

Cutting the edge in fighting cybercrime: reverse-engineering a search language to cross-compile it to PySpark



Jude Ken-Kwofie

Principal Engineer, HSBC Cybersecurity Sciences & Analytics



Abigail Shriver Lead Cybersecurity Engineer, HSBC Cybersecurity Sciences & Analytics

ORGANIZED BY Sdatabricks





Serge Smertin Senior Resident Solutions Architect, Databricks



226k+ Employees services company.

HSBC is a multinational investment and financial

40m Customers

~64 Countries

~1.3m Assets



\$4.24m Average Breach Cost

\$401m Mega Breach Cost



4 Main Cloud Providers



The scale of cybersecurity at HSBC

1,880+ Global cybersecurity team members

16 countries and territories with dedicated cybersecurity centers

200+ Global organisations receive intelligence feeds from HSBC

75+ Global cybersecurity education events held for employees each year





What do we mean by Advanced Data Analytics?

Using advanced algorithms on massive, complex data sets to derive actionable insights at pace for cyber





Introductions

Jude Ken-Kwofie

Principal Software Engineer Cybersecurity Sciences & Analytics HSBC

Jude Ken-Kwofie is a Principal Software Engineer at HSBC. He has over 15 years of experience across engineering, data systems, and security. He is the author of several telecom battery standards and led projects for large-scale financial services and adtech data solutions. At HSBC, Jude's role is to help engineer next-gen Cybersec analytic capabilities. When not working, Jude is busy with his Business PhD researching alternative small business financing and decision support systems.



Abigail Shriver

Lead Cybersecurity Engineer Cybersecurity Sciences & Analytics HSBC

Shriver, Lead Cyber Security Abigail Engineer at HSBC's Cybersecurity Sciences and Analytics division, develops in-house data-driven endpoint, cloud, & network security analytics. Previously, she developed the largest centralized data lake within Capital One and leveraged the PBs of data within the lake to identify anomalies & remediate vulnerabilities. She attended The George Washington University studying CS and IA concentrating in Cyber Security & Security Policy.



Serge Smertin



Senior Resident Solutions Architect

In his over 15 years of career, Serge has solutions. been dealing with data cybersecurity, and heterogeneous system integration. His track record got novel ideas from whiteboard to operating them in production vears. like for large-scale malware forensic analysis for the cyber-threat intelligence, or real-time data science platform as the basis for anomaly detection and decision support systems for industry-leading payments service an provider.







- 1. Intro
- 2. Traditional SIEMs for Cyber
- 3. Migration in Practice
- 4. Business Outcomes
- Search Processing Language Transpiler
 Questions







Traditional SIEMs for Cyber



Jude Ken-Kwofie Principal Engineer

HSBC Cybersecurity Sciences & Analytics

ORGANIZED BY 😂 databricks

Cybersecurity is a massive big data problem

100-200 TB/day x 13 months = 38-79 PB



Everything is an asset



Storing 10x the contents of every book in the US Library of Congress, every day





Traditional Security Information and Event Management (SIEM)

Cyber Analysts prefer SIEMs and Search Processing Languages

SIEM Cybersec Analysis



Cybersec Analysts / Threat Hunter







Traditional SIEMs Engineered for Processing Recent Data Traditional SIEM Processing > 100TB+ / day is Challenging

- SIEMs struggle with processing large data volumes, >30 TB/day
- Cybersecurity Data Volumes > 100-200 TB / day
- \$\$\$\$
- Joins and Aggregations
- HSBC Cyber Lakehouse SIEM developed to:
 - Overcome challenges
 - Empower Cyber analysts and Threat Hunters

Lift-n-Shift is Challenging







So much data – so hard to get insight



Lift-n-Shift Threat Detections to PySpark Benefits of Approach

- Process > 100TB / Day
- Empower SIEM economically
- Use advanced analytics & ML
- Near real-time detections
- Extensible
 - Search & Processing

Language Analytic

index=*dns*
|sourcetype=*cloud*
|stats sum(bytes_out)
|rename sum(bytes_out) AS "Total Bytes"
|table domain_name ...





