

# Straightforward stream processing with



[Materials](#)

Adam Glustein

PyData NYC



[Point72/csp](#)

November 7<sup>th</sup>, 2024

# What you will learn

- Basics of stream processing
- Challenges posed by **real-time** data pipelines
  - Debugging and unit testing
  - Performance tuning
  - Interactive execution (ipython)
- How **CSP** removes some of these obstacles



# The main goal of this talk

- Add a new tool to your data science toolbox

## The law of the instrument (corollary)

If the only tool you have is *pandas*, then all data looks like a *DataFrame*.

- Know when to identify situations where stream processing makes sense for you

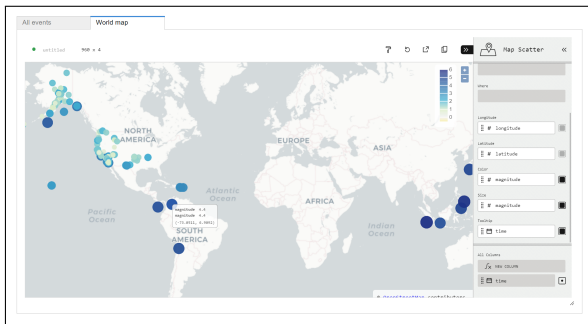
# A bit about myself

- Quantitative Developer at **Point72** in New York, USA
- Contributor to csp since 2022



# Composable Stream Processing (CSP)

- First open-source project from Point72
- Many internal (and now external!) contributors
- Primary goal: improve **developer experience** when building real-time data applications



# What is developer experience?

## Definition

**(Developer Experience)** How easy it is to build a robust application using a tool or framework

- Inherently subjective!
- Common criteria that *most* people agree on
  - 1 How much code do I have to write?
  - 2 How certain can I be in my app's behavior?
  - 3 How easy it is to integrate with other tools?

# Static vs. dynamic data

- Two different types of *data*

## Static Data

- All entries known up-front



## Dynamic Data

- Data updates during application runtime



# Batch vs. stream processing

- Two different methods of *data processing*

## Batch Processing

- Periodically process *batches* of data
- Simpler and less compute

ADJUSTING JOURNAL									
Data last entry date: 4/18/2024									
Date	Description	Supplier/Job	Description	Cost	Account Name	Debit	Credit	Credit	
4/18/2024	AP/PS	Revenue from bank interest	1200	Revenue Charge Income				200.00	
4/18/2024	AP/PS	Revenue from bank interest	1200	Revenue Nonrevenue		200.00			
4/18/2024	AP/PS	Advertising expenses	1200	Exp - Other		1,000.00			
4/18/2024	AP/PS	Advertising expenses	1200	Exp - Trade			1,000.00		
4/18/2024	AP/PS	Insurance expenses	1400	Exp - Insurance		2,000.00			
4/18/2024	AP/PS	Insurance expenses	1400	Exp/PSAD - Insurance			2,000.00		
4/18/2024	AP/PS	Supplies expenses	1400	Exp - Office Supplies		1,000.00			
4/18/2024	AP/PS	Supplies expenses	1400	Exp/PSAD - Insurance			1,000.00		
4/18/2024	AP/PS	Revenue from rental services	1400	Rental Income				1,000.00	
4/18/2024	AP/PS	Revenue from rental services	1400	Rental Income				1,000.00	
4/18/2024	AP/PS	Revenue from rental services	1400	Exp - Depreciation Expense		200.00			
4/18/2024	AP/PS	Revenue from rental services	1400	Exp/PSAD - Computer Equipment			200.00		

## Stream Processing

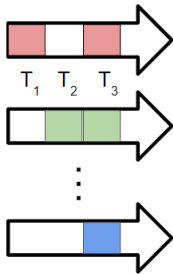
- Process data immediately
- Provides immediate results





