# Software contribution from MIPT: Development of Event Metadata System and Monitoring & High-Availability Service

Peter Klimai

# Current Projects Summary
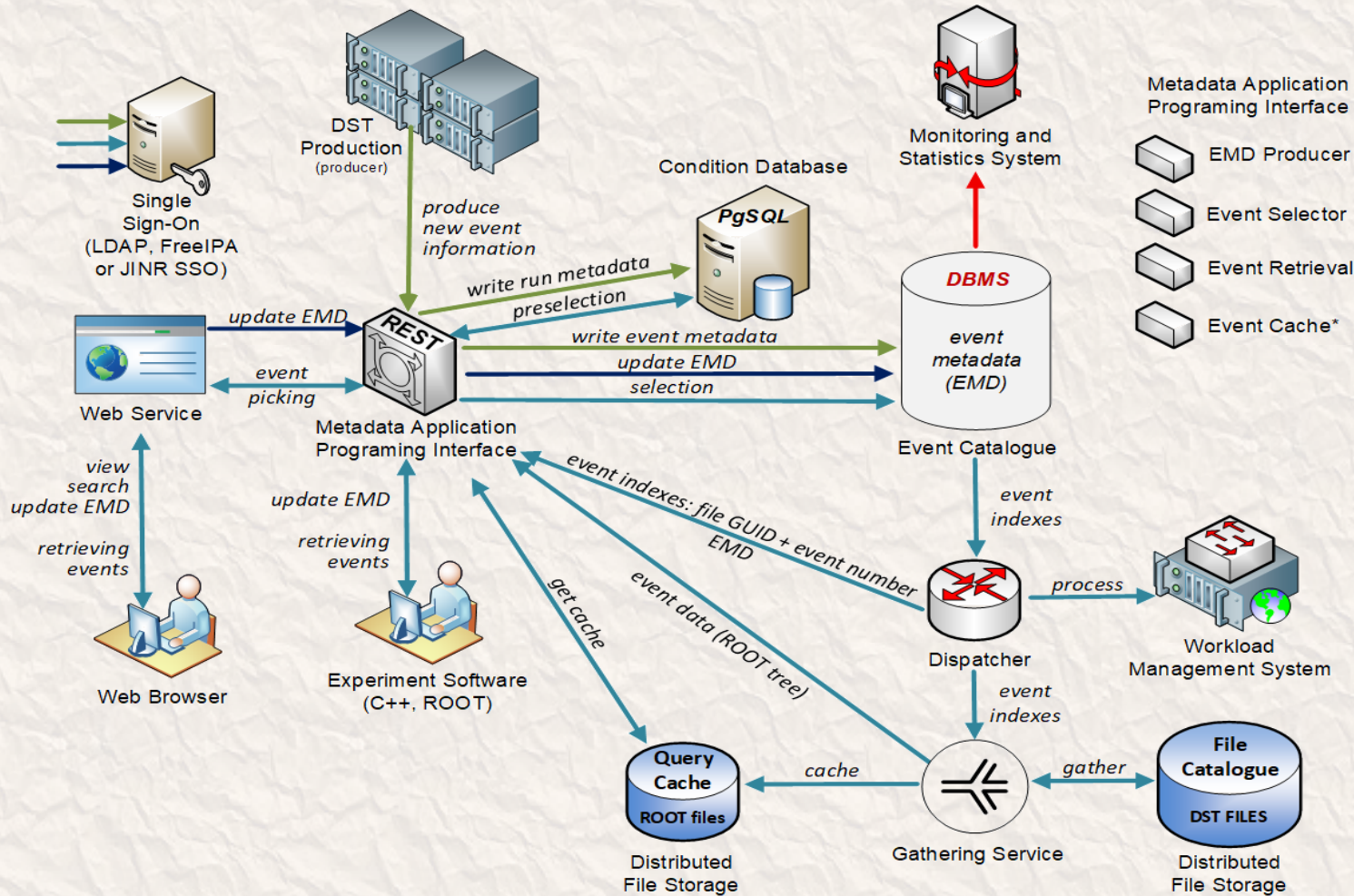
| Project | URL | Notes |
|---|---|---|
| **Event Metadata System** | https://git.jinr.ru/nica_db/emd | First version deployed; updates discussed in this talk |
| Deployment scripts for EMS | | This talk |
| High Availability design for EMS | | This talk |
| Statistics collection and visualization for EMS | | WIP |
| **Next-generation event display** | https://github.com/SciProgCentre/visionforge | See talk by A. Nozik |
| **Monitoring service** | https://mon-service.jinr.ru https://git.jinr.ru/nica/bmnroot/-/tree/dev/services/is_monitor | In production; updates planned. This talk |
| Slow control system viewer | https://bmn-tango.jinr.ru | Needs update to match new SCS and its database |

# Event Metadata System – Update

# EMS Architecture and Features



- Event Metadata System
  - Event Catalogue is based on PostgreSQL
  - Integrates with BM@N Condition database
  - REST API and Web UI developed based on Kotlin multiplatform
  - Configurable to support different metadata
  - ROOT macro to write BM@N events in the catalogue
  - Role-based access control implemented
  - Monitoring

For more details:
E. Alexandrov, I. Alexandrov, A. Chebotov, A. Degtyarev, I. Filozova, K. Gertsenberger, P. Klimai and A. Yakovlev, "Implementation of the Event Metadata System for physics analysis in the NICA experiments", J. Phys.: Conf. Ser. 2438, 012046 (2023).

