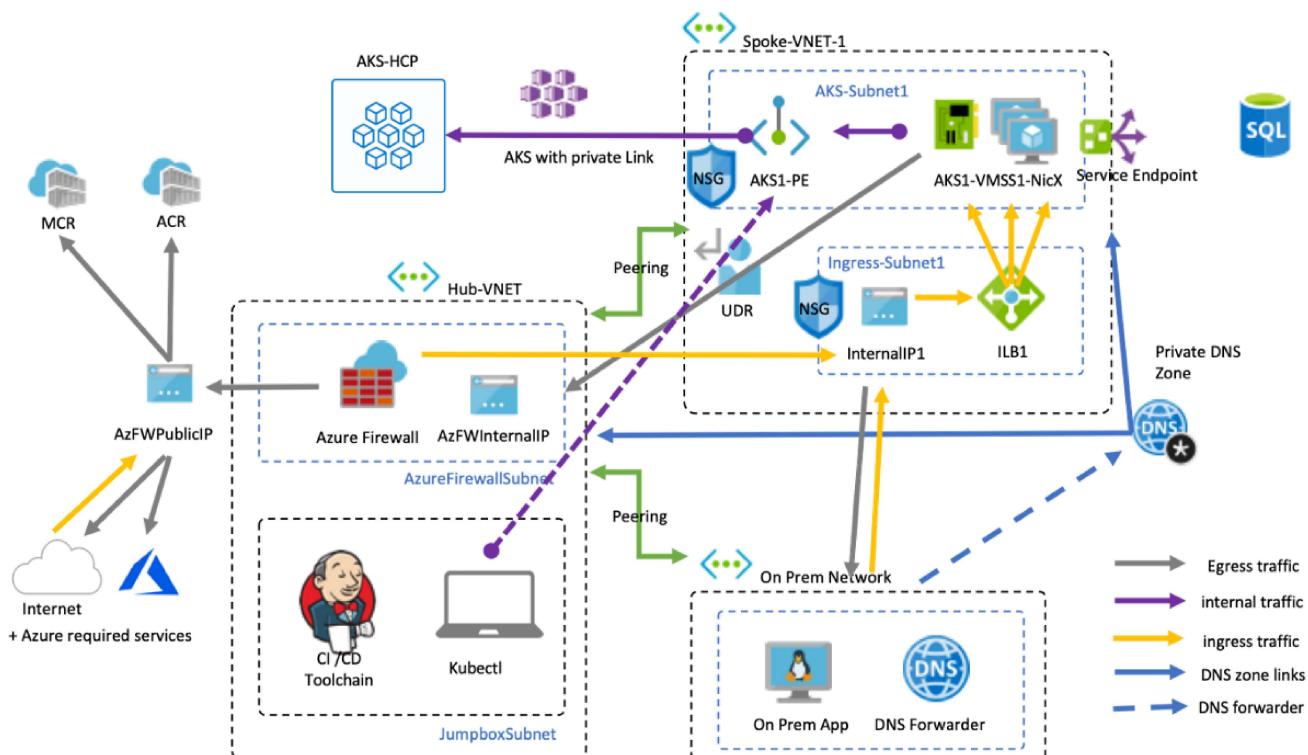


# Create a “Fully private” AKS infrastructure with Terraform



Pavel Tuzov

Jun 24, 2020 · 7 min read ★



Infrastructure from <https://medium.com/@denniszielke/fully-private-aks-clusters-without-any-public-ips-finally-7f5688411184>

In Azure you can create a private AKS cluster, in which the traffic between the node pools and the API server does not leave the private network. But often, making only the communication between nodes and the control plane private is not enough for your security needs. In a *fully* private cluster you also want to not expose and use public IPs.

Such an architecture is described in details in this article:

Fully private AKS clusters — without any public ips — finally!

One year after my last post on leveraging an azure firewall I want to

medium.com

Here I will describe the Terraform script which will manage the exact infrastructure as in the above post. I will explain only occasionally why we need the different Azure resources, so you probably want to read Dennis's article first.

I find this project also a very good example to practice different Terraform features:

- organizing the project into modules
- defining terraform and provider versions
- type constraints
- each and for\_each expressions
- functions
- defining explicit dependency
- and even using provisioners!

