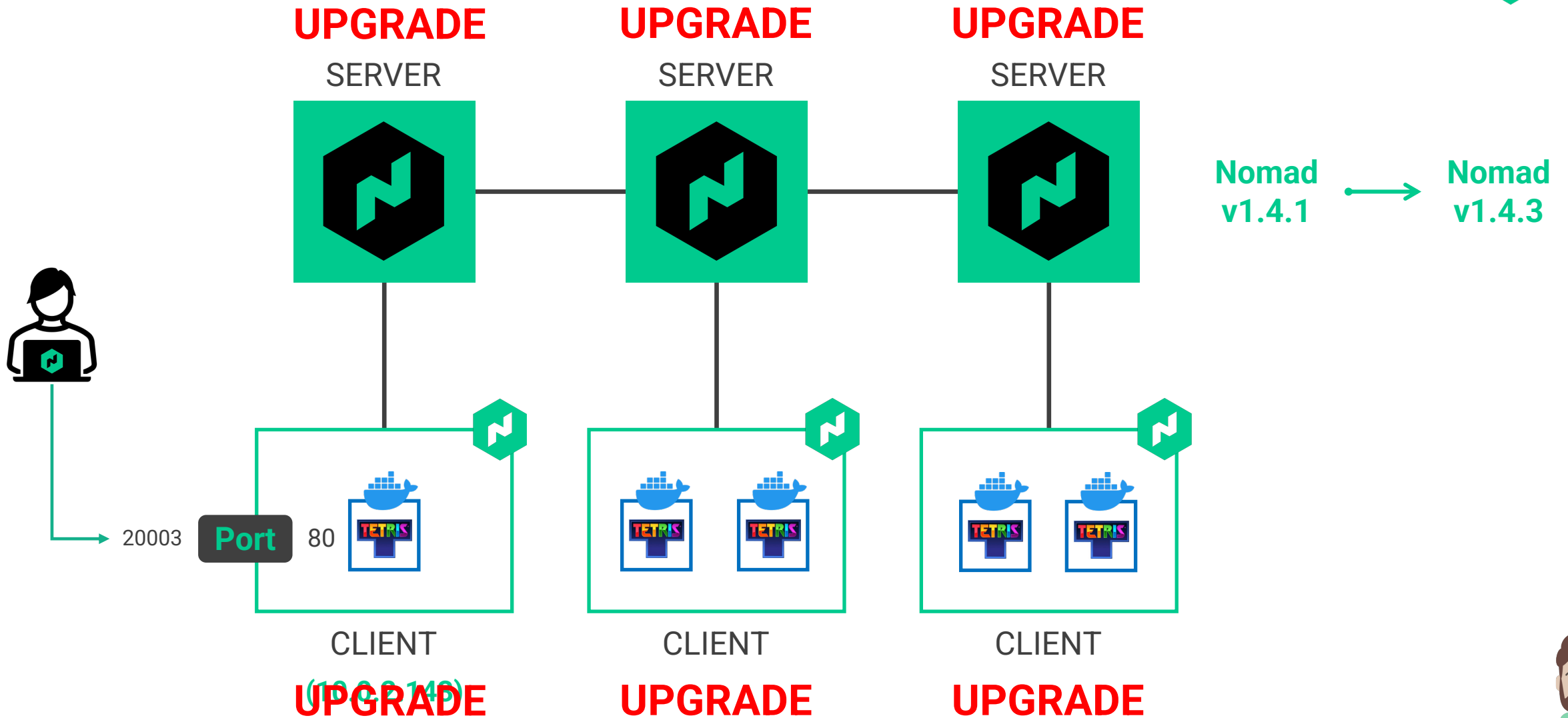




Expanding on Nomad Jobs



Expanding on Nomad Jobs

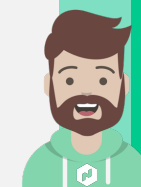
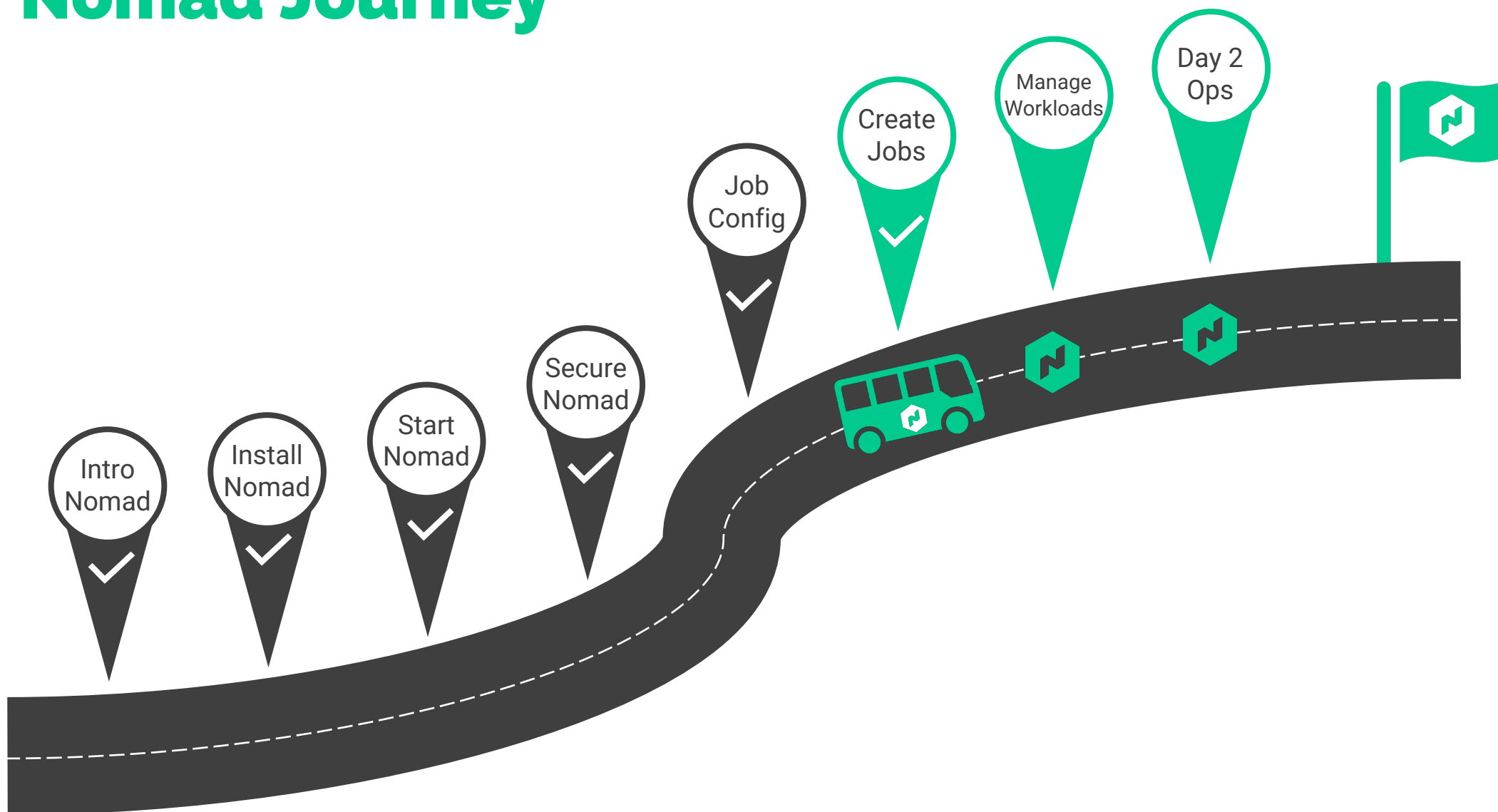


How Can We Improve Our Environment?

- ❑ Spread our application across nodes for high availability
- ❑ Monitor our application logs
- ❑ Upgrading the version of our application
- ❑ Improve the way we access our application
- ❑ Scale our application for dynamic workloads
- ❑ Use variables to limit hardcoding values
- ❑ Integrate Nomad with Vault and Consul



Nomad Journey





Job Placement



Job Placement

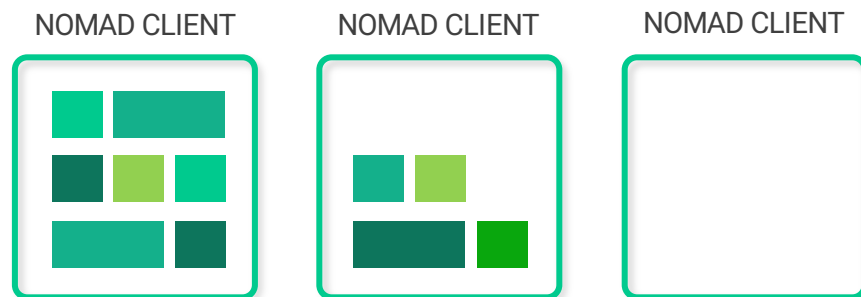


- When deploying applications, understanding where they will run is essential to ensure the maximum uptime and high availability for consumers
- By default, Nomad will use the **"binpack"** algorithm for allocations on available client nodes
- Bin packing can save costs by **maximizing the resources** of Nomad clients
- However, bin packing **can introduce risk** to your application because it may not be deployed across multiple client nodes
- You can configure Nomad scheduling to use a **"spread"** algorithm for the entire cluster, or you can customize it per job within the job specification

