DEVICE MANAGER    Controller.Run()    ①

(1)    (2) ∞ fix nodesInfo

InitNodesInfo() (ONCE)

pods = dummy pods (label "kubeshare/role" = "dummyPod")
sharepods = all SharePods on the apiserver
- for pod in pods (dummypods) : {
    GPUID = labels ("kubeshare/GPUID")
    if NodesInfo[nodeName] DOESN'T EXIST {
        make nodesInfo[nodeName] with
            empty fields
Ⓐ    [  make nodesInfo[NodeName].GPUID2GPU[GPUID]
           and set default fields : UUID = " "
                                      Usage = 0.0
                                      Mem = 0
                                      PodList = empty list
    }
    else (if nodeInfo[nodeName] EXISTS, just do Ⓐ.
for sharepod in sharepods : {
    check Annotations (gpu-request, limit, mem) make sense (<1.0, non-negative, etc)
    if nodesInfo[nodeName].GPUID2GPU[GPUID] exists : {
        add PodRequest to PodList of GPUID

fix processDummyPodLater (this is in case device
                          manager is restarted
                          while Pods are still running)

nodesInfo = map[nodeName] → NodeInfo

NodeInfo = struct {
    GPUID2GPU = map[GPUID] → GPUInfo
        struct {UUID, Usage, Mem,
                              PodList}⊕
    UUID2Port = map[UUID] → string

    PodManagerPortBitMap    bitmap
    PodIP
}
⊕ PodList = List of struct PodRequest

struct PodRequest : {
    Key : string
    Request : float
    Limit : float
    Memory : int
    PodManagerPort : int
}

Controller.Run() ②

Start ConfigManager()

Server that talks to config-client on the nodes:

(1) receives list of UUIDs of ALL GPUs on the Node from config-client (establishes a UUID Database)

(2) update Annotations of NodeStatus on the apiserver
"Kubeshare/gpu-info" = GPU-UUID : mem(in MiB)

(3) update (UUID2Port) in nodesInfo[nodeName]
 + PodIp

(4) syncConfig() ; send PodList of that Node to config-client for every GPUID that belongs to that node

config-client :

check notebook Notes !!

- sends hostname, list of UUIDs

- waits for "Requests", writes PodList to special files.

handle-object()

For Pods (dummy).
extracts UUID from their logs and updates

nodesInfo[nodeName].
GPUID2GPU[UPUID].UUI[

Controller.Run() ③

processNextWorkItem()   [for sharepods on the apiserver]

↓

SyncHandler()

Input is a sharepod that has been scheduled (nodeName ≠ " ")

if node.GPUID2GPU[GPUID] exists:
add PodRequest {key, Req, Limit, Memory, Port}
to the GPU's PodList
and call syncConfig()

if it doesn't exist (UUID within

- add PodRequest to gpu-PodList
- create DummyPod (nodeName, GPUIDS)
- when dummypod gets created:
  (getobject & updates UUID of nodesInfo)
  the key of PodRequest is used
  to requeue it for
  processing

# Kubeshare-scheduler

Controller.Run()
    (WATCHES FOR WORK ITEMS!)
    ↓
processWorkItem(key)
    ↓

**bindSharePodToNode()**
[updates Spec.NodeName and Annotations "kubeshare/GPUID"] (7)

(8) → SyncHandler(key)

(1) (2) sharepod →
(3)
(4) isGPUPod=true
(5)
(6) returns { schedNode, sched GPUID }

**get sharepod based on key.**

**examine Annotations**
" kubeshare/{ gpu-request, gpu-limit, gpu-mem }"

**get nodeList, PodList, SharePodList**
**scheduleSharePod** (gpu-request, gpu-mem, sharepod, nodelist, podList)

(1) ↓

(2) return nodeResources

**nodeResources:= syncClusterResources (nodeList, podList, sharePodList)** [sync-resources.go]

**SyncNodeResources(NodeList)**

nodeResources is a map [nodeName] → struct NodeResource

struct NodeResource {
    CpuTotal    int64
    MemTotal    int64
    GpuTotal    int
    GpuMemTotal    int64
    CpuFree    int64
    MemFree    int64
    GpuFreeCount    int
    GpuFree    map[GPUID] → GPUInfo
}

↳ struct { GPUFreeReq
    GPUFreeMem
    GPUAffinityTags []
    GPUAntiAffinityTags []
    GPUExclusionTags [] }

fills in (for every node without the "NoSchedule" taint)
the NodeResources[nodeName]. GpuTotal = allocatable ("nvidia.com") GpuMemTotal is obtained from the "kubeshare/gpuinfo" Annotation
~~GpuFree~~ it doesn't touch GpuFree map

↳ device-manager creates this after getting info from config-client.

