# Scheduler Extender

Learn how the Kubernetes scheduler extender works.

## What is the scheduler extender?

A scheduler extender works as an external service. It isn't another scheduler. Just as the name suggests, it simply extends the default Kubernetes scheduler with more customized features or implementations. That doesn't mean we can customize everything in the default Kubernetes scheduler. But there are some extension points that we can include in our customized implementations.

```
1   // https://github.com/kubernetes/kubernetes/blob/master/pkg/scheduler/framework/extender.go#L24-L72
2   // Extender is an interface for external processes to influence scheduling
3   // decisions made by Kubernetes. This is typically needed for resources not directly
4   // managed by Kubernetes.
5   type Extender interface {
6       // Name() returns a unique name that identifies the extender.
7       Name() string
8
9       // Filter based on extender-implemented predicate functions. The filtered list is
10      // expected to be a subset of the supplied list.
11      // The failedNodes and failedAndUnresolvableNodes optionally contains the list
12      // of failed nodes and failure reasons, except nodes in the latter are
13      // unresolvable.
14      Filter(pod *v1.Pod, nodes []*v1.Node) (filteredNodes []*v1.Node, failedNodesMap extenderv1.FailedNod
15      // Prioritize based on extender-implemented priority functions
16      // The returned scores and weight are used to compute the weighted score for an extender.
17      // The weighted scores are added t the scores computed by Kubernetes scheduler. The total
18      // scores are used to do the host selection.
19
20
21      Prioritize(pod *v1.Pod, nodes []*v1.Node) (hostPriorities *extenderv1.HostPriorityList, weight int64
22      // Bind() delegates the action of binding a Pod to a node to the extender.
23      Bind(binding *v1.Binding) error
24
25      // IsBinder() returns whether this extender is configured for the Bind() method.
26      IsBinder() bool
27
28      // sInterested() returns true if at least one extended resource requested by this
29      // Pod is managed by this extender.
```

Scheduler extender interface

Above is the extender interface, where we can add our extension implementations during the `Filter`, `Preempt`, `Prioritize`, and `Bind` phases. The Kubernetes scheduler extender works like a webhook. When Pod scheduling comes into these four phases, the `kube-scheduler` can call the extension logic registered by the extender.

Now, let's take a closer look at how the scheduler extender works.

# How the scheduler extender works

With the scheduler extender, we can change how the `kube-scheduler` makes decisions on Pod scheduling to meet our business needs.

## Register scheduler extenders

Firstly, let's see how we register our extenders to the `kube-scheduler`. The flag `--config` in the `kube-scheduler` is the entry where we register our external extenders. Here, we can register a list of extenders when needed. In that configuration file, we specify the parameters used by the `kube-scheduler` to communicate with the extenders, such as the opt-in `Filter` phase and extender service endpoints. Below is a sample configuration file:

```
 1  apiVersion: kubescheduler.config.k8s.io/v1
 2  kind: KubeSchedulerConfiguration
 3  clientConnection:
 4    kubeconfig: /etc/kubernetes/scheduler.conf
 5  extenders:
 6  - urlPrefix: "http://127.0.0.1:8080/"
 7    filterVerb: filter
 8    bindVerb: bind
 9    enableHTTPS: false
10    nodeCacheCapable: true
11    managedResources:
12    - name: resource.mycompany.com/gpu
13      ignoredByScheduler: false
14    ignorable: false
```

Sample scheduler configuration

In the file above, we specify an insecure extender running at `http://127.0.0.1:8080/`. In this extender, we specify the extended resources that are managed by this extender. Only the Pods request at the extended resource `resource.mycompany.com/gpu` will be sent to this extender during the `Filter` phase. The field `managedResources` is optional, which can help us improve the scheduling latency by calling

extenders when necessary. Below is a completed struct of an `Extender`, where we can add more extender verbs, such as `Filter` and `Preempt`:

```
 1   // https://github.com/kubernetes/kubernetes/blob/master/staging/src/k8s.io/kube-scheduler/config/v1/type
 2   // Extender holds the parameters used to communicate with the extender. If a verb is unspecified/empty,
 3   // it is assumed that the extender chose not to provide that extension.
 4   type Extender struct {
 5       // URLPrefix at which the extender is available
 6       URLPrefix string `json:"urlPrefix"`
 7       // Verb for the filter call, empty if not supported. This verb is appended to the URLPrefix when iss
 8       FilterVerb string `json:"filterVerb,omitempty"`
 9       // Verb for the preempt call, empty if not supported. This verb is appended to the URLPrefix when is
10       PreemptVerb string `json:"preemptVerb,omitempty"`
11       // Verb for the prioritize call, empty if not supported. This verb is appended to the URLPrefix when
12       PrioritizeVerb string `json:"prioritizeVerb,omitempty"`
13       // The numeric multiplier for the node scores that the prioritize call generates.
14       // The weight should be a positive integer
15       Weight int64 `json:"weight,omitempty"`
16       // Verb for the bind call, empty if not supported. This verb is appended to the URLPrefix when issui
17       // If this method is implemented by the extender, it is the extender's responsibility to bind the po
18       // can implement this function.
19       BindVerb string `json:"bindVerb,omitempty"`
20       // EnableHTTPS specifies whether https should be used to communicate with the extender
21       EnableHTTPS bool `json:"enableHTTPS,omitempty"`
22       // TLSConfig specifies the transport layer security config
23       TLSConfig *ExtenderTLSConfig `json:"tlsConfig,omitempty"`
24       // HTTPTimeout specifies the timeout duration for a call to the extender. Filter timeout fails the s
25       // timeout is ignored, k8s/other extenders priorities are used to select the node.
26       HTTPTimeout metav1.Duration `json:"httpTimeout,omitempty"`
27       // NodeCacheCapable specifies that the extender is capable of caching node information,
28       // so the scheduler should only send minimal information about the eligible nodes
29       // assuming that the extender already cached full details of all nodes in the cluster
30       NodeCacheCapable bool `json:"nodeCacheCapable,omitempty"`
```