Placement { job → spread
job → group → spread

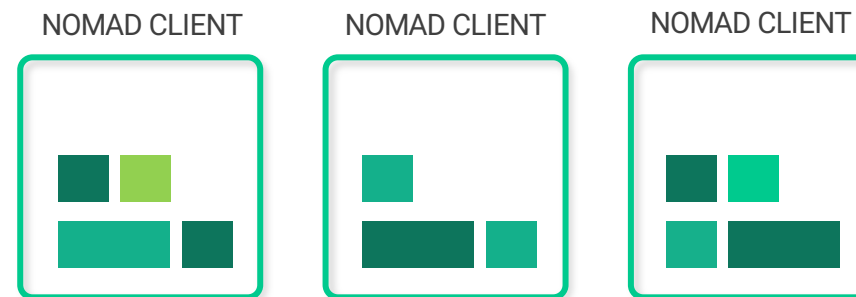


Bin Packing



Bin Packing (default)

- Utilize all of a node's resources before deploying applications on a different node
- Maximizes resource utilization and density
- Minimize infrastructure cost



Spread-Based Scheduling

- Evenly deploy applications by prioritizing based on nodes that are least utilized
- Distribute risk and ensures faster response and recovery in case of a node outage
- Optimized for performance



Customizable Scheduling



- Choose “binpack” or “spread”
- Customizable as one simple string in Nomad’s configuration
- Scheduling algorithm applies to all applications deployed on the cluster

```
1  # Server & Raft configuration
2  server {
3      enabled          = true
4      bootstrap_expect = 3
5      encrypt          = "Do7GerAsNtzK527dxRZJwpJANdS2NTFbKJIxIod84u0="
6      license_path     = "/etc/nomad.d/nomad.hclic"
7      server_join {
8          retry_join = ["10.4.23.44", "10.4.54.112", "10.4.56.33"]
9      }
10     default_scheduler_config {
11         scheduler_algorithm = "spread"
12     }
13 }
```



Job Placement



- When using spread, the scheduler will attempt to **place allocations equally** among the available values of the given target
- Spread can be used at the job level **and/or** the group level
- Spread can distribute tasks across datacenters if you have federated datacenters

```
1 job "tetris" {
2   datacenters = ["dc1"]
3
4   group "games" {
5     count = 5
6
7     spread {
8       attribute = "${node.datacenter}"
9       target "dc1" {
10        percent = 100
11      }
12    }
13  }
```



Job Scheduling



```
1 job "tetris" {
2   datacenters = ["dc1", "dc2"]
3
4   group "games" {
5     count = 5
6
7     spread {
8       attribute = "${node.datacenter}"
9       target "dc1" {
10        percent = 70
11      }
12      target "dc2" {
13        percent = 30
14      }
15    }
16  }
```

Schedule 70% of allocations to DC1 and 30% to DC2

```
1 job "tetris" {
2   datacenters = ["nyc", "sfo"]
3
4   group "games" {
5     count = 5
6
7     spread {
8       attribute = "${meta.dc}"
9       target "nyc-prod" {
10        percent = 50
11      }
12      target "sfo-dr" {
13        percent = 50
14      }
15    }
16  }
```

Schedule 50% of allocations to the primary datacenter (NYC) and 50% to the DR datacenter (SFO) using user-defined metadata



Job Scheduling



```
1 job "tetris" {
2   datacenters = ["dc1", "dc2"]
3
4   spread {
5     attribute = "${node.datacenter}"
6     target "dc1" {
7       percent = 70
8     }
9     target "dc2" {
10      percent = 30
11    }
12  }
13
14  group "frontend" {
15    # ...
16    task "webapp" {
17      # ...
18    }
19
20    group "backend" {
21      # ...
22      task "data" {
23        # ...
24      }
25    }
26  }
27 }
```

Spread is now under the job stanza

Job → Spread

This applies to **all groups and tasks** in the job specification

Multiple groups in the job specification



Job Scheduling

Multiple Spread Configurations

Spread all groups included in the job specification across both the **nyc** and **sfo** datacenters

Within each datacenter, spread allocations for the **frontend** group across Nomad clients that have user-defined metadata for **prod1** and **prod2**

Note: If there is a conflict, the **group level spread** will take priority

```
1  job "tetris" {
2    datacenters = ["nyc", "sfo"]
3
4    spread {
5      attribute = "${node.datacenter}"
6      target "nyc" {
7        percent = 70
8      }
9      target "sfo" {
10       percent = 30
11     }
12   }
13
14   group "frontend" {
15     spread {
16       attribute = "${meta.env}"
17       target "prod1" {
18         percent = 50
19       }
20       target "prod2" {
21         percent = 50
22       }
23     }
24
25     group "webapp" {
26       # ...
27     }
28   }
29 }
```





DEMO

Job Placement





Job Constraints



Job Constraints

- Constraints are requirements Nomad must **evaluate about the client**, such as the operating system, architecture, kernel version, and more before allocations are made....
- Constraint requirements are specified at the **job**, **group**, or **task** level
- Examples:
 - Client must be **Linux**
 - Client must be running **x64** architecture
 - Client must be running on an underlying AWS instance that is **m5.8xlarge**
 - Client must have specific **metadata**