Migration via Transpiler





Transpiling Search Processing Language Code Example

```
import spl.Transpiler
println(Transpiler.toPython("""
index=main
(sourcetype::ProcessRollup2* OR sourcetype::CommandHistory*)
(TERM(ProcessRollup2) OR TERM(CommandHistory))
TERM(TODO_BACKSLASHES)
   CommandHistory="*\\XXX\\XXXX\\XXXX\\XXXX\\XXXX\\XXXX
   OR
   CommandLine="*\\XXX\\XXXX\\XXXX\\XXXX\\XXXX\\XXXX\
   earliest=-1h
|| rex max_match=0 field=_raw "(?i)(?<IOC>(YYYYY)YYYYYYYYY))"
|| stats count earliest(_time) as earliest latest(_time) as latest
   values(ComputerName) as ComputerName by aid
|| eval x = strftime(earliest, "askjdhlakshdjashd")
11 table aid
          ComputerName
          LocalAddressIP4
          FileName
          event_simpleName
          IOC
          count
          CommandLine
          earliest
          latest
          earliest_h
          latest_h
·····))
                                                                  scala
```

DATA+AI SUMMIT 2022







Migration in practice

How it went



Abigail Shriver

Lead Cybersecurity Engineer HSBC Cybersecurity Sciences & Analytics

ORGANIZED BY 😂 databricks



Support Several Functions/Commands

Commands		Functions		Extensions	
<pre>search eval table rex regex where lookup bin rename join fields convert collec stats sort return mvcombine map inputlookup</pre>	head format fillnull eventstats dedup makeresults mvexpand streamstats addtotals multisearch	<pre>strftime() values() latest() earliest() if() mvcount() coalesce() mvindex() mvappend() null() min() round() max() substr() isnotnull() sum() mvfilter() len() count()</pre>	<pre>memk() rmunit() rmcomma() ctime() auto() num() replace() lower() none() CIDR search</pre>	term()	





How to Build the Transpiler?

Built using maven (3.6.x)

- mvn -DskipTests Plocal package
- A (fat) jar will be created containing the code & all the dependencies required
- Created by default in target folder in the root directory of the repository







Attach to a Databricks Cluster

Attach the created jar through the Libraries tab in the cluster configuration page in Ul.

Configuration	Notebooks (0)	Libraries	Event log	Spark UI	Driver logs	Metrics	Apps	Spark cluster UI - Master -
Committee (7	Non laterid							
Name			Туре	Statua		Source		
🗆 spark_spl	_0_4.jer		JAR	3		dbfs:/FileStore/jars/837b2217_042c_4d32_b275_9e3f0edfee3c-spark_spl_0_4-96e6e.jar		837b2217_o42o_4d32_b275_9e3f0edfee3c-spark_spl_0_4-96a6e.jar

(Or) Use Databricks Cluster init-script (required for extensions like TERM())

<pre>dbutils.fs.put(*/databricks/scripts/spl-transpiler-install.sh*, #1/bin/bash cp /dbfs/tmp/spl/spark_spl_8_3.jar /databricks/jars/spark_spl_8_3.jar ***, True)</pre>	Spark Tags Logging Init Scripts Permissions Init scripts @ Type File path Init Scripts				
	DBFS V dbfs/databricks/soripts/spi-transpiler-install.sh Add				





Cybersecurity Lakehouse Overview

Endpoint Custom Detections

100s of queries to migrate over to
 PySpark on Databricks

 Endpoint detection data is already ingested into the cybersecurity lakehouse and registered in the metastore







Search Processing Language to PySpark







Data in Search Processing Platform != Format in Databricks







Data in Search Processing Platform != Format in Databricks







How to Inject the Enrichment Logic







How to Inject the Enrichment Logic

return (df .where(F.expr("sourcetype LIKE 'ProcessRollup2%'")) .where(F.expr("term('ProcessRollup2')")) .where(F.expr("tern('TODO_BACKSLASHES')")) .where(F.expr("CommandLine LIKE '%\\\\XXX\\\\XXXX\\\\XXXX\\\\XXXX\\\\ .where(F.col('_time') ≥ (F.expr('now()') + F.expr(*INTERVAL .withColumn('IOC' F.regexp_extract(F.col('_raw'), '(?i)(?<IOC>(YYYYYY|YYYYYYYYYYY))', 1) def do_aggregation(df: DataFrame) → DataFrame: return (df .orderBy(F.col('_time').asc()) .groupBy('aid') .agg(F.count(F.lit(1)).alias('count'), F.first(F.col('_time'), True).alias('earliest'), F.last(F.col('_time'), True).alias('latest'), F.collect set(F.col('ComputerName')).alias('ComputerName')) .withColumn('x', F.date_format(F.col('earliest'), 'askjdhlakshdj .select(F.col('aid'), F.col('ComputerName'), F.col('LocalAddressIP4'), F.col('FileName'), F.col('event_simpleName'), F.col('IOC'); F.col('count'), F.col('CommandLine'), F.col('earliest'), F.col('latest'), F.col('earliest_h'), F.col('latest_h'))

def do_as_much_filtering_as_possible(df: DataFrame) → DataFrame:



processrollup2') as_much_filtering_as_possible) h_aidmaster) h_ipv4) aggregation)







