# Introducing STAC-ML

**Bishop Brock**
**Head of Research, STAC**

bishop.brock@STACresearch.com

**Peter Nabicht**
**President, STAC**

peter.nabicht@STACresearch.com

# History

- Driven by financial firms
  - Motivation: market making, hedging, customer pricing, etc.

- STAC-ML working group has refined the original POC idea into a finished benchmark specification

- Tech vendors provided crucial input

- But control ultimately rests with users – i.e., those who must deliver business value from technology in the real world
  - Like all STAC Benchmarks

**STAC** ®
SECURITIES TECHNOLOGY ANALYSIS CENTER
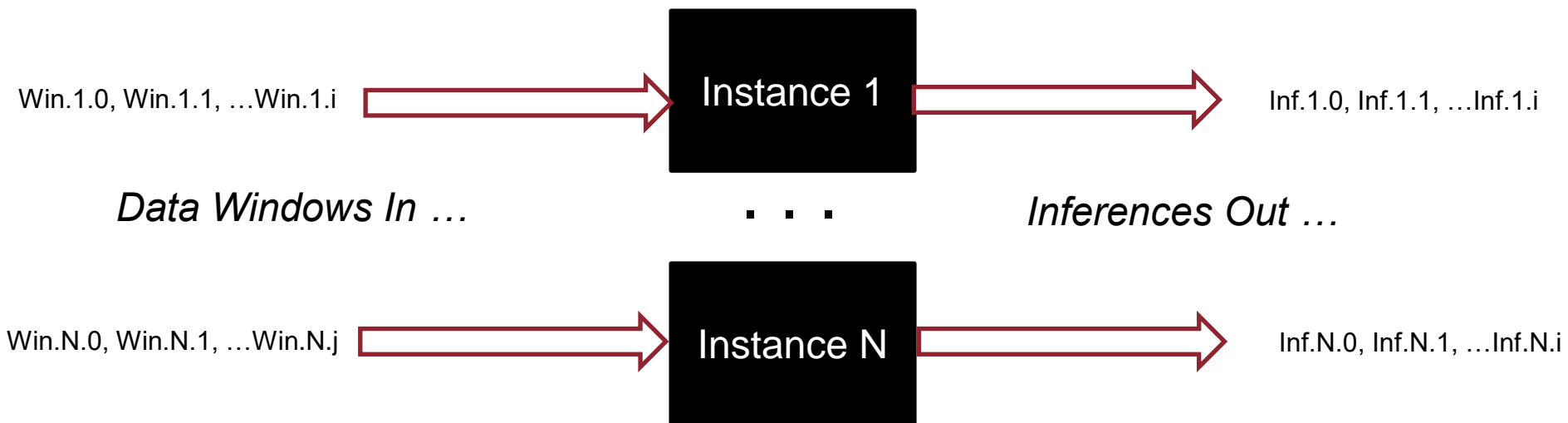
# Latest Status

- Off to the races!

- The benchmark specifications, test harness, reference implementation, and documentation are released

- 5 vendor implementations are currently underway

- Test Harness engineered to allow end-users to "mark their own stacks to market"

**http://www.STACresearch.com/ml**

STAC®
SECURITIES TECHNOLOGY ANALYSIS CENTER

# Basics

- LSTM models that simulate real models derived from market data

- Goal: isolate <u>inference</u> performance
  - Inference engine software
  - Underlying processors, memory, accelerators, etc.
  - Anything required to optimally use the former with the latter (e.g., data transfer to processor memory)

- Metrics:
  - Latency, throughput, power efficiency, space efficiency, error

- Benchmarks allow any level of precision (including mixed-precision)

**S T A C**
SECURITIES TECHNOLOGY ANALYSIS CENTER

# Scale Dimensions

- ## Model size
  - Three are currently specified
  - Input data window scales with model size

- ## Number of Model Instances running in parallel
  - As specified by the SUT provider