# EMS Updates

- Recent EMS Updates (discussed next):
  - New unified REST API scheme
  - Simplified to support only one metadata table per EMS instance
  - OpenAPI documentation (aka Swagger) now available
  - Database performance improvement studies (indexes)
  - High Availability solution
  - Deployment scripts (Ansible based)

# New scheme for REST API

- The new scheme is unified for different BM@N Information Systems

GET
POST
DELETE

**https://bmn-event.jinr.ru/event_api/v1/event?**

run_number=3950:4000&beam_particle=Ar&target_particle=Al
energy=3.16:3.18&target_particle=SRC%20Lead

HOSTNAME / SERVICE / VERSION / ENTITY?parameter_set

*HOSTNAME=https://bmn-[SYSNAME].jinr.ru*

*SERVICE=[SYSNAME]_api*

*VERSION=v1 (v2...)*

*ENTITY=tablename without last '_' (if present)*

*parameters are separated by '&'*
*ranges: min:max → >=min AND <=max*
*min: → >=min      :max → <=max*

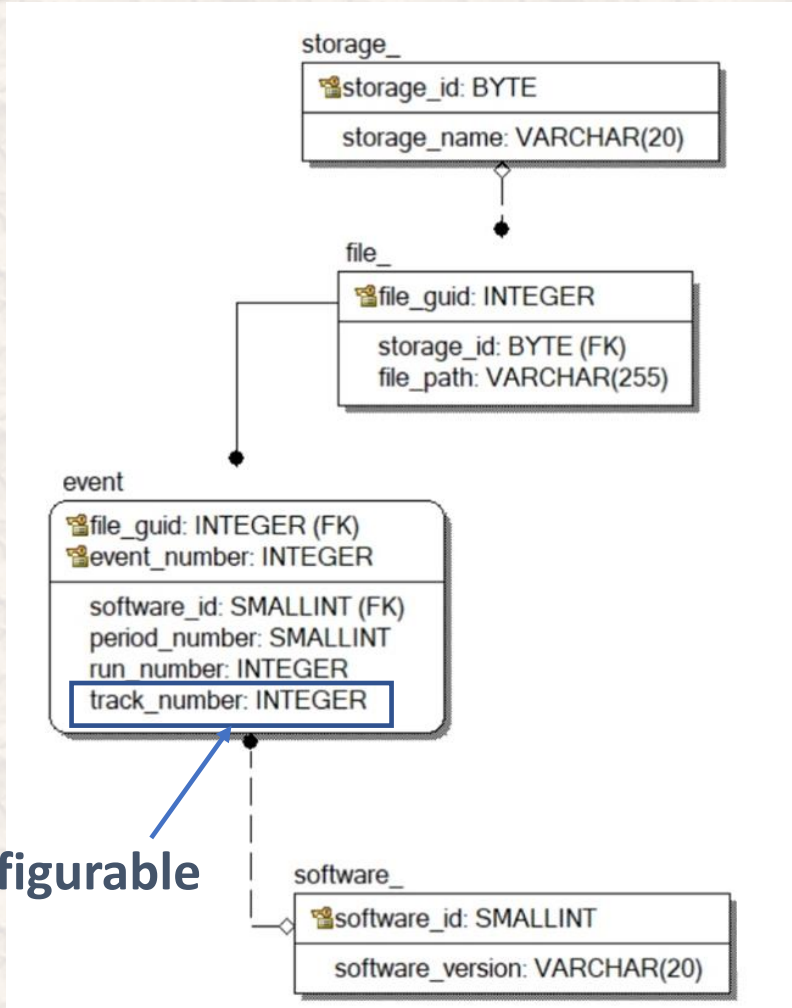For the Unified Condition Database (UniConDa), SYSNAME = uniconda
For the Event Metadata System (EMS), SYSNAME = event
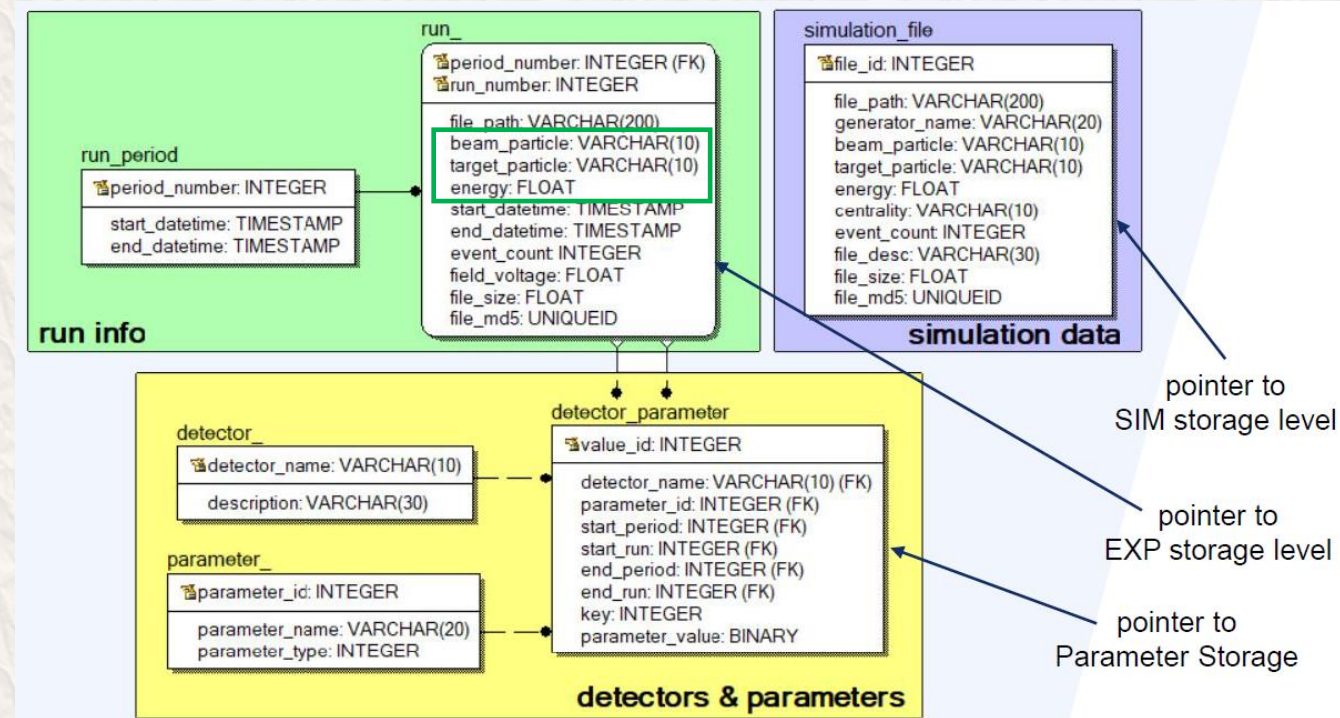
# Web UI Main Page



