# Building cloud native Machine learning based applications

Mr. Bhushan Hiralal Patil

Mr. Mohan Nikam

Mr. Arunvel A

M.Tech, Department of Computing Science and Engineering,

Sandip University, Nashik

CDAC ACTS - PUNE

*Abstract-* **Cloud computing has an emerged as animportant relevence to improve resource utilization, efficiency, flexibility, However, cloud platforms cause performance degradations due to their virtualization layer and may not be appropriate for the requirements of high performance applications, such as machine learning. The Research tackles the problem of improving network performance in container based cloud instances to create a viable alternative to run network intensive machine learning applications.**

**Our approach consists of deploying container based machine learning based apps via LXC (Linux Container) cloud instances to increase the available bandwidth & performance. so to evaluate the efficiency of this approach and the overhead added by the container-based cloud environment, we ran a set of experiments to measure throughput, latency, bandwidth utilization, and completion times. The outcomes proves that this approach adds minimum overhead in cloud environment as well as increases throughput and reduces latency.**

**Containers are a lightweight virtualization solution to replace virtual machines for deploying cloud applications as they are less resource and time consuming. Easy deployment of applications on containers is done using kubernetes which helps user to analyze applications with various network topologies improving transfer rate of large data sets which are required for machine learning models**.

## I. INTRODUCTION

Today, cloud computing is rapidly developing and large scale application. As a kind of network-based computing, cloud computing can provide shared software and hardware resource to users by using network, which realizes the reasonable distribution of information and resource. Meanwhile, cloud computing brings along the changes of the traditional data center, which produces a new generation of data center: cloud data center. By building computing resources, storage resources and network resources into dynamic virtual resource pool using virtualization technology But now cloud data center faces with many problems such as low resource utilization, application and platform can not be decoupled, application

runtime environment limitations are strong, and operational staff control decrease, which limit the development of it.

Virtual machines (VMs) are an infrastructure as a service (IaaS) focusing on hardware virtualization whose techniques have been used at the infrastructure layer. To achieve sharing and elasticity of resources, the cloud makes use of such virtualization techniques for administering scheduling, provisioning and security. However, VMs suffer from slow start up time and exhibit lower densities even when full. Applications occupy entire host operating system (OS) in VMs and there are also several limitations on storage and Containers offer as a lightweight virtualization solution to deploying applications across various domains and sectors. They allow infrastructure and platform to be shared in a secure and portable manner, along with application packaging and management. They facilitate faster deployment of applications and have faster start up times. They are less resource and time consuming and can be scaled up or down providing higher density levels than full VMs. They facilitate easier and portable across infrastructure deployment of applications in an interoperable way. Using containers will accelerate agile application development of distributed applications providing an additional layer of protection by isolating applications and the host, without using incremental resources. They also allow easy updates to applications. Fig. 1 shows the difference between traditional hypervisor and container-based architectures
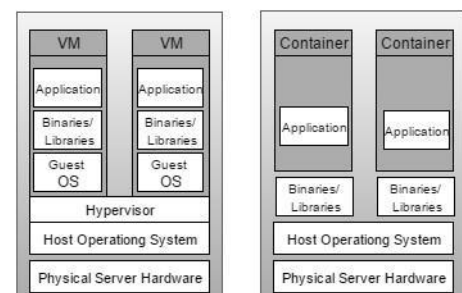


Fig.1 hypervisor architecture and a container-based architecture

Container technology is a virtualization technology at the operating system level, it can provide a separate file system, network and process context space for each running server without modifying the host operating environment. So that,