Win.1.0, Win.1.1, …Win.1.i → **Instance 1** → Inf.1.0, Inf.1.1, …Inf.1.i

*Data Windows In …*  . . .  *Inferences Out …*

Win.N.0, Win.N.1, …Win.N.j → **Instance N** → Inf.N.0, Inf.N.1, …Inf.N.i

STAC®
SECURITIES TECHNOLOGY ANALYSIS CENTER
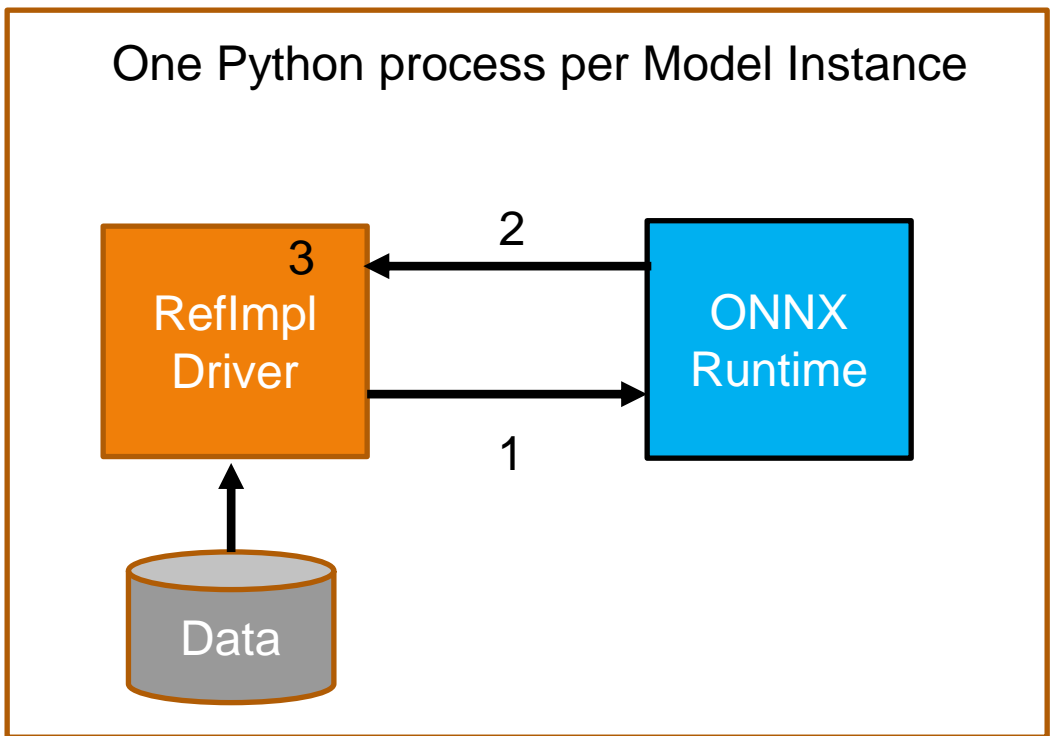
# Use Cases and Optimizations

- Different Use Cases:
    - Trading – Latency Optimization
    - Backtesting – Throughput Optimization

- Optimization tradeoffs (latency vs throughput vs efficiency vs error) are up to the SUT provider
    - The tests collect all metrics every time, no matter the optimization goal
    - Any quantization scheme allowed, if used consistently

