

Bitstamp

Towards Modern Technology Stack

User Activity History

Application Platform Engineering

Jure Žvelc, January 2024



User Activity History

- One of the oldest features
- Important customer actions
 - Login
 - 2FA reset
 - Deposits
 - Withdrawals
 - Etc.
- Available to every customer
 - Web view
 - Export
- Integrated in various business processes

Bitstamp S.A. Basic Pro

Dashboard Markets Deposit Withdrawal Earn Blog Learn

Buy / Sell

Find asset

Main account



Transaction history

Activity history

Activity history

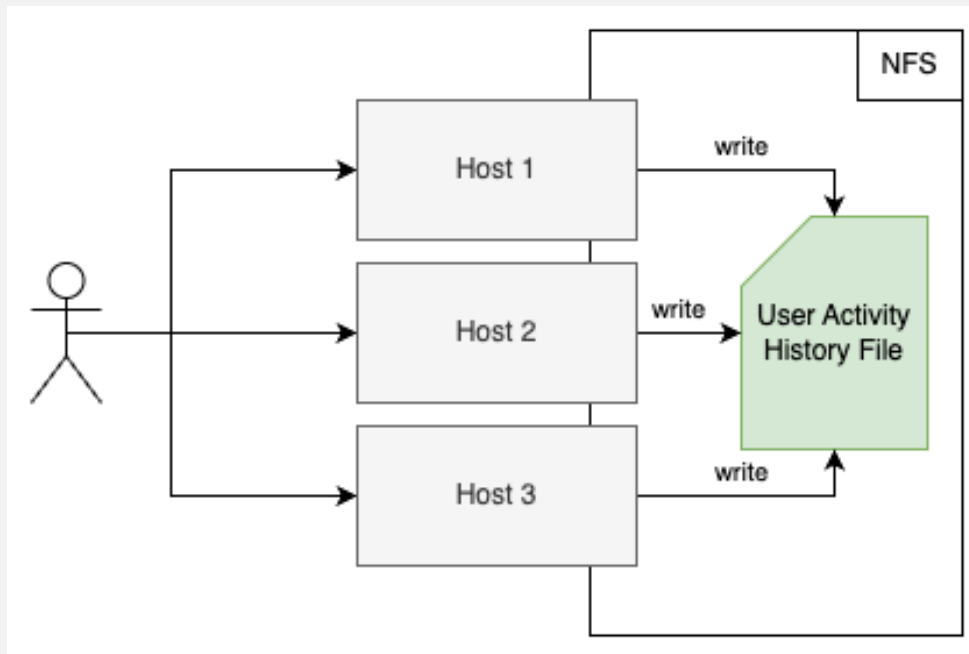
Date and time	IP address	Action
Jan. 25, 2024 09:19 AM	127.0.0.1	BTC withdrawal request: email was sent to user
Jan. 25, 2024 09:19 AM	127.0.0.1	Opened BTC withdrawal request for 2.00000000 BTC (Bitcoin) to 34zVs9QIRZhkQD4xyZXfr6tV3nYdSD7Kh from Web Browser
Jan. 25, 2024 09:18 AM	127.0.0.1	Opened new Sub Account with id: 198057 (Test 2) for trading on Spot Market
Jan. 25, 2024 09:17 AM	127.0.0.1	Logged in using two-factor authentication
Jan. 25, 2024 09:17 AM	127.0.0.1	User 5ad87b06 logged in with 2FA from IP address '127.0.0.1'
Jan. 25, 2024 09:02 AM	127.0.0.1	Opened new Sub Account with id: 57927968 (Test) for trading on Spot Market
Jan. 25, 2024 09:00 AM	127.0.0.1	BTC withdrawal request: email was sent to user
Jan. 25, 2024 09:00 AM	127.0.0.1	Opened BTC withdrawal request for 1.00000000 BTC (Bitcoin) to 34zVs9QIRZhkQD4xyZXfr6tV3nYdSD7Kh from Web Browser
Jan. 25, 2024 08:59 AM	127.0.0.1	Logged in using two-factor authentication

/activity/v2/

Initial Implementation

B

- User history activity is appended to a per user log file as CSV
- Log files are stored on NFS



Initial Implementation



Problems

- Parallel writes to files
- Reading directly from NFS
- NFS is a single point of failure
- NFS upgrade requires a scheduled maintenance (downtime)
- Any NFS write latency spikes blocking hot path
- Only full user activity exports available
- Large exports handled through support tickets
- Extreme caution when working on NFS server