each service can arbitrarily change the contents of the operating system files, configuration of network and user in their own virtual running space, without damage to real host system files. Although an instance of the operating system is shared between containers, they run independently and without interference. Besides, do not relaying on the integrity of the operating system makes container more lightweight and fast start than virtual machine, which meets the needs of cloud data center. Cloud computing infrastructures support rapid resource provisioning, which in turn is a suitable option to deploy a container environment. Primarily due to the advances of cloud platforms, there are possibilities to develop new business models based on the pay-per-use billing structure. Also, there are cloud computing technologies (specifically those focused on network and storage) that provide lower performance degradation. However, the use of cloud environments for processing distributed applications has traditionally been avoided when it requires high bandwidth and throughput as well as low latency. Traditional cloud providers, which use virtualization technologies that are not optimized for the execution of distributed applications, add significant overheads. Our main challenge in this paper is to improve network performance in the container-based cloud environment for applications. We chose to increase the network performance by using the IEEE 802.3ad link aggregation standard. It gives and uses a standard method to combine multiple physical links that can be used as a single logical link. The standard is a layer 2 control protocol that can be used to automatically detect, configure, and manage a single logical link with multiple physical links between two adjacent enabled devices. Thus, link aggregation provides higher availability and capacity while network performance improvements are obtained using existing hardware (IEEE 802.3ad requires support in a network switch). Also, our goal is to target a high-performance cloud environment by deploying container-based instances, where applications may run. We chose container-based technologies to enable multiple isolated Linux systems to run on a single host through Namespaces (providing isolated user environments in the form of containers) and cgroups (providing resource management and accounting). LXC (Linux Container) technology is an important part of this cloud infrastructure, because it is a free software that provides a powerful set of user space tools and utilities to manage Linux containers

## II. RELATED WORK

In this section, we provide a brief background about machines placement, Docker containers, Docker Container Networking. We also described related work reported in the literature with regards to performance of Docker containers when compared to the performance of VMs. ALSO we should focus on container network models. There are two mainstream models called Container Network Model (CNM) and Container Network Interface (CNI).

### A.        Virtual machine placement:

Nowadays, cloud computing has been widely used.The Virtual Machines (VMs) are created on servers in cloud computing. With the rapid development of Cloud computing technologies, the management of Cloud resources has become crucial to Cloud Service Providers (CSPs). CSPs provide on demand services to the end users through provisioning Virtual Machines (VMs) in their available

Physical Machines (PM)s. In order to reduce the Total Cost of Ownership(TCO), CSPs consolidate multiple VMs in each of the available PMs. In this context, significant number of research work have already been done which indicate the benefit of running multiple VMs in a PM. The VM scheduling on servers for energy saving in the cloud computing has been studied. The Virtual Machine Scheduling Algorithm (VSA) is proposed to schedule VMs in cluster environments. However, it is not effective and has high time complexity. The VM scheduling on servers for energy saving inthe cloud computing has been studied. The Virtual Machine Scheduling Algorithm (VSA) is we propose a new scheduling method called Energy-aware Virtual Machine Placement (EVP)method to schedule VMs that can reduce power consumption complexity. Power consumption of a PM depends on the number of VMs deployed, reserved resources of VMs, and the state ofthe processes running within the VMs. There are significant number of works that try to minimize the the overall power consumption of a data-center by reducing the number of active PMs. They classify this problem of mapping the VMs to minimum number of PMs satisfying the constraints of physical resources and the requirements of the VMs as the VM Placement Problem

### B.        Docker Container:

Docker containers are lightweight, highly portable and scalable. Such features become attractive to develop containerized services (or microservices). In the literature, container-based services are reported to always outperform VM-based services. For example, it was shown in that containers outperform virtual machines in terms of execution time, latency, throughput, power consumption, CPU utilization and memory usage. However, according to Amazon AWS documentation. it was reported that Docker containers are deployed on top of virtual machines (VMs), and not on bare-metal hardware. This contradicts the common practice of container deployment which is widely adopted in the literature of deploying container on bare-metal hardware. This was the primary motivation for our research work in which we aim to investigate and assess the performance difference resulting from deploying services using VMs (virtual machines) and using Docker containers
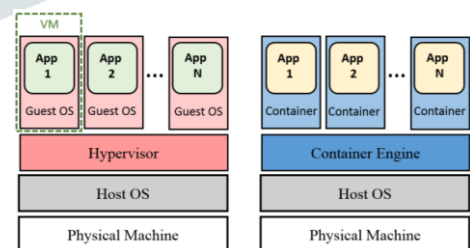


Fig. 2 Common approach of hosting of VMs vs. Containers

The Services deployed using containers are expected to take less execution time thus resulting in less latency over virtual machines as in. In addition, containers were shown to introduce lower power consumption over virtualized technologies.