**S T A C** ®
SECURITIES TECHNOLOGY ANALYSIS CENTER
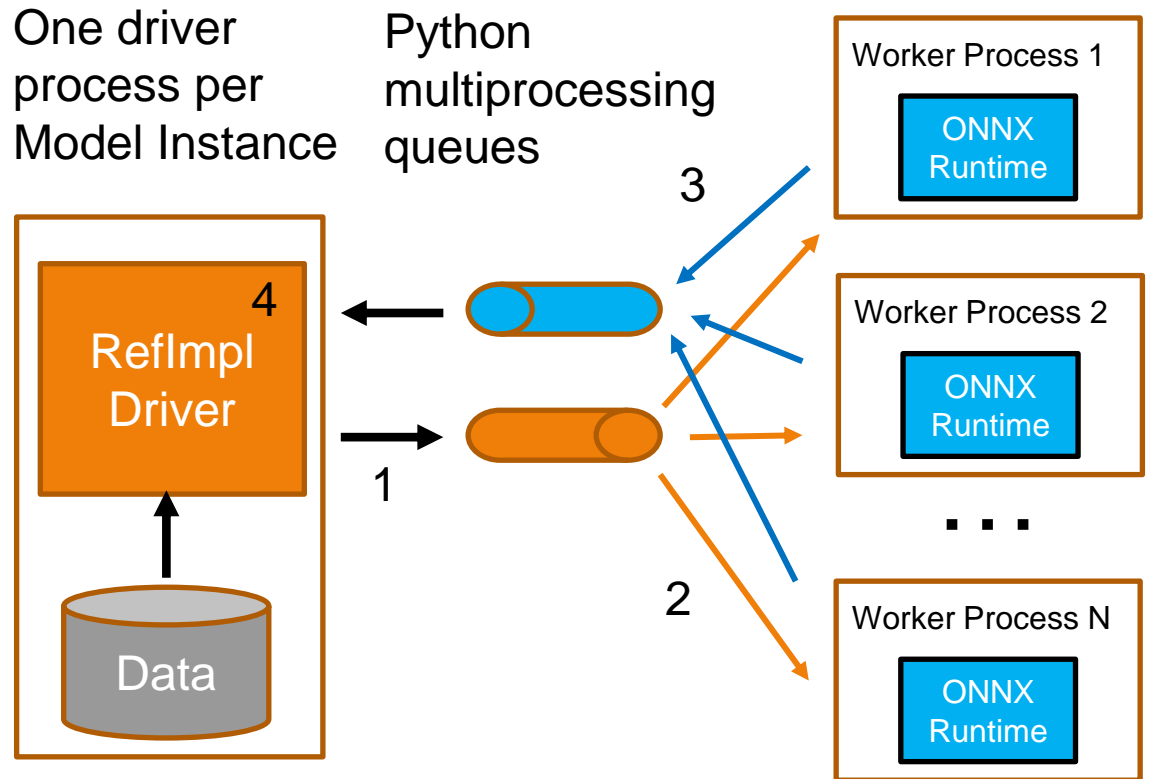
# Speaking of Tradeoffs…

- STAC has just published the first audit reports, representing internal research

- 2 SUTs:
  - Latency Optimized
  - Throughput Optimized

- Same Software (different tuning):
  - Pure Python Implementation
  - STAC-ML Markets (Inference) Naive Implementation
  - Unmodified ONNX 1.11.0 runtime

- Same Hardware:
  - 60-vCPU @ 3.1Ghz, Sole-Tenant cloud instance, 240GiB Memory

- Models: Standard benchmark models at default FP32 precision

**STAC®**
SECURITIES TECHNOLOGY ANALYSIS CENTER

# Serial Implementation | Parallel Implementation



One Python process per Model Instance

```
3                    2
RefImpl  ←──────────   ONNX
Driver   ──────────→   Runtime
              1
   ↑
 Data
```