**Special script is collecting this statistics on the backend (WIP)**

# Main search page



Selection based on standard parameters

Preselection based on Condition DB

Selection based on configured parameters

Limit and offset

# OpenAPI pages for EMS

127.0.0.1:8080/openapi#api-Default-eventGet

**API SUMMARY**

**API METHODS - DEFAULT**

eventGet
eventPost
softwareGet
softwarePost
storageGet
storagePost

## eventGet

Returns event metadata

**GET**

/event

**Usage and SDK Samples**

Curl | Java | Android | Obj-C | JavaScript | C# | PHP | Perl | Python

```
curl -X GET\
 -H "Authorization: Basic [[basicHash]]"\
 -H "Accept: application/json"\
 "http://127.0.0.1:8080/event_api/v1/event?limit=&offset=&software_version=&period_number=&run_number="
```

## Parameters

Query parameters

| Name | Description |
|------|-------------|
| limit | Integer <br> *Limit on the number of events* |
| offset | Integer <br> *Offset for obtained events (typically used with limit)* |
| software_version | String <br> *Software version for events* |
| period_number | String <br> *Period number (possibly range) for requested events* |
| run_number | String <br> *Run number (possibly range) for requested events* |

**Responses**

Status: 200 - A JSON array of event objects
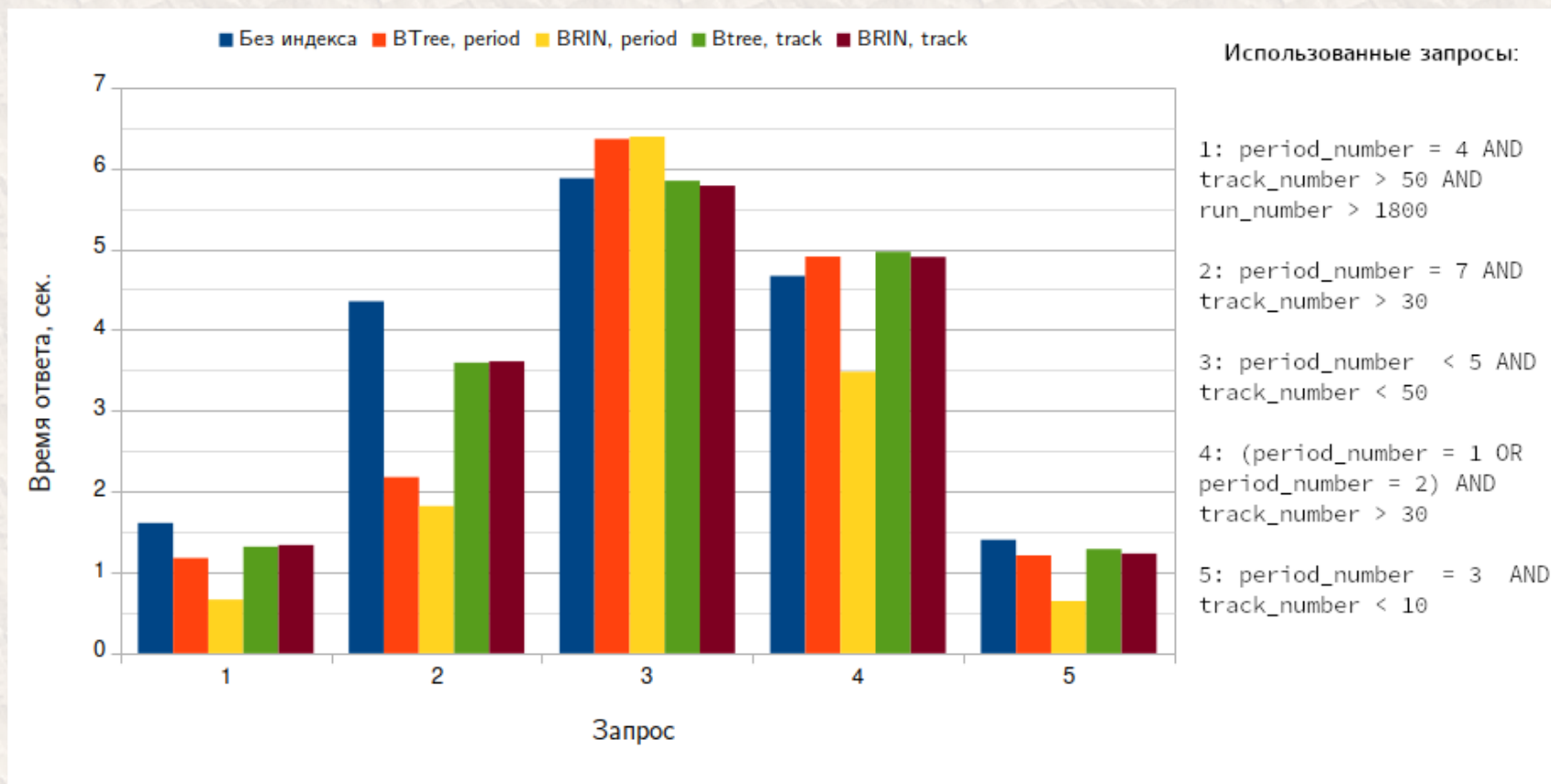
Schema

```
▼ [
  ▼ {
      Single event metadata
      Required: parameters,period_number,reference,run_number,software_version
      reference:        ▼ {
                          Reference to event's storage/file/event_number
                          Required: event_number,file_path,storage_name
                          storage_name: string
                                        example: data1
                          file_path:    string
                                        example: /var/tmp/file1.root
                          event_number: integer
                                        example: 100000
                        }
      software_version: string
                        example: 20.1
      period_number:    integer
                        example: 8
      run_number:       integer
                        example: 5000
      parameters:       ▼ {
                          Map of optional parameters key/values, according to EMS config
                          track_number: integer
                                        example: 30
                        }
    }
]
```
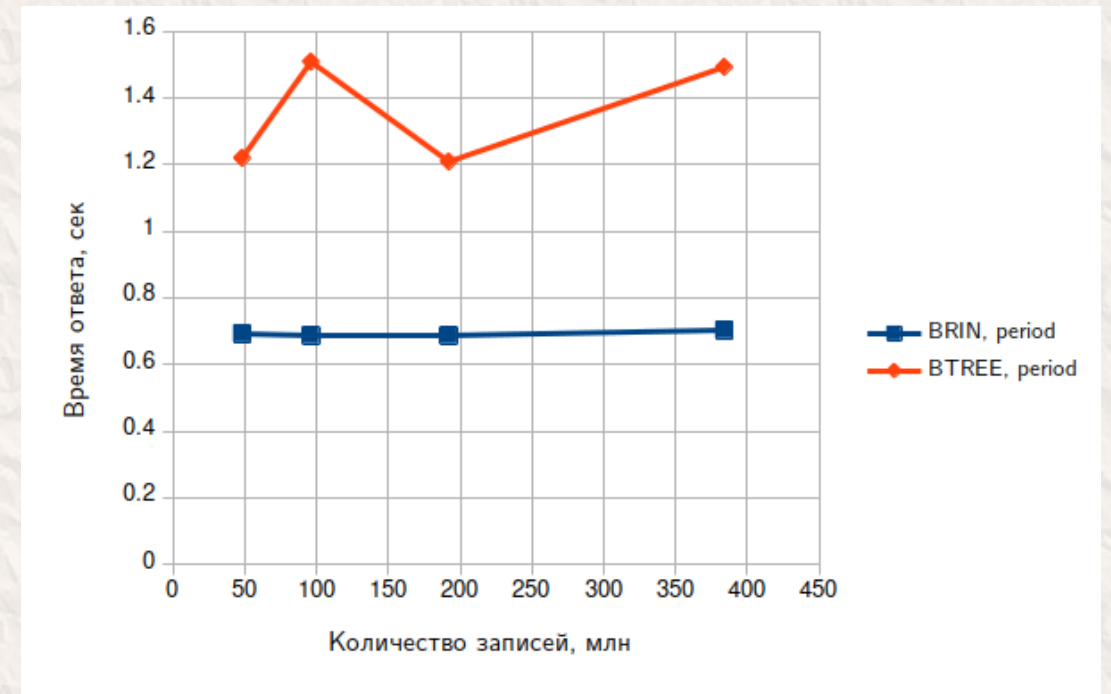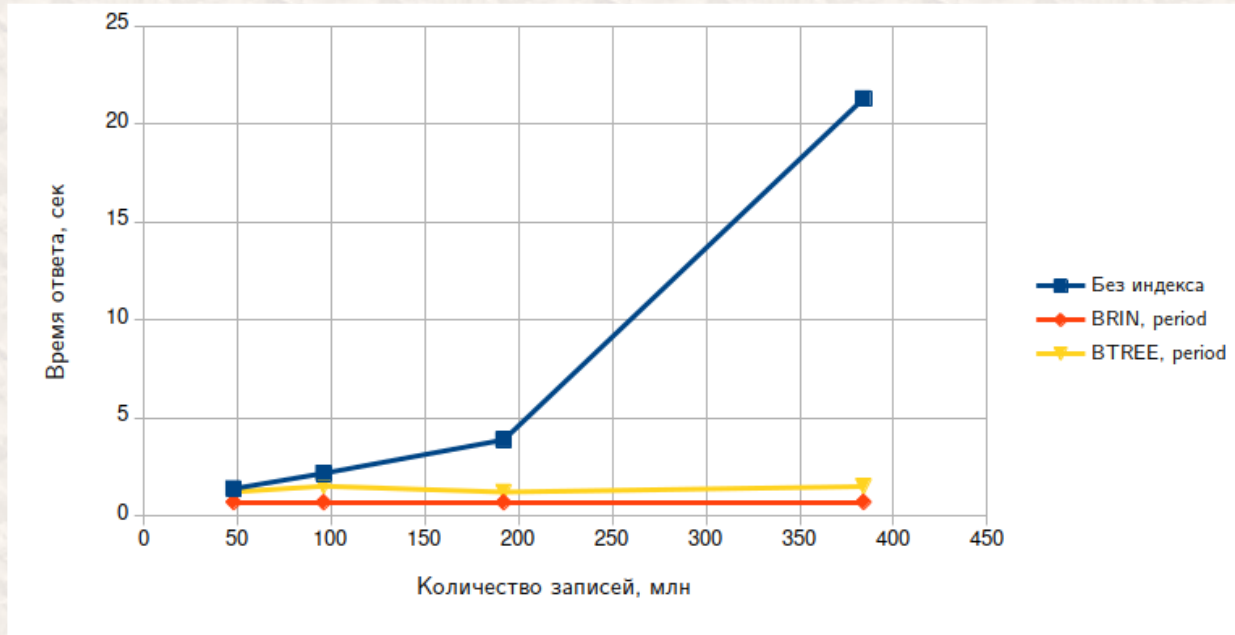
10