The sources can be found in this repository: <https://github.com/patuzov/terraform-private-aks>

## Resource groups

There will be two separate resource groups for the hub and spoke networks. So, here they are in the *main.tf*, along with the *terraform* and *provider* blocks.

*Note, we will need azurerm provider starting with version 2.5, as this is when the outbound\_type argument was introduced for AKS cluster resource.*

```

1  terraform {
2      required_version = ">= 0.12"
3  }
4
5  provider "azurerm" {
6      version = "~>2.5"

```

```

7   features {}
8 }
9
10 resource "azurerm_resource_group" "vnet" {
11   name     = var.vnet_resource_group_name
12   location = var.location
13 }
14
15 resource "azurerm_resource_group" "kube" {
16   name     = var.kube_resource_group_name
17   location = var.location
18 }

```

~~

`terraform init`  
`terraform apply`

...and we can see our two resource groups in the Azure Portal.

## Networking

### VNETs, Subnets

As for the networks, we will need a VNET and one or more SUBNETs for both: hub and spoke. This is a great opportunity to use **Terraform Modules**. Using modules, you can hide complexity and, most importantly, reuse code in Terraform.

We will create a *modules* folder and a *vnet* folder inside it, with all the usual files for a module: *main.tf*, *variables.tf*, *outputs.tf*.

There is nothing fancy for the virtual network resource. Just set the values provided via the input variables:

```

1 resource "azurerm_virtual_network" "vnet" {
2   name          = var.vnet_name
3   address_space = var.address_space
4   location      = var.location
5   resource_group_name = var.resource_group_name
6 }
```

main.tf hosted with ❤ by GitHub

[view raw](#)

Now we might want to create a list of subnets in this vnet and we need to receive following input values for each subnet: name and the IP address range. We could receive as input a list of objects containing this information:

```

1 variable subnets {
2   description = "Subnets configuration"
3   type        = list(object({
4     name          = string
5     address_prefixes = list(string)
6   }))
7 }
```

[variables.tf](#) hosted with ❤ by GitHub

[view raw](#)

...which, in my opinion, makes sense for the *interface* of this module.

For each such object a subnet resource can be created, using the *for each* expression. But this expression expects either a list of strings or a map. So we'll have to transform our list of objects to a map. *for* expression comes to the rescue:

```

1 resource "azurerm_subnet" "subnet" {
2   for_each = { for subnet in var.subnets : subnet.name => subnet.address_prefixes }
3
4   name          = each.key
5   resource_group_name = var.resource_group_name
6   virtual_network_name = azurerm_virtual_network.vnet.name
7   address_prefixes = each.value
8 }
```

[main.tf](#) hosted with ❤ by GitHub

[view raw](#)

The *for* expression will create a map with *subnet.name* as keys and *subnet.address\_prefixes* as corresponding values. For each such key/value pair, an *azurerm\_subnet* resource will be created.

From this module, we will need vnet and subnet IDs. We define this in the *outputs.tf*. The same *technique* with the *for* expression is used to get back the subnet IDs as a map:

```

1 output vnet_id {
2   description = "Generated VNET ID"
3   value       = azurerm_virtual_network.vnet.id
4 }
5
6 output subnet_ids {
```

```

7   description = "Generated subnet IDs map"
8   value        = { for subnet in azurerm_subnet.subnet : subnet.name => subnet.id }
9 }
```

outputs.tf hosted with ❤ by GitHub

[view raw](#)

Finally, we make use of this module in our *main.tf*:

```

1 module "hub_network" {
2   source          = "./modules/vnet"
3   resource_group_name = azurerm_resource_group.vnet.name
4   location        = var.location
5   vnet_name       = var.hub_vnet_name
6   address_space    = ["10.0.0.0/22"]
7   subnets = [
8     {
9       name : "AzureFirewallSubnet"
10      address_prefixes : ["10.0.0.0/24"]
11    },
12    {
13      name : "jumpbox-subnet"
14      address_prefixes : ["10.0.1.0/24"]
15    }
16  ]
17 }
18
19 module "kube_network" {
20   source          = "./modules/vnet"
21   resource_group_name = azurerm_resource_group.kube.name
22   location        = var.location
23   vnet_name       = var.kube_vnet_name
24   address_space    = ["10.0.4.0/22"]
25   subnets = [
26     {
27       name : "aks-subnet"
28       address_prefixes : ["10.0.5.0/24"]
29     }
30   ]
31 }
```

You need to run `terraform init` again to initialize the new module. After that you can `apply` and the VNETs are in the portal as well, under the corresponding resource groups.