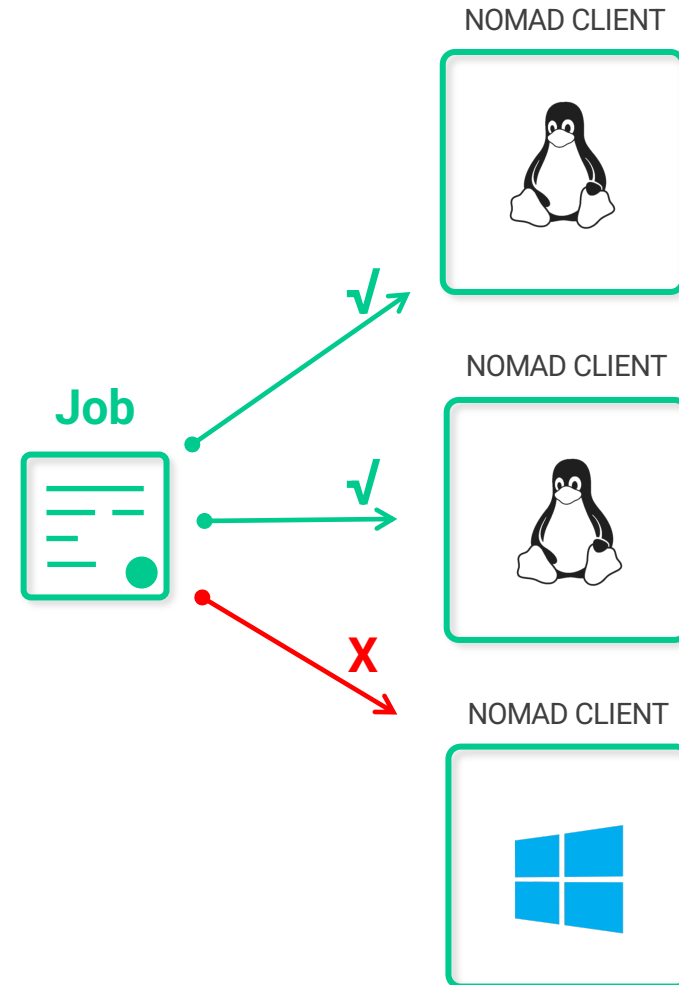
Placement {
job → constraint
job → group → constraint
job → group → task → constraint



Job Constraints



```
1 job "tetris" {  
2   datacenters = ["nyc", "sfo"]  
3  
4   constraint {  
5     attribute = "${attr.kernel.name}"  
6     value     = "linux"  
7   }  
}
```



Job Constraints



```
1 job "tetriss" {  
2   datacenters = ["nyc", "sfo"]  
3  
4   constraint {  
5     attribute = "${meta.env}"  
6     value     = "prod1"  
7   }  
}
```

Client Agent Configuration

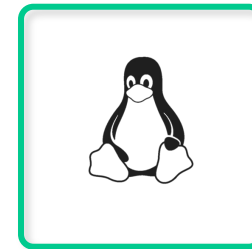
```
client {  
  enabled = true  
  
  meta {  
    env       = "prod1"  
    rack      = "rack-12"  
    hardware  = "cisco"  
    instructor = "krausen"  
  }  
}
```



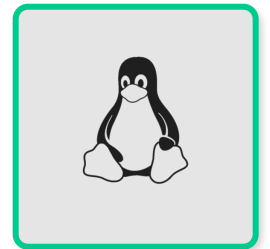
prod1

X

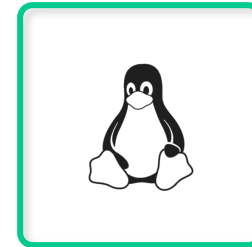
prod2



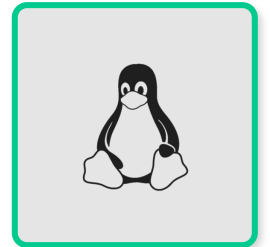
NOMAD CLIENT



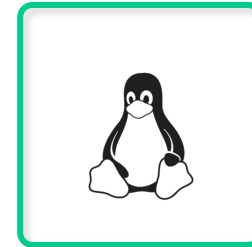
NOMAD CLIENT



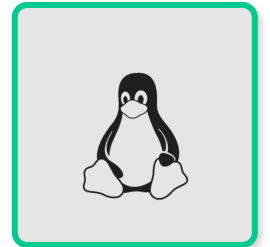
NOMAD CLIENT



NOMAD CLIENT



NOMAD CLIENT



NOMAD CLIENT



Job Constraints

Additional Examples



Distinct Hosts

Don't co-locate any tasks on the same client



```
1 group "frontend" {
2   constraint {
3     operator = "distinct_hosts"
4     value   = "true"
5   }
6
7   task {
8     driver = "docker"
9     #...
10  }
```