# Index Selection (Type and Columns)

- Measurements with test database instance are shown (50M events)

# Response time

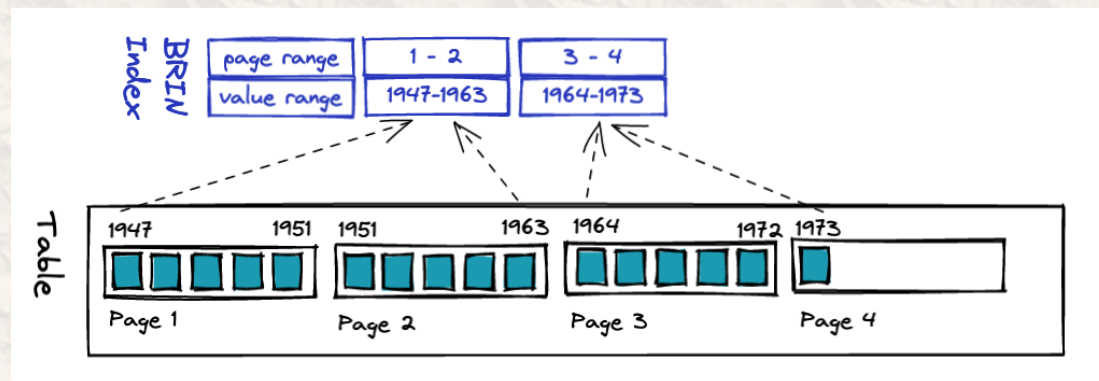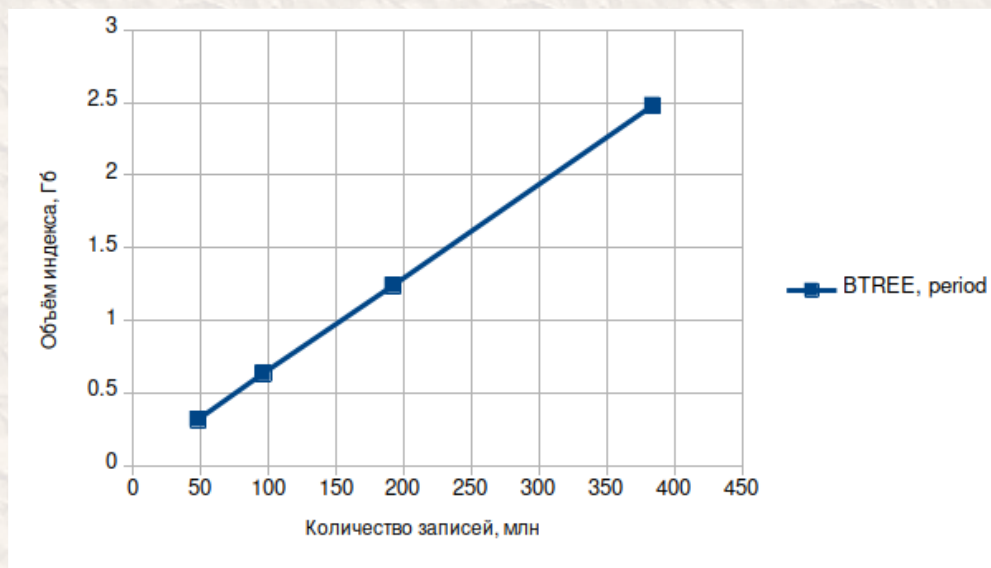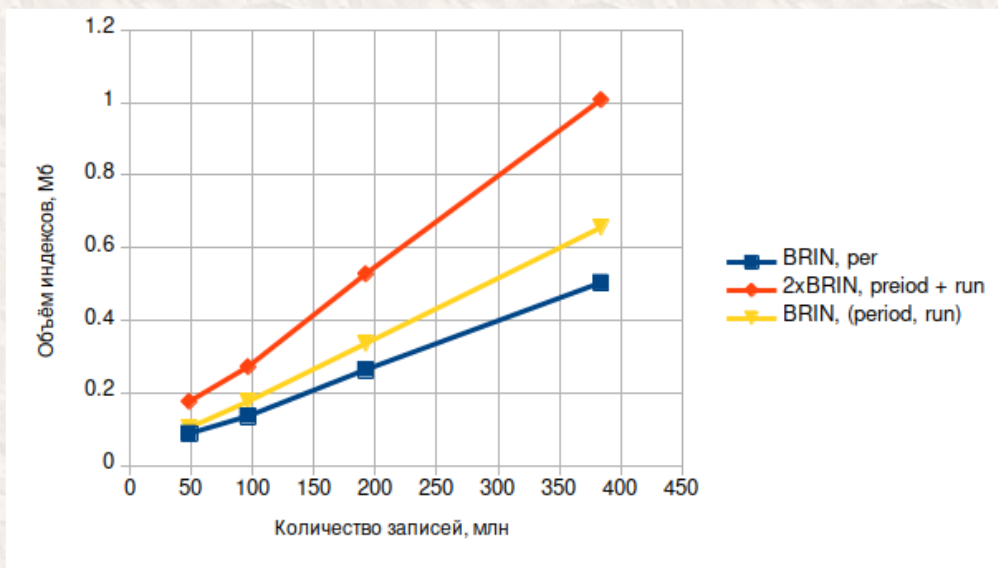- Adding more periods to test database

# Index size on disk

- BRIN vs. BTREE
  - Overall, BRIN (Block Range Index) works better for indexing columns having some natural correlation with their physical location within the table