## Peering

Now we need to peer the VNETs. The *azurerm* provider has the *azurerm\_virtual\_network\_peering* resource. But, as the peering will have to be done both ways, it's better to abstract this in a new module.

Here is the very simple module's *main.tf* code:

```

1 resource "azurerm_virtual_network_peering" "peering" {
2   name           = var.peering_name_1_to_2
3   resource_group_name = var.resource_group
4   virtual_network_name = var.vnet_1_name
5   remote_virtual_network_id = var.vnet_2_id
6 }
7
8 resource "azurerm_virtual_network_peering" "peering-back" {
9   name           = var.peering_name_2_to_1
10  resource_group_name = var.resource_group
11  virtual_network_name = var.vnet_2_name
12  remote_virtual_network_id = var.vnet_1_id
13 }
```

[main tf hosted with ❤ by GitHub](#)

[view raw](#)

And this is how we use it:

```

1 module "vnet_peering" {
2   source      = "./modules/vnet_peering"
3   vnet_1_name = var.hub_vnet_name
4   vnet_1_id   = module.hub_network.vnet.id
5   vnet_1_rg   = azurerm_resource_group.vnet.name
6   vnet_2_name = var.kube_vnet_name
7   vnet_2_id   = module.kube_network.vnet.id
8   vnet_2_rg   = azurerm_resource_group.kube.name
9   peering_name_1_to_2 = "HubToSpoke1"
10  peering_name_2_to_1 = "Spoke1ToHub"
11 }
```

[main.tf hosted with ❤ by GitHub](#)

[view raw](#)

*Note how we use outputs of one module as inputs for the other.*

Don't forget to `terraform init`, in order to initialize the new module.

If you apply this configuration, you will see the *Peerings* settings under both VNETs.

## Firewall

An Azure Firewall will need a public IP and we'll also configure the exceptions for all the *external access* AKS needs in order to function properly (like pulling images or OS updates). So why not hide it all into yet another module?

Here is the code for the public IP, firewall itself and an example of a network rule for the *aks basics* (there are more rule collections defined — see the [git repo](#)):

```

1 resource "azurerm_public_ip" "pip" {
2   name          = var.pip_name
3   resource_group_name = var.resource_group
4   location      = var.location
5   allocation_method = "Static"
6   sku           = "Standard"
7 }
8
9 resource "azurerm_firewall" "fw" {
10  name          = var.fw_name
11  location      = var.location
12  resource_group_name = var.resource_group
13
14  ip_configuration {
15    name          = "fw_ip_config"
16    subnet_id     = var.subnet_id
17    public_ip_address_id = azurerm_public_ip.pip.id
18  }
19 }
20
21 resource "azurerm_firewall_application_rule_collection" "aksbasics" {
22  name          = "aksbasics"
23  azure_firewall_name = azurerm_firewall.fw.name
24  resource_group_name = var.resource_group
25  priority      = 101
26  action         = "Allow"
27
28  rule {
29    name          = "allow network"
30    source_addresses = ["*"]
31
32    target_fqdns = [
33      "*.cdn.mscr.io",
34      "mcr.microsoft.com",
35      "*.data.mcr.microsoft.com",
36      "management.azure.com",
37      "login.microsoftonline.com",
38      "acs-mirror.azureedge.net",
39      "dc.services.visualstudio.com",

```

```

40     "*.opinsights.azure.com",
41     "*.oms.opinsights.azure.com",
42     "*.microsoftonline.com",
43     "*.monitoring.azure.com",
44   ]
45
46   protocol {
47     port = "80"
48     type = "Http"
49   }
50
51   protocol {
52     port = "443"
53     type = "Https"
54   }
55 }
```

```

1 module "firewall" {
2   source      = "./modules/firewall"
3   resource_group = azurerm_resource_group.vnet.name
4   location     = var.location
5   pip_name     = "azureFirewalls-ip"
6   fw_name      = "kubenetfw"
7   subnet_id    = module.hub_network.subnet_ids["AzureFirewallSubnet"]
8 }
```

main.tf hosted with ❤ by GitHub

[view raw](#)

Note how we get the Azure Firewall Subnet ID from the network module's output.