1. Driver calls ONNX runtime with data window; waits for response
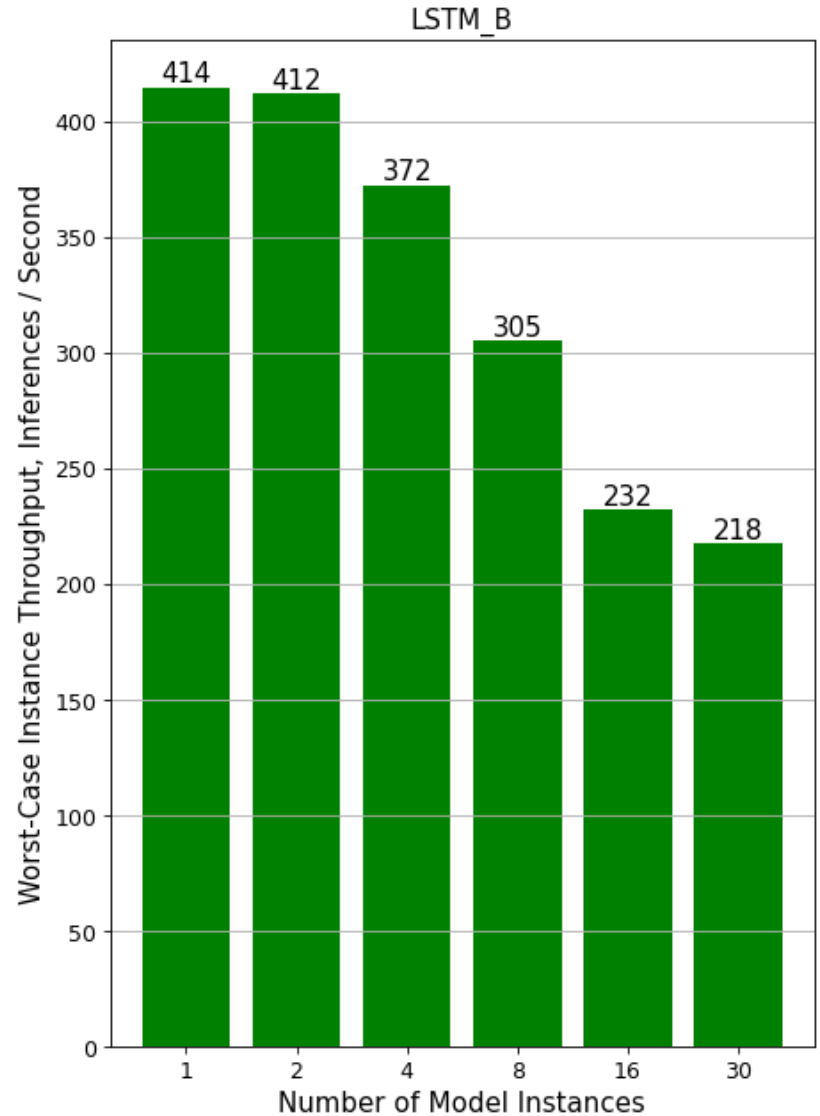2. ONNX returns inference value
3. Driver stores value in memory

Repeat…

One driver process per Model Instance

Python multiprocessing queues

Worker Process 1
ONNX Runtime

Worker Process 2
ONNX Runtime

. . .

Worker Process N