# High Availability – Task
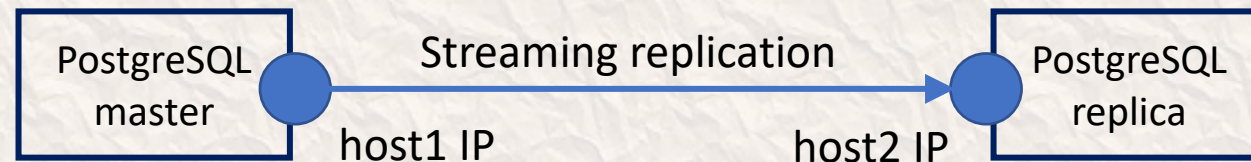
- Need for HA
  - EMS as well as other IS are essential for timely obtaining physical results of the experiment
  - From client point of view, connection must be initiated to single IP / domain name
    - We do not want to ask client to keep several addresses like primary/secondary ones
  - Considering 2 to 1, active/passive redundancy
  - Need to avoid split brain and no brain scenarios

# High Availability – Solution

- Base for HA solution
  - PostgreSQL supports streaming replication out of the box (one master to one/many replica servers)
    - https://www.postgresql.org/docs/current/warm-standby.html#STREAMING-REPLICATION
  - Completely synchronous replication is also available (at a performance price)
    - https://www.postgresql.org/docs/current/warm-standby.html#SYNCHRONOUS-REPLICATION

PostgreSQL
master

Streaming replication
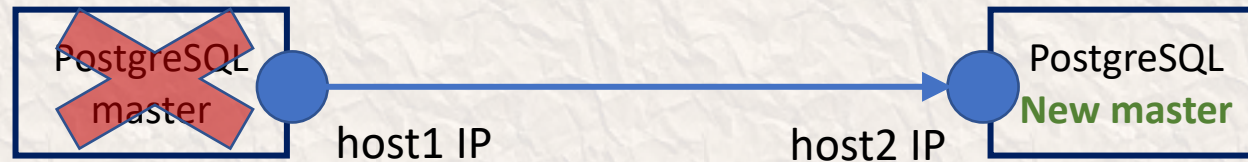
host1 IP

host2 IP

PostgreSQL
replica

# Switchover to new master

- Switchover
  - One command on replica - `pg_ctl promote`
  - Old master must be turned off to avoid split brain
  - Monitoring system can perform switchover (WIP), or it can be done manually
  - It works, but the big question is – where does a client connect?



PostgreSQL master

host1 IP

host2 IP

PostgreSQL
New master

# Solution based on VRRP (single L2 domain)

- Keepalived provides virtual IP address for client connection
- This works when both servers are in the same L2 (broadcast) domain

# Avoiding single point of failure

- VRRP-based solution can be considered final if:
  - L2-segment is built with redundancy (both for links and switches)
  - VIP's network is announced from at least two routers
  - Not possible to implement without access to network infrastructure



master     replica

**Net/mask metric=10**

**Net/mask metric=20**

# Solution based on DNS

- Solution details (WIP)
  - PostgreSQL replication unmodified
  - Client connection to host/domain name (needs DNS settings)
  - Monitoring system performs switchover
    - Change DNS record
    - Perform `pg_ctl promote`
  - Switchover time determined by DNS TTL settings

# EMS Automated Deployment

- Why automated deployment?
  - Manual deployment of a distributed system is slow and error-prone
  - Automation increases speed and predictability
  - Avoids issue of "forgotten step" in documentation
  - EMS instance may be deployed by other NICA experiments

- Main components of solution
  - Ansible
  - Docker

- Inputs
  - EMS configuration as YAML template
  - Deployment configuration as Ansible variables in hosts file
    - To be replaced by unified JSON config (WIP)

```
(env) [lab@alma1 ems-deploy]$ cat deploy-pgsql.pb.yaml
---
- name: Deploy PostgreSQL on Event Catalogue hosts
  hosts: event_catalogue
  become: yes

  tasks:

  - name: Install packages
    dnf: "name={{ item }} state=present"
    with_items:
      - postgresql
      - postgresql-server

  - name: Install Python packages
    pip: "name={{ item }}  state=present"
    with_items:
      - psycopg2-binary
…
```

```
…

  - name: Check if PostgreSQL is initialized
    ansible.builtin.stat:
      path: "/var/lib/pgsql/data/pg_hba.conf"
    register: postgres_data

  - name: Initialize PostgreSQL
    command: "postgresql-setup initdb"
    when: not postgres_data.stat.exists

  - name: Start and enable services
    service: "name={{ item }} state=started enabled=yes"
    with_items:
      - postgresql

…
```