### C.        Docker Container Networking:

As an important part of container technology, container network mainly solves the problem of communication between containers in the case of ensure container isolation. Since container is a new technology, container network is

not yet mature, and many container network solutions were proposed to solve the problems of container network. Each solution has its own characteristics, and how to choose appropriate solution according to different scenarios is not only very important for the efficiency of communication between containers but also of great significance in improving cloud data center performance. In order to provide experiences and references for choosing and deploying appropriate container network, we introduce three kinds of mainstream container network solutions (Flannel, Docker Swarm Overlay, Calico), and design experiments to measure and evaluate the performance of them

### III. LITERATURE REVIEW

This section describes the primary related works aimed at improving the network performance and deploying Container Based instances or Services. Our first scenario consists of evaluating network performance and then second scenario consist deploying Container based ML services The Authors Wattanasomboon, P., & Somchit, Y. (2015, October)[1] proposes the method called EVP(Energy-Aware virtual machine placement ) to method to schedule VMs that can reduce power consumption. We also formulate power consumption model to evaluate the performance of the EVP method. Finally, we evaluate the EVP method by simulation. On The basis of experimental results show that the EVP method has better performance. There are many researches of VM scheduling for energy saving. Befor using EVP method for our approach, we have to check performance of our VMs and containers and choose which platform is better to to deploy the ML(machine learning )Services.the Salah, Tasneem, et al.[5] stated the performace comparison between VMs and Containers and gave theirs results as Container always outperforms the VM in terms of the parameters like Cpu Utilization, latency, Throughput, Memory usage, Execution Time etc. So Services deployed using containers are expected to take less execution time thus resulting in less latency over virtual machines as in
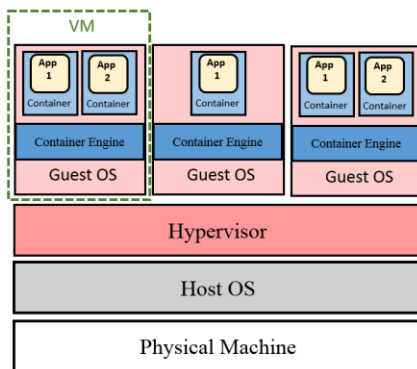


Fig,3 Hosting of Containers in Amazon cloud

Furthermore, according to Amazon , Docker containers are deployed on top of virtual machine, which is against the common practice of deployment as depicted in Fig.A&B and Fig.3 shows the container deployment approach adopted by Amazon cloud.

### IV. PROPOSED METHODOLOGY

This approach involves the accessing of services that are stored in containers and checking the maintenance of computing nodes. These services are designed to perform two types of tasks: first accessing the all configuration,

dependancies of machine learning services and second is deploying container based ML services in cloud Enviornment.
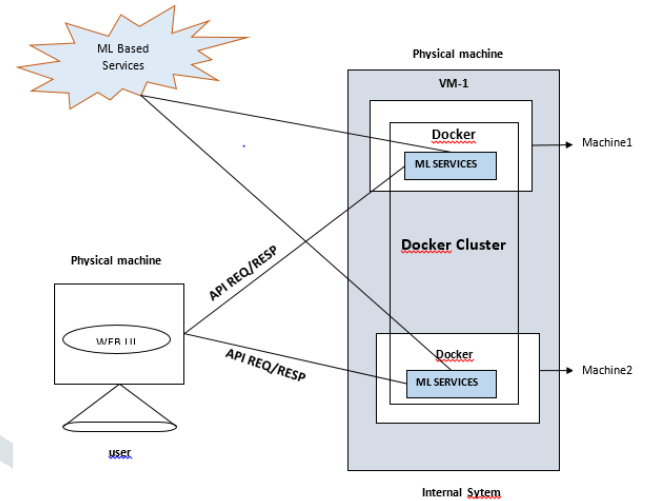


Fig.4 Proposed Architecture

The above Figure shows proposed architecture ,in that we make the use of VMs over PMs and Configured all the necessary packeges related to Docker. after docker installed, we create the cluster of docker running the command - docker Swarm init on one of machine and then token is generated for adding or joining workers/slaves to master,so we used that token for forming a docker cluster. so this is how cluster is created. once cluster is formed, we are trying to build containers which contains the set of ML services.so on top of docker ,we deploy container-based ML services on cloud using deployment technology (eg: using kubernetes). for deployment we served container-based ML services as a API using Keras and Flask.We then put that application inside of a Docker container, uploaded the container to Docker Hub, and deployed it with Kubernetes.this is how the workflow and implementation of proposed methodology is done So the machine learning services or applications that I'm builded and deployed is nothing but the object detections in a images. It finds out object in the given image that where is the object in images and what is ratio of object and maps it. The ML services is built on top of python 3.6 (docker) and has these packages built into it. The packeges are:
• Tensorflow — 1.7
• Tensorflow Object detection
• Keras
• Scipy
• Matplotlib
• Numpy
• Pandas
• Sklearn
• Open CV
• Scikit Image
• Pillow
• Jupyter Notebook