ONNX Runtime

```
4
RefImpl   ←──────  [queue]
Driver    ──────→  [queue]
   ↑          1
 Data         2
```
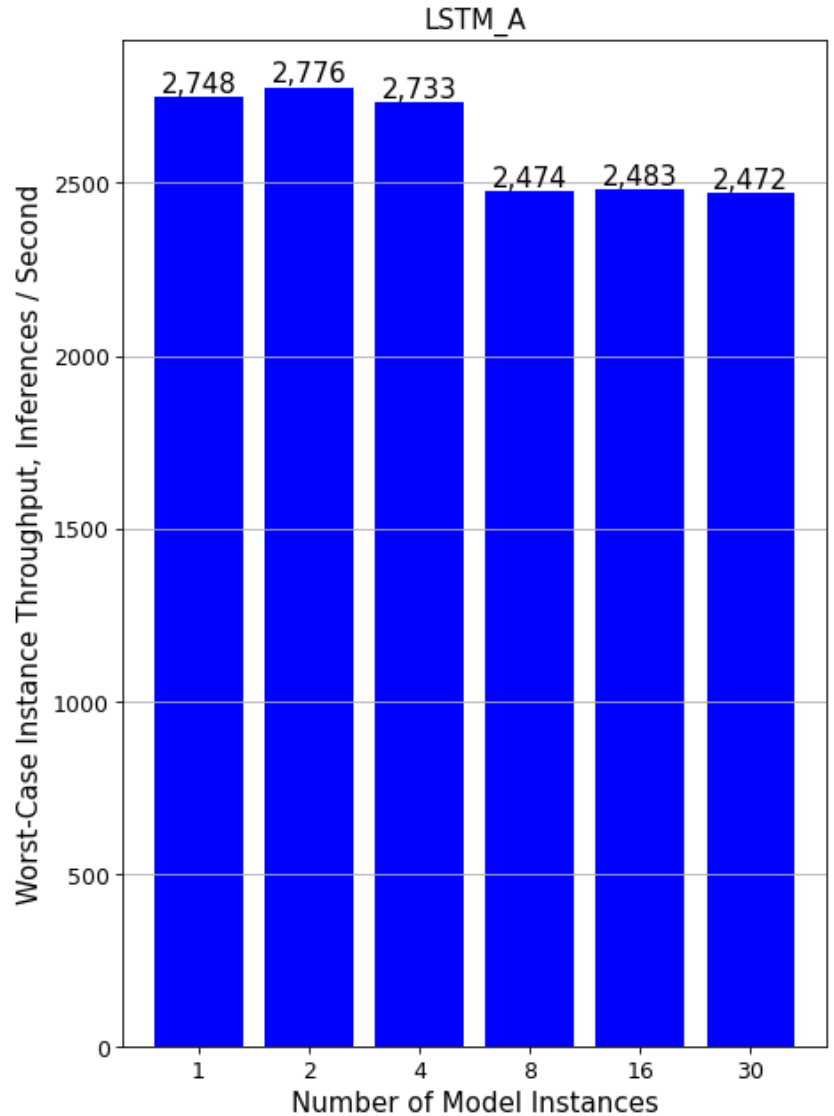
1. Driver inserts data window(s) into queue; waits for response
2. Workers take windows from queue
3. Workers store results in queue
4. Driver gets result from queue and orders it in memory