Cloud Metadata

Prevent the task to only run on nodes of size m4.xlarge



```
1 group "frontend" {
2   constraint {
3     attribute = "${attr.platform.aws.instance-type}"
4     value     = "m4.xlarge"
5   }
6
7   task {
8     driver = "docker"
9     #...
10  }
```

Client OS

Restrict the task to running on clients running Ubuntu



```
1 group "frontend" {
2   constraint {
3     attribute = "${attr.os.name}"
4     value     = "ubuntu"
5   }
6
7   task {
8     driver = "docker"
9     #...
10  }
```





DEMO

Job Constraints





Networking





Networking

- When deploying applications, networking is a critical component to ensure users can access the running applications
- Networking is defined at the group level using the **network** stanza
- Options for networking include:
 - **bridge**: group will have an isolated network namespace with an interface bridged with the host
 - **host**: (default) - each task will join the host network namespace – a shared network namespace is not created
 - **cni/<network>**: task group will have an isolated network namespace with the CNI network
 - **none**: each task will have an isolated network without any network interfaces

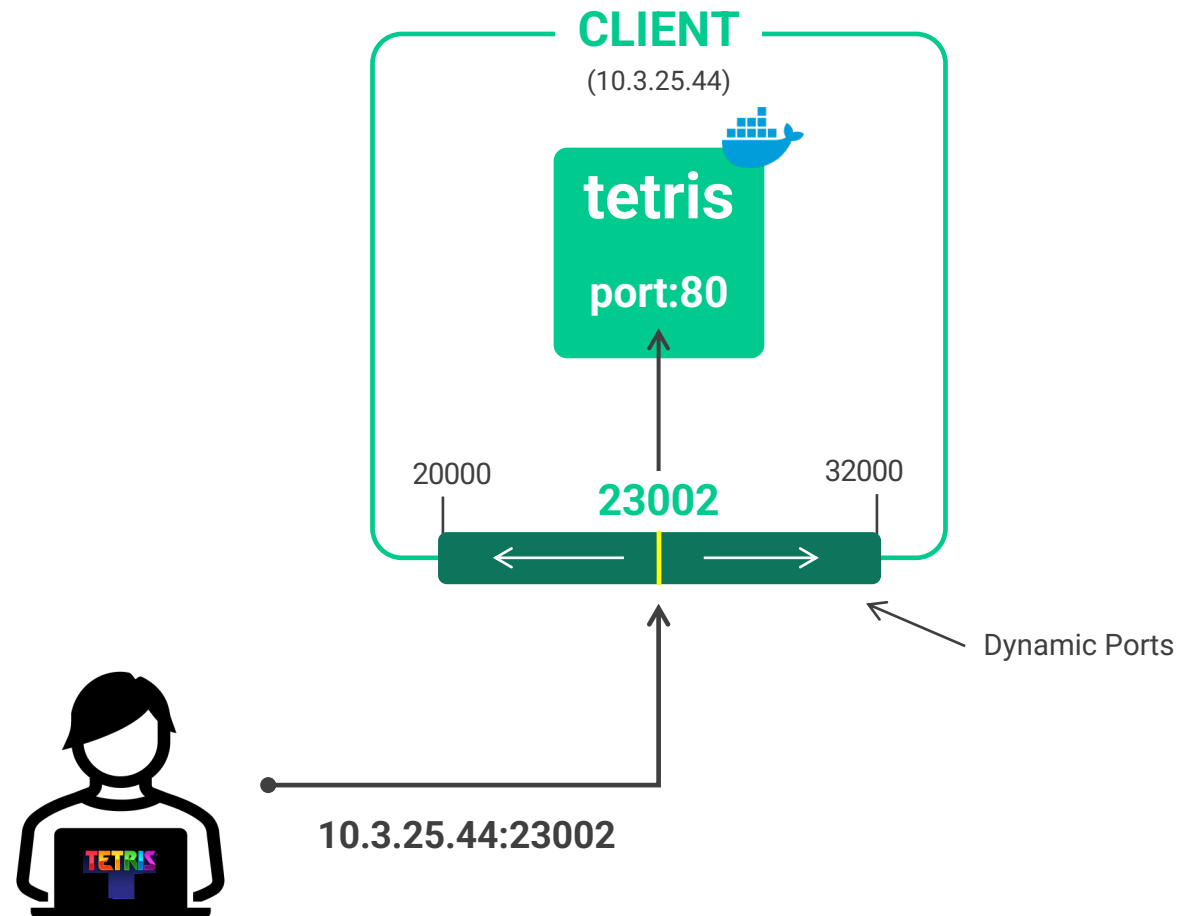
Placement { job → group → network



Networking

Host Mode

- Access via host IP address and dynamically allocated port
- Default port range is 20000 – 32000
- Relies on the task drivers to implement port mapping



