

BD|CESGA

Anomaly Detection at Scale

Javier López Cacheiro
Orquídea Seijas Salinas
Samuel Soutullo Sobral



Measurement is the first step that leads to control and eventually to improvement. If you can't **measure** something, you can't understand it. If you can't **understand** it, you can't control it. If you can't **control** it, you can't **improve** it.

H. James Harrington

Anomaly Detection

- Phase 1: Measure → Collect & Store
- Phase 2: Understand → Analyze&Visualize
- Phase 3: Control → Monitoring
- Phase 4: Improve → Anomaly Detection

Phase 1: Measure

Collect & Store

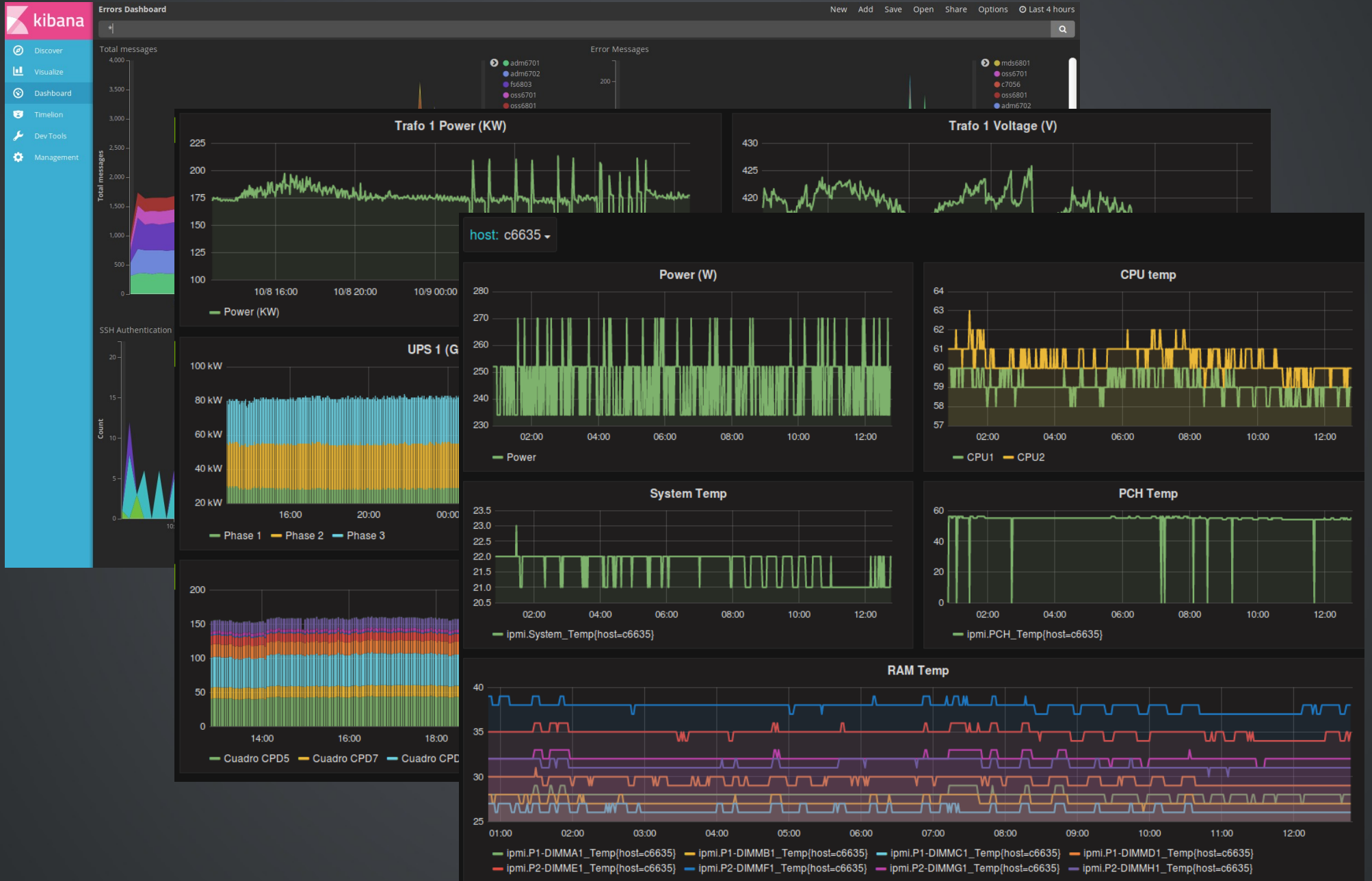
- System logs
- Process accounting
- Server metrics
- IPMI sensors
- SMART disk info
- Chillers & AHUs
- UPSs
- Power Meters
- Network traffic
- Application module usage
- SLURM Accounting
- Temperature/Humidity sensors

