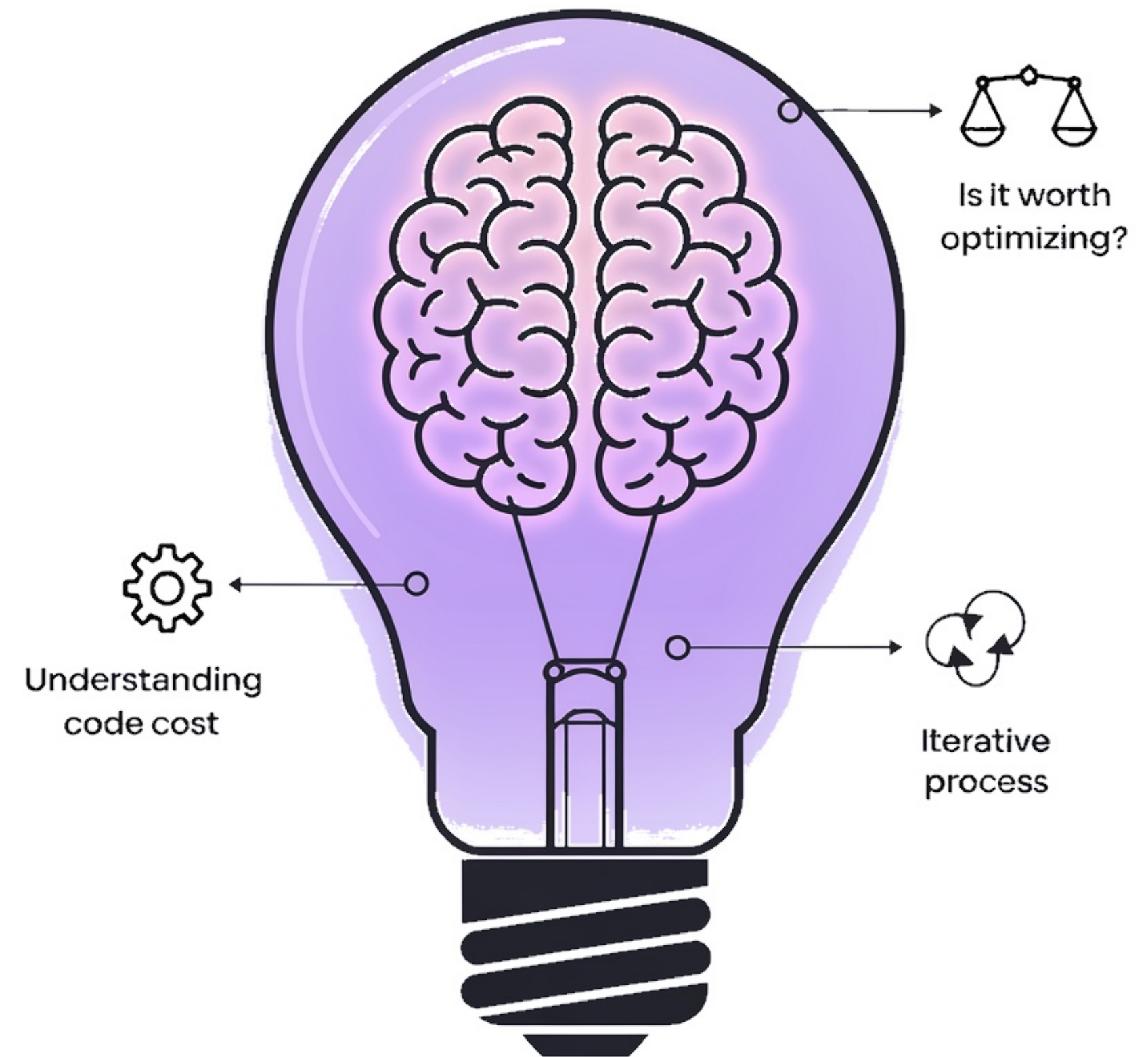
**100 Million Rows**  
🕒 **2.46 sec** ↘ 98.7%

**1 Billion Rows**  
🕒 **24.71 sec** ↘ 98.0%



# Some takeaways

- Developing a performance mindset.
- Understanding the true cost of code and its performance implications.
- Ask: “is this code worth optimising”?
- Optimisation is an iterative process.
- Understanding trade-offs
- Using profiling tools effectively





# Thanks!

rajrajhans/  
**elixir\_1brc**



The One Billion Rows Challenge in Elixir



@\_rajrajhans

# Thanks!



@\_rajrajhans

## rajrajhans/ elixir\_1brc

The One Billion Rows Challenge in Elixir