Introduction of Database Outbox

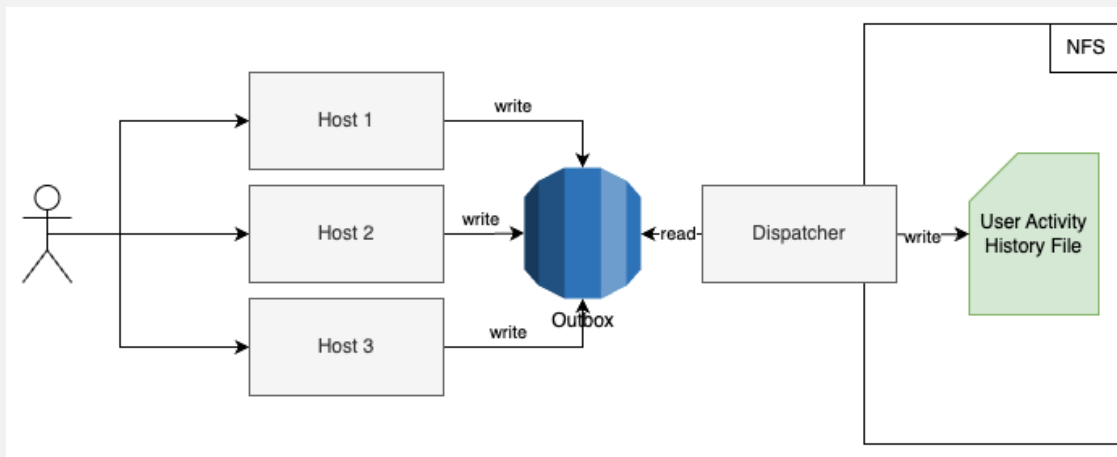


Goals

- Unblock hot path by removing NFS write dependency
- Centralize user history writes in a single service
- Datalake integration via S3
- Support for current business use cases

Introduction of Database Outbox

B



Solution

- All user activity history writes redirected to a database table
- A separate process (Dispatcher) reading from the outbox table and writing to NFS
- Migrate from file per user to per day user files
- Switch from CSV to JSON
- Upload per day files to S3 for datalake integration

Introduction of Database Outbox



Results

- Resolved latency spikes in hot path
- Increased average write latencies
- Shared database across services
- Very high load on database requiring large instances
- Huge database costs

Race with Time

Approaching the NFS max capacity!

B



In Search for a New Solution



Requirements

- Cost effective
- Performant
 - Fast propagation (agents need up to date data)
 - Has to handle 4000 writes per second with avg. line size of 0.5KB (168GB of data per day)
- Scalable and highly available
- Use AWS managed infrastructure

In Search for a New Solution

B

Storage



- Replace NFS with EFS
 - Still requires volume mounts
 - Expensive
 - Slower than NFS



- S3
 - Cheap storage and read operations
 - Scalable
 - Does not support append



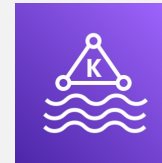
- NoSQL solutions too expensive to permanently store all the data

In Search for a New Solution

Transport

Kafka

- Available as an AWS managed solution (MSK)
- Ideal for data streaming



Vector

- Tool built in Rust to get data from point A to point B
- Lots of available sources and sinks (File, Kafka, S3,...)
- Easy to configure!



We already used Kafka for log shipping

- Vector ships logs from files to Kafka
- Vector consumes logs from Kafka and stores them on S3