33487 metrics

10 Million time series



Phase 2: Understand

Dashboards






Phase 3: Control

Monitoring

Current Incidents Overdue Muted  

Service Problems

CRITICAL 12:28	c7235: Log Alerts Singularity: sexec (U=12529, P=22457)> Retval = 255	! 
CRITICAL 12:28	c7236: Log Alerts Singularity: sexec (U=12529, P=28244)> Retval = 255	! 
CRITICAL 1d 22h	c0511: ssh Server answer:	!
CRITICAL Sep 27	puertos2.cesga.es: Radares 3, KO, 20170731, ftp_off, PHY_RADARH, ID1	!
CRITICAL Aug 3	gestion-sis.cesga.es: Partition /mnt/EMC/Store_uscfm DISK CRITICAL - free space: /mnt/EMC/Store_uscfm 115354 MB (0.71% inode=100%):	! 

Recently Recovered Services

OK 1m 6s	c0621: Load OK - load average: 32.00, 31.97, 31.70
OK 20m 28s	c1102: ping4 PING OK - Packet loss = 0%, RTA = 0.44
OK 46m 3s	c7101: IB Status (0x00000000000000000000000000000000 00000000000000000000000000000000 00000000000000000000000000000000 00000000000000000000000000000000 0000)
	c7102: IB Status (0x00000000000000000000000000000000

Phase 4: Improve

Anomaly Detection

- Types of anomalies in time series
 - Outliers
 - Change points
 - Anomalous time series
- Generic Anomaly Detection Systems