# Deployment example (abbreviated)

```
[lab@alma1 ems-deploy]$ source env/bin/activate
(env) [lab@alma1 ems-deploy]$ ansible-playbook deploy-pgsql.pb.yaml

PLAY [Deploy PostgreSQL on Event Catalogue hosts] ****************************************

TASK [Gathering Facts] ******************************************************************
ok: [ems2]
ok: [ems1]

TASK [Install packages] *****************************************************************
ok: [ems1] => (item=postgresql)
ok: [ems2] => (item=postgresql)
ok: [ems1] => (item=postgresql-server)
ok: [ems2] => (item=postgresql-server)

…

TASK [Apply SQL schema file] ************************************************************
changed: [ems1]

PLAY RECAP ******************************************************************************
ems1                       : ok=13    changed=1    unreachable=0    failed=0    skipped=2    rescued=0    ignored=0
ems2                       : ok=16    changed=4    unreachable=0    failed=0    skipped=2    rescued=0    ignored=0

(env) [lab@alma1 ems-deploy]$ ansible-playbook deploy-vrrp.pb.yaml
(env) [lab@alma1 ems-deploy]$ ansible-playbook deploy-web-api-docker.pb.yaml
```
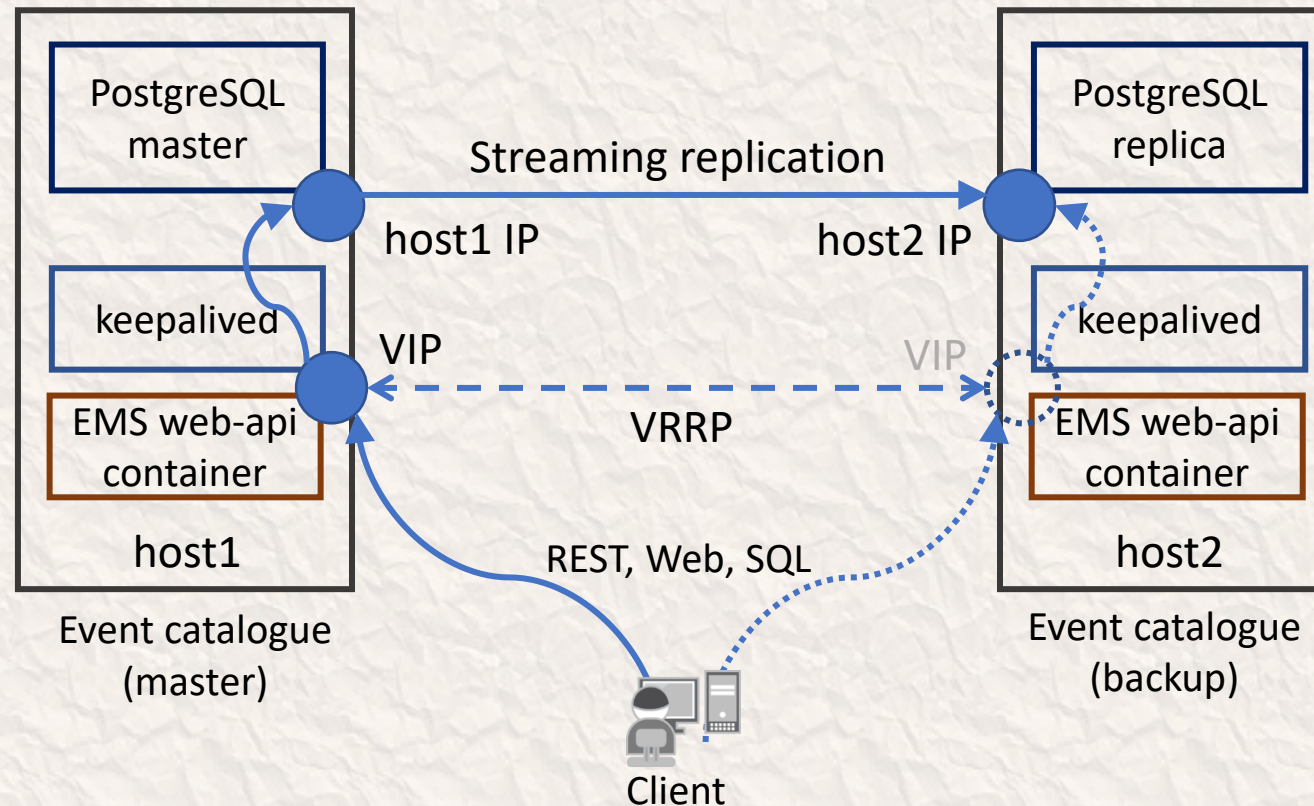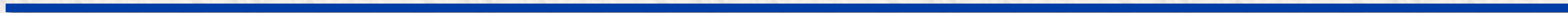
# The Result

- After running the three playbooks:
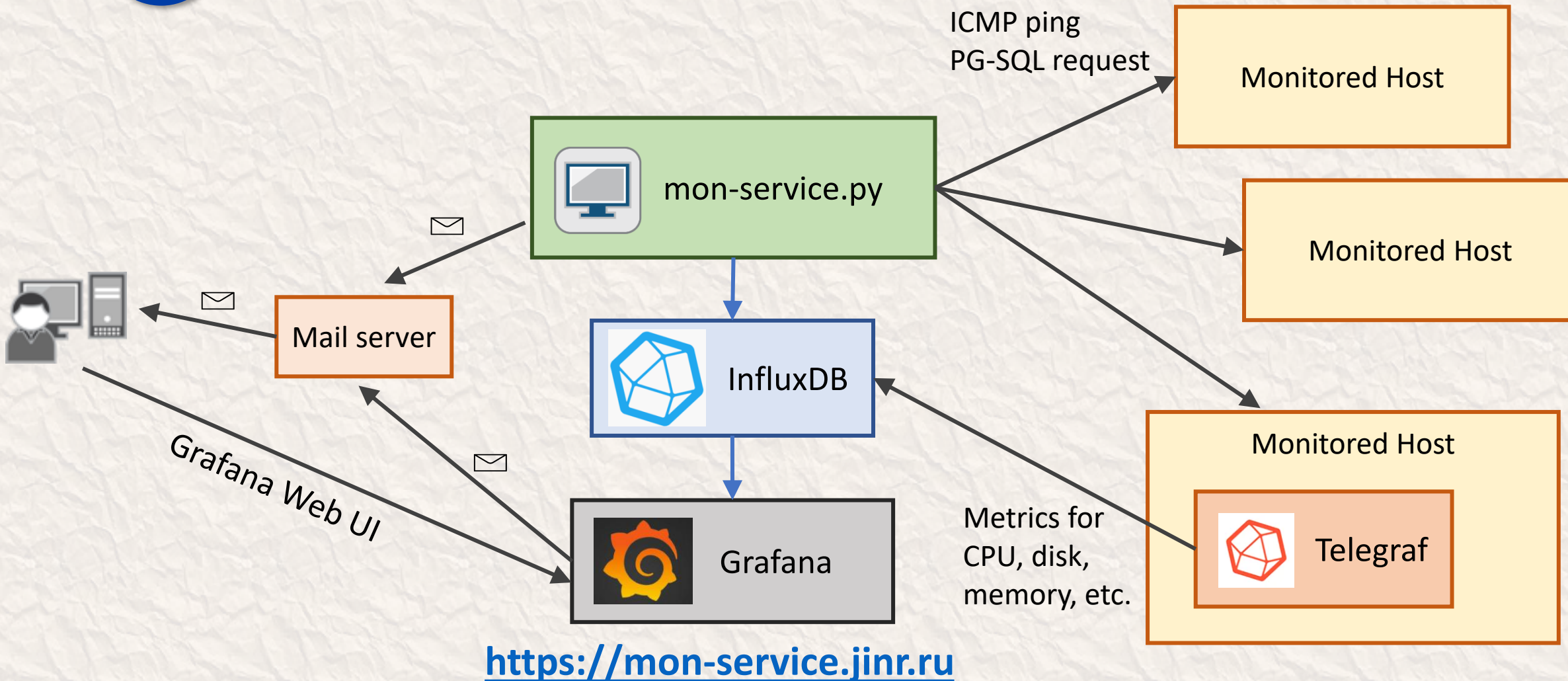
# Monitoring Service Overview