---

**apply filters (exclusion, affinity, anti-affinity)**

(3)

**ScheduleAlgorithmBestFit()**
args = (isGPUPod, gpu-request, gpu-mem, sharepod, nodeResources)

(4) delete from NodeResources the incompatible GPUs

(1) calculate cpuReqTotal, memReqTotal from the containerSpecs within the podspec.
(2) obtain gpu-request-millivalue from SharePodSpec
(3) for every node in NodeResources: call **ScheduleNode(nodeName, node...)**

(5)
(6) RETURN NodeName, GPUID

(1) check cpu, mem fit
(2) if it's a GPUPod {
    for every GPU(ID) in the nodeRes. GpuFree {
    check if req + mem fit
    find best fit (least remaining req on GPUID)
    if none found and GpuFree > 0, create a new GPUID value (just ~~an~~ a random string)

**syncPodResources()** → (that is not owned by a SharePod and is not a dummy pod)

For every Pod that is on a node in NodeResources[] and is Running: {
    update the remaining NodeResources[pod.nodeName]
    (cpu, mem, GpuFreeCount)
    pending or → that means it is already running scheduled

For every SharePod that is pending or running: {
    update the remaining NodeResources[nodeName]
(1) based on the PodSpec that is contained within the SharePodSpec [sync-resources.go : 85]
(2) check if it has the kubeshare annotations (if it is a GPUPod)
    if isGPUPod {
    check if nodeRes[nodeName]. GpuFree[GPUID] exists
    if it doesn't ↳ create GPUInfo and update
    if it exists: update GpuInfo. GPUFreeReq
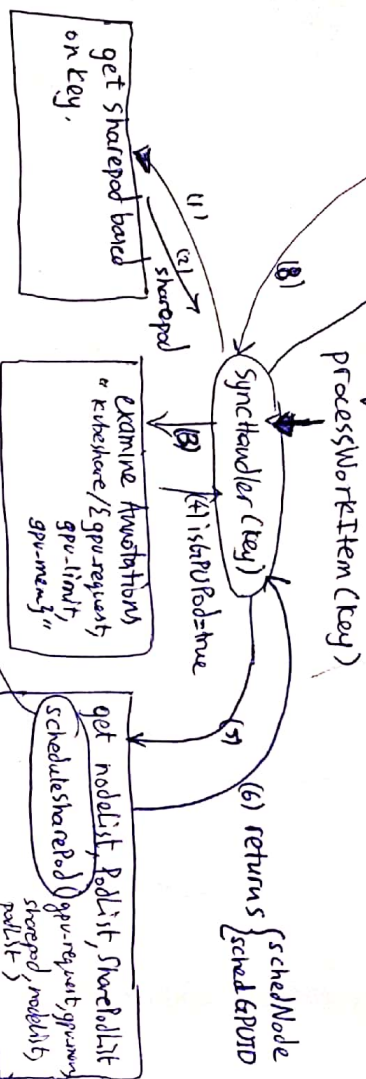    + Affinity    GpuFreeMem

Κάτι θα ξέρουμε για scheduling, αν δεν θα ξέρουμε σε κάποιο από τα υπόλοιπα
GPUID και υπάρχει GpuFreeCount > 0 (πόσοι "nvidia.com/gpu") έχουν ξεκαμφθεί
... ... σε GPUID( (...) ... θέλει να check1 θα μας update τα nodeRes.

# Kubeshare-scheduler

**controller.Run()** (WATCHES FOR WORK ITEMS!)

**findShareInfoToNode()**
[updates Spec.NodeName and Annotation "kubeshare/GPUID"]

(7) → **processWorkItem(key)**

(8) → **SyncHandler(key)**

(1') (a) sharepod → **get sharepod based on key.**

(B) → **examine Annotations**
"kubeshare/£ gpu-request,
gpu-limit,
gpu-mem £"

[4] isGPUPod=true

(5) → **get nodelist, PodList, SharePodList**
schedulesharedpod (gpu-request,gpu-mem,
sharepod, nodelist,
podlist)

(6) returns £ schedNode
£ schedGPUID

nodeResources := **SyncClusterResources** (nodelist, podlist, sharePodList) [sync-resources go]

(1) → 

(2) return nodeResources

---

Remember: a SharePodSpec, contains a normal PodSpec within it.

(3) → **apply filters** (exclusion,affinity, anti-affinity)

→ **ScheduleAlgorithmBestFit()**
args = £isGPUPod, gpu-request,gpu-mem,
sharepod, nodeResources£

(1) calculate cpuRegTotal, memRegTotal
from the containerspec within the podspec.
(2) obtain gpu-request-millivalue from
sharePodSpec.