Caution: Alert overload



Server Anomalous Performance

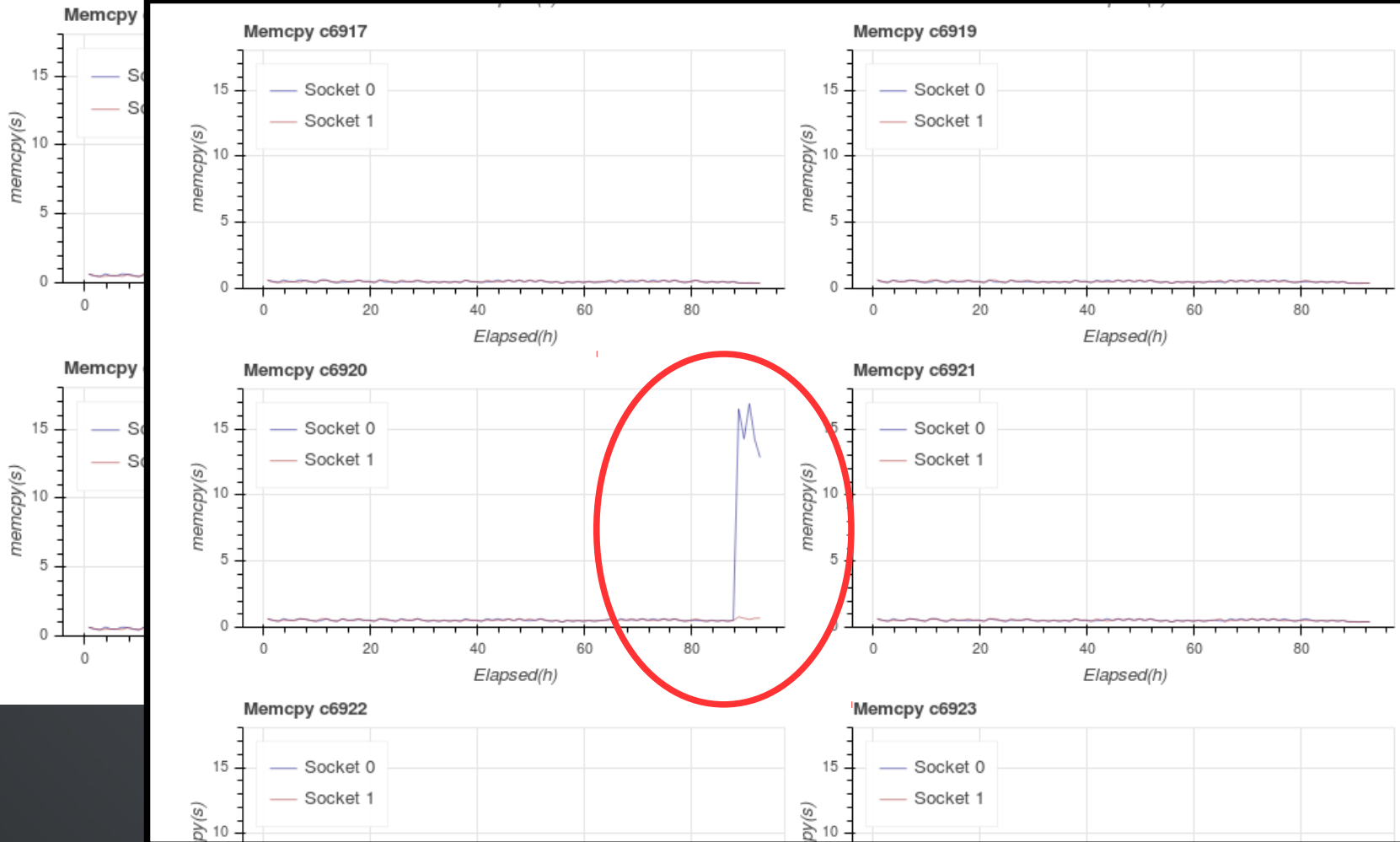
- **Problem:** Parallel jobs are cancelled because some of the nodes have poor performance. Computation is lost.
- **Detection:** Analyze & visualize server metrics to spot the anomalous node
- **Objective:** Automatically detect low performance nodes