Repeat…

Multiple worker processes per Model Instance

S T A C
SECURITIES TECHNOLOGY ANALYSIS CENTER

# Optimization 'Knobs'

- ## Number of Model Instances
  - We looked at 1, 2, 4, 8, 16 and 30 instances

- ## Number of ONNX threads
  - ONNX can (sometimes) effectively utilize multiple threads per model

- ## Number of Parallel Instances
  - Note: 1 Parallel Instance == Serial Model
  - We have a choice of allocating a HW thread to either an ONNX thread or a Python process
  - The optimal choice again varies by model and optimization goal

- ## A Research Note describing the optimization experiment is available in the STAC Vault.

**S T A C** ®
SECURITIES TECHNOLOGY ANALYSIS CENTER

# Worst-Case Instance Throughput Comparison, LSTM_A vs LSTM_B

LSTM_A vs. LSTM_B: Instance Throughput of Latency-Optimized Configurations
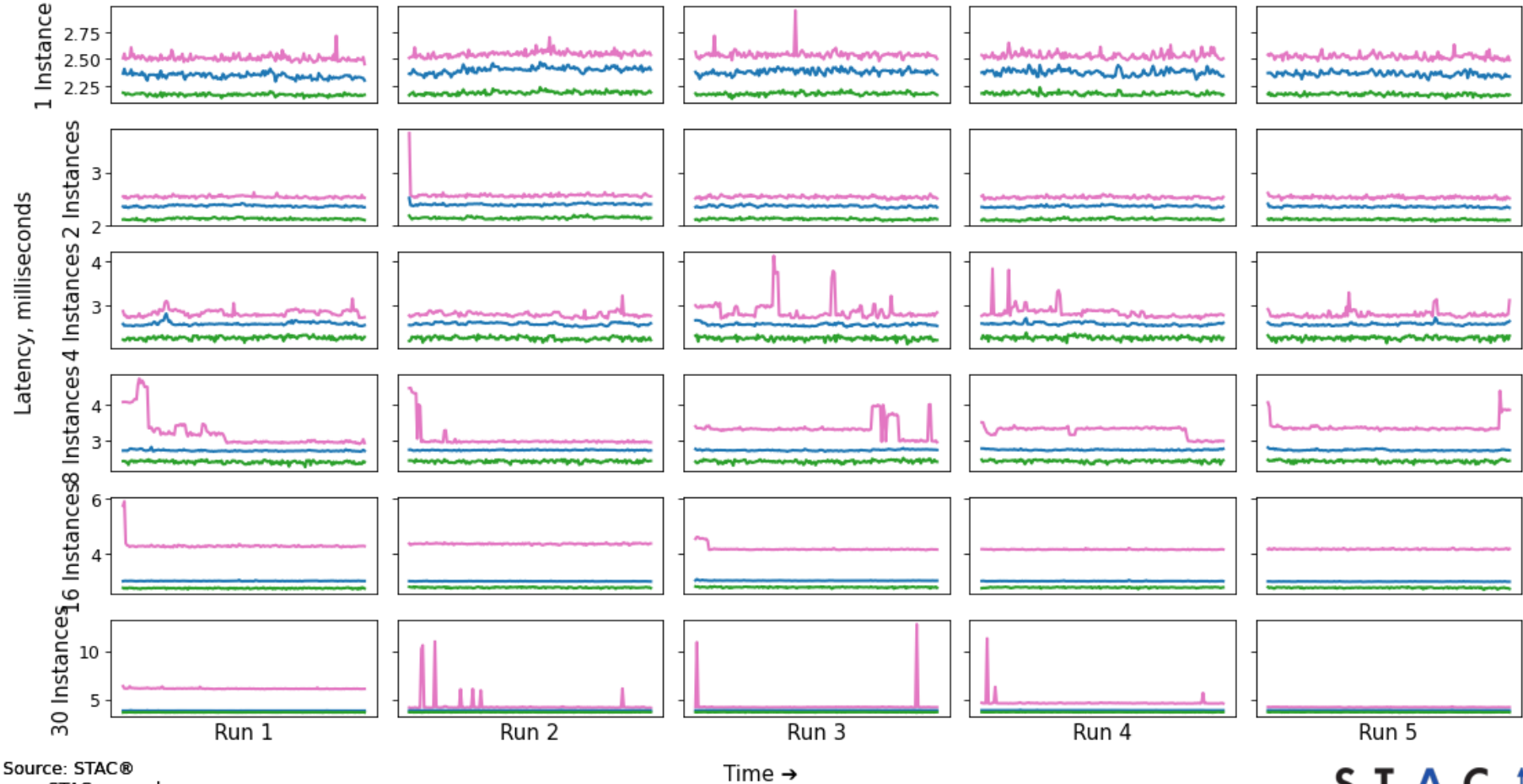
# Latency Trends



STAC-ML™ Markets (Inference)

STAC-ML™ Markets (Inference) Naïve Implementation
on a 60-vCPU Sole-Tenant Cloud Node with 240 GiB Memory
Latency-Optimized Configuration
SUT ID: STAC220503