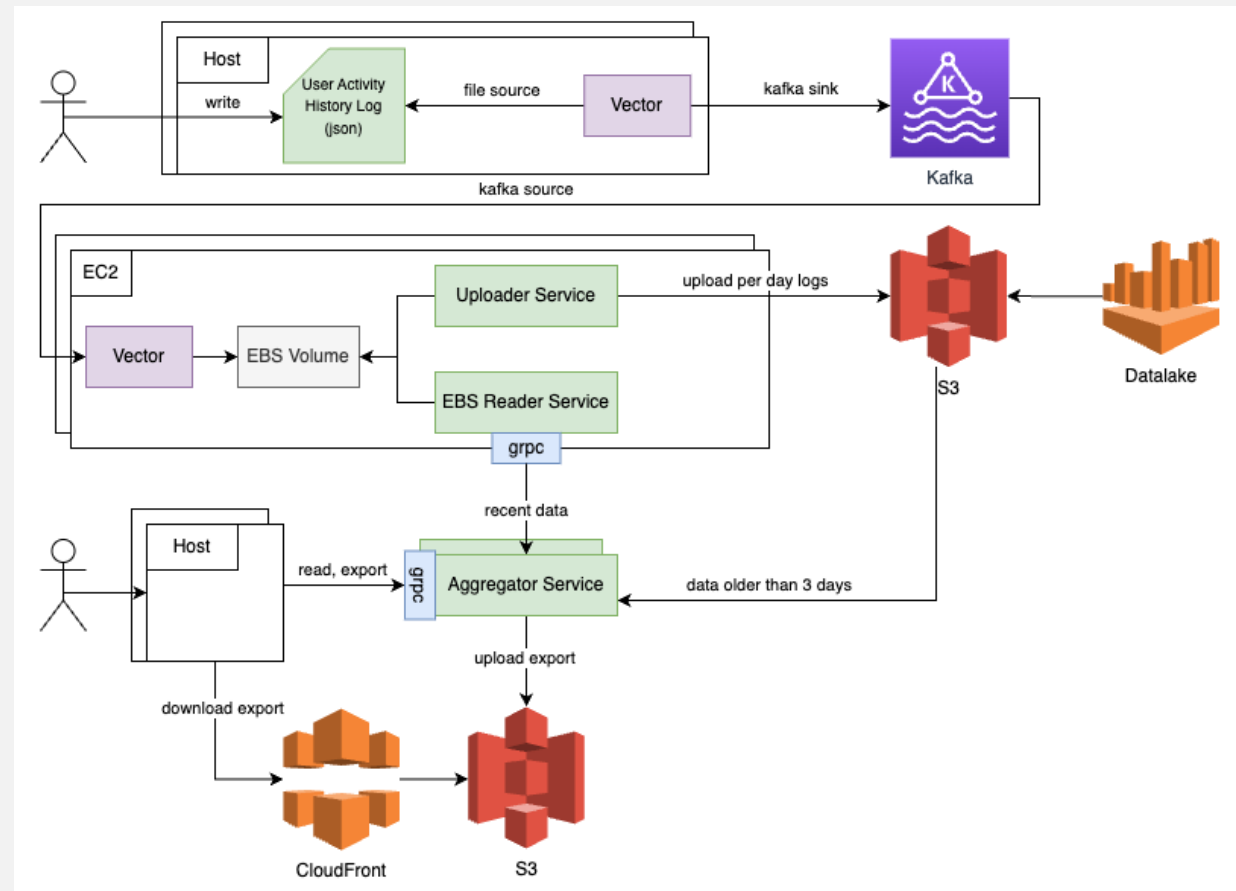
Idea



- Reuse Kafka & Vector to ship user activity history logs to some service which handles ordering and deduplication
- Use S3 for permanent storage
 - Data needs to be stable before uploading to S3
 - Per day user activity history files
 - Data ordered and deduplicated
- Use something for fast lookup of data up to 3 days old

Solution

- Each Vector on EC2 instance forms a separate Kafka consumer group
- **Vector (EC2)**
 - Parses JSON logs consumed from Kafka and writes them in a hierarchical directory structure on EBS volume
- **Uploader service (EC2)**
 - Sorts and deduplicates data from EBS vol.
 - Uploads daily logs
 - Master election through S3
 - Deletes old logs
- **EBS reader service (EC2)**
 - Sorts and deduplicates data from EBS vol.
 - Serves 3 days of data through gRPC
- **Aggregator service**
 - Serve data from EBS reader and S3 with cursor based pagination through gRPC
 - Processes and converts JSON to CSV
 - Background workers for exports with S3 based locking



Highlights



- Solution is highly available and scalable
- More than 6x cost reduction! 🍺
 - Kafka easily handles huge amounts of logs
 - S3 storage and GET/SELECT operations are very cheap
- Services written in Golang with a single external dependency (S3)
- Improvements
 - Very fast user activity history logging (decreased latencies)
 - Customizable date for exports & support for large exports
 - Fast export downloads through CloudFront
 - Compatible with existing business processes

Challenges



- Migration from old solution
- Sort and deduplication algorithms
- S3 master election and queueing solution
- S3 user key catalogs
 - Contains keys for all files with sizes
 - Reduces number of S3 LIST operations
- S3 prefetching

Challenges



- Cursor based navigation and reading in chunks
- Security
 - STS authentication
 - Audit logs for requestes
 - CloudFront presigned urls valid for 1 minute
 - Secret rotations
- Tracing and metrics



Q & A