### V. IMPLEMENTATION
Following steps are listed and certain tools are used for developing this project.The Tools That are required to build and deploy container based ML services are follows

A.      Docker Cluster:

Docker is a tool designed to make it easier to create, deploy, and run applications by using containers.In a way, Docker is a bit like a virtual machine. But unlike a virtual machine, rather than creating a whole virtual operating system, Docker allows applications to use the same Linux kernel as the system that they're running on and only requires applications be shipped with things not already running on the host computer. This gives a significant performance boost and reduces the size of the application.So, to gain boost in performance we create the cluster of docker. crating cluster of docker is done by using Docker swarm. Docker Swarm is a clustering and scheduling tool for Docker containers. With Swarm, IT administrators and developers can establish and manage a cluster of Docker nodes as a single virtual system.you will also need to install the Docker engine on both nodes. By default, Docker is not available in the Ubuntu 16.04 repository. Requirements:

- Hypervisor
- Docker engine
- Docker swarm

First, you will need to initialize the cluster with the IP address, so your node acts as a Manager node. On the Manager Node, run the following command for advertising IP address:

*~docker swarm init --advertise-addr 192.168.222.138*

You should see output like this:



The token shown in the above output will be used to add worker nodes to the cluster in next step. The Docker Engine joins the swarm depending on the join-token you provide to the docker swarm join command. The node only uses the token at join time. If you subsequently rotate the token, it doesn't affect existing swarm nodes. Now, check the status of the Manager Node with the following command: docker node ls If everything is fine, you should see the following output:



You can also check the status of the Docker Swarm Cluster: code>*docker info*



Manager Node (master) is now configured properly, it's time to add Worker Node(slave) to the Swarm Cluster. First, copy the output of the "swarm init" command from the previous step, then paste that output on the Worker Node to join the Swarm Cluster: *docker swarm join --token SWMTKN-1-5p5f6p6tv1cmjzq9ntx3zmck9kpgt355qq0uaqoj2ple629dl458 80qso8jio78djpx5mzbqcfu192.168.222.138:2377*

You should see the following output:



This node joined a swarm as a worker.

Now, on the Manager Node(master), run the following command to list the Worker Node(slave):

*~docker node ls*

You should see the Worker Node(slave) in the following output



Fig A. Docker Cluster

Docker Swarm Cluster is now up and running, it's time to launch the web service inside Docker Swarm Mode. On the Manager Node, run the commands to deploy webservice

B.      Configuring Kubernetes:

Kubernetes (commonly stylized as k8s) is an open-source container orchestration system for automating application deployment, scaling, and management. It was originally designed by Google, and is now maintained by the Cloud Native Computing Foundation. It aims to provide a platform for automating deployment, scaling, and operations of application containers across clusters of hosts. It works with a range of container tools, including Docker. Many cloud services offer a Kubernetes-based platform or infrastructure as a service (PaaS or IaaS) on which Kubernetes can be deployed as a platform-providing service.

C.      Installing Kubernetes:

After Docker, we install kubernetes on top of docker to deploy Containerized applications,kuberenetes is configured is shown by the readiness state.the below images shown that

Fig C. kubernetes installation

**D.　　Kubernetes Components:**

Below image shows that kubernetes components configured via command.



Fig D. Kubernetes Components

**E.　　Dashboard Creation:**



After Creating dashboard it shown as below.



Fig E. Dashboard Creation

**F.　　Token Generation:**

we can get access to dashboard by generating token, for that we used token generation command and execute it and used it to get access to dashboard.The token generation by command is shown below:



Fig F. Token Generation

**G.　　Kubernetes Cluster:**

After getting access to dashboard we have create cluster by joining nodes to master.so we have to deploy container based ML services using kubernetes. The cluster of kubernetes is shown below and there deployements ,workloads replica sets are also shown below.