LSTM_B: Latency Trends Across Runs

Latency Statistics Over
1 Second Intervals
— 99th Percentile Latency
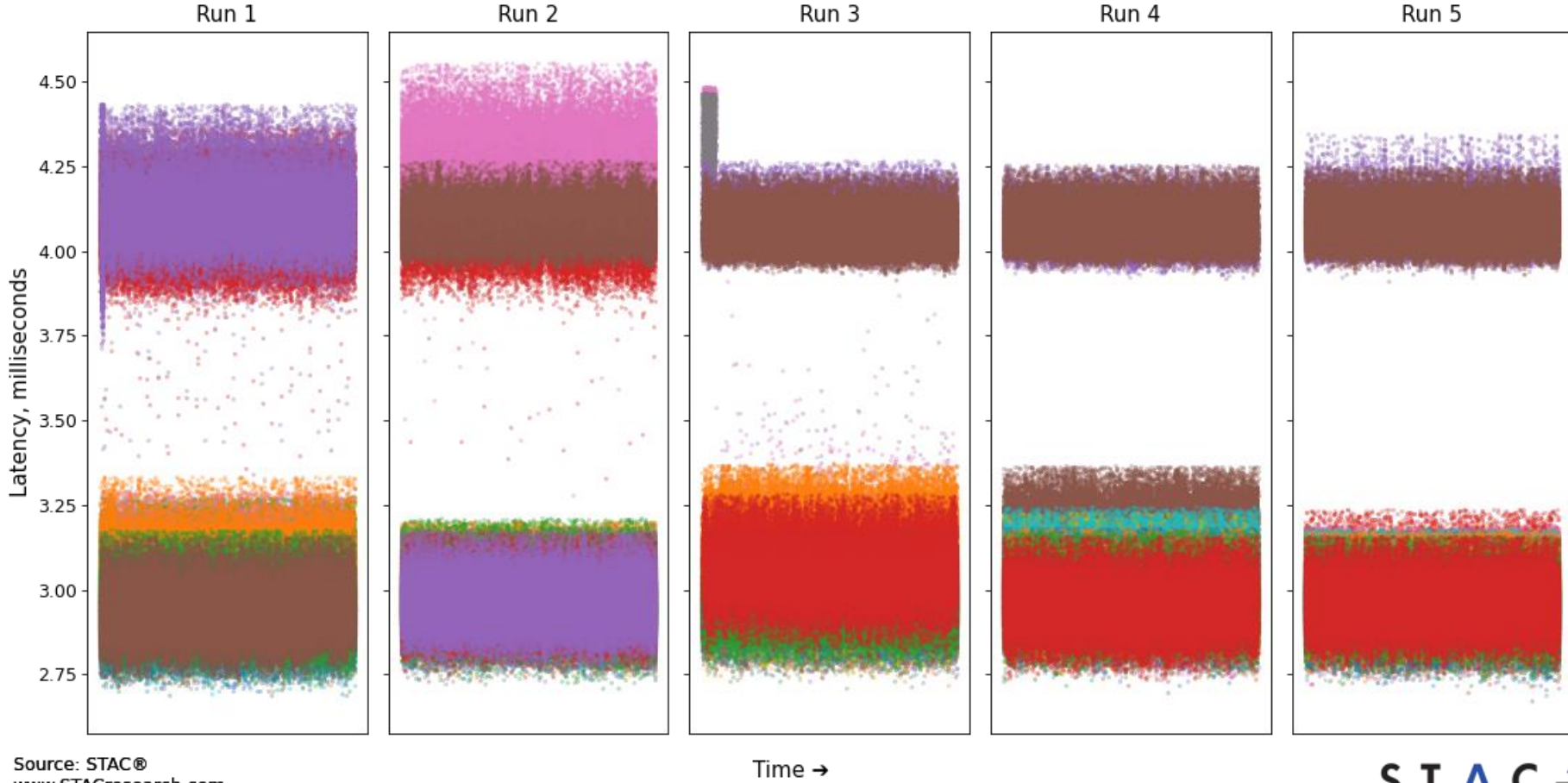— Median Latency
— Min Latency

# Latency Details



STAC-ML™ Markets (Inference)

STAC-ML™ Markets (Inference) Naïve Implementation
on a 60-vCPU Sole-Tenant Cloud Node with 240 GiB Memory
Latency-Optimized Configuration
SUT ID: STAC220503

LSTM_B, 16 Model Instances
Detailed Latency over Time, Samples Above the 99.0th% Quantiles per Model Instance/Run Omitted

# How to get Access

- All STAC subscribers can access

  - Public STAC Reports

- Premium Subscribers can access

  - Benchmark Specifications

  - Highly detailed configuration information

  - Extensive, detailed visualizations and tables on Performance, Efficiency, and Error

