# Supporting CREDO project with scalable data acquisition and processing infrastructure

Maciej Pawlik, Krzysztof Oziomek, Marek Magryś, Patryk Lasoń, Piotr Homola

## Presentation plan

- 1. Introduction to CREDO, technical details
- 2. Contribution / challenges
- 3. CREDO ecosystem components
- 4. Backend infrastructure
  - a. Software
  - b. Hardware
- 5. Conclusions and references

# **CREDO** project

Cosmic Ray Extremely Distributed Observatory

aims to:

- Expand our knowledge about the universe, understand dark matter
- Detect super-preshower phenomena particle showers
  - Use mobile phones as a network of detectors covering a large area
- multiple other applications, including:
  - detecting changes in earth's magnetic field
  - predicting earthquakes
  - influence on human physiology
  - education, community involvement
- Depends on community involvement, community sourced data



### Some statistics

- 700K+ visible detections (2M+ overall)
- 1M+ device pings (sums up to 50 years of looking for particles)
- 3K+ users with at least 1 detection
- 5K+ devices
- 1K+ user teams
- 10s of GBs used for storage of data, metrics and backups

# Contribution / challenges

- Backend infrastructure
- Provide computing and storage resources
- Manage ingestion of data
- Provide data accessibility
- Integrate into existing ecosystem
- Encourage community involvement

this implies:

- Developing software for gathering and storage information about detections
- User management
- Providing means for extensibility in multiple areas

# Non functional design goals

- Open Source everything:
  - o API's
  - Server application
  - Detectors (software and hardware!)
  - Tools
- Provide documentation
- Apply single purpose principle, KISS, etc.





• Provide means for growth -> improve science reproducibility!



#### credo-webapp

Credo web application



#### https://github.com/credo-science

# Software components of CREDO project

- credo-webapp (Server application)
- credo-api-tools (Utilities)
- credo-detector-android
- CREDO-monitor-TimeClusteringAlgo
- Credo-Desktop-Detector
- CREDO-PC-Windows

. . .

#### Component diagram



# Credo Server application

- <u>https://github.com/credo-science/credo-webapp</u>
- Django app running under Apache and mod\_wsgi
- Display real time detection information
  - basic on-line analysis of data
  - filtering based on provided criteria
- Provides API for other components
  - versioned APIs!
- Manages user accounts
- Provides data export facilities







CRED : THE QUEST FOR THE UNEXPECTED					
Main page					
Detections (702821) Users (3711)	Teams ( 1503 )				
네 Top users		لط Rec	ently registered user	s	
‡ Login	♦ Detections ♦ Login		Detections		
Barti	32992	Izka		0	
Wzorzasty	23901	m.jurczak		0	
smph-kitkat	22626	Mateusz		0	
Natalia	17389	Patrycjaboryczko		0	
Emanuele Maria Latorre	17129	Nivv		0	
Last 20 detections					
≑ date	\$ login		≑team		\$ img
2018-10-22 09:57:50.135	Patryk Olszowski		CREDO ASP OXFORD		
2018-10-22 09:56:07.544	smph-kitkat		IFJ		
2018-10-22 09:55:44.510	Patryk Olszowski		CREDO ASP OXFORD		
2018-10-22 09:52:49.200	Marcin Kuzniak				

#### https://api.credo.science

# Credo Api Tools

- <u>https://github.com/credo-science/credo-api-tools</u>
- Collection of tools designed to simplify process of exporting and analyzing available data
- "data-exporter" handles making API calls to authenticate user, request and download data export incrementally
- "data-processor" handles incremental data processing and provides simple plugin interface for scientists to write their own code

### Hardware infrastructure

- Virtualized environment
- Hosted at Cyfronet's cloud
- Runs on 8 VMs
- Periodic multi tier backups
- Automated deployment, through Ansible
- Available resources allow for scaling if required

# Performance and availability monitoring

- Application level
- Operating system level
- Time series DBs: django-influxdb-metrics + InfluxDB + Grafana
  - request count
  - latency
- Event and log monitoring: Raven + Sentry
  - warnings/errors
- Elasticsearch + Kibana:
  - operating metrics and utilization (Metricbeat)
  - log shipping (Filebeat)
  - high-level overview of gathered data (heatmaps, rates)











[CREDO] Top users (by on time)

Ping count 🗘

13,868

13,457

9,139

4,288

4,708

302,462 347,922

-	Detection count $\ddagger$	User ID 🗢	On time 🗘
	32,992	4980	3 months
	23,956	4125	3 months
	17,389	157	2 months
	14,114	1163	2 months
	11,462	7358	a month
er	234,060	Other	5 years
	333,973		6 years





### Conclusions

- Implemented system delivered the required functionality
  - Implementation methodology encourages community involvement
- APIs proved to be reusable, multiple integrated components
- Methodology used to develop software was a success, multiple contributions

remaining challenges:

• Traffic can be unpredictable - ongoing work

#### References

- 1. Homola, P., Bhatta, G., Bratek, Ł., Bretz, T., Cheminant, K. A., Castillo, D. A., ... & Jarvis, J. F. (2018). Search for Extensive Photon Cascades with the Cosmic-Ray Extremely Distributed Observatory. arXiv preprint arXiv:1804.05614.
- 2. Conrad, C. C., & Hilchey, K. G. (2011). A review of citizen science and community-based environmental monitoring: issues and opportunities. Environmental monitoring and assessment, 176(1-4), 273-291.
- 3. CREDO's first light: The global particle detector begins its collection of scientific data: https://www.eurekalert.org/pub\_releases/2018-10/thni-cfl100418.php (accessed 4.10.2018)