## Route Table

Finally, we'll create a user defined route to force all traffic (0.0.0.0/0) from the AKS subnet to the firewall (specifically, to its private IP). For this, we need to create the route table with the corresponding route and associate it to the AKS subnet. And, yes, we'll create a new module for this:

```

1 resource "azurerm_route_table" "rt" {
2   name        = var.rt_name
3   location    = var.location
4   resource_group_name = var.resource_group
5
6   route {
7     name          = "kubenetfw_fw_r"
```

```

8     address_prefix      = "0.0.0.0/0"
9     next_hop_type       = "VirtualAppliance"
10    next_hop_in_ip_address = var.firewal_private_ip
11  }
12 }
13
14 resource "azurerm_subnet_route_table_association" "aks_subnet_association" {
15   subnet_id      = var.subnet_id
16   route_table_id = azurerm_route_table.rt.id
17 }
```

And use the module:

```

1 module "routetable" {
2   source          = "./modules/route_table"
3   resource_group = azurerm_resource_group.vnet.name
4   location        = var.location
5   rt_name         = "kubenetfw_fw_rt"
6   r_name          = "kubenetfw_fw_r"
7   firewal_private_ip = module.firewall.fw_private_ip
8   subnet_id       = module.kube_network.subnet_ids["aks-subnet"]
9 }
```

main.tf hosted with ❤ by GitHub

[view raw](#)

## AKS

We are finally ready to define the AKS cluster!

```

1 resource "azurerm_kubernetes_cluster" "privateaks" {
2   name          = "private-aks"
3   location      = var.location
4   kubernetes_version = "1.16.9"
5   resource_group_name = var.kube_resource_group_name
6   dns_prefix     = "private-aks"
7   private_cluster_enabled = true
8
9   default_node_pool {
10     name          = "default"
11     node_count    = var.nodepool_nodes_count
12     vm_size       = var.nodepool_vm_size
13     vnet_subnet_id = module.kube_network.subnet_ids["aks-subnet"]
14   }
15
16   identity {
17     type = "SystemAssigned"
```

```

18 }
19
20 network_profile {
21   docker_bridge_cidr = var.network_docker_bridge_cidr
22   dns_service_ip     = var.network_dns_service_ip
23   network_plugin      = "azure"
24   outbound_type        = "userDefinedRouting"
25   service_cidr        = var.network_service_cidr
26 }
27
28 depends_on = [module.routetable]
29 }
30
31 # https://github.com/Azure/AKS/issues/1557
32 resource "azurerm_role_assignment" "vmcontributor" {
33   role_definition_name = "Virtual Machine Contributor"
34   scope                 = azurerm_resource_group.vnet.id
35   principal_id          = azurerm_kubernetes_cluster.privateaks.identity[0].principal_id
36 }

```

~~There are 2 important things to mention here.~~

- Enabling the actual private AKS cluster feature: `private_cluster_enabled = true`. This will also create a *Private DNS Zone* resource, which will have an *A Record* pointing to the **private** IP address of the API, allowing *everyone* from the AKS VNET (Spoke) to resolve the control plane’s address.
- Setting *outbound type* to *userDefinedRouting*, which will prevent creating a public IP for the cluster’s egress traffic and make use of the “user defined” Route Table:  
`outbound_type = "userDefinedRouting"`.

Also, did you notice the `depends_on` argument? Terraform can identify if a resource depends on another. For example, in order to deploy this AKS cluster in the “aks-subnet” subnet, Terraform knows it has to create the vnet and subnet first. These are the *implicit* dependencies. If Terraform doesn’t see dependencies between resources, it will [luckily] create them in parallel.

But in our case, as we set the *userDefinedRouting* outbound type, AKS will expect to have a route table associated with its subnet, so it actually depends on the route table resource. We can let Terraform know about this dependency by *explicitly* defining it with the `depends_on` argument.

Also, because of the still unresolved issue, we need to allow the created Service Principal (by the managed identity) to edit the user defined routes in the route table from the *vnet* resource group, by defining the *azure role assignment*.