Fig G. Kubernetes Cluster

Once an applications is deployed ,we can see the how in the given image the object is detected and the outcome or results shows clearly the two given images in which we see the objects are detected.we can see those objects are also detected that we can see clearly by our eyes. So you can imagine how much powerful the machine learning algorithms and how efficient it is. By using kubernetes you can see cpu usage, memory usage of deployed applicaton on kubernetes dashboard

**H.　　Machine Learning Services:**

Here, the machine learning services are nothing but the applications that we builded and deployed. there are many kind of services we deployed using docker and kubernetes, we just need to packaged the source code/scripts inside

container and deployed it using one the deployement technology. But it is not that easy job, we also require lots of libraries , configuration files and dependancies. One most important things we need to install is machine learning packages, if we want to build and deploy the machine learning services or applications. Following are the list of packages that are necessary for any machine learning application and we also need to configure these packages for our application which is machine learning based. The packeges are:
• Tensorflow Object detection
• Keras
• Scipy
• Matplotlib
• Numpy
• Pandas
• Sklearn
• Open CV
• Scikit Image
• Pillow
• Jupyter Notebook

The machine learning based services which is being build is nothing but the object detection in image. It means, build an machine learning application which detects the object in given image and gives accurate measurement related to object or entity. It gives you ratio about object means it scales where is your object and at how big in size. So I'm currently working on doing some changes like i said above. So here's below is the code or snippet about my application

a.       Imports:
Once the packages for application is installed and configured then its time to import them


Fig a. Imports

b.       Environment Setup:
Enviornmental setup is important part of application To display an image where the result located.So we used a library called matplotlib to display an images and plot the detected object on that images.


Fig b. Enviornment Setup

Enviroment setup is important step because once your foundation is strong then you Have much more opportunites to prepare and train your model effectively accurate.

c.       Model Preparation and Download
Next step is to prepare the model which is being builded and deployed after that we download that model for future improvement in application


Fig c. Model Preparation and Download

d.       Detection:
Finally, the last part come i.e to detect object in image that's our aim of application or machine learning based service. We have using two images for simplicity to test out application to check either it is working correctly or not. We added the path of images that we gonna used for testing purpose to check whether application is correctly able to detect an object on given images or not So that's our final step or final script which is the detection script to detects the object in a given images. It finds or detects the object in images that we specifically used two images that we seen in result of application.


Fig d. Detection

e.       Results
As you can see below the outcome of our application that detect the object in a given image.. The results generated by the application is shown below in two images. First result shows the detection of being and birds. The second images shows the detection of dog.
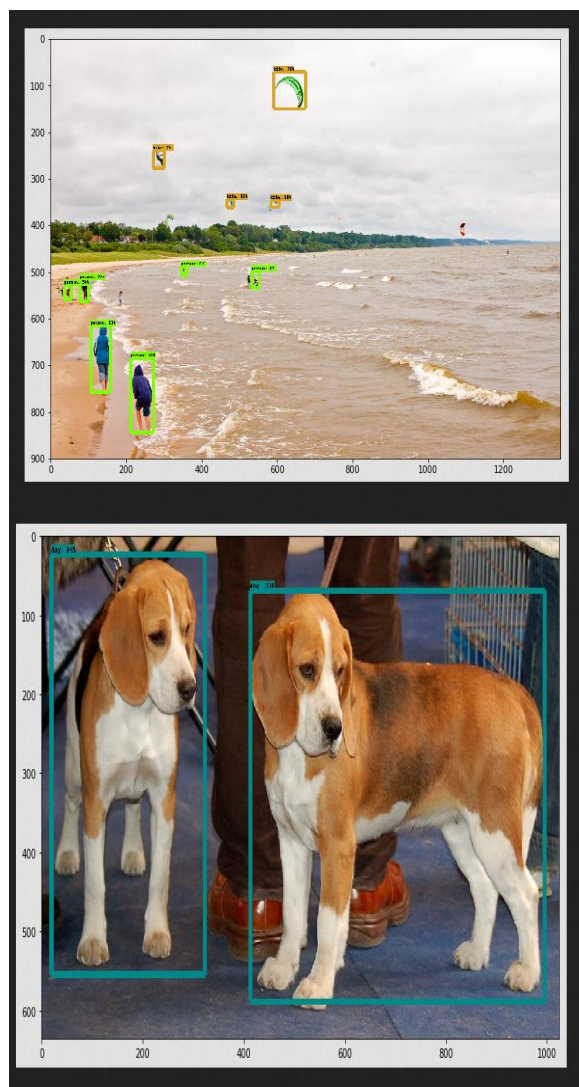
Fig e. Results