Networking Mode - Host

```
1 group "games" {
2     count = 5
3
4     network {
5         port "web" {
6             to = 80
7         }
8     }
9
10    task "tetris" {
11        driver = "docker"
12
13        config {
14            image      = "bsord/tetris"
15            ports      = ["web"]
16            auth_soft_fail = true
17        }
18    }
```

This configuration is used when your application listens on a fixed port and you want to map to a dynamic port on the host

- **port** – used to specify both dynamic and reserved ports. This example uses a reserved port using the **to** parameter
- uses a label called **"web"**
- Note that we do not need a **mode = host** here because that is the default configuration

Reference the port label inside the task



Networking Mode - Host

This configuration is used when your application listens on a fixed port, and you need to map it to a static port on the host

- **port** – used to specify both dynamic and reserved ports. This example specifies the **static** port to allocate on the host (not recommended)

Reference the port label inside the task

```
1  group "games" {
2    count = 5
3
4    network {
5      port "web" {
6        static = 80
7        to = 80
8      }
9    }
10
11   task "tetris" {
12     driver = "docker"
13
14     config {
15       image      = "bsord/tetris"
16       ports      = ["web"]
17       auth_soft_fail = true
18     }
19   }
```





Networking Mode - Host

```
1 group "games" {
2   count = 5
3
4   network {
5     port "web" {}
6   }
7
8   task "tetris" {
9     driver = "docker"
10
11     config {
12       image      = "bsord/tetris"
13       ports      = ["web"]
14       auth_soft_fail = true
15     }
16   }
17 }
```

- No configuration specifies a dynamic port allocation for the task
- The application/service would need to read an environment variable to know the port to bind to at startup
- **NOMAD_PORT_web** would be passed to the task so it can know the dynamic port
- Again, no need to specify **mode = host** because it's the default configuration

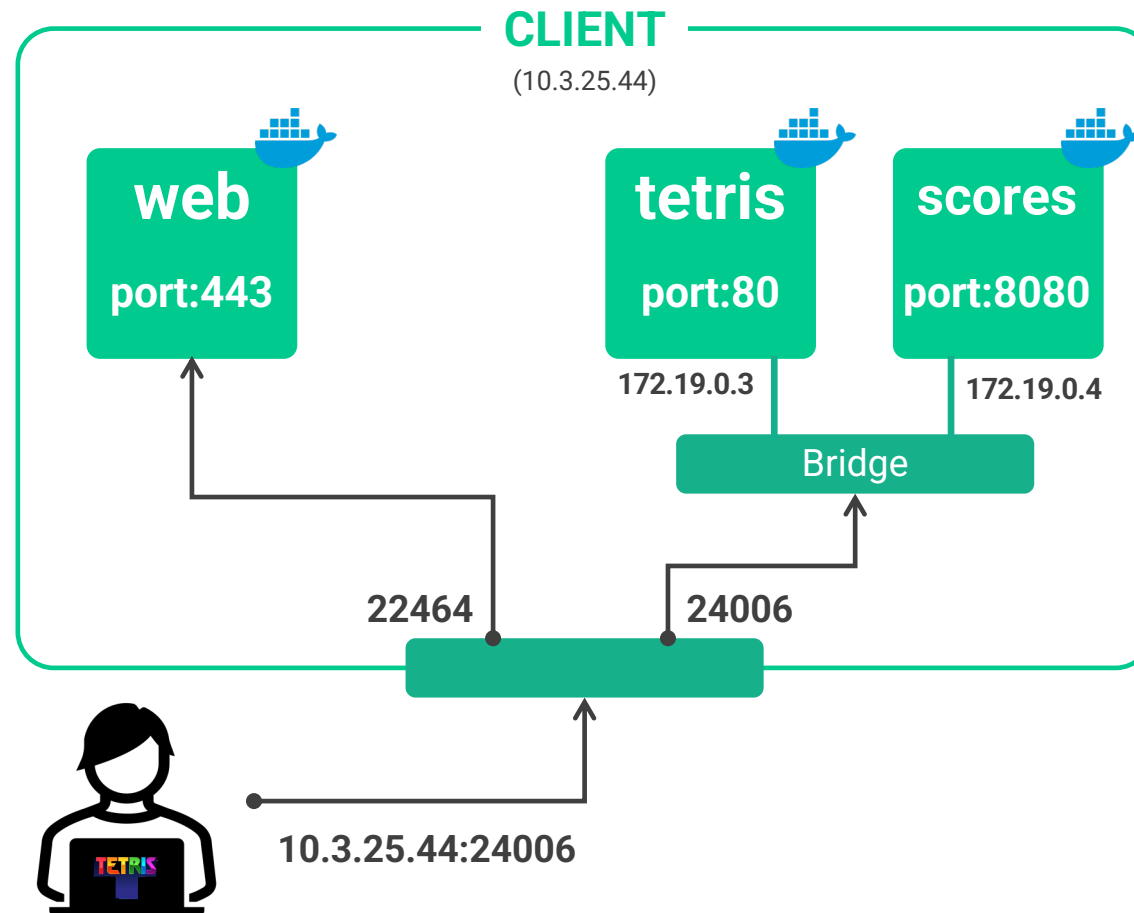


Networking

Bridge Mode



- Access via host IP address and dynamically allocated port
- Default port range is 20000 - 32000



