

# Processing and Analyzing Increasing Amounts of Traffic Data at eMAG

Laurentiu Matei  
Budapest Data Forum  
June 24th 2021

**Hello**

# Who Am I?

---

- **Faculty of Automatic Control Bucharest 2003**
- **Developed software for various industries**
- **Joined eMAG in 2011**
- **Leading people and teams in software development after 2013**
- **Business Intelligence Director starting February 2020**

# BI Teams

---

- **100 employees**
- **Roles: PO, BA, BI Dev, DWH Eng, Big Data Eng, Data Scientist, Web Dev,**  
**Team Leads/Managers**
- **Teams: 6 BI, 2 DWH, 1 Big Data, 1 Data Science, 1 Core (Platform Team)**

# Traffic in the past

# Flow

---

- **NGINX logs -> Vector.dev -> Kafka -> PHP -> MySQL -> QlikView**
- **Writing to MySQL limited to 40k rows/s**
- **Reading in QlikView from MySQL – 2 processes max 23k rows/s each**

# Issues

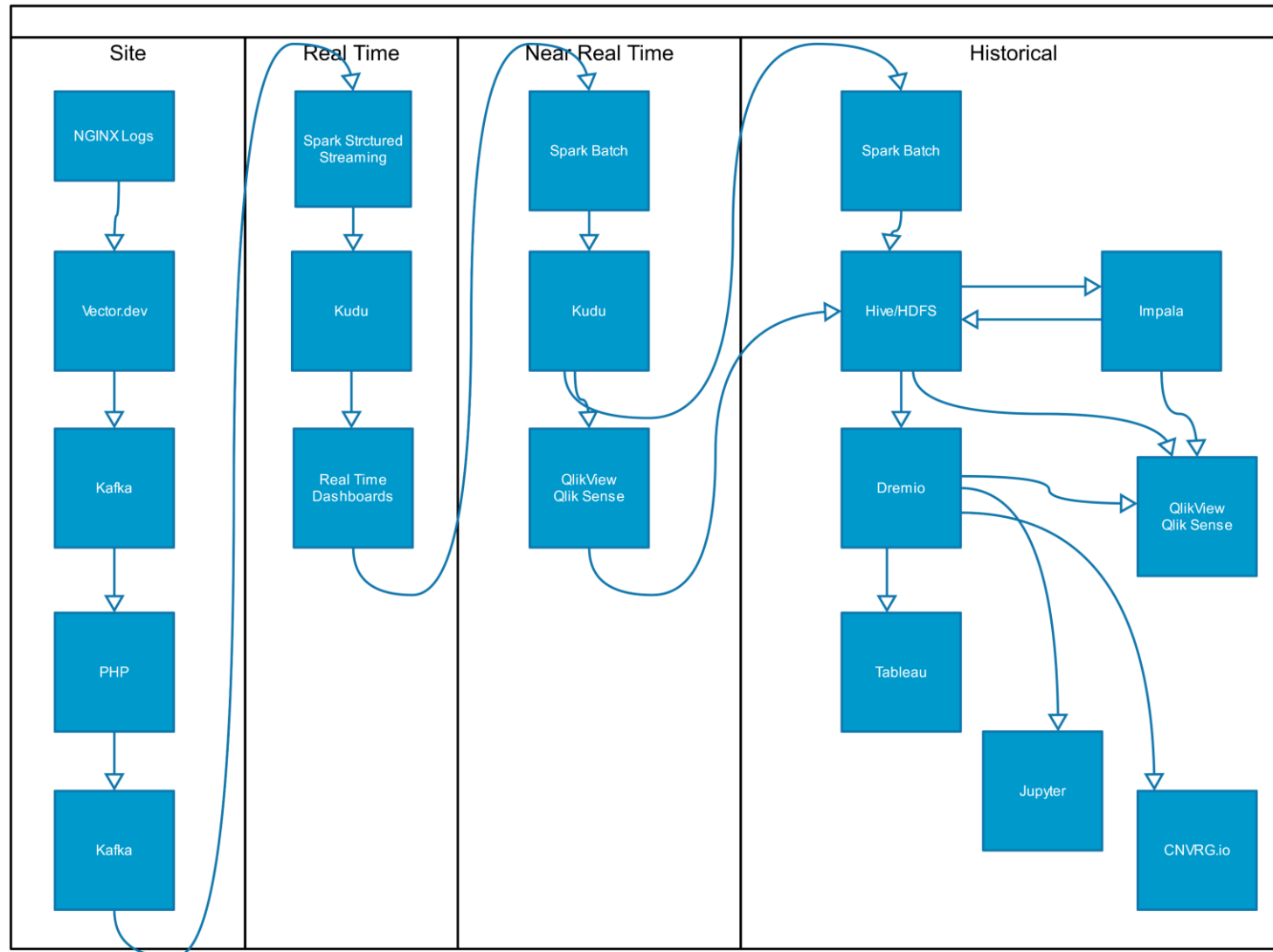
---

- **Peak traffic in BF over 40k reqs/s in BF**
- **No HA for stressed MySQL – risk of failure**
- **Need to add dedicated infrastructure in BF**
- **Processing lag**
- **Time consuming to process data in QlikView (no paralel processing)**
- **Storing processed data in QVD files, ad-hoc analysis was time consuming**

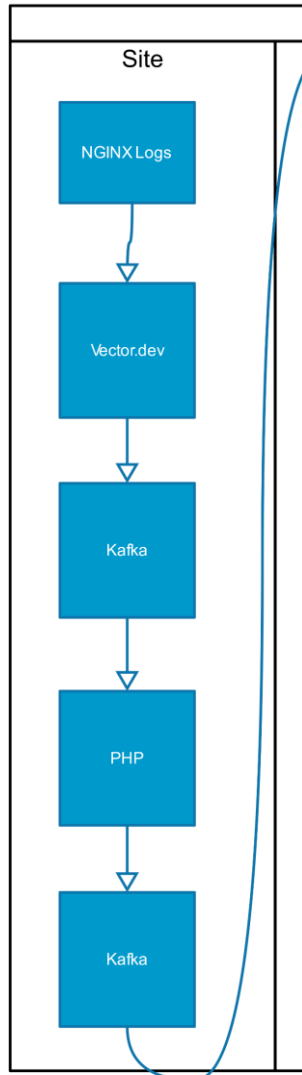
**Traffic now**



# Flow

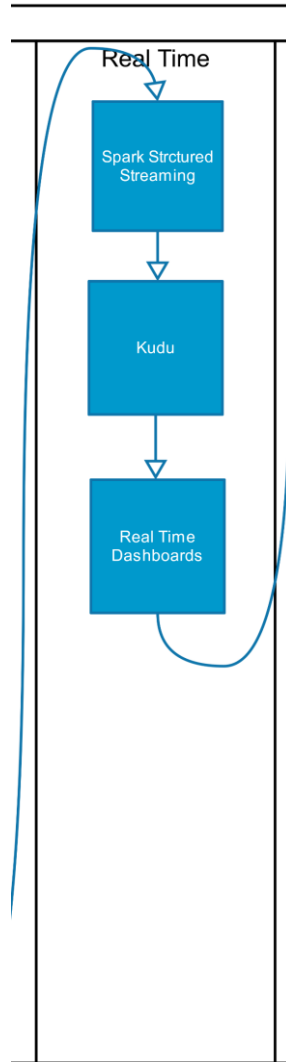


# Site



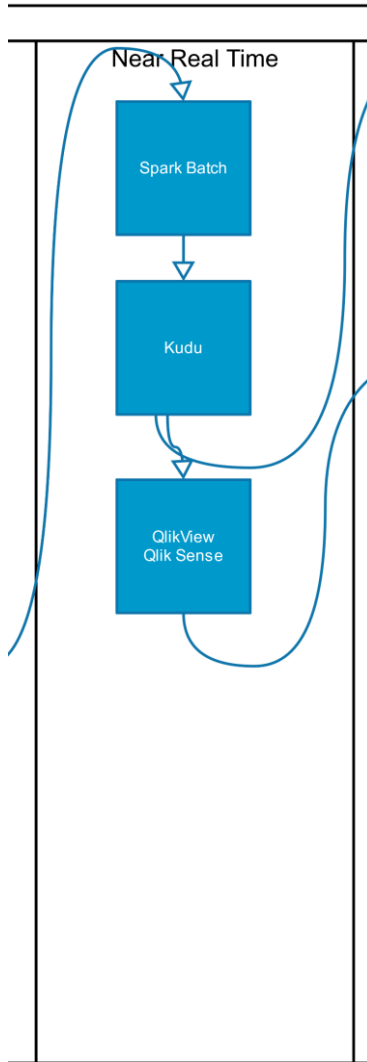
- All NGINX access logs are sent to Kafka via Vector.dev
- PHP cleans and enriches logs then sends to another Kafka topic
- Minimal useful info added, such as category ids

