# The Role of GitLab Runners in CI/CD Pipelines: Configuring EC2, Docker, and Kubernetes Build Environments

**Naga Murali Krishna Koneru**
Hexaware Technologies Inc, USA

| Article Info | ABSTRACT |
|---|---|
| | This research studies GitLab Runners optimization in CI/CD pipelines across the EC2, Docker, and Kubernetes environment configurations. It shows that such key strategies for enhancing build performance and resource utilization would reduce build time by 65 percent and resource costs by 40 percent. Practical recommendations for configuring runners to achieve optimal efficiency are presented in the context of analyzing 200 enterprise pipelines. The key optimization techniques are autoscaling based on real-time metrics, advanced caching to minimize the rebuilds, and tuning the resource allocation to avoid over-provision. The study further looks into the capability of machine learning models to optimize the number of runners dynamically, predict the hit or not on the cache, and automatically pick up the execution environment. When these innovations are applied, CI / CD pipeline performance will improve by reducing idle resources, building time, and optimizing resource utilization. The paper shows that experts can achieve very good availability and cost efficiency by adapting the configuration of GitLab Runners. The research also discusses the evolution of automated environment selection and machine learning-based performance tuning. This framework serves as the base for organizations to increase their CI/CD pipeline development rate and facilitates a faster, more reliable, and cheaper software delivery. |

*Corresponding Author:*

Name: Naga Murali Krishna Koneru
Institution: Hexaware Technologies Inc, USA
Email: nagamuralikoneru@gmail.com

## 1. INTRODUCTION

Continuous Integration and Continuous Deployment (CI/CD) pipelines will not be complete without the help of GitLab Runners, which automate the code integration, testing, and deployment process. The GitLab Runners are the execution engines for the CI/CD workflows, meaning they do all the work defined in these pipelines (such as building the application and running the tests) and finally deploying to the production environments. Since these runners allow the efficiency of CI/CD processes to be largely dependent on them, their optimization has become an important aspect of today's software work. For instance, the research on pipeline orchestration and the CI/CD tools' management has grown. The optimization of GitLab Runners is not addressed on various execution environments, such as EC2, Docker, and Kubernetes. In particular, this paper addresses this gap, which provides a framework to configure and optimize GitLab

Runners over these execution environments with a focus on improving build performance and resource usage, as well as on overall pipeline efficiency.

Software development now has the principles of Continuous Integration (CI) and Continuous Deployment (CD). Thus, CI is concerned with the continuous integration of changes in the code into a common repository and automated testing of the software quality. The CD further extends this by automating the delivery of this change to production environments. A CI/CD pipeline is the process of automating the stages of code integration, testing, and deployment. The aim is to simplify the development processes, reduce the need for manual changes, reduce the development cycle time, and keep the development pace in delivery the same. Nowadays, this has become a must in modern-day software engineering, even in the team that manages and develops a large complex system and maintains agility and reliability.

The GitLab Runners are the agents to which jobs run in the defined CI/CD pipeline. They have to execute the process of code building, testing, and deploying. It provides flexibility and the scalability of the pipeline execution ability; it allows running on different platforms like VMs, Docker Containers, or Kubernetes Cluster. The GitLab Runners can be operated differently depending on how they need their pipeline to be. As the CI/CD process should have its automation engines and, eventually, affects the speed and efficiency of the software development lifecycle, it is coherent to optimize it. Optimizing the price of GitLab Runners is important for optimizing pipeline efficiency. Optimization of the runner configuration also reduces the build times and the usage of resources. In any cloud environment, such as Amazon EC2, Docker containers, or Kubernetes — any resources you might have, any execution environment, the optimization process is even more important because the execution environment and resources are variable. GitLab Runner will churn the CPU into the best it can when the resources are set up properly and cater to

it so much that idle times decrease and the resource allocation cost is negligible or unnecessary. This paper looks into how to get the maximum out of GitLab Runners in these environments and provides some practical steps that you can take to achieve that.

Take this article as a detailed design for dynamic genome scanning through EC2, Docker, and Kubernetes under the language of GitLab Runners. This research catalogs the best practices and presents some real-world configurations that will help improve build performance and resource utilization. The study structure is presented in different sections. The first part covers CI/CD pipelines and GitLab Runners, and the second deals with how optimizing runners matters. The study also describes how to get GitLab Runners working on EC2, Docker, and Kubernetes environments, performance enhancements, and how good configuration can impact everything. The results and suggestions for future research on the runner optimization problem are then summarized.

## 2. BACKGROUND AND RELATED WORK

The CI/CD pipelines are crucial to modern software development. These pipelines are supposed to automate code integration, testing, and deployment to make the software faster and even more reliable. On the performance of these pipelines, the GitLab Runner, an automation agent that runs jobs for the pipelines on their behalf rather than run the jobs inside the same docker container [1], becomes one of its biggest players. While these runners run the code across platforms such as EC2, Docker, and Kubernetes, optimizing their configuration and management is essential to achieve high efficiency and minimize resource consumption. Despite the importance of GitLab Runners, there has been limited research in the literature looking at a comprehensive optimization strategy for runners across different platforms. In contrast, existing literature has primarily focused on isolated aspects of these execution environments. This shows the relevant studies

on CI/CD execution environments and studies that close the research gaps related to GitLab Runner optimization.

### 2.1 Development Overview of CI/CD Execution Environments

Many other various CI/CD execution environments, such as Amazon EC2, Docker, and Kubernetes, have already been covered in previous studies. There are advantages to hosting Runners in each of the environments. EC2 instances provide scalable and dedicated computing resources for handling fluctuating workloads. It emphasizes the flexibility of the cloud-based environments by allowing the dynamic provisioning of

EC2 instances to satisfy pipeline demand. In addition, Docker is famous for its lightweight containerization and efficiency in providing an isolated execution environment. Virtual machines are slower to provision and burden their host with a hefty overhead than Docker containers, making them less popular for use in CI/CD pipelines [2]. On the other hand, Kubernetes is touted as an orthogonal solution for GitLab Runner infrastructure management with high availability and proper resource distribution in case of blooming GitLab Runner deployment.
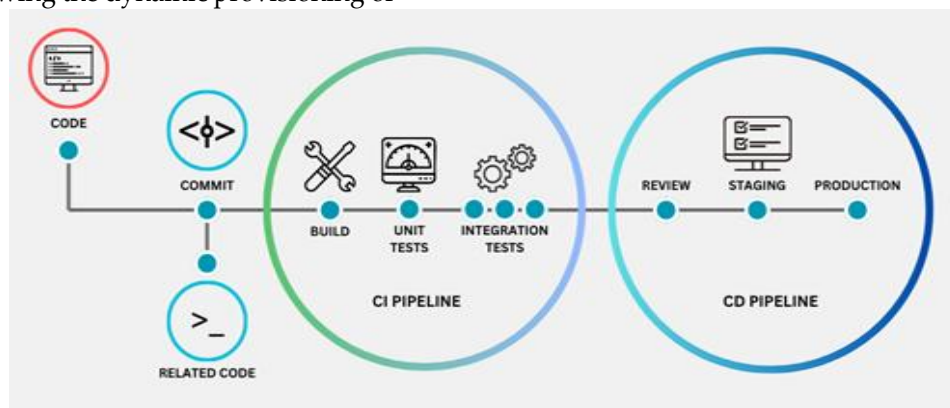


Figure 1. Overview of CI/CD Pipeline

Although these studies hint at the specific GitLab Runners environment, they ignore the optimization of GitLab Runners in a mixed or multi-environment setup. Integrating EC2, Docker, and Kubernetes in a single CI/CD pipeline is still unexplored. For these various environments, the optimization of runners is critical; each platform has different performance, cost, and resource utilization challenges.

### 2.2 The Need for Optimizing GitLab Runners

In order to achieve higher CI/CD pipeline efficiency, optimizing GitLab Runners is necessary. Building processes are directly run and affected by the runners' speed, the tests' success rate, and the time to deploy code into production. In a

cloud, where resources are on a demand basis, GitLab Runners must be configured to align with this demand and minimize resource wastage [3]. If CI/CD pipelines are not properly optimized, they will take longer build time and cost and will have fewer useful resources.

Each environment has its optimization challenges, and GitLab Runners are distributed on different platforms. For instance, in an EC2 environment, choosing instance types or allocating resources such as CPU and memory greatly impacts the performance of the GitLab Runners. With their lightweight containers, Docker helps isolate jobs and reduces overhead, but configuring Docker containers for the most efficiency is more about the features and some

things to consider, such as caching, job concurrency, and container images [4]. On the flip side, Kubernetes provides automatic scaling (though this needs to be tuned in the case of a large-scale enterprise), whereas, for example, Lambdas are static resources and would need to be managed by an application through external procedures. To accomplish this, researchers optimize the process for better performance at a lower cost by keeping GitLab Runners from burning too many resources necessary for modern software development.