Analyze & Visualize

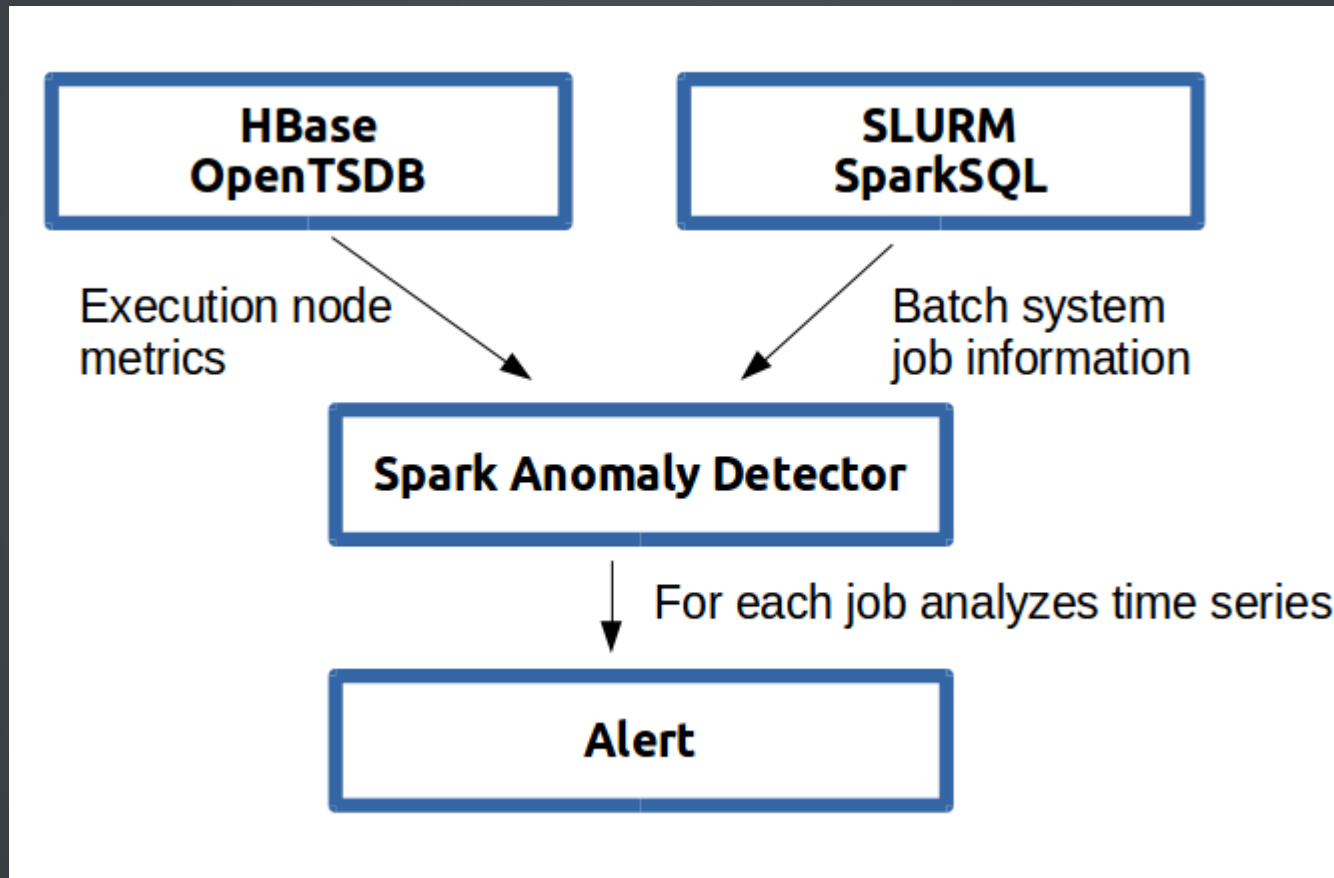
jupyter memcpy_results_for_slurm_job Last Checkpoint: 10/02/2017 (autosaved)

File Edit View Insert Cell Kernel Help Trusted

Code



Anomalous Performance Detection



Results

- 6 months
- 11965 jobs
- 22 anomalies detected
 - Precision: 100%
 - Recall: 96%
 - F-score: 0.98