## VI. CONCLSUTION & FUTURE SCOPE STUDY

This project concludes the developing a machine learning application is common and might be easy and we call such applications as a user-specific application but bringing a machine learning application by converting its whole scenario in cloud enviornment is something different. so in our project, we achieve this feet by using new cloud technologies. The purpose of this project is to develop an application using emerging trending technologies. the tools and technologies which are used in this project are rarely know to people because this are the upcoming technologies that i have used in project. so people are not much aware about this technologies and i build the project in such technologies might be great and nice experience for me. The future step might be doing the performance comparison of deployed services using Prometheus & ISTIO and also improve their performance as its our primary concern..we trying to the monitor performance of containers based ml services by using the prometheus which is a open-source monitoring and altering toolkit.

## REFERENCES

[1] Ernst &Young."Virtual machine placement method for energy saving in cloud computing." Information Technology and Electrical Engineering (ICITEE), 2015 7th International Conference on. IEEE, 2015.

[2] Su, Nan, et al. "Research on virtual machine placement in the cloud based on improved simulated annealing algorithm." World Automation Congress (WAC), 2016. IEEE, 2016.

[3] Karmakar, Kamalesh, SunirmalKhatua, and Rajib K. Das. "Efficient virtual machine placement in cloud environment." Advances in Computing, Communications and Informatics (ICACCI), 2017 International Conference on. IEEE, 2017.

[4] Yang, Song, et al. "Reliable virtual machine placement and routing in clouds." arXiv preprint arXiv:1701.06005 (2017).

[5] Salah, Tasneem, et al. "Performance comparison between container-based and VM-based services." Innovations in Clouds, Internet and Networks (ICIN), 2017 20th Conference on. IEEE, 2017.

[6] Rista, Cassiano, et al. "Improving the Network Performance of a Container-Based Cloud Environment for Distributed Systems." High Performance Computing & Simulation (HPCS), 2017 International Conference on. IEEE, 2017.

[7] Xu, Pengfei, Shaohuai Shi, and Xiaowen Chu. "Performance Evaluation of Deep Learning Tools in Docker Containers." Big Data Computing and Communications (BIGCOM), 2017 3rd International Conference on. IEEE, 2017.

[8] Zeng, Hao, et al. "Measurement and Evaluation for Docker Container Networking." Cyber-Enabled Distributed Computing and Knowledge Discovery (CyberC), 2017 International Conference on. IEEE, 2017

[9] Mao, Ying, et al. "Draps: Dynamic and resource-aware placement scheme for docker containers in a heterogeneous cluster." Performance Computing and Communications Conference (IPCCC), 2017 IEEE 36th International. IEEE, 2017

[10] Model, Víctor, et al. "Modelling performance & resource management in kubernetes." 2016 IEEE/ACM 9th International Conference on Utility and Cloud Computing (UCC). IEEE, 2016.

[11] , Víctor, et al. "Adaptive application scheduling under interference in kubernetes." Utility and Cloud Computing (UCC), 2016 IEEE/ACM 9th International Conference on. IEEE, 2016

[12] Tsai, Pei-Hsuan, et al. "Distributed analytics in fog computing platforms using Tensorflow and Kubernetes." Network Operations and Management Symposium (APNOMS), 2017 19th Asia-Pacific. IEEE, 2017.

[13] Chang, Chia-Chen, et al. "A Kubernetes-Based Monitoring Platform for Dynamic Cloud Resource Provisioning." GLOBECOM 2017-2017 IEEE Global Communications Conference. IEEE, 2017. [14] Xie, Xiao-Lan, Peng Wang, and Qi Wang. "The performance analysis of Docker and rkt based on Kubernetes." 2017 13th International Conference on Natural Computation, Fuzzy Systems and Knowledge Discovery (ICNC-FSKD). IEEE, 2017.

[15] Rovnyagin, Mikhail M., et al. "Using the ML-based architecture for adaptive containerized system creation." Young Researchers in Electrical and Electronic Engineering (EIConRus), 2018 IEEE Conference of Russian. IEEE, 2018.