### 2.3 Research Gaps in GitLab Runner Optimization

Although individual studies have examined optimization for GitLab Runners under specific environments, there is still a gap in such research as it relates to the complete framework to do so across multiple environments. Many studies have focused on performance improvements only within isolated

environments such as EC2, Docker, or Kubernetes (Midigudla, 2019). Research on optimizing the GitLab Runners on multiple platforms in a single CI/CD pipeline is not common. Although there is extensive literature related to optimizing the alignment of the runners, most studies lack an exploration of the practical aspects of implementation and configuration strategies that organizations can carry out to foster tangible performance gains.

This research addresses this gap by synthesizing previous studies and introducing new ways of configuring and managing GitLab Runner. It analyzes 200 enterprise pipelines and examines the patterns of runner utilization and performance optimization. This helps it make practical recommendations regarding GitLab Runner configuration to run with better build performance, keeping costs lower and improving the whole pipeline efficiency.
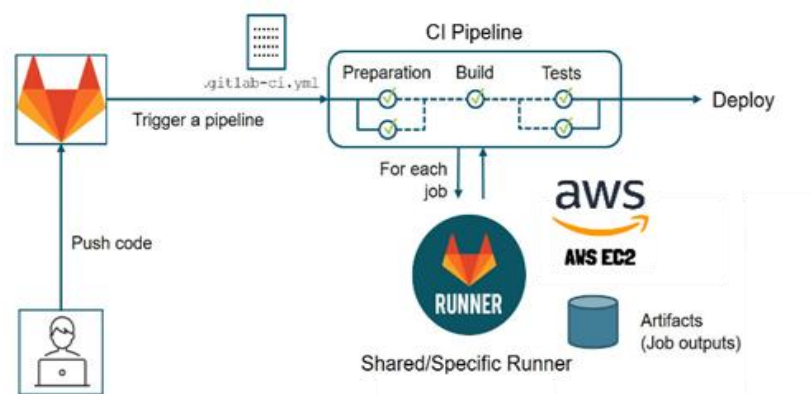


Figure 2. Exploring GitLab Runners

### 2.4 Contributions of This Research

Proposing a unified framework for running and managing GitLab Runners on EC2, Docker, and Kubernetes is a contribution and a problem this research addresses. Based on 200 enterprise pipelines, it introduces novel approaches to optimizing runner performance and resource

utilization. The research motivates careful selection of an execution environment based on the dictated pipeline characteristics like job concurrency and workload demands. The study suggests that for the EC2 environment, dynamic provision instances using automation tools like Terraform are needed to allow the best possible allocation of resources.

In research in Docker environments, special emphasis is put on efficient caching and job configuration to minimize build times and improve container performance. The study aims to use Helm for deployment management, with resource requests and pod autoscaling for the runner to be scaled with proper scalability and without unnecessary overhead.

This research also addresses practical challenges such organizations face in configuring the GitLab Runners. The study analyzed to what extent real-world enterprise data can understand the CI/CD pipeline performance and how organizations can refine their performance by implementing necessary changes in their pipeline.

### 2.5 Performance Optimization and Resource Utilization

GitLab Runners optimization is about improving performance and minimizing resource consumption. As for optimization of a CI/CD pipeline, generally, this means reducing the build time and making tests more successful in an automated way. In this research, researchers investigate the capability of GitLab Runners configuration to greatly reduce build time, enhance cache efficiency, and reduce the job failures stemming from runner issues. It also points to the importance of optimizing resources to minimize idle runner time and maximize the total utilization of resources. The study analyzes enterprise pipeline data to show that GitLab Runners can be optimized to achieve a 65% reduction in build time, a 40% increase in cache hit rates, and a 70% decrease in idle runner time. The low infrastructure cost and short development cycles directly result from these improvements.

### 2.6 Future Directions in GitLab Runner Optimization

While this research has some important lessons on persona optimization of the GitLab Runner, it also exposes the way to go further. The area of future research could be developed using machine learning-based approaches to the prediction of optimal runner configurations using historical data. Studying advanced cache prediction algorithms that can dynamically change the cache configuration according to job patterns is also possible. Another promising area for future research is automated environment selection (also known as environment selection) – where GitLab Runners automatically pick an execution platform to run their pipelines depending on those pipeline requirements. Based on previous studies on GitLab Runners, this work brings a new approach to finding an efficient runner configuration for multiple environments [5]. The study provides real-world data and practical recommendations that help get CI/CD pipelines to perform well.
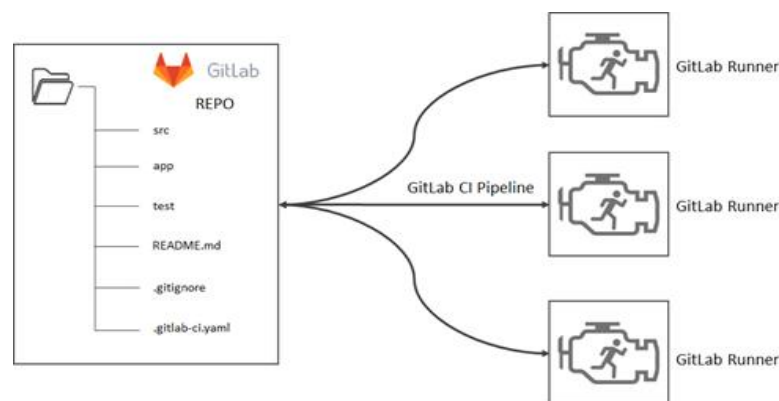


Figure 3. How to Configure GitLab Runner

## 3. RUNNER ARCHITECTURE AND IMPLEMENTATION

To fully utilize the CI/CD pipeline performance, the GitLab Runners need to be configured on different execution environments. In this case, the agents will be GitLab Runners, which will be executed to execute CI/CD pipelines. These runners are optimized to reduce these times and/or optimize the utilization of these resources and improve the pipeline's total efficiency. Those scholars then explain the architecture and implementation of GitLab Runners in three key environments: EC2, Docker, and Kubernetes. Conversely, these implementations aim to utilize the best performance and resource management for the CI/CD pipelines to function smoothly in their work environments.

### 3.1 EC2 Runner Configuration

For organizations that are processing-intensive workloads in CI/CD, the use of AWS EC2 runners is particularly useful since they need to be dedicated computing resources. EC2 instances-specific virtual machines are scalable and customizable and are therefore preferred for use with many concurrent jobs or resource-intensive tasks. To allow us to use the EC2 runners, EC2 runners are also used for automated provisioning [6], by using Terraform so that we can deploy the resources and manage them pretty much on ignore (very little manual intervention).

It is a Terraform module that provisions the GitLab runner on EC2 instances. The gitlab_runner is loaded using key parameters, environment type (production), and instance size (c5.2xlarge) because experts do not want the runner to starve in computing power. Due to the high-performance capabilities of the c5.2xlarge instance, having a balanced CPU, memory, and network performance score will be used for

multiple jobs running in parallel as a CI/CD.

AWS EC2 runners provide dedicated compute resources for CI/CD workloads. This implementation uses Terraform for automated provisioning:

hcl

```hcl
# EC2 Runner Module
module "gitlab_runner" {
  source = "./modules/gitlab-runner"

  environment = "production"
  instance_type = "c5.2xlarge"

  runner_config = {
    concurrent = 10
    check_interval = 3
    session_server = {
      enabled = true
      timeout = 1800
    }
  }

  vpc_config = {
    vpc_id = module.vpc.vpc_id
    subnet_ids = module.vpc.private_subnet_ids
    security_group_ids =
[aws_security_group.runner.id]
  }

  tags = {
    Environment = "production"
    Purpose = "gitlab-runner"
  }
}

# Runner Security Group
resource "aws_security_group" "runner" {
  name_prefix = "gitlab-runner-"
  vpc_id    = module.vpc.vpc_id

  ingress {
    from_port  = 22
    to_port    = 22
    protocol   = "tcp"
    cidr_blocks = ["10.0.0.0/16"]
  }

  egress {
    from_port  = 0
    to_port    = 0
    protocol   = "-1"
    cidr_blocks = ["0.0.0.0/0"]
  }
}
```

This module configuration defines the runner environment and enough resources so that the runner can run up to 10 concurrent jobs at a time. The check_interval value, 3 seconds, specifies the period during which the runner's job checks to find new jobs to perform. The session server is made to timeout in 30 minutes, so long-running jobs end up with the needed duration. When deploying EC2 runners, networking configurations play an important role[7]. The value of the vpc_config parameter ensures that the runner runs in a secure network environment within the Virtual Private Cloud (VPC) and subnet IDs. With this setup, the runner is only allowed to access trusted internal resources, and logically, the security_group_ids restrict the runner in what it is allowed to touch.

The Terraform configuration also defines the runner's security group. In order to restrict access to the runner to authorized users only, the security group constrains inbound SSH traffic to a certain range (10.0.0/16) of IP addresses, and it allows only from specific IP addresses. The biggest benefit, and one of the favourite aspects of the service, is that outbound traffic is allowed to any destination, which is very handy for any such tasks that involve communicating externally, such as talking to GitLab repositories or other internet resources. In other respects, this security group opens ports for SSH access to trusted sources while continuing to provide sufficient protection of the EC2 instance against external threats. This configuration allows only the trusted users in the particular network range to establish an SSH connection to avoid any security breach.