Conclusions

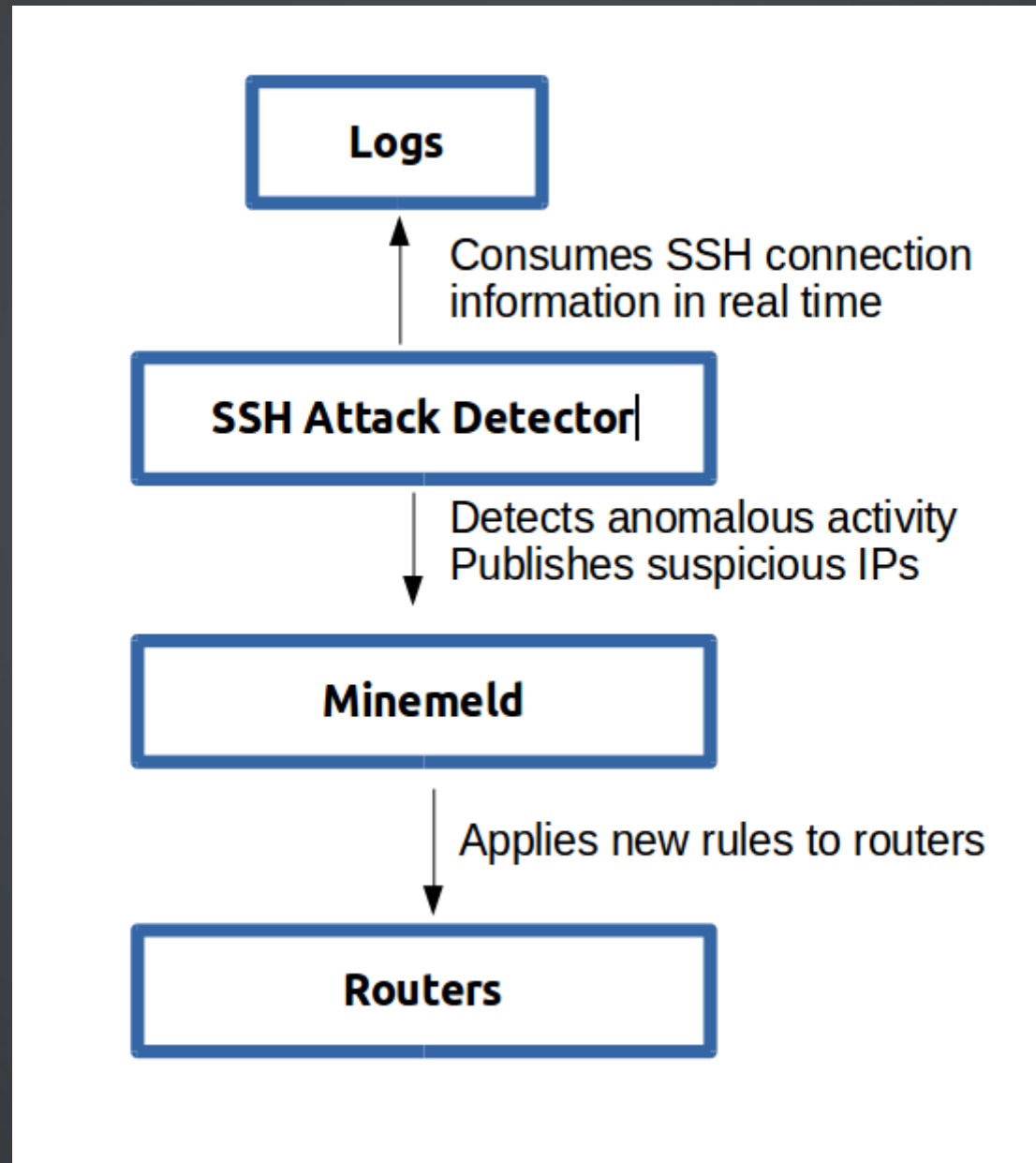
- No longer needed to delete old data
- Understand your data
- Generic anomaly detection systems generate too many alerts
- Target specific use cases
- Maintain number of alerts low

Thanks!

SSH Attack Detection

- **Problem:** Daily our public servers are scanned and attacked
- **Detection:** Correlate real-time SSH connection information to detect attacks
- **Objective:** Automatically update router configuration to stop the attacks

SSH Attack Detection




Results

STATISTICS

24 hours

METRIC	CURRENT	HISTORY (LAST 24H)
--------	---------	--------------------

INDICATORS	59	
------------	----	--

METRIC	SINCE ENGINE START	HISTORY (LAST 24H)
--------	--------------------	--------------------

ADDED	609	
-------	-----	--

AGED_OUT	691	
----------	-----	--

Results

STATISTICS

7 days