# Real Time



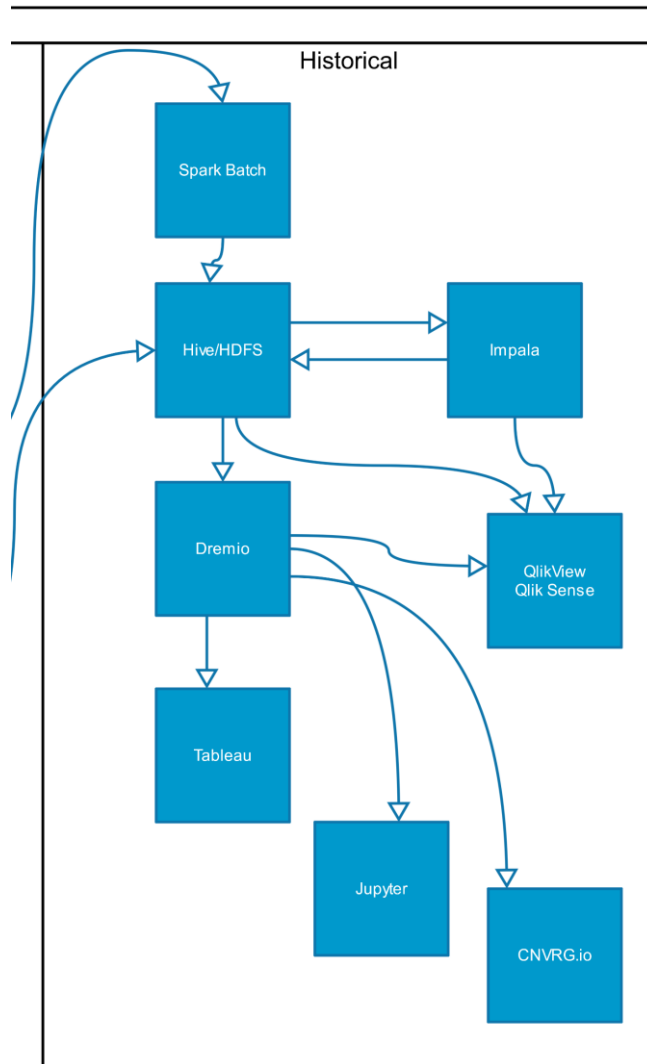
- **Reading from Kafka with Spark Structured Streaming (Scala)**
- **Enriching with minimal info and writing to Kudu**
- **Raw traffic data is available via Impala in real time to other apps**
- **160k events/s**

# Near Real Time



- **Spark (Scala) Batch process reads from Kudu**
- **User matching**
- **Eliminate unwanted events (irrelevant, spiders)**
- **Compute time on page, session duration, other KPIs**
- **Write back to Kudu**
- **Data read via Impala ODBC mostly by Qlik reports**
- **60k events/s**

# Historical Analysis



- **Spark (Scala) Batch process reads from Kudu and writes to Hive/HDFS once a day**
- **Cleans data older the one week from Kudu**
- **Impala used for complex cohort calculations**
- **Dremio used to expose traffic from HDFS for ad-hoc analysis in Tableau and data science**