Extender struct

These extender verbs will be appended to the `URLPrefix` when the `kube-scheduler` calls the extenders. The verbs for the extender call should exactly match the HTTP handler paths.

## A simple extender

The scheduler extender runs as an HTTP(S) service and exposes an accessible endpoint. We can write our implementations in any language we're familiar with. Below is a simple Golang snippet for reference. We can insert our own business logic on **line 49**.

```
 1   package main
 2
 3   import (
 4       "bytes"
 5       "encoding/json"
 6       "io"
 7       "log"
 8       "math/rand"
 9       "net/http"
10       "time"
11
12       v1 "k8s.io/api/core/v1"
13       extenderapi "k8s.io/kube-scheduler/extender/v1"
14   )
```

```
15
16  func init() {
17      rand.Seed(time.Now().UTC().UnixNano())
18  }
19
20  func Filter(w http.ResponseWriter, r *http.Request) {
21      var buf bytes.Buffer
22      body := io.TeeReader(r.Body, &buf)
23      var extenderArgs extenderapi.ExtenderArgs
24      var extenderFilterResult *extenderapi.ExtenderFilterResult
25      if err := json.NewDecoder(body).Decode(&extenderArgs); err != nil {
26          extenderFilterResult = &extenderapi.ExtenderFilterResult{
27              Error: err.Error(),
28          }
29      } else {
30          extenderFilterResult = filter(extenderArgs)
```

Code snippet of the scheduler extender

All the extender functions accept an input with the type `schedulerapi.ExtenderArgs` and return with `*schedulerapi.ExtenderFilterResult`. Within each function, we can do further implementations to satisfy our particular business needs.

The `Filter` and `Prioritize` functions are the most important ones to extend the default Kubernetes scheduler. We can append the `preempt` and `bind` functions as well when needed.

**Interactions between the `kube-scheduler` and extenders**

Now, let's see the completed interactions between the `kube-scheduler` and scheduler extenders.

The default scheduler calls the registered extenders at certain extension points by sending an HTTP request. Below, the code snippet describes how the request is sent out.

```
 1  // Helper function to send messages to the extender
 2  func (h *HTTPExtender) send(action string, args interface{}, result interface{}) error {
 3      out, err := json.Marshal(args)
 4      if err != nil {
 5          return err
 6      }
 7
 8      url := strings.TrimRight(h.extenderURL, "/") + "/" + action
 9
10      req, err := http.NewRequest("POST", url, bytes.NewReader(out))
11      if err != nil {
12          return err
13      }
14
15      req.Header.Set("Content-Type", "application/json")
16
17      resp, err := h.client.Do(req)
18      if err != nil {
19          return err
20      }
21      defer resp.Body.Close()
22
23      if resp.StatusCode != http.StatusOK {
24          return fmt.Errorf("failed %v with extender at URL %v, code %v", action, url, resp.StatusCode)
25      }
26
```

```
27        return json.NewDecoder(resp.Body).Decode(result)
28  }
```

How the scheduler calls the extender

It's worth noting that the default Kubernetes scheduler keeps rescheduling failed Pods periodically. As a result, the extenders will be called again and again until the Pods get scheduled.

## Limitations

The scheduler extender could be a good option for some scenarios, but it does have some limitations, as outline in the list below, which are why we don't implement a completed extender):

- **Large communication cost and downgraded performance:** The extender is called by HTTP(S). All the data is transferred between them. If the cluster gets lots of nodes, the request payload will have a big size as well. Moreover, all the data needs to be serialized and deserialized, which is time consuming. This will greatly affect the performance of the default scheduler.
- **Limited extension points:** From above, we can see we have only a few options for customizations.
- **Cache sharing:** In the default scheduler, it maintains a local cache to store the status of the whole cluster. Additionally, the caching can't be shared with external extenders. This brings trouble to scheduling results, especially for concurrent scheduling.

## Conclusion

Extending the default scheduler with extenders is a good choice for small clusters with low scheduling efficiency requirements. However, for large or high-performing clusters, it isn't a good idea to use scheduler extenders.

To address these limitations, a better way is provided by the Kubernetes community, known as the scheduler framework. It resolves the pain points above and is the officially recommended way to extend Kubernetes scheduling.