Business Outcomes

Why it matters



Abigail Shriver

Lead Cybersecurity Engineer HSBC Cybersecurity Sciences & Analytics

ORGANIZED BY 😂 databricks



Automation at Scale



Scheduled or executed one time	Pass paramete notebooks	rs through	Executi Microsoft Azure Dat	ON tabricks	
$\gg _$ <u>Idd</u> $\sim - \times$			G. Home	😂 Azure Data	bricks
	Widget Demo (Pytho	n)	₩uvikapace © Paccents	{ } •	
Task name * @ Order_ETL	- Detached		Deta Ge	Explore the Quickstart Tutorial Spin up a cluster, run queries on preloaded dat display results in 5 minutes.	a, and
Type * Source * @	ene Dottaoned		Clusters	and all second as a community	
Notebook Git (production) Path * @ order-etl/notebooks/ Cluster * @	year 2014	~	Jobs 52 Models Search	Common Tasks	۲]]]
Job_cluster (274.50 GB 36 Cores DBR 10.4 LTS Spark 3.2.1 Scala Parameters UI JSON Add Advanced options >					R















Historical and Scalable Metrics









Incident Response









Learning is Good

Educational



Onboarding and usability









Anatomy of a translating compiler

... or how to use debugger all the time.



Serge Smertin Senior Resident Solutions Architect

Databricks

ORGANIZED BY 🗟 databricks





Apache Spark query execution recap





code IN(4*, 5*)



code IN(4*, 5*)

```
def argu[_: P]: P[Expr] = termCall | call | constant
def parens[_: P]: P[Expr] = "(" ~ expr ~ ")"
def primary[_: P]: P[Expr] = unaryOf(expr) | fieldIn | parens | argu
def expr[_: P]: P[Expr] = binaryOf(primary, ALL)
```

```
def impliedSearch[_: P]: P[SearchCommand] =
    "search".? ~ expr.rep(max = 100) map(_.reduce((a, b) => Binαry(a, And, b))) map SearchCommand
```



code IN(4*, 5*)

```
def argu[_: P]: P[Expr] = termCall | call | constant
def parens[_: P]: P[Expr] = "(" ~ expr ~ ")"
def primary[_: P]: P[Expr] = unaryOf(expr) | fieldIn | parens | argu
def expr[_: P]: P[Expr] = binaryOf(primary, ALL)
```

```
def impliedSearch[_: P]: P[SearchCommand] =
    "search".? ~ expr.rep(max = 100) map(_.reduce((a, b) => Binary(a, And, b))) map SearchCommand
```

```
test( testName = "code IN(4*, 5*)") {
    p(impliedSearch(_), SearchCommand(
        FieldIn("code", Seq(
            Wildcard("4*"),
            Wildcard("5*")
            ))))
```



```
SearchCommand(
  FieldIn("code", Seq(
    Wildcard("4*"),
    Wildcard("5*")
)))
```



```
SearchCommand(
  FieldIn("code", Seq(
    Wildcard("4*"),
    Wildcard("5*")
)))
```



Or(

Like(UnresolvedAttribute("code"), Literal.create("4%"), '\\'), Like(UnresolvedAttribute("code"), Literal.create("5%"), '\\'))





```
SearchCommand(
  FieldIn("code", Seq(
    Wildcard("4*"),
    Wildcard("5*")
)))
```



Or(

Like(UnresolvedAttribute("code"), Literal.create("4%"), '\\'), Like(UnresolvedAttribute("code"), Literal.create("5%"), '\\'))



display(spark.table('main')
.where((F.col('code').like('4%') |
 F.col('code').like('5%'))))





Thank you

Jude Ken-Kwofie



Principal Software Engineer, Cybersecurity Sciences & Analytics HSBC

jude.ken-kwofie@us.hsbc.com

Abigail Shriver

Lead Cybersecurity Engineer, Cybersecurity Sciences & Analytics HSBC

abigail.shriver@us.hsbc.com

cybersecurity@hsbc.com



Serge Smertin



Senior Resident Solutions Architect Databricks

serge.smertin@databricks.com





DATA+AI SUMMIT 2022



https://databricks.com/customers/hsbc

Thank you