  - Code for test harness, generating post-test visualizations, and STAC Packs

  - Additional reports and research in the confidential STAC Vault™

**http://www.STACresearch.com/ml**

**S T A C**®
SECURITIES TECHNOLOGY ANALYSIS CENTER
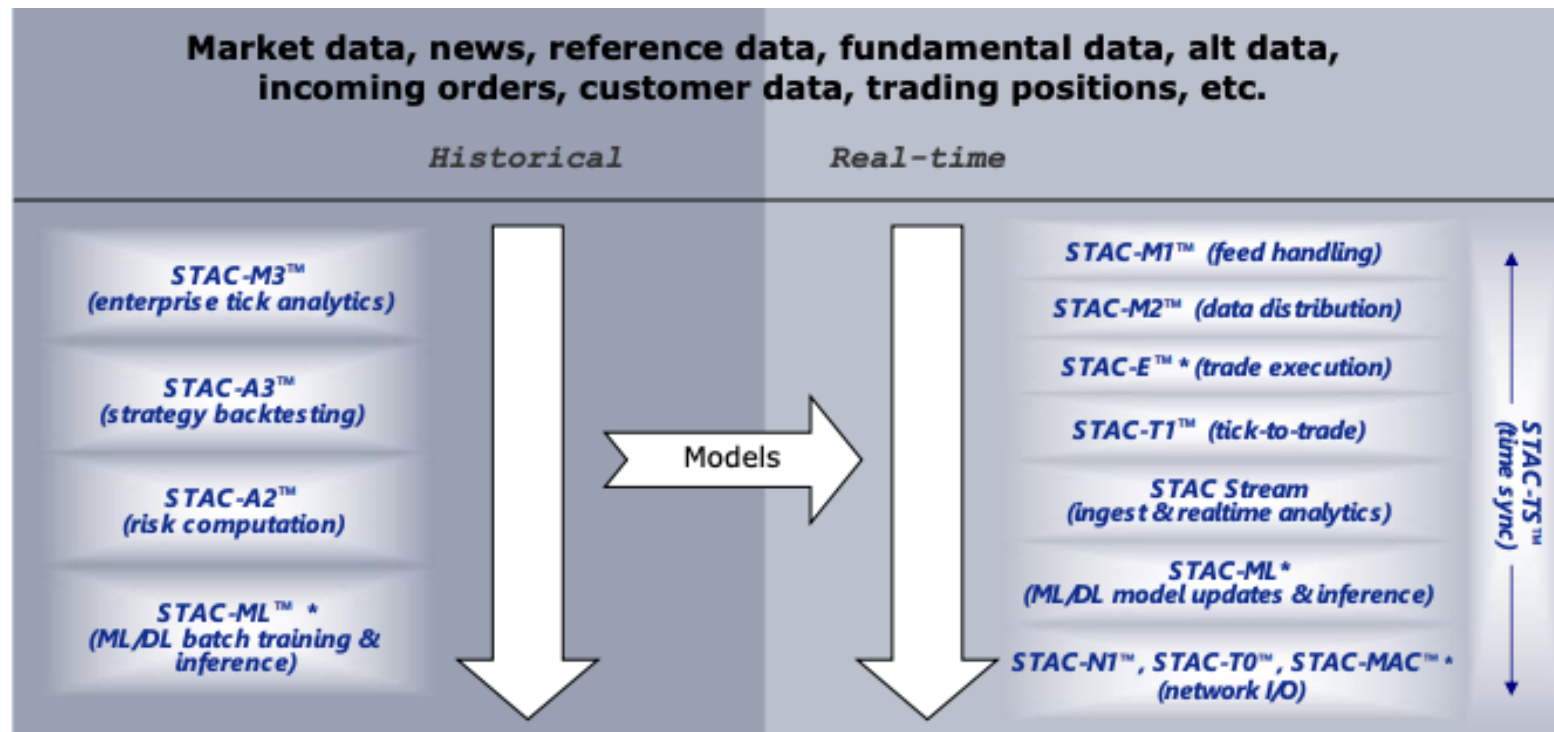
# STAC Tracks

Analytics STAC Track          Tradeflow STAC Track

STAC-ML Markets (Inference) is included in both previously existing STAC Tracks

No action needed by current subscribers to those

# Machine Learning STAC Track

- New STAC Track that includes
  - STAC-ML Markets (Inference)
  - Future STAC-ML benchmarks

- Free trial for the remainder of 2022
  - For those responsible for ML research and infrastructure
  - Full access, including STAC Vault content
  - To request the trial:

**council@STACresearch.com**

**STAC**®
SECURITIES TECHNOLOGY ANALYSIS CENTER