Networking Mode - Bridge

This configuration is used when your application listens on a fixed port, and you want to connect to a bridged network on the host

- `mode = bridge` required since it's not the default configuration

Reference the port label inside the task

```
1  group "games" {
2    count = 5
3
4    network {
5      mode = "bridge"
6      port "http" {
7        to = 8080
8      }
9    }
10
11   task "tetris" {
12     driver = "docker"
13
14     config {
15       image      = "bsord/tetris"
16       ports      = ["http"]
17       auth_soft_fail = true
18     }
19   }
```





DEMO

Networking





Working with Volumes





Working with Volumes

- For workloads that require storage for stateful workloads, Nomad can provide access to host volumes or CSI volumes
- Nomad is aware of host volumes during the scheduling process, allowing it to make scheduling decisions based on the availability of host volumes on a specific client
- However, Nomad is **NOT** aware of Docker volumes because they are managed outside of Nomad. Therefore, the scheduler cannot make decisions based on availability
- Of course, volumes are available to other resources beyond containers, such as **exec** and **java** apps

Placement { job → group → volume





Working with Volumes

- A **host** volume is essentially a path on the client that is made available to jobs. The data is stored locally on the Nomad client
- A **csi-volume** is exposed to jobs using a CSI plugin to consume externally created storage volumes. Examples of these include AWS EBS volumes, GCP persistent disks, Ceph, vSphere, and more

Placement { job → group → volume



Host Volumes

Host volumes must be first configured using the **client** block in the Nomad agent configuration file:

```
1  # Client Configuration - Node can be Server & Client
2  client {
3    enabled = true
4    server_join {
5      retry_join = ["provider=aws tag_key=nomad_cluster_id tag_value=us-east-1"]
6    }
7    host_volume "database-primary" {
8      path = "/opt/nomad/volumes/db-pri"
9      read_only = false
10   }
11   host_volume "web-temp-data-01" {
12     path = "/opt/nomad/volumes/web-temp"
13     read_only = false
14   }
15 }
```

Name of Volume

Path on Nomad client



Host Volumes

```
1  group "games" {
2    count = 5
3
4    # ...
5    volume "data-01" {
6      type      = "host"
7      source     = "web-temp-data-01"
8      read_only  = false
9    }
10
11   task "tetris" {
12     driver = "docker"
13
14     volume_mount {
15       volume      = "data-01"
16       destination = "/var/lib/http"
17       read_only   = false
18     }
19
20     config {
21       image      = "bsord/tetris"
22       ports      = ["http"]
23       auth_soft_fail = true
24     }
25
26     resources {
27       cpu      = 50
28       memory   = 256
29     }
30   }
31 }
```

Specify the Required Volume for the games Group

- Once the volume is ready to use, it is referenced *within the group* using the **volume** parameter

Specifies where a group volume should be mounted

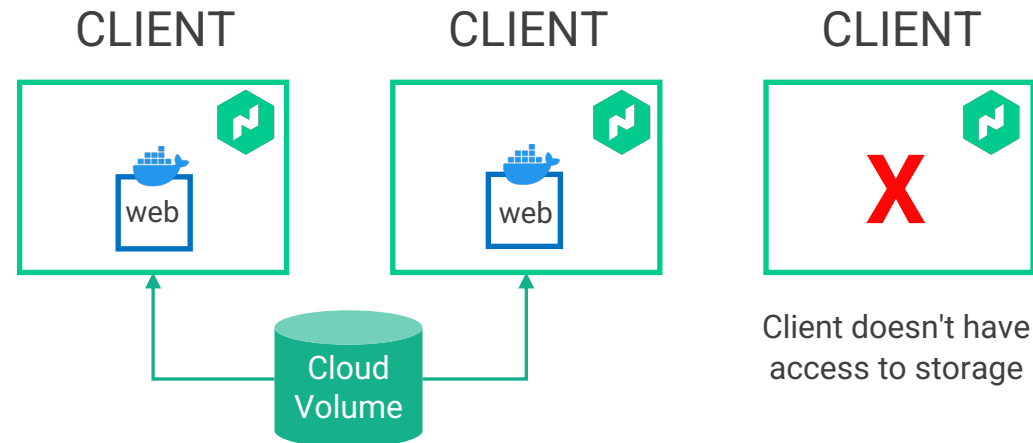
- The task will reference the volume *required for the task to run* using the **volume_mount** parameter





CSI Volumes

- CSI volumes are dynamically **mounted by CSI plugins**.
- These plugins actually **run as a system job** and can mount volumes created by cloud providers or storage platforms. No updates to the client agent configuration file needed*!
- Again, Nomad is aware of CSI volumes, which allows Nomad to schedule your workloads **based on the availability of volumes** on a client