# CDP Cluster Layout

---

## 3x Coordinators

40 cores - 80 vCores  
256GB Ram  
1TB Storage

## 2x Gateways

40 cores - 80 vCores  
256GB Ram  
1TB Storage

## 2x Utility

40 cores - 80 vCores  
256GB Ram  
1TB Storage

## 2x Processing nodes

40 cores - 80 vCores  
768G B Ram  
1TB Storage

## 11x Data nodes

48 cores - 96 vCores  
768GB Ram  
70TB Storage

## 5x Data nodes

48 cores - 96 vCores  
512GB Ram  
28TB Storage

# Traffic in the future

# Next Steps

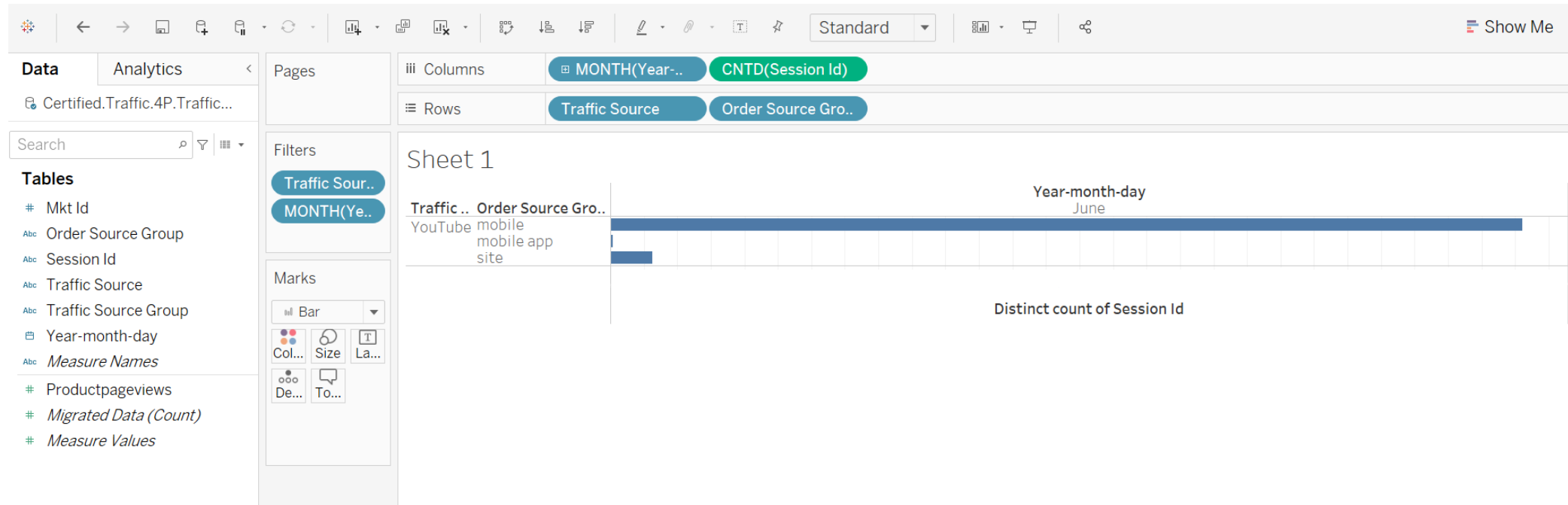
---

- **Keep up with growth (what works for real time today may not work tomorrow)**
- **Analyze on-prem data in cloud**
- **Build a graph model from traffic data**
- **Increase self service for ad-hoc analysis**



# Use Cases


# Traffic From YouTube in June 2019



# Customer Cohorts

Year 2021 YearMonth Jun-2021 YearWeek

## Customized report

 Date	Geolocation Label	Visits
<b>Total</b>		
6/21/2021	1. Bucharest & Ilfov	
	2. ...	0
	3. ...	3
	4. ...	6
	5. ...	1
	6. ...	0
	no info	1

# Daily KPIs

☰ Daily Snapshot
 Analyze Sheet
Narrate Storytelling
🔖
📄 Duplicate
Daily
⏪
⏩
📄 Selections