And we're done! We now have an AKS cluster, in which the API doesn't have a public endpoint and communicates with the nodes only through the Azure Backbone. The cluster also doesn't expose and doesn't use any public IPs (except those that we allowed it to).

## Jumpbox

The problem is, since the control plane doesn't expose any public endpoint, you will not be able to manage your cluster, say, with `kubectl` from your local machine, or even from the Azure DevOps pipelines (if you're using Microsoft-hosted agents).

Remember the *Private DNS Zone* from the previous step? Only the Spoke VNET can resolve the API. But we can link other VNETs too. So, let's create a (Linux) Virtual Machine in the Hub VNET, install `kubectl` on it and link the VNET in the Private DNS Zone.

No surprises here: I'll *hide* it all in a separate module.

I'll skip all the declarations for public IP, network security group, etc. You can find them in the github repo under *jumpbox* module. Here are the most important things:

```

1  resource "azurerm_linux_virtual_machine" "jumpbox" {
2      name                  = "jumpboxvm"
3      location              = var.location
4      resource_group_name   = var.resource_group
5      network_interface_ids = [azurerm_network_interface.vm_nic.id]
6      size                  = "Standard_DS1_v2"
7      computer_name         = "jumpboxvm"
8      admin_username        = var.vm_user
9      admin_password         = random_string.adminpassword.result
10     disable_password_authentication = false
11
12     os_disk {
13         name          = "jumpboxOsDisk"
14         caching       = "ReadWrite"
15         storage_account_type = "Premium_LRS"
16     }
17
18     source_image_reference {
19         publisher     = "Canonical"
20         offer         = "UbuntuServer"
21         sku           = "18.04-LTS"
22         version       = "latest"
23     }
24
25     tags = {
26         environment = "Development"
27     }
28 }
```

```
19     provider = Canonical
20     offer     = "UbuntuServer"
21     sku       = "16.04.0-LTS"
22     version   = "latest"
23 }
24
25 provisioner "remote-exec" {
26   connection {
27     host     = self.public_ip_address
28     type     = "ssh"
29     user     = var.vm_user
30     password = random_string.adminpassword.result
31   }
32
33   inline = [
34     "sudo apt-get update && sudo apt-get install -y apt-transport-https gnupg2",
35     "curl -s https://packages.cloud.google.com/apt/doc/apt-key.gpg | sudo apt-key add -",
36     "echo 'deb https://apt.kubernetes.io/ kubernetes-xenial main' | sudo tee -a /etc/apt/sourc
37     "sudo apt-get update",
38     "sudo apt-get install -y kubectl",
39     "curl -sL https://aka.ms/InstallAzureCLIDeb | sudo bash"
40   ]
41 }
42 }
43
44 resource "azurerm_private_dns_zone_virtual_network_link" "hublink" {
45   name           = "hubnetdnsconfig"
46   resource_group_name = var.dns_zone_resource_group
47   private_dns_zone_name = var.dns_zone_name
48   virtual_network_id    = var.vnet_id
```

*random\_string* resource) and show it in the terraform output, so that I can later ssh to the VM.

Probably you noticed a block in the VM resource which is a *taboo* in Terraform: **remote-exec provisioner**. Terraform provides them and kindly asks you not to use them. But, in this case, it fits perfectly! We're spinning up this VM just to test that we can connect to the control plane with kubectl from our private network. And, being lazy to do it manually, I'll let the provisioner automatically install kubectl and azure cli when it creates the VM.

As mentioned previously, we also need to link the hub VNET to the Private DNS Zone (which is automatically created with the AKS cluster). We can easily achieve this with

the `azurerm_private_dns_zone_virtual_network_link` resource.