# Monitoring Service - Task

- Monitoring Service Features
  - Ping and PG-SQL request to check database server status
  - Configurable via JSON file
  - Email notifications
  - Response time stored in InfluxDB
  - Use Grafana for visualization and additional alerting
  - Monitor server parameters such as Disk, CPU, Memory, etc.
- Planned new features:
  - Web-services monitoring
  - API endpoint monitoring
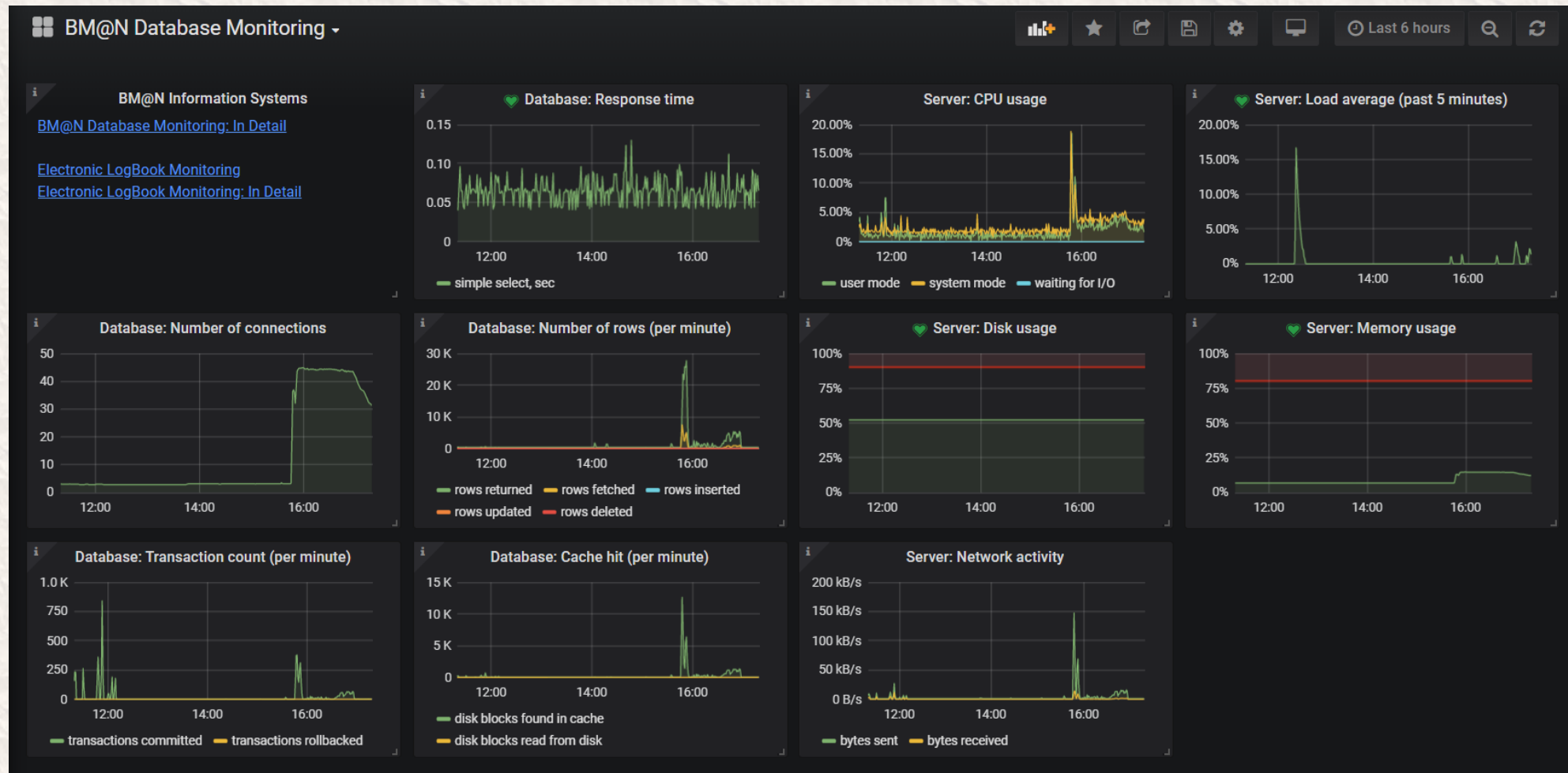  - HA switchover functionality

# Monitoring Service - Components



ICMP ping
PG-SQL request

Monitored Host

mon-service.py

Monitored Host

Mail server

InfluxDB

Monitored Host

Grafana Web UI

Grafana

Metrics for CPU, disk, memory, etc.

Telegraf

**https://mon-service.jinr.ru**

# Monitoring Service View Example

# Thank You!