* ok, maybe one small one but not per volume





CSI Volumes

- Since the CSI plugins are written by the storage vendors, any CSI plugin that supports Kubernetes should work with Nomad
- CSI volumes are supported by a large number of storage platforms and cloud providers, including:
 - Alicloud Disk/NAS/OSS
 - AWS EBS
 - AWS EFS
 - AWS FSx
 - Azure Blob/Disk/File
 - CephFS
 - Cisco Hyperflex
 - Dell EMC PowerMax/PowerScale/Unity/etc
 - Google Cloud Filestore/Storage
 - IBM
 - NFS
 - Portworx (Pure)
 - SMB
 - Synology
 - vSphere



Required Steps for CSI Volumes



1

Create the volume on your platform of choice

2

Enable Privileged Docker jobs

3

Run the plugin job for the specific driver needed

4

Register the volume with Nomad

5

Deploy the application that will use the volume



Enable Privileged Docker Jobs

- CSI plugins run as a **privileged Docker job** since they use bidirectional mount propagation to mount disks to the host
- The default configuration does not allow privileged Docker jobs, so you must **update your client configuration** to permit them

```
1  # Client Configuration - Node can be Server & Client
2  client {
3      enabled = true
4      server_join {
5          retry_join = ["10.3.2.16", "10.4.33.66", "10.3.2.18"]
6      }
7  }
8
9  plugin "docker" {
10     config {
11         allow_privileged = true
12     }
13 }
```



Create the Plugin Job

- Each CSI plugin supports one or more types of jobs – **Controller**, **Nodes**, or **Monolith**
- Most CSI plugins require that you run both a controller and a nodes job
- Node plugins are usually run as **system** jobs
 - These jobs use the **csi_plugin** stanza in the job spec file



Create the Plugin Job

this is a **system** job

location of the driver image

Run as a privileged job

Specifies this task provides a CSI plugin to the cluster

```
1  job "plugin-aws-ebs-nodes" {
2    datacenters = ["dc1"]
3
4    type = "system"
5
6    group "nodes" {
7      task "plugin" {
8        driver = "docker"
9
10       config {
11         image = "amazon/aws-ebs-csi-driver:v0.10.1"
12
13         args = [
14           "node",
15           "--endpoint=unix://csi/csi.sock",
16           "--logtostderr",
17           "--v=5",
18         ]
19
20       privileged = true
21     }
22
23     csi_plugin {
24       id      = "aws-ebs0"
25       type    = "node"
26       mount_dir = "/csi"
27     }
28
29     resources {
30       cpu    = 500
31       memory = 256
32     }
33   }
34 }
35 }
```



Deploy the Plugin Jobs to Nomad



Once you have the plugin job specifications ready, you can create the jobs in Nomad:

```

$ nomad job run csi_plugin.hcl

==> 2023-01-26T17:20:29-05:00: Monitoring evaluation "b66aaf83"
    2023-01-26T17:20:29-05:00: Evaluation triggered by job "plugin-aws-ebs-nodes"
    2023-01-26T17:20:30-05:00: Allocation "f329afdd" created: node "7ff357a0", group "nodes"
    2023-01-26T17:20:30-05:00: Evaluation status changed: "pending" -> "complete"
==> 2023-01-26T17:20:30-05:00: Evaluation "b66aaf83" finished with status "complete"

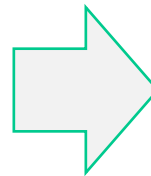
```



Register the Volume

Each volume needs to be registered with Nomad to ensure the CSI plugins know about each volume

```
1 # volume registration
2 type      = "csi"
3 id        = "web-data"
4 name      = "mysql"
5 external_id = "vol-5c699e29779e11fa7"
6 plugin_id  = "aws-ebs"
7
8 capability {
9   access_mode      = "single-node-writer"
10  attachment_mode = "file-system"
11 }
```



```
$ nomad volume register ebs_volume.hcl
```

```
ID = web-data
Name = web-data
External ID = vol-5c699e29779e11fa7
Plugin ID = aws-ebs
Provider = ebs.csi.aws.com
Version = v0.10.1
...
```



Create the Job to Use the Volume

And....finally. Create the job for the task(s) that will consume the new CSI volume

```
1  job "tetris" {
2    datacenters = ["dc1"]
3    type        = "service"
4
5    group "tetris" {
6      count = 1
7
8      volume "web-data-temp" {
9        type          = "csi"
10       read_only      = false
11       source          = "web-data"
12       access_mode     = "single-node-writer"
13       attachment_mode = "file-system"
14     }
15
16     task "mysql-server" {
17       driver = "docker"
18
19       volume_mount {
20         volume       = "web-data-temp"
21         destination  = "/pdata"
22         read_only    = false
23       }
24     }
25   }
```





DEMO

Nomad Volumes