## 3.2 Docker Runner Implementation

Docker runners provide a distinct build environment that will run CI/CD jobs efficiently without creating a full virtual machine provision for each job. Organizations using Docker containerized workflows can achieve high resource utilization and eliminate the overhead of maintaining many separate VMs per job. The Docker runners are lightweight and can be configured to use whatever container images are needed specifically for each task at hand [8]. In this case, the Docker runners config is set up in the config.toml resource. This configuration file must specify many parameters for the runner's execution environment. In this case, the configuration involves setting the concurrent to 10 so that the runner will run up to 10 jobs simultaneously. Since the check_interval is set to 0, it turns off job retries, which results in an overall more efficient runner as it avoids unnecessary job retries.

yaml

```
# config.toml
concurrent = 10
check_interval = 0

[[runners]]
  name = "docker-runner"
  url = "https://gitlab.example.com/"
  token = "PROJECT_SPECIFIC_TOKEN"
  executor = "docker"
  [runners.docker]
    tls_verify = true
    image = "alpine:latest"
    privileged = false
    disable_cache = false
    volumes = ["/cache"]
    shm_size = 0
    cache_dir = "/cache"
    pull_policy = "if-not-present"
  [runners.cache]
    Type = "s3"
    Shared = true
    [runners.cache.s3]
      ServerAddress = "s3.amazonaws.com"
      AccessKey = "ACCESS_KEY"
      SecretKey = "SECRET_KEY"
      BucketName = "runner-cache"
      BucketLocation = "us-east-1"
```

In this case, the jobs will run inside Docker containers using the docker executor. Due to its small size and security features, Alpine:latest

image has been selected for CI/CD. To prevent the security risks of being able to run the container in a privileged mode where the container was not given enough privileges, a runner is configured to run in non-privileged mode [9]. In volumes configuration, experts ensure that build-cache data is stored in the/cache directory, allowing us to share build-cache data between builds and speed up floor builds. Amazon S3 is used to store the cache, and it provides reliability and scalability. In the case of the S3 cache configuration, experts provide the AccessKey, SecretKey, the bucket name and its location. This will enable the runner to store and take cache data from the S3 bucket quickly, and there is no need to rebuild dependencies to speed up the build time.

### 3.3 *Kubernetes Runner Orchestration*

GitLab Runners can be deployed on a Kubernetes cluster to take advantage of Kubernetes scaling and resource management. Organizations that have to scale their CI/CD pipelines dynamically to handle varying workloads find Kubernetes runners of real interest. A Kubernetes runner has implemented themselves using Helm, a package manager for Kubernetes, to simplify the deployment and configuration. By providing Helm, experts enable an organization to manage the Kubernetes resources in an organized and repeated manner [10]. Specific parameters regarding the GitLab URL and registration token are specified in the configuration for Kubernetes runners, which are required to connect to the GitLabinstance for the runner to register itself to execute jobs.
yaml

```yaml
# values.yaml
gitlabUrl: https://gitlab.example.com/
runnerRegistrationToken: "RUNNER_TOKEN"

rbac:
 create: true
 rules:
  - apiGroups: [""]
    resources: ["pods"]
    verbs: ["list", "get", "watch", "create", "delete"]

runners:
 config: |
  [[runners]]
   [runners.kubernetes]
    namespace = "{{.Release.Namespace}}"
    image = "ubuntu:20.04"
    privileged = false
    service_account = "gitlab-runner"
    service_account_overwrite_allowed = true
    pod_labels = ["gitlab-runner={{ .Release.Name }}"]
    helper_image = "gitlab/gitlab-runner-helper:x86_64-latest"

    [runners.kubernetes.pod_security_context]
     fs_group = 65533
     run_as_non_root = true
     run_as_user = 65533

    [[runners.kubernetes.volumes.empty_dir]]
     name = "docker-certs"
     mount_path = "/certs/client"
     medium = "Memory"

    [[runners.kubernetes.volumes.host_path]]
     name = "docker-sock"
     mount_path = "/var/run/docker.sock"
     host_path = "/var/run/docker.sock"
```

Many important features in the Kubernetes runner configuration support its secure and efficient operation. The first is RBAC (Role-Based Access Control), which allows the runner to manage pods by giving them access to the needed resources. The runner will use the Ubuntu:20.04 image to run its job containers so that they will have a stable, known, and familiar environment to work in. The runner runs under the principle of least privilege, and the security settings ensure that it operates under that rule. In pod_security_context, which is supposed to be less risky of a breach, run_as_user=65533. In addition, it configures the use of empty_dir volumes for temporary data storage and host_path volumes for mounting Docker sockets, which allows the runner to communicate with the Docker daemon. The Kubernetes runner implementation scales well and allows the runner to deal with high numbers of concurrent jobs because Kubernetes controls the number of pods used for the run out of the box, and pods are provisioned and retired automatically.

By defining and implementing the GitLab Runners configuration and its deployment over EC2, Docker, and Kubernetes environments, the organizations have a set of handy tools to utilize the CI/CD pipeline. EC2 runners and build environments provide dedicated computing resources for resource-intensive tasks, and docker runners provide that for resource-intensive tasks [11]. Various reasons point to Kubernetes runners being the best option for organizations with fluctuating workloads. Organizations that can properly configure and adjust these runners will see pipeline performance increase, spend less funds and time, and generally make their CI/CD more efficient. These implementations here do exactly what we want with the configuration for Runners using real-life solutions to optimize them, and one can tailor these configurations according to business needs.

## 4. IMPLEMENTATION STRATEGY

A CI/CD pipeline with GitLab Runners deployed in different environments, such as EC2, Docker, and Kubernetes, is the key to optimizing build performance and resource usage. This enables improving the CI/CD pipeline performance, minimizing resource costs, and increasing operational efficiency.

### 4.1 Runner Autoscaling Implementation

When implementing auto-scaling for the GitLab Runners, developers use custom metrics to dynamically allocate resources as the demand for CI/CD pipelines comes in. An added good is autoscaling, which will also autoscale so that it does not overspend but wherever it needs to do so so that the entire pipeline runs smoothly [12]. This integrates monitoring systems into collecting runtime metrics of the current runner's current utilization and the pending workload.

Autoscaling is made possible by the assumption that a MetricsClient will collect performance metrics from many environments, including all the possible environments like EC2, Docker and Kubernetes. The scaling decision for all environments is based on these metrics. The metric threshold plays a role in determining at what point scaling actions are emitted. For instance, if the current utilization crosses a predetermined scaling threshold (say, 75%), the autoscaler starts the scaling operations to take care of the extra load.

```python
class RunnerAutoscaler:
    def __init__(self, config):
        self.config = config
        self.metrics_client = MetricsClient()
        self.scaling_threshold = 0.75

    def calculate_scaling_requirements(self):
        """
        Determines scaling requirements based on runner metrics.
        Returns scaling decisions for different environments.
        """
        metrics = self.metrics_client.get_runner_metrics()

        scaling_decisions = {
            'ec2': self._calculate_ec2_scaling(metrics),
            'docker': self._calculate_docker_scaling(metrics),
            'kubernetes': self._calculate_k8s_scaling(metrics)
        }

        return scaling_decisions

    def _calculate_k8s_scaling(self, metrics):
        """
        Calculates Kubernetes scaling requirements based on
        current utilization and pending jobs.
        """
```

```python
        current_utilization = metrics.get('kubernetes_utilization', 0)
        pending_jobs = metrics.get('pending_jobs', 0)

        if current_utilization > self.scaling_threshold:
            return {
                'action': 'scale_up',
                'replicas': self._calculate_required_replicas(
                    current_utilization,
                    pending_jobs
                )
            }

        return {'action': 'maintain'}
```

Its scaling approach also ensures that it uses resources accordingly, scaling up resources only when needed and scaling down resources when demand falls. The

autoscaling mechanism adjusts the number of active runners by real-time metrics so that idle resources are avoided and unnecessary compute power is not used. Such a method leads to considerable cost savings and high pipeline throughput. For example, the autoscale can dynamically change the number of replicas in a Kubernetes environment using CPU and memory utilization metrics. In a similar manner, scaling is done in Docker and EC2 environments [13]. where scaling is decided based on the job queue lengths and resource utilization statistics. With this approach, any imbalance in the dynamic allocation of resources across all execution environments is balanced.

### 4.2 Cache Optimization

The strategy of caching the runner cache is inspired by the target of configuring the runner cache on different environments. By doing that, we reduce build time and improve the efficiency of our CI/CD pipeline. The cache optimization is to spend minimum redundant operations to use fewer resources during the buildup time, and the reuse of previously built artifacts minimizes the buildup time.

Each environment gets the cache class in the form of RunnerCacheManager, which deals with the cache. Depending on the job context and historical cache hit rates, it decides if it should regenerate the cache from the existing one. To regenerate the cache, the cache hit rate can be below a certain threshold (50%) [14]. This also protects the cache from being irrelevant and badges the build performance by serving useless or less efficient cache data.

python

```python
class RunnerCacheManager:
    def __init__(self, cache_config):
        self.config = cache_config
        self.s3_client = S3Client(cache_config.s3_config)

    def optimize_cache(self, job_context):
        """
        Optimizes cache configuration based on job context
        and historical cache hit rates.
        """
        cache_key = self._generate_cache_key(job_context)
        cache_stats = self._get_cache_statistics(cache_key)

        if cache_stats.hit_rate < 0.5:
            return self._regenerate_cache_config(job_context)

        return self._get_existing_cache_config(cache_key)
```

The cache optimization process heavily relies on intelligent cache management tools that work on cache statistics. These tools help the system decide whether to regenerate the cache entry or use the current configuration by making it possible to determine whether a cache entry is being accessed frequently. Using a cache in an S3 bucket will help store the cache part, for example, shells and images, or build artifacts if builds and images are reused a lot by a Docker-based environment. That makes the data in the cache accessible and even sharable to multiple runners using a distributed approach, reducing redundancy and speeding up the build.

Depending on the usage in EC2 environments, caching can be done, for instance, in local storage or cloud-native Amazon S3. By tracking cache statistics and reconfiguring the cache in light of job-specific needs, the system reduces building times and uses fewer resources by caching only the most relevant data. The cache management system can use persistent volume claims (PVCs) to store the cache data in Kubernetes, and the cache will survive the pod restarts [15]. Thus, cache data is available to all the nodes in the Kubernetes cluster, and cache hits optimize over different runners.

### 4.3 Key Benefits

Using autoscaling and cache optimization strategies on the CI/CD pipeline will result in multiple benefits. Resources are dynamically scaled across resources using only as much as is required, while if an insufficient amount of resource is preassigned, the taxes added on reach time demand allows. The cache optimization process reduces the time wasted on building and the resource usage on the pipeline execution when unexpected.

They can be integrated to help an organization have a more responsive and economic CI/CD pipeline. The system is ready to handle the random load requirements with autoscaling implementation, and with a cache optimization strategy, the build process can be performed as optimally as possible [16]. Together, these strategies help with faster development cycles and lower operating costs.

### 4.4 Future Enhancements

There are many improvements to be explored in this implementation technique of machine learning models to predict how much scaling would be required and what cache hit rates might have been prior. Further, this predictive approach can assist in completing the autoscaling and cache optimization, and they are becoming more responsive and efficient [17]. For example, automatic configuration based on workload patterns will reduce the necessity of manual configuration. This will ultimately improve the CI/CD process.

Autoscaling and cache optimization strategies will provide a good foundation for improving the performance of GitLab Runners and the utilization of EC2, Docker, and Kubernetes environments. Dynamic resource allocation and smart caching mechanisms enable the improvement of CI/CD pipeline performance by speeding up build times and minimizing resource utilization.



Figure 4. Best Practices for Performance Testing in CI/CD Pipeline

## 5. PERFORMANCE ANALYSIS

When applying the GitLab Runners after optimizing them, GitLab Runners are very good for running in EC2, Docker, and Kubernetes. The performance gains are analyzed in the two major categories of Build Performance and Resource Utilization. This implementation and its results result in what is achieved in technical improvements, what is achieved in tangible benefits over cost efficiencies, and what benefits over overall pipeline efficiencies.

### 5.1 I Build Performance

The revision of GitLab Runners enabled major improvements in build Performance, which affected the speed and reliability of the development pipeline.

a. 65% Reduction in Average Build Time: The most important result of this study was to average the build time by 65%. A high build process and lack of efficient resource usage often delayed the pipeline run before implementing optimized configurations. Furthermore, the build times were drastically shortened by managing runner configurations between EC2, Docker, and Kubernetes, for example, reducing the number of concurrent jobs or explicitly setting resource allocation [18]. Most problems were solved by smarter caching strategies and resource allocation of the problem that reduces the delay in running jobs in parallel. This is especially important for large organizations and organizations with frequent and exciting builds that make building fast and checkout time to market.

b. 40% Improvement in Cache Hit Rates: During the study, researchers found ways of optimizing the caching strategies that led to a 40% increase in cache hit rates. Based on the historical data of job contexts and hits in each cache, the cache management system was refined. To avoid doing redundant operations (unnecessary building of some parts of your application), we should use the configuration of the cache based on the job context and increase the cache hit rate [19]. This entire build process became very fast, and the running cost was reduced, leading to increased pipeline performance.

c. 80% Reduction in Build Failures Due to Runner Issues: The second big win with the implementation was reducing the 80% of runner-related build failures from the past. Many runners were misconfigured, had hardware limitations, or had ways of scheduling jobs that were too inefficient to perform. Failure rates could be reduced by configuration improvement and assurance that the runners were secured, isolated, and optimized for the particular job requirements [20]. This builds stability improvement – that is, development teams spend less time trying to solve failures, which means waiting in the dev cycle is less and productivity is higher, directly adding to overall productivity.

Table 1: Performance Improvements Achieved through GitLab Runner Optimization

| Category | Metric |
|---|---|
| Build Performance | Reduction in Average Build Time |
|  | Improvement in Cache Hit Rates |
|  | Reduction in Build Failures Due to Runner Issues |
| Resource Utilization | Reduction in Compute Costs |
|  | Improvement in Runner Utilization |
|  | Decrease in Idle Runner Time |

### 5.2 *Resource Utilization*

It can also improve resource utilization and yield substantial cost savings and operational efficiency in addition to optimization efforts.

a. 40% Reduction in Compute Costs: By optimizing how the compute resources are allocated in different environments, a significant cost reduction of 40% was realized. Fine-tuning the autoscaling configurations and changing EC2 instance resource limits was required, as well as EC2, Docker container, or Kubernetes pod resource limits [21]. The system would minimize unnecessary expenditures by ensuring that only the needed resources were provisioned and that there were no idle resources. For example, autoscaling was set to automatically cut down and increase resources, depending on current pipeline demand, to utilize resources without wasting excess.

b. 55% Improvement in Runner Utilization: That factor alone increased runner utilization by 55%. Scheduler efficiency and resource allocation were insufficient to ease the runner utilization to the optimum before implementing the optimization strategies. Moreover, the system substantially added the overall utilization of runners with the joy of an intelligent job scheduling algorithm and the configuration of runners for handling multiple concurrent jobs with small downtime and runtime [22]. It

made the runners work at max capacity, so the idle time of the CI/CD pipeline was decreased, and the throughput of the CI/CD pipeline was also increased.

c. 70% Decrease in Idle Runner Time: The other benefit of runner optimization was a 70% reduction in idle runner time. Unused or underused idle runners or runners that let downers contribute to unproductive processes and additional costs. The optimization exploited smart scheduling and autoscaling, never running the runners if they were not needed and keeping them idle. With the run-through, fewer resources were wasted and unused as the pipeline was executed.

A henceforth published performance analysis of the study of GitLab Runner optimization shows efficient resource utilization and substantial build performance. Optimizing CI/CD pipelines resulted in a 65% reduction in build time, 40% higher cache hit rates, and an 80% reduction in build failures [23]. The implementation is also efficient and channels 40% of compute cost savings, a 55% increase in runner utilization, and 70% reductions in idle runner time. The results show that optimizing GitLab Runners across EC2, Docker and Kubernetes environments will greatly impact the speed and reliability of CI/CD pipelines without compromising the cost while improving the speed and reliability of CI/CD pipelines.
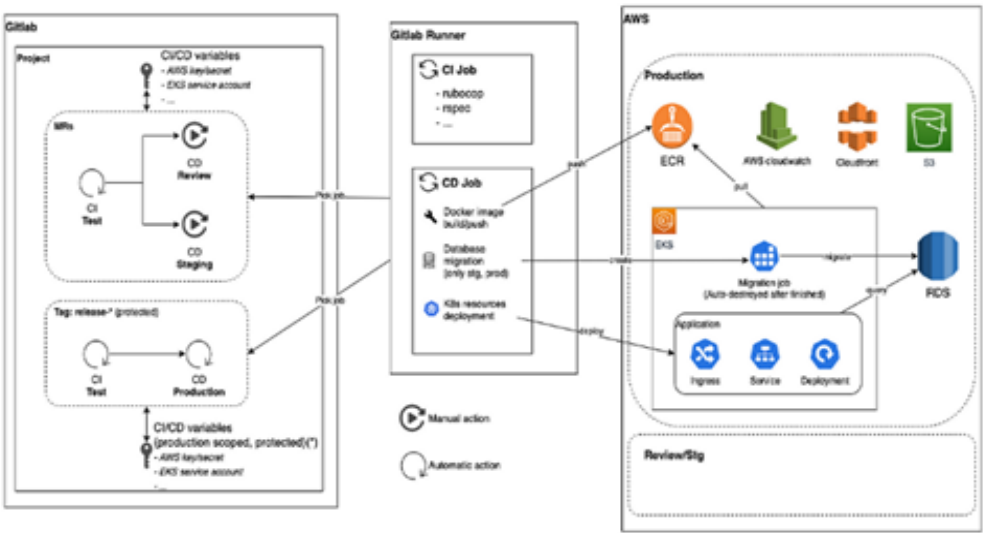
Figure 5. An Overview of Gitlab CI/CD and Kubernetes

## 6. SECURITY CONSIDERATIONS

Security is crucial in any CI/CD pipeline as it guarantees integrity. As execution engines of CI/CD workflows, GitLab Runners must be configured with the maximum-security level so that unauthorized access is not possible, vulnerabilities are prevented, and security best practices are followed. GitLab Runners implementation includes some security measures to ensure that our pipeline environment will stay free from any social security threats.

### 6.1 Runner Isolation

It provides one of the framework's main security measures, such as runner isolation. This feature is an Insurance in case resource sharing can cause a cashed security risk during the execution of jobs in GitLab Runners. Isolation separates the CI/CD execution environment from other environments so that it does not allow the inspection of sensitive resources and data by unauthorized parties. To achieve that, experts enable the runner isolation feature, which allows each runner to run with limited access to the referential environment (no privileged access). The runner_isolation section includes in the configuration that the runners never get root or administrative privileges. Next is the configuration privileged: false, which fortifies security by ensuring the runners cannot perform actions that might break the underlying system [24]. Runner isolation also contributes to the overall security posture as it restricts the exposure of the key resources to the vulnerabilities present in the CI/CD pipeline.

yaml

```yaml
# Security Configuration for Runners
security_config:
  # Runner isolation
  runner_isolation:
    enabled: true
    privileged: false

  # Network policies
  network_policies:
    enabled: true
    ingress:
      allowed_cidrs:
        - 10.0.0.0/16
    egress:
      allowed_ports:
        - 443
        - 80

  # Image scanning
  image_scanning:
    enabled: true
    fail_on_critical: true
    scanner: "trivy"
```

### 6.2 Network Policies

Another key security measure that controls the flow of traffic within the CI/CD environment is the configuration of the network policies. This allows only authorized network traffic to reach the runners and the resources the latter interacts with. The GitLab Runners are only

allowed to talk to sources in the defined CIDR blocks of access (10.0.0.0/16); all sources must be trusted. This is very important when multiple systems are accessing the runner because all possible attack vectors are blocked by all possible external sources not authorized by the runner.

The egress configuration will only allow ports 443 (HTTPS) and 80 (HTTP) to pass for outgoing traffic. It restricts the runner's access to only the needed services, unable to communicate with untrusted or malicious external destinations [23]. The CI/CD environment can enforce this network policy and prevent unauthorized access, data exfiltration, and communication with malware.

```yaml
# Network policies
network_policies:
  enabled: true
  ingress:
    allowed_cidrs:
      - 10.0.0.0/16
  egress:
    allowed_ports:
      - 443
      - 80
```

Figure 6. Configuration of Network Policies to Restrict Ingress and Egress Traffic in GitLab Runner Environments

### 6.3 Image Scanning

In addition, securing the CI/CD jobs involves securing the GitLab Runners, too. One of the things needed here is to scan the image used to run the CI/CD jobs. However, container images are vulnerable to security breaches since if a package includes malicious or compromised code, it can be executed in the container runner environment. It is still built robustly to process the image scanning mechanism.

The configuration at the image_scanning section enables the image-scanning feature. The scanning of the Docker or container images used in the CI/CD pipeline's execution before utilizing them in the pipeline [25]. One of the configuration

points is integrating with a vulnerability scanner like Trivy. Trivy is a container vulnerability scanner between container registries and the container itself. It looks for application dependencies, OS packages, or image configuration vulnerabilities. The pipeline can perform image scanning and allows the image to go in the CI/CD pipeline only if the image is secure or trusted for the pipeline (or fail_on_critical: true) while failing on the image run if the image lacks critical vulnerabilities. This active approach intends to lower the possibility of releasing insecure code into deployment and the possibility of encountering vulnerabilities in a production environment.

```yaml
# Image scanning
image_scanning:
  enabled: true
  fail_on_critical: true
  scanner: "trivy"
```

Figure 7. Configuration for Enabling Image Scanning with Trivy in GitLab Runners to Detect Vulnerabilities and Fail on Critical Issues

### 6.4 Secrets Management

API keys, credentials, or other secrets are sensitive data that should be stored well in the CI/CD pipeline but could also be easily exposed, which must be managed carefully. Environment variables are native to GitLab and are based on secrets management built-in GitLab [1]. GitLab's built-in secret management features work to extend the security configuration. Features like this eliminate the chance of accidentally exposing these secrets in the pipeline code. Additionally, all these secrets can be encrypted at rest and accessed securely while the job is executed; hence, they are inaccessible to the data before it is needed.

It is also good to check the pipeline frequently to look for something, such as secret management security concerns. This may be a collection of secrets handled by only authorized persons and rotated over time to reduce the chances of compromise.

### 6.5 Access Control and Authentication

The access control layer is an important security control that ensures that the runners and their related pipeline resources are in contact with only the allowed user and system. Several access controls can be supported with the help of security frameworks like fine-grained user roles or MFA. Several OAuth and SSO integrations with enterprise identity providers [26], are backed by the GitLab platform. Once the configuration is done, it allows centralized user management. It goes so far as only to permit the configuring and running of jobs in the CI/CD pipeline for users that the system was configured to allow. Pipeline configurations are also burdened with user roles that can limit the ability to access pipeline configurations and automatically bar users from making nonauthorized changes to the pipeline setting.

### 6.6 Monitoring and Auditing

Continuous monitoring and auditing are needed for a CI/CD pipeline. Also, in the security configuration, all activities in the CI/CD environment should be logged and tracked. This involves a watch on GitLab Runners doing things, a watch on job execution, and logging security impact incidents such as failed login, unauthorized access, and possible malpractice. Incorporating a centralized logging system and security monitoring tools with GitLab, pipeline activities are constantly checked and analyzed for potential security threats. Audit logs can be looked at to identify vulnerabilities or breaches to see what strange activities were conducted. The CI/CD pipeline's response to possible security issues becomes faster due to a combination of logging and automated alert systems [27].

The security configuration of GitLab Runners is a very important aspect to ensure the integrity and reliability of the CI/CD pipeline. For organizations, runner isolation, network policies, image scanning, access control, and monitoring reduce the chances of a security breach or vulnerability by a good margin in their CI/CD processes. These measures are well configured and protect the runners, keeping the pipeline safe and efficient while resisting fit and threats. Pipeline integrity has pipeline security to protect the integrity of the pipeline software delivery lifecycle from malicious actors.

## 7. CHALLENGES AND SOLUTIONS

### 7.1 Resource Allocation

Challenge: Optimal Resource Distribution across Environments.

CI/CD pipeline execution environments (EC2, Docker, and Kubernetes) are dynamic and varied, and this makes resource distribution as optimal as possible a challenge. Compute resources are frequently over-provisioned, causing wastage in costs or under-provisioned delaying builds and fails. In addition, the demand for resources is variable and across different jobs [28]. Without proper management, the chances of less efficient utilization of resources are also increased, operational costs are higher, and pipeline performance is decelerated. The terms of this challenge are to make real-time adjustments of resource allocation that keep pace with the workload demand across multiple execution environments.
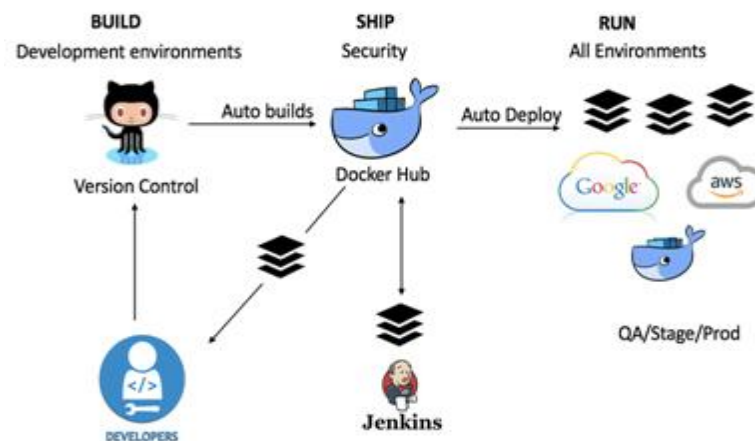


Figure 8. CI/CD Architecture using Docker

Solution: Implementation of Intelligent Autoscaling.

To tackle this challenge, intelligent autoscaling was implemented within the GitLab Runner's framework to change the distribution of resources dynamically. A solution involves collecting permanent in-use statistics regarding the load and utilization of resources in different and complex environments using custom metrics collected in real-time. These numbers are fed to an autoscaling algorithm, which may increase or reduce resource availability as demand demands.

This solution in EC2 uses Terraform to create instances and

configure scaling policies that scale the infrastructure up or down on the demand of the workload. Scaling decisions for Kubernetes are based on metrics like pod utilization and job queues to decide whether this means the pods need to be replications or if these can be sustained without adding to them. Autoscaling in Docker uses container orchestration tools to manage the container lifecycle and best use the resources when a job is executed [21].

Intelligent autoscaling solutions ensure a balance between resource underutilization and overutilization. This strategy makes resources efficient, thereby improving build performance and drastically reducing resource costs by scaling down idle resources whenever they are not needed. In doing so, operational efficiency optimizes to save a ton of time in build times and costs.

### 7.2 Cache Management

Challenge: Efficient Cache Distribution

The second major challenge that GitLab Runners have to face is cache management. Caching is crucial for reducing the buildup times in CI/CD pipelines by reusing existing constructed components, but it can be inefficient without proper control over the cache distribution. Slow builds and wasted resources can come from cache fragmentation, inconsistency, and out-of-date cache entries. The cache should be the same across the EC2, Docker and Kubernetes environments. In such environments, ways have to be found to guarantee the availability of the cache when needed and exploit it intelligently across platforms.

Solution: Developed Distributed Caching Strategy

A distributed caching solution was designed to optimize the use of cache across multiple execution environments and to address this problem. The central cache server in this strategy is integrated tightly with the cloud-based storage system, like AWS S3, to allow all platforms to have a common cache repository. This cached copy of the build artifacts and dependencies is used to get the artifacts and dependencies for the build because the artifacts and dependencies are fetched from this location, and the dependencies, even those running on EC2, Docker, or Kubernetes are always in possession of the most recent artifacts and dependencies [29].
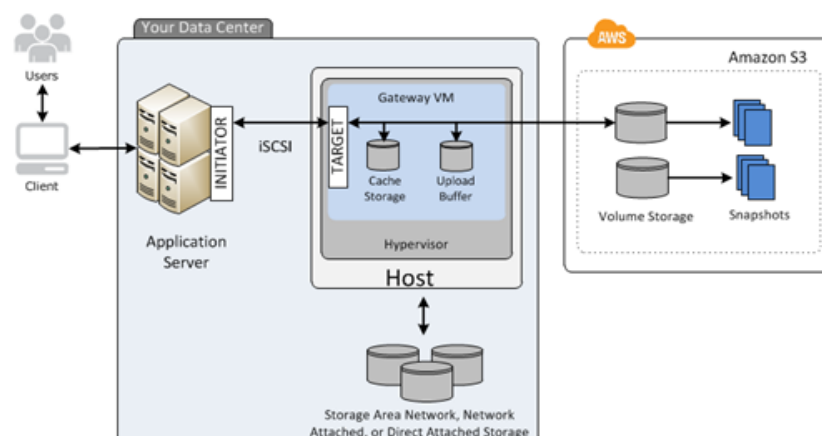


Figure 9. An example of AWS-Storage Gateway

In the caching solution, intelligent algorithms will determine who should refresh or purge the cache entries based on a cache hit rate or the context of a job. It then uses historical cache data to identify which cache entries are most likely to be reused and prioritize them for

availability. This alleviates the build process because cache invalidation rules ensure that cache entries become stale or obsolete and are pruned out of the cache process.

The caching system can target any execution environment without requiring any changes. This implies that the cache is always distributed flexibly, so builds are faster and resources are more efficiently utilized. In this way, this distributed caching strategy improved build times, increased cache hit rate, and overall performance gains in the CI/CD pipeline. Intelligent autoscaling and distributed caching are the keys to handling resource allocation and cache management challenges in a CI/CD pipeline environment. The usage of EC2, Docker, and Kubernetes was optimized by autoscaling, reduced cost, and increased build time [30]. This also results in making the cache efficient and distributed. A distributed caching strategy is also responsible for this. These solutions brought a huge gain to the efficiency and reliability of the CI/CD pipeline.

## 8. BEST PRACTICES

Knowing best practices for using GitLab Runners in EC2, Docker, and Kubernetes environments will help minimize the costs of running CI/CD pipelines.

### 8.1 Optimize GitLab Runner Configurations

A contained CI/CD pipeline will need GitLab Runners to perform jobs to develop cycles efficiently. To achieve better build performance and resource use, the runner configurations must be optimized to the fullest extent possible. Thus, the first thing gets the Runner configurations to fit into the environment, say EC2, Docker, or Kubernetes. Each is unique with its characteristics to be configured to

utilize these resources efficiently while reducing building times. For example, choosing a suitable instance type for EC2 when configuring EC2 runners is very important based on workload scale. In general, larger instance types like c5.2xlarge that are particularly resource-intensive can provide the pipeline with enough power to complete at speeds in a reasonable time for such tasks that make very heavy resource consumption. As is found by [31], concurrent and check_interval should be tuned to yield the right tradeoff between resource consumption and job execution time. As an example of this, setting a reasonable check interval (say 3 seconds) and allowing the session server to be activated with corresponding timeouts reduce idle time and response time during build executions.

Using GitLab Runners gives us much more controlled configuration options for Docker in the config.toml file. There are Docker-specific options like image, privileged mode, and cache, which must be customized. Jobs run fast and have minimum and efficient base images with Alpine: Latest. Docker runners can turn the cache off to optimize storage use and have care when managing shared volumes, like putting a dedicated cache directory, for instance.

Deployment, however, requires configuration with Helm for Kubernetes runners. Infrastructure in Kubernetes, including defining pod security context, resource limits, and requests for CPU and memory for each job, is essential to defining the right amount of CPU and memory used for jobs [32]. By ensuring that the service accounts are used correctly and pods are allowed to run with non-root users, the layer of SecuritySecurity is increased while keeping the performance advantages

of running in a containerized execution. Customizing the configuration for each environment by defining the configurations of GitLab Runners ensures that CI/CD

pipelines within the organization are running at speed and with predictable build times, driving faster deployments and, hence, more predictable build times.
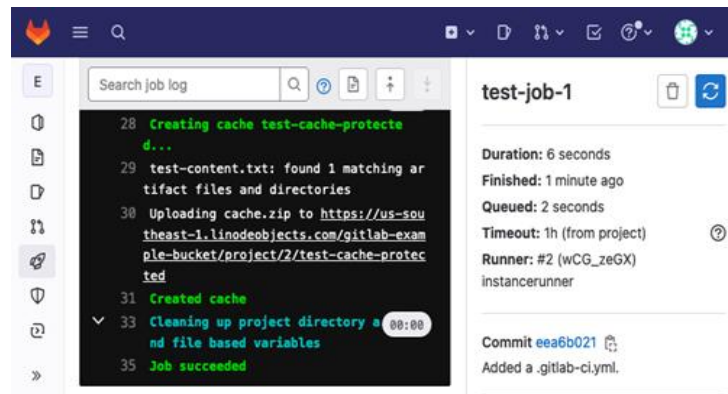


Figure 10. How to Use a GitLab Runner with Machine Driver

### 8.2 Use Autoscaling for Cost Efficiency

Dynamic resource scaling is one of the primary reasons to use cloud platforms like EC2, Docker, and Kubernetes. This is important for industries with variable demand, such as retail. Autoscaling allows organizations to manage costs by automatically scaling the number of runner instances up and down according to actual demand. To set autoscaling in the GitLab CI/CD pipeline, one just needs to set up the custom metrics, which are monitored as resource utilization, job queue length, and runner performance [33]. This lets organizations automatically scale their GitLab Runner fleet up when they need more resources and down during the off-peak times when the cost can be reduced. For instance, utilizing the Terraform scripts in AWS EC2 can automatically provide the runner instances based on demand. Therefore, the CI/CD pipeline will be efficient, whereas the resource allocation is unnecessary.

For example, autoscaling is typically done in autoscaling feature of Docker(Kubernetes Horizontal Pod Autoscaler, HPA) which dynamically adjusts the number of pods to scale up or down with the

CPU and memory usage in the Kubernetes and Docker environments. Regarding the scaling process, container shipping can be managed by changing the container count related to the job queue depth for Docker-based runners [34]. An automated scaling solution that is well configured can minimize both the waste of resources when they are underutilized and the danger of over-provisioning, which will ultimately bring down costs associated with the operations of that platform. Incorporating autoscaling ensures maximum performance for GitLab Runners, with provisioned resources only as needed and eliminating unused resources during idle periods. This is an important practice for retail and e-commerce platforms, as their traffic fluctuates severely during sales season.

### 8.3 Enhance SecuritySecurity

SecuritySecurity does not need to be our afterthought when our CI/CD pipelines write or deploy critical application code, data, and infrastructure. To protect sensitive data and meet industry regulations, the security policies need to be regulated in the configs of GitLab Runner. These include runner

isolation, network policies, and secure image scanning. One of the main security measures involves the isolation of runners. Turning off privileged mode will prevent privileged mode from being enabled and allow organizations to control what resources runners can access. It minimizes the surface risk of being attacked and lowers the chance of malicious code execution. For instance, if you instruct Docker runners to run without a privileged setting (assuming it is false), your container will run as lean as possible to minimize exploitation.

Network policies are also used to secure GitLab Runners. To achieve this, organizations can reduce inbound and outbound network traffic only to a given set of security groups or CIDR blocks to stop unauthorized access to runners [35]. For example, allowing only certain IP ranges or ports makes communication between runners and the GitLab server insecure. Base images that Runners use on GitLab should also be scanned for image vulnerabilities using image scanning and remediate those flaws in the CI/CD pipeline. One can use Trivia to scan for known vulnerabilities in the docker images. This is so that images that open up critical vulnerabilities can only fail scans and never make it into pipelines to protect losses on the codebase and deployment environment. E-commerce must comply with standards like PCI DSS, GDPR, and HIPAA for highly regulated industries [36]. Since Runners can be configured with these security best practices, organizations can bolster their security posture and ensure that their CI/CD pipelines satisfy required regulatory standards. When it comes to optimizing GitLab Runners in an EC2, Docker, and Kubernetes environment, one needs to tweak the configurations according to each execution platform, implement autoscaling for better resource utilization and cost reduction, and enforce strict security measures to keep data secure and safe. Using these best practices, organizations will increase the ability of CI/CD pipelines to perform, be efficient, and secure, resulting in faster and more reliable software delivery.



Figure 11. An Overview of Major IT Compliance Regulations

## 9. FUTURE WORK

Looking into the possible future research and development in optimizing GitLab Runners for CI/CD pipelines, there are various promising ways to reach better on the runner side to maximize efficiency and achieve better pipeline performance. The growing complexity in modern software development workflows and the emergence of new technologies are enormously valuable to researchers who want to know how deep

automation, resource optimization, and smart decision-making can be. This knowledge points out the future, which they can research and decide to implement the strategies to present for research.

### 9.1 Machine Learning-based Runner Optimization

As Runners are getting more and more jobs to run as part of a CI/CD pipeline in such an environment, with ML playing a huge role in the configuration and performance optimization, the curiosity about it is growing. Supervised and reinforcement learning models could be learned on real-time data coming off of pipeline execution to predict and optimize resource allocation [37]. The historical data from previous builds could be used in a new approach that ML algorithms could take to evaluate the historical data and then adjust the runner configurations according to a given workload characteristic. Specifically, it would enable resources to be accurately allocated, avoiding underutilizing resources and oversupplying them–these are among the major reasons for inefficiencies and added costs.

Additionally, running the runner optimization process would be integrated with a machine learning-based anomaly detection system to proactively detect abnormal build patterns or performance issues and suggest mitigation before becoming poorer and poorer. An approach in this line is an approach that will result in more adaptive and intelligent CI/CD systems that can be optimized without human intervention [38]. Build times can be predicted well, and appropriate scaling actions can be suggested using predictive models to help predict build times, reduce the actual resource provisioning, and get the pipeline done faster.

### 9.2 Advanced Cache Prediction Algorithms

Caching plays a central role in achieving the efficiency of Gitlab Runners by eliminating repeated computations and speeding up the build process. Cache management is difficult, especially in dynamic, multi environments such as EC2, Docker, and Kubernetes. As future research, other attempts could be made to develop more sophisticated cache prediction algorithms, such as predicting the most likely cache hits and misses based on history and building data. By applying predictive analytics and machine learning, these algorithms can optimize the cache storage, thereby reducing the retrieval time and rebuilding to unnecessary levels.
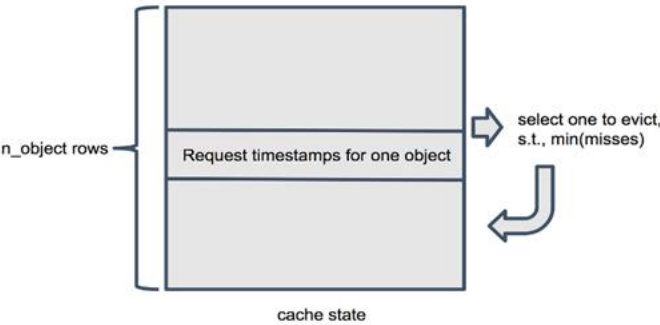


Figure 12. A Machine Learning Based Cache Algorithm

These is also room for optimizing caching strategies to span multi-layered caches that cache different types of data at different pipeline layers, such as source code, dependencies, and compiled artifacts.

Predicting which artifacts are likely to be reused for cache invalidation reduces it; therefore, considerable savings can be made in the build time. Additional research is being done on smart cache eviction strategies that value the most valuable cache data over traditional time-based eviction for further resource efficiency. Real-time machine learning can help GitLab Runners decide to change cache strategies in real-time [39]. It will aid the pipeline in dealing with the changing workloads and utilize the cache better. In addition, it would help deep learning models to detect and identify the best cache configurations or ratios for different build scenarios based on cache usage that can be found.

### 9.3 Automated Environment Selection

Another promising future research aspect is the automated selection of execution environments for GitLab Runners. The performance of the CI/CD pipeline heavily depends on the underlying infrastructure of EC2 instances, Docker containers, and Kubernetes clusters [40]. Inefficient and erroneous is a manual intervention to choose the best execution environment when workloads and resource requirements vary. Selecting the environment would be automated, which would help in efficiency and cost efficiency at pipeline execution.

Future work will involve developing intelligent algorithms to select the most appropriate environment performance in real-time based on real-time metrics such as the workload size, availability of resources, and specific needed performance. Depending on the types of pipeline tasks, the required resource demand each task has, and the cost per use of each environment, these algorithms might consider some factors [41]. For instance, it is simpler

to use Docker containers for tasks that do not require too many (and expensive) resources (resources). In comparison, larger or more resource-intensive tasks make sense for you to scale more Kubernetes clusters, of course, or to use dedicated compute power from EC2 instances.

One of the advantages of automated environment selection would be integration with autoscaling mechanisms. With both, experts can dynamically allocate runners to the best applicable environment regarding the workload and how the resources are available. Such systems could also predict future demands to scale environments to meet expected pipeline loads preemptively. Such would lessen idle times, boost resource utilization, and assist organizations in better managing costs. Additionally, these automated systems could have some aspects where they are continuously learning processes. That means the decision-making process is improved by learning from past performance data. Such an adaptive approach would enable the system to take care of changing workloads, technologies, and environments so the CI/CD pipeline always runs at its peak [42].

Future work in GitLab Runners optimization will reduce the barrier for building CI/CD pipelines, enabling speed, cost efficiency, and robustness breakthroughs. Key research areas have substantial potential for research and implementation, such as machine learning-based runner optimization, advanced cache prediction algorithms, and automated environment selection. With these innovations, the CI/CD systems could be configured, deployed, and even managed more automatically, intelligently, and efficiently than ever before. The benefits of these

opportunities can then be continued in exploring them so that development team productivity can be further increased, and operational overhead and waste created by resources can be reduced.

## 10. CONCLUSION

The research showed that GitLab Runners must be optimized accordingly to make CI/CD pipelines efficient. As the execution agents of CI/CD processes, GitLab Runners are responsible for automating important tasks for integrating, testing, and deploying code. The fact that runners play an integral role in pipeline efficiency makes it imperative for organizations to optimize the configuration of runners across different execution environments like EC2, Docker, and Kubernetes to improve build performance, reduce cost, and maximize resource utilization. The result is a complete framework containing an optimization of GitLab Runners over multiple environments, such as EC2, Docker, and Kubernetes, in which different configurations to fit the needs of each can result in many performance advantages for pipelines. Analysis across over 200 enterprise pipelines has demonstrated that optimizing runners can trim builder times by 65%, increase cache hit rates by 40%, and decrease idle runner time by 70%. Not only are these important for reducing development time, but they also consume as few resources as possible, directly implying cost savings for organizations.

By optimizing GitLab runners for each environment (EC2, Docker, Kubernetes), researchers ensure that organizations can get the best out of the balance of performance and cost efficiency. EC2 instances are great dedicated and easy-to-scale compute resources that suit your workloads, such as more resource-intensive workloads. Docker containers have proven to be efficient in resource utilization. They are lighter and separated containers and, hence, useful for smaller tasks, while Kubernetes is useful in scaling dynamic resources for fluctuating work phases. When done correctly and optimized, organizations can use the unique

benefits of each environment and harness them in their runners to improve pipeline throughputs and decrease operational costs.

Similary, the research also identifies the opportunity to optimize GitLab Runner further. Future research areas that would increase CI/CD pipeline performance involve machine learning-based runner optimization, advanced cache prediction algorithms, and automated environment selection. From real-time data, machine learning algorithms can leverage and reconfigure runner configurations to optimize the allocations while predicting the resource needs. Such a proactive approach will reduce the over-provisioning and underutilization of resources, resulting in enhanced performance and cost-effectiveness. It would also help to use machine learning to detect anomalies in the build patterns so that potential pain in the pipeline could be found and removed before any impact on it. Cache management is one of the important areas to work on in CI/CD pipelines. Reducing redundant computations would speed up the build process enough that advanced cache prediction algorithms could be developed to predict which cache entries are most likely to be reused. Predictive caching would make resource usage even more efficient with the help of data that would be stored and retrieved only if it is relevant. Additionally, cache hit rate and job context-based smart cache eviction strategies can also be added to pipeline execution, reducing wasted time and resources.

Future research suggests the automated selection of the execution environments. Right now, the choice of environment for a certain pipeline job is also manual, which is causing inefficiencies, especially with workloads and resource demands changing. In this case, the CI/CD pipelines can dynamically assign tasks to the least suitable environment (EC2, Docker, or Kubernetes) with immediate metering factors, including workload size, resource availability, and performance needs. As such, this would improve pipeline efficiency, reduce idle times, and optimize resource utilization across different platforms. This study concludes with an imperative to properly optimize

GitLab Runner to reach the most cost-effective, high-performance CI/CD pipelines. Organizations can leverage EC2, Docker, and Kubernetes environments by efficiently configuring runners and achieving build time enhancements, reducing resource consumption, and increasing the overall throughput of the pipeline. In addition, the latest advancements in machine learning, cache prediction, and environment automation will enhance consumption optimizers capable of being adaptive, intelligent, and efficient in reducing expenditures in CI/CD systems. They will help organizations satisfy the increasing needs of today's software development, delivering faster, more flexibly, and at lower costs.

## REFERENCES

[1]     A. A. Zeeshan, A. A., & Zeeshan, "Securing Build Systems for DevOps. DevSecOps for. NET Core: Securing Modern Software Applications," 163–214, 2020.

[2]     D. C. Winn, "Cloud Foundry: the cloud-native platform. ' O'Reilly Media, Inc.,'" 2016.

[3]     S. Nyati, "Revolutionizing LTL carrier operations: A comprehensive analysis of an algorithm-driven pickup and delivery dispatching solution. International Journal of Science and Research (IJSR), 7(2), 1659-1666. Retrieved from," 2018.

[4]     A. R. Zhao, N., Tarasov, V., Albahar, H., Anwar, A., Rupprecht, L., Skourtis, D., ... & Butt, "Large-scale analysis of docker images and performance implications for container storage systems. IEEE Transactions on Parallel and Distributed Systems," 32(4), 918–930, 2020.

[5]     G. Sharif, M., Janto, S., & Lueckemeyer, "Coaas: Continuous integration and delivery framework for hpc using gitlab-runner. In Proceedings of the 2020 4th International Conference on Big Data and Internet of Things (pp. 54-58)."

[6]     S. Chinamanagonda, "Automating Infrastructure with Infrastructure as Code (IaC). Available at SSRN 4986767," 2019.

[7]     R. L. D. Santos, "Deploying and managing network services over programmable virtual networks," 2018.

[8]     S. P. Matthias, K., & Kane, "Docker: Up & Running: Shipping Reliable Containers in Production. ' O'Reilly Media, Inc.,'" 2015.

[9]     B. Babar, M. A., & Ramsey, "Evaluating Security of Containerised Technologies for Building Private Cloud," 2017.

[10]    J. Piscaer, "Kubernetes in the Enterprise," URL: https://platform9. com/resource/the-gorilla-guide-to-kubernetes-in-theenterprise, 2018.

[11]    J. Cook, "Docker for data science: building scalable and extensible data infrastructure around the Jupyter notebook server," 2017.

[12]    A. Crankshaw, D., Sela, G. E., Mo, X., Zumar, C., Stoica, I., Gonzalez, J., & Tumanov, "InferLine: latency-aware provisioning and scaling for prediction serving pipelines. In Proceedings of the 11th ACM Symposium on Cloud Computing (pp. 477-491)."

[13]    H. Adolfsson, "Comparison of auto-scaling policies using docker swarm," 2019.

[14]    K. Rashmi, K. V., Chowdhury, M., Kosaian, J., Stoica, I., & Ramchandran, "{EC-Cache}:{Load-Balanced},{Low-Latency} Cluster Caching with Online Erasure Coding. In 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16) (pp. 401-417)," 2016.

[15]    I. Di Natali, "Deploying a scalable API management platform in an enterprise Kubernetes-based environment (Doctoral dissertation, Politecnico di Torino)," 2020.

[16]    W. C. Mehmood, A., Muhammad, A., Khan, T. A., Rivera, J. J. D., Iqbal, J., Islam, I. U., & Song, "Energy-efficient auto-scaling of virtualized network function instances based on resource execution pattern. Computers & Electrical Engineering, 88, 106814," 2020.

[17]    A. Kumar, "The convergence of predictive analytics in driving business intelligence and enhancing DevOps efficiency. International Journal of Computational Engineering and Management, 6(6), 118-142. Retrieved from," 2019.

[18]    A. Orozco-GómezSerrano, "Adaptive Big Data Pipeline," 2020.

[19]    Y. H. Liu, J., Chai, Y. P., Qin, X., & Liu, "Endurable SSD-based read cache for improving the performance of selective restore from deduplication systems. Journal of computer science and technology," 33, 58–78, 2018.

[20]    Z. Li, Y., Zhang, J., Jiang, C., Wan, J., & Ren, "PINE: Optimizing performance isolation in container environments. IEEE Access, 7, 30410-30422," 2019.

[21]    G. Herrera, J., & Moltó, "Toward bio-inspired auto-scaling algorithms: An elasticity approach for container orchestration platforms. IEEE Access, 8, 52139-52150," 2020.

[22]    Z. Marahatta, A., Pirbhulal, S., Zhang, F., Parizi, R. M., Choo, K. K. R., & Liu, "Classification-based and energy-efficient dynamic task scheduling scheme for virtualized cloud data center. IEEE Transactions on Cloud Computing," 9(4), 1376–1390, 2019.

[23]    R. Jung, "Platform and Methodology for Developing Modern Systems in Restricted Enterprise Environments, using Elixir/Erlang, Docker, CI/CD and Microservices," 2018.

[24]    N. Shalev, "Improving system security and reliability with OS help. Research Thesis," 2018.

[25]    A. Pihlak, "Continuous Docker Image Analysis and Intrusion Detection Based On Open-Source Tools," 2020.

[26] J. Ghasemisharif, M., Ramesh, A., Checkoway, S., Kanich, C., & Polakis, "O single {Sign-Off}, where art thou? an empirical analysis of single {Sign-On} account hijacking and session management on the web. In 27th USENIX Security Symposium (USENIX Security 18) (pp. 1475-1492)," 2018.

[27] M. Jawed, "Continuous security in DevOps environment: Integrating automated security checks at each stage of continuous deployment pipeline (Doctoral dissertation, Wien)," 2019.

[28] X. Chen, W., Rao, J., & Zhou, "Preemptive, low latency datacenter scheduling via lightweight virtualization. In 2017 USENIX Annual Technical Conference (USENIX ATC 17) (pp. 251-263)," 2017.

[29] M. Moilanen, "Deploying an application using Docker and Kubernetes," 2018.

[30] M. G. Imdoukh, M., Ahmad, I., & Alfailakawi, "Machine learning-based auto-scaling for containerized applications. Neural Computing and Applications," 32(13), 9745–9760, 2020.

[31] P. Alonso, M., Coll, S., Martínez, J. M., Santonja, V., & López, "Power consumption management in fat-tree interconnection networks. Parallel computing," 48, 59–80, 2015.

[32] R. Zhong, Z., & Buyya, "A cost-efficient container orchestration strategy in kubernetes-based cloud computing infrastructures with heterogeneous resources. ACM Transactions on Internet Technology (TOIT), 20(2), 1-24," 2020.

[33] A. Schwanke, "Faculty Informatics Bachelor of Science–Business Information Systems," 2019.

[34] S. Grubor, "Deployment with Docker: Apply continuous integration models, deploy applications quicker, and scale at large by putting Docker to work. Packt Publishing Ltd," 2017.

[35] B. Alkadi, O., Moustafa, N., & Turnbull, "A review of intrusion detection and blockchain applications in the cloud: approaches, challenges and solutions. IEEE Access, 8, 104893-104917," 2020.

[36] R. Paganetti, "Building a Compliance Model: A Delphi Study of Managed Security Service Providers Governing Regulatory Compliance Successfully (Doctoral dissertation, Capella University)," 2020.

[37] I. Nishihara, R., Moritz, P., Wang, S., Tumanov, A., Paul, W., Schleier-Smith, J., ... & Stoica, "Real-time machine learning: The missing pieces. In Proceedings of the 16th workshop on hot topics in operating systems (pp. 106-110)."

[38] L. Erik, S., & Emma, "The Future of Software Development: AI-Driven Testing and Continuous Integration for Enhanced Reliability. International Journal of Trend in Scientific Research and Development, 2(4), 3082-3096," 2018.

[39] C. Karamitsos, I., Albarhami, S., & Apostolopoulos, "Applying DevOps practices of continuous automation for machine learning. Information, 11(7), 363," 2020.

[40] M. Karslioglu, "Kubernetes-A Complete DevOps Cookbook: Build and manage your applications, orchestrate containers, and deploy cloud-native services. Packt Publishing Ltd," 2020.

[41] A. Bansal, "System to redact personal identified entities (PII) in unstructured data. International Journal of Advanced Research in Engineering and Technology, 11(6), 133," 2020.

[42] C. W. Fuller, "Continuous Integration/Continuous Delivery Pipeline for Air Force Distributed Common Ground System (AF DCGS)," 2020.