From this module, we will need the following information:

```

1  output "jumpbox_ip" {
2      description = "Jumpbox VM IP"
3      value       = azurerm_linux_virtual_machine.jumpbox.public_ip_address
4  }
5
6  output "jumpbox_username" {
7      description = "Jumpbox VM username"
8      value       = var.vm_user
9  }
10
11 output "jumpbox_password" {
12     description = "Jumpbox VM admin password"
13     value       = random_string.adminpassword.result
14 }
```

[outputs.tf](#) hosted with ❤ by GitHub

[View raw](#)

To use this module, we have a challenge: getting the Private DNS Zone's name. It is generated automatically by Azure when the private AKS cluster is created (it's not directly maintained by Terraform) and it even contains a GUID in it:

*9c3df69f-b281-42ef-9996-89941b484906.privatelink.westeurope.azmk8s.io*

Fortunately, it's not that random. If you look at the private FQDN of the API, you will see this:

*private-aks-7b099028.9c3df69f-b281-42ef-9996-89941b484906.privatelink.westeurope.azmk8s.io*

And this information is known by Terraform. We just need to extract the private DNS Zone's name from it, using various Terraform functions: *join*, *slice*, *split* and *length*:

```
join(".", slice(split(".", azurerm_kubernetes_cluster.privateaks.private_fqdn), 1,
length(split(".", azurerm_kubernetes_cluster.privateaks.private_fqdn))))
```

This is how the module is used:

```

1 module "jumpbox" {
2   source          = "./modules/jumpbox"
3   location        = var.location
4   resource_group = azurerm_resource_group.vnet.name
5   vnet_id         = module.hub_network.vnet_id
6   subnet_id       = module.hub_network.subnet_ids["jumpbox-subnet"]
7   dns_zone_name  = join(".", slice(split(".", azurerm_kubernetes_cluster.privateaks.pr
8   dns_zone_resource_group = azurerm_kubernetes_cluster.privateaks.node_resource_group
9 }

```

main.tf hosted with ❤ by GitHub

[view raw](#)

Lastly, we want to define our outputs:

```

1 output "ssh_command" {
2   value = "ssh ${module.jumpbox.jumpbox_username}@${module.jumpbox.jumpbox_ip}"
3 }
4
5 output "jumpbox_password" {
6   description = "Jumpbox Admin Password"
7   value        = module.jumpbox.jumpbox_password
8 }

```

outputs.tf hosted with ❤ by GitHub

[view raw](#)

After applying the entire configuration, we get an output like this:

## Outputs:

```

jumpbox_password = (nSq7qMxbJ
ssh_command = ssh azureuser@52.157.213.177

```

Terraform output

Now we can ssh into the Linux VM with kubectl and azure cli already installed. All what's left is authenticating with Azure, getting the configuration for kubectl and we can connect to the API:

```

azureuser@jumpboxvm:~$ kubectl get nodes
NAME                  STATUS    ROLES     AGE      VERSION
aks-default-75302467-vmss000000   Ready    agent    11h      v1.16.9

```

### Using kubectl with a private AKS cluster

The traffic from the pods will be redirected to Azure Firewall and corresponding rules will be applied. Let's check it.

In the *firewall* module, a *test* rule is defined, which will allow traffic to *bing.com*:

```

1 resource "azurerm_firewall_application_rule_collection" "test" {
2   name          = "test"
3   azure_firewall_name = azurerm_firewall.fw.name
4   resource_group_name = var.resource_group
5   priority       = 104
6   action         = "Allow"
7
8   rule {
9     name          = "allow network"
10    source_addresses = ["*"]
11
12    target_fqdns = [
13      "*.bing.com"
14    ]
15
16    protocol {
17      port = "80"
18      type = "Http"
19    }
20
21    protocol {
22      port = "443"
23      type = "Https"
24    }
25  }
26}

```

Let's create a pod and try to make a couple of requests from its container:

```

azureuser@jumpboxvm:~$ kubectl run busy --image=busybox -- sleep "3600"
pod/busy created
azureuser@jumpboxvm:~$ kubectl get pods
NAME    READY    STATUS    RESTARTS   AGE
busy    1/1     Running   0          9s
azureuser@jumpboxvm:~$ kubectl exec -it busy -- sh
/ # wget http://www.google.com
Connecting to www.google.com (172.217.17.68:80)
wget: server returned error: HTTP/1.1 470 status code 470
/ #

```

```
/ "#  
/ # wget http://www.bing.com  
Connecting to www.bing.com (204.79.197.200:80)  
saving to 'index.html'  
index.html          100% |*****  
'index.html' saved  
/ #
```

Making external http requests from a container

Exactly what we wanted to achieve: external traffic is blocked, unless we have explicit rules to allow it.

Terraform    Kubernetes    Azure Kubernetes Service    Azure

About    Help    Legal

Get the Medium app