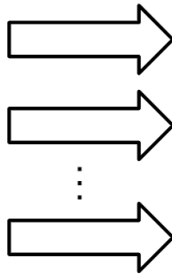
# Stream processing

Asynchronous  
data sources



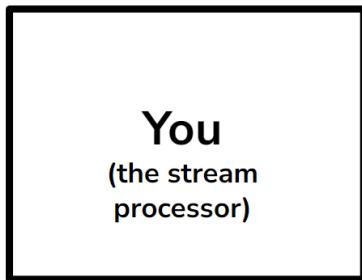
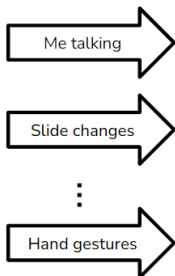
Event stream  
processing  
engine

Outputs

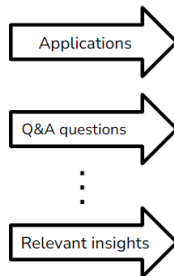


# Stream processing

Asynchronous  
data sources



Outputs



# A semi-complete data toolbox

- Well-known tools for static data



- Distributed batch processing tools




- Unified (batch and stream) platforms



- Pure stream processing tools



# What about CSP?

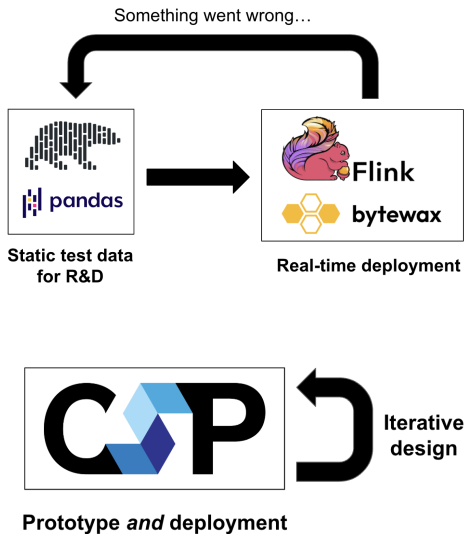
 is a stream processing library for *both* static and dynamic data

Key feature: seamless transition between historical and real-time workflows

# Why CSP integrates static dataflows

- Streaming data is inherently dynamic
- However, static data is much easier to work with
  - Quick prototyping
  - Backtesting on historical data
  - Verification before deployment
  - Stable CI tests and benchmarks

# The CSP development cycle



# Why use a different framework for research?

- Static data libraries have a **great** developer experience!
- Streaming data libraries are focused on deployment, not development

## Example 1: Reading a CSV file

### pyflink

```
1 t_env.create_temporary_table(  
2     "source",  
3     TableDescriptor.for_connector("filesystem")  
4     .schema(Schema.new_builder()  
5     .column("word", DataTypes.STRING())  
6     .build())  
7     .option("path", input_path)  
8     .format("csv")  
9     .build())  
10 tab = t_env.from_path("source")
```

### pandas

```
1 df = pd.read_csv(input_path)
```



## Example 2: Writing to stdout

### pyflink

```
1 t_env.create_temporary_table(  
2     "sink",  
3     TableDescriptor.for_connector("print")  
4     .schema(Schema.new_builder()  
5     .column("word", DataTypes.STRING())  
6     .column("count", DataTypes.BIGINT())  
7     .build())  
8     .build())
```

### pandas

```
1 print(df)
```

- Make a *streaming data* library with a similar developer experience to a *static data* library

## Reading a CSV file

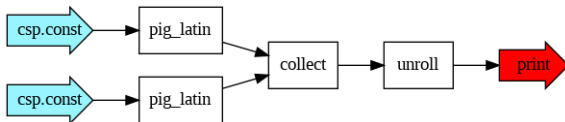
```
1 reader = CSVReader("data.csv", time_format)
2 stream = reader.subscribe_all(str)
```

## Writing to stdout

```
1 csp.print(stream)
```

# Brief walkthrough of CSP

- CSP is a **computation graph** framework



- 1 **Input adapters** bring data into the graph
  - 2 **Nodes** transform data
  - 3 **Output adapters** send data out of the graph
- `csp.graph`: collection of these 3 elements

# Hello World in CSP

```
1 @csp.graph
2 def hello_world():
3     hw = csp.const("hElLo wOrLd") # input adapter
4     fmt = csp.apply(hw, lambda x: x.title(), str) # node
5     csp.print("Message", fmt) # output adapter
6
7 csp.show_graph(hello_world)
```



# Hello World in CSP

```
1 from datetime import datetime, timedelta
2
3 csp.run(hello_world, starttime=datetime.utcnow(),
4         endtime=timedelta(seconds=1), realtime=True)
```

## Output

```
1 2024-10-27 13:10:33.160102 Message:Hello World
```

# Hello World in CSP

- CSP keeps its own internal *engine time*
- Seamless transition to playback mode

```
1 csp.run(hello_world, starttime=datetime(2022,1,1),  
2         endtime=timedelta(seconds=1), realtime=False)
```

## Output

```
1 2022-01-01 00:00:00 Message:Hello World
```

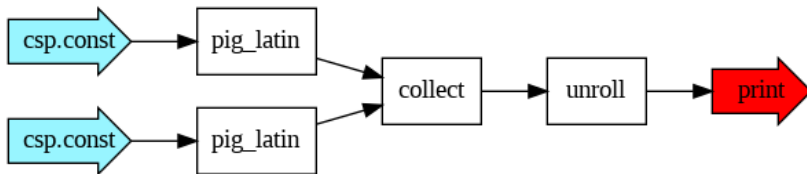
# Writing your own CSP nodes

- Nodes transform **time series** data

```
1 from csp import ts
2
3 @csp.node
4 def pig_latin(text: ts[str]) -> ts[str]:
5     is_vowel = lambda c : c in "aeiou"
6     assert len(text) > 2
7     if not is_vowel(text[0]):
8         if is_vowel(text[1]):
9             pl = text[1:]+text[0]
10        else:
11            pl = text[2:]+text[:2]
12    else:
13        pl = text+"w"
14
15    return pl+"ay"
```

# Writing your own CSP nodes

```
1 @csp.graph
2 def ellohay_orldway():
3     hello = csp.const("hello")
4     world = csp.const("world")
5     latin = [pig_latin(hello), pig_latin(world)]
6     csp.print("Pig Latin", csp.flatten(latin))
7
8 csp.show_graph(ellohay_orldway)
```

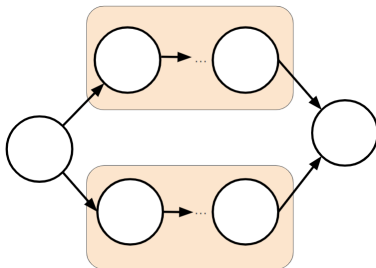




# Nodes and graphs

- Nodes are **atomic** units of computation
- Graphs are collections of nodes

```
1 @csp.graph
2 def subgraph(x: ts[T]):
3     ...
4
5 @csp.graph
6 def main_graph():
7     for source in
8         data_sources:
9             subgraph(source)
```



# Pre-written nodes and adapters

- CSP engine is written in **C++** for performance
- Includes an extensive library of optimized nodes (control flow, statistics etc.)
- Input and output adapters for common formats (Parquet, Kafka, Perspective)

# Current limitations of CSP

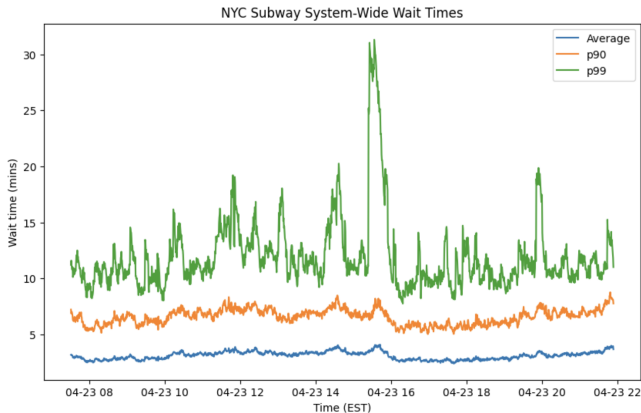
- Fewer pre-built adapters than some other frameworks
- Developer API is only in Python
  - Flink, Beam available in multiple languages
- Can execute graphs in **parallel** but not nodes
  - Node-level parallelization is often too fine grained

# Areas where CSP can enhance developer experience

- 1 Research before deployment
- 2 Debugging
- 3 Unit testing
- 4 Performance tuning
- 5 Interactive execution (jupyter)

# Research before deployment

- Playback **time-series** data for research
- CSP handles event ordering and synchronization



# Research before deployment

- **Look-ahead bias:** using future data in historical analysis
- Example: knowing delays in transit route planning
- CSP removes look-ahead bias in research

```
>> python e_01_nyct_subway.py 635:456 L03:L R20:NQRW --num_trains 5

2024-04-20 00:27:24.166629 Departure Board:
  At station 14 St-Union Sq

Uptown L train to 8 Av in 1 minutes
Downtown L train to Canarsie-Rockaway Pkwy in 1 minutes
Uptown L train to 8 Av in 2 minutes
Uptown L train to 8 Av in 5 minutes
Downtown L train to Canarsie-Rockaway Pkwy in 5 minutes
```

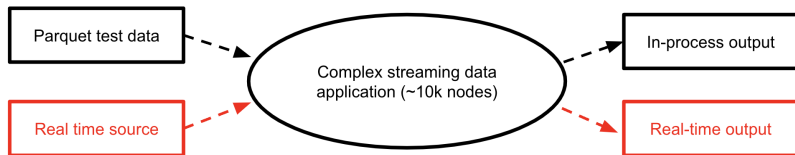
- Single application for research and deployment
- Replay event stream to replicate a crash or error



**Prototype *and* deployment**

# Unit testing

- Real-time data is unstable in CI tests
- Static data allows for more comprehensive tests
- In CSP, simply swap out your adapters



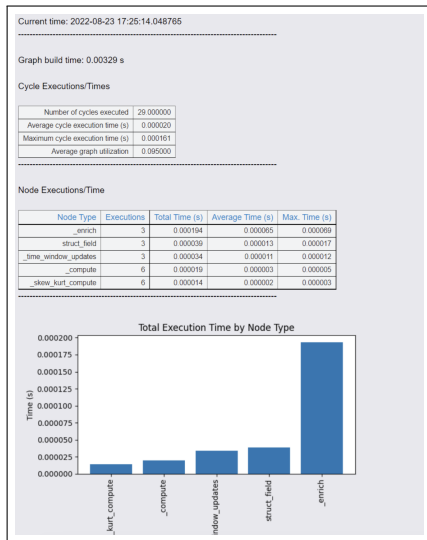


- *"Premature optimization is the root of all evil"* - Knuth
- CSP provides a **live profiler** to identify bottlenecks in your application
- Requires 1 line of code to configure!

```
1 with profiler.Profiler(http_port=8888) as p:  
2     # run the graph normally  
3     results = csp.run(graph, starttime=st, endtime=et)
```

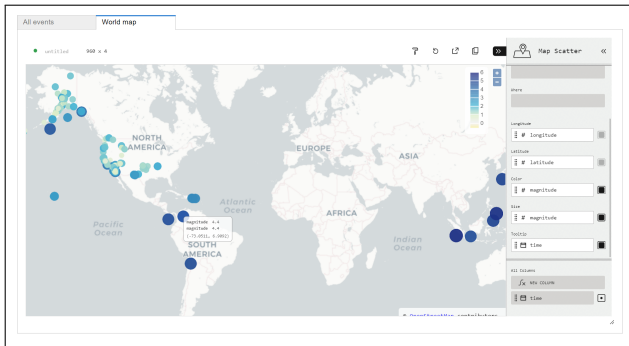
# Performance tuning

- Live timing metrics
- Periodic memory snapshots



# Interactive execution

- CSP is fully functional in IPython/Jupyter
- Useful for prototyping and interactive demos



Live earthquake tracker running in a Jupyter notebook

# Stream processing is hard, but it doesn't have to be!

- CSP abstracts away the complexity of stream processing
- Over 30 examples at [github.com/Point72/csp](https://github.com/Point72/csp)
- User projects at [github.com/csp-community](https://github.com/csp-community)

```
pip install csp
```