(4) delete from NodeResources the
incompatible GPUs

(5) → (6)RETURN NodeName,GPUID

→ **ScheduleNode (nodeName,nodes)**

(1) check cpu,mem fit
(2) if fit a GPUPod £
for every GPUID in the nodes.GpuFree £
check if req + mem fit
find best fit (least remaining req on GPUID)
if none found and GpuFree >0, create a
new GPUID value (just a random string)

(3) for every node in NodeResources:
call ScheduleNode (nodeName,nodes)

---

**nodeResources is a map[nodeName]→struct NodeResource**

struct NodeResource £
| | |
|---|---|
| CpuTotal | int64 |
| MemTotal | int64 |
| GpuTotal | int |
| GpuMemTotal | int64 |
| CpuFree | int64 |
| MemFree | int64 |
| GpuFreeCount | int |
| GpuFree map[GPUID]→GPUInfo | |

↳ struct £ GPUFreeReq
GPUFreeMem
GPUAffinityTags []
GPUAntiAffinityTags []
GPUExclusionTags [] £

---

**fills in** (for every node without the "kubeshare" taint)
the NodeResources[nodeName]. GpuTotal = allocatable (unallocator)
GPUFreeCount
GpuMemTotal is obtained from
"kubeshare/gpu" Annotation

↳ **CPUFree**
it doesn't touch GpuFree map

↳ **GPUFree**
it doesn't touch GpuFree map
↳ opens-manager creates that affect
getting info from config-client

---

**SyncPodResources()** → (that's not owned by a
sharePod and is in a
dummy pod)

For every Pod that is
on a node in NodeResources []
and is Running : £
update the remaining
NodeResources[pod.nodeName]
(cpu, mem, GpuFreeCount)

For every SharePod that is running :
pending or
(1) based on the remaining NodeResources[nodeName] → that means if it's already
update the SharePec that is controlled with £ scheduled

(2) check if it has the kubeshare annotations (if it is a GPUPod)
if isGPUPod £
check if nodeList[nodeName].GpuFree[GPUID] exists
if it doesn't → create GPUID into and update it
if it exists, update GpuInto.GPUFreeReg
+ Affinity GPUFreeMem

↳ Kata in SharePod gpu scheduling , ku £a
GPUID Kα unclear GpuFreeCount > 0 (nodan "milticonfgpu") £gov £spryon
£gov £ον ερεδη το sharePod

# Kubeshare-scheduler

Remember: a SharePodSpec, contains a normal PodSpec within it.

**Controller.Run()**
(WATCHES FOR WORK ITEMS!)
↓
**processWorkItem(key)**
↓
**SyncHandler(key)**

**bindSharePodToNode()**
[updates Spec.NodeName and Annotations "kubeshare/GPUID"]
(7)
(8)

(1) (2) sharepod
**get sharepod based on key.**

(3) [4] isGPUPod=true

**examine Annotations**
"kubeshare/{gpu-request, gpu-limit, gpu-mem}"

(5) (6) returns {schedNode, schedGPUID}

**get nodeList, PodList, SharePodList**
**scheduleSharePod** (gpu-request, gpu-mem, sharepod, nodelist, podList)

**apply filters** (exclusion, affinity, anti-affinity)
(3)

(4) delete from NodeResources the incompatible GPUs

**ScheduleAlgorithmBestFit()**
args = (isGPUPod, gpu-request, gpu-mem, sharepod, nodeResources)

(1) calculate cpuReqTotal, memReqTotal from the container specs within the podspec.
(2) obtain gpu-request-millivalue from SharePodSpec.
(3) for every node in NodeResources: call ScheduleNode(nodeName, nodeRes)

(5) (6) RETURN NodeName, GPUID

(1) check cpu, mem fit
(2) if it's a GPUPod {
for every GPUID in the nodeRes. GpuFree {
check if req + mem fit
find best fit (least remaining req on GPUID)
if none found and GpuFree > 0, create a new GPUID value (just a random string)

(1)
**nodeResources := syncClusterResources(nodeList, PodList, sharePodList)** [sync-resources.go]
(2) return nodeResources

**SyncNodeResources(NodeList)**

nodeResources is a map [nodeName] → struct NodeResource

struct NodeResource {
    CpuTotal    int64
    MemTotal    int64
    GpuTotal    int
    GpuMemTotal int64
    CpuFree     int64
    MemFree     int64
    GpuFreeCount int
    GpuFree  map[GPUID] → GPUInfo
}

↳ struct { GPUFreeReq
           GPUFreeMem
           GPUAffinityTags []
           GPUAntiAffinityTags []
           GPUExclusionTags []

fills in (for every node without the "NoSchedule" taint) the NodeResources[nodeName]. GpuTotal = allocatable ("nvidia.com") GpuMemTotal is obtained from the "kubeshare/gpuinfo" Annotation
it doesn't touch GpuFree map
↳ device-manager creates this after getting info from config-client.

