# DRP: Dynamic Resource Provisioning

*Ioan Raicu, Ian Foster*

## 1.    Overview

This short report on dynamic resource provisioning is in part the result of the CEDPS meeting on October 24th-25th by the Scalable Services breakout session where the following people were present: Ian Foster, Carl Kesselman, Kate Keahey, Keith Jackson, David Konerding, Tim Freeman, Ravi Madduri, and Ioan Raicu (if I missed anyone, please feel free to add them).

After discussing three different systems (AstroPortal [1], GridFTP [2], and Condor Brick [3]) that all needed to do dynamic resource provisioning, we came to the conclusion that there are several things in common across all these three systems, and that other feature are desirable:

- Typical system components and other related components:
    - User(s) to produce work
    - Work Dispatcher
    - Resource Monitor (observes state information at the work dispatcher)
    - Resource Manager (resource provisioner based on information from Resource Monitor)
    - Worker Starter (used to bootstrap worker code by the Resource Manager)
- Coarse grained resource provisioning (i.e. GRAM, etc); typical resources allocated for hours to days
- Fine grained work dispatching (i.e. WS calls, Condor glidein, proprietary protocol based on UDP/TCP/IP, etc); typical work completion in seconds to minutes
- Variable resource pool size; increasing/decrease the number of allocated resources based on load and/or other metrics
- It would be desirable to have a generic system that could be customizable per application to provide the dynamic resource provisioning; the customization needed should be kept to a minimum, and mostly involve only configuration of the resource monitor (which state information to monitor), and the configuration of the resource manager (what policies for resource provisioning it should use depending on the resource monitor observations).

The main motivations behind having the capability for dynamic resource provisioning are:

- Allows for finer grained resource management, including the control of priorities and usage policies
- Optimize for the grid user's perspective: reduces delays on per job scheduling by utilizing pre-reserved resources
- If a metric of success for the Grid VO is resource utilization, then using a DRP will likely increase the resource utilization; note that the VO does not know of the actual state of the allocated resources (idel vs busy), so once they have been allocated by a user, they are by all practical means busy
- Opens the possibility to customize the resource scheduler per application basis, including the use of both data resource management and compute  resource management information for more efficient scheduling
- Reduced complexity to the application developer as the details of the dynamic resource provisioning are abstracted away behind a web service interface

This document has the following organization:

1) an overview of the requirements of a dynamic resource provisioning system
2) a description of the proposed architecture DRP: Dynamic Resource Provisioning (if