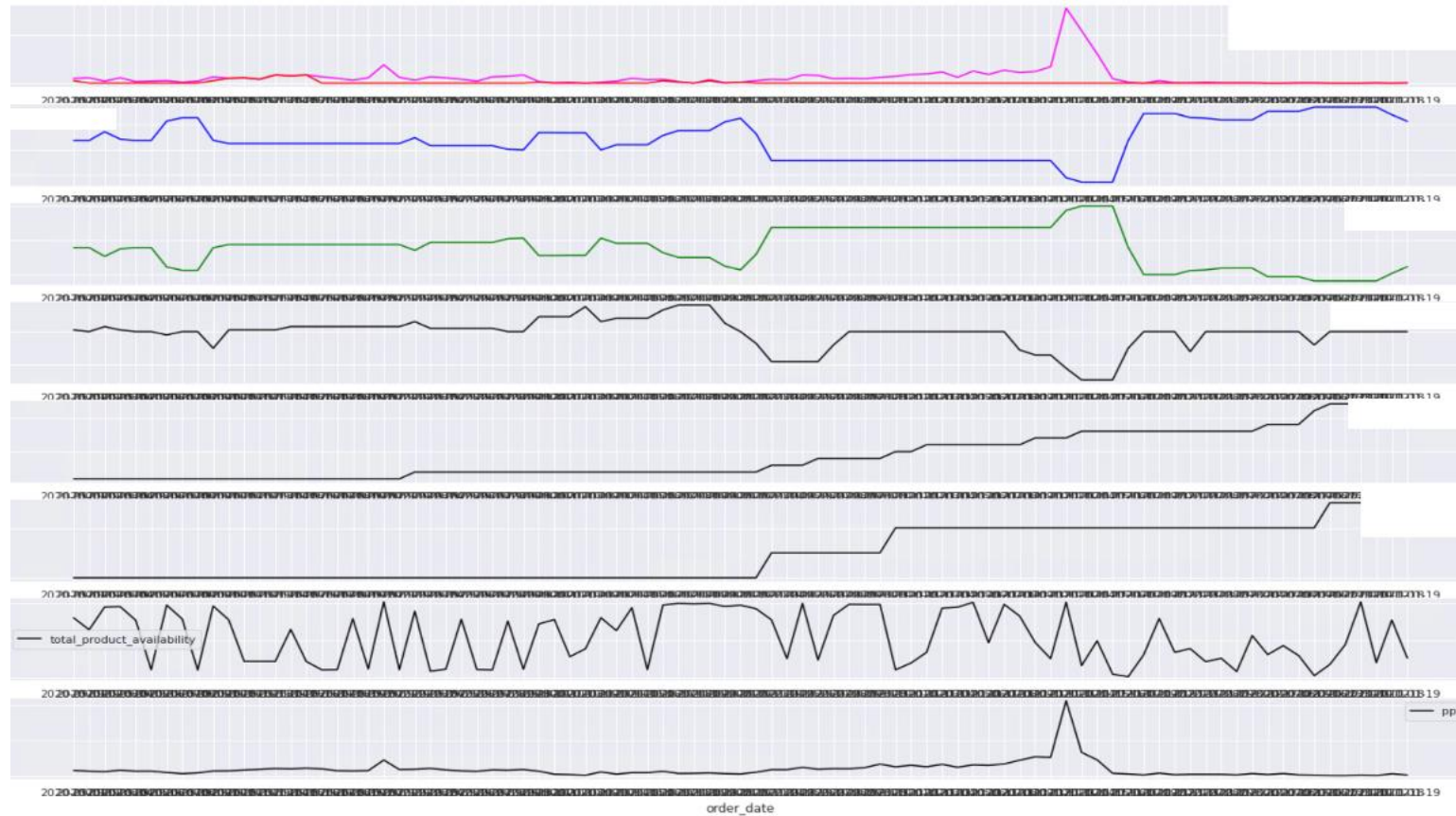
Daily

★★★★★ Send Rating

Last Reload Time is 6/23/2021 12:39:33 PM

Values	Date	Weekday		
	6/9/2021	6/10/2021		
	Wed	Thu		
Visits 4P Target			18	21
Traffic vs Target				
Visits 4P actual			2	1
Visits 1P Actual			6	10
Visits 3P Actual			15	7
PPV 4P actual			10	7
PPV 1P Actual			11	9
PPV 3P Actual			9	8
Visits 4P LY			6	5
Visits 1P LY			5	1
Visits 3P LY			13	2
PPV 4P LY			.5	9
PPV 1P LY			16	5
PPV 3P LY			175	4
Visits 4P YoY				
PPV 4P YoY				
Visits YoY 7d				

# Data Science



# Real Time Dashboard

EMAG		Black Friday 2020 - Main Dashboard HU					Select another dashboard		Last Update 03:04:51
	Visitors	Visits	Pages per Visit	AVG Time On	No of Orders	Conversion Rate	Total Visits	Total Conversion Rate	AVG Visits / Logged Account
emag.hu		38	,4		4	%			
m.emag.hu		0	,1		59	,9%			
app-emag		0	,2	:29	)	30,6%			

# Lessons Learned

# Lessons

---

- **Replace only bottlenecks, not the whole setup at once**
- **Use technologies that you can combine as needed**
- **Leverage existing knowledge for fast results (Impala/SQL)**
- **Implementing new technologies takes time, be patient but don't give up**



**Q&A**

# Feedback

Conference Slack

laurentiu.matei@emag.ro

<https://twitter.com/lmatei80>

<https://www.linkedin.com/in/laurentiumatei/>