**syncPodResources()** → that is not owned by a sharePod and is not a dummy pod
For every Pod that is on a node in NodeResources[] and is Running: {
    update the remaining NodeResources[pod.nodeName]
    (cpu, mem, GpuFreeCount) —— pending or → that means it is already scheduled
For every SharePod that is running: {
    update the remaining NodeResources[nodeName]
(1) based on the PodSpec that is contained within the SharePodSpec [sync-resources.go:85]
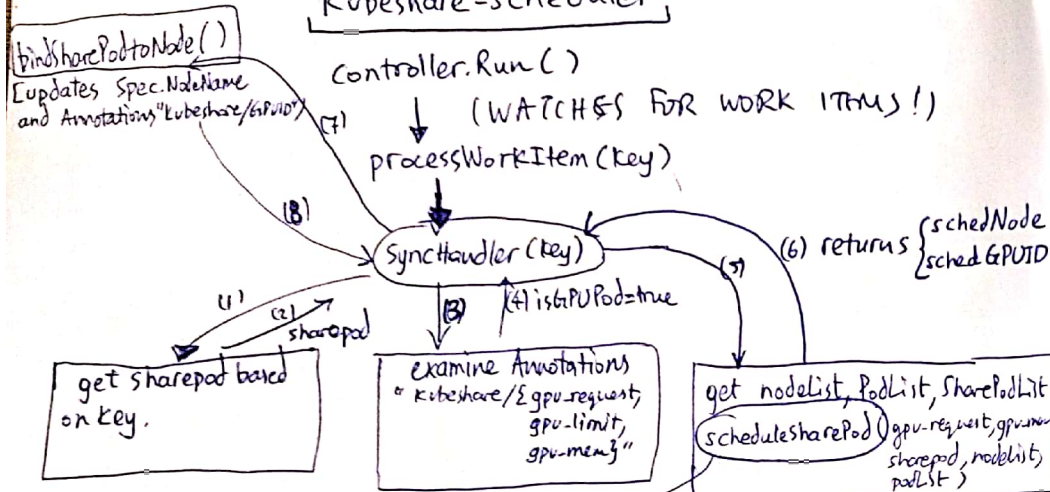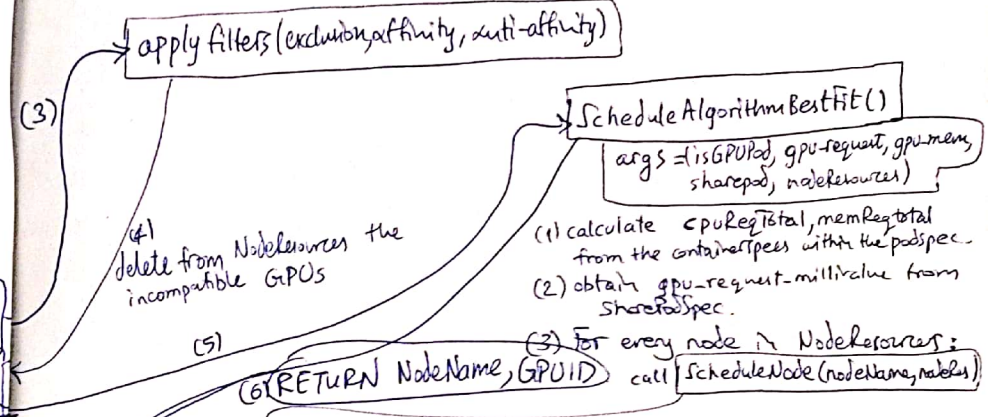(2) check if it has the kubeshare annotations (if it is a GPUPod)
    if isGPUPod {
        check if nodeRes[nodeName]. GpuFree[GPUID] exists
        if it doesn't → create GPUInfo and update it
        if it exists: update GpuInfo. GPUFreeReq
                                      GPUFreeMem
                      + Affinity

- Κάτα τη διάρκεια του scheduling, αν δω ότι δεν χωράει σε κανένα από τα υπάρχον GPUID και υπάρχει GpuFreeCount > 0 (πόσοι "nvidia.com/gpu") έχουν εκκρεμάσει τότε δημιουργώ νέο GPUID! (Μετά όταν κρύψω το sharePod θα κάνω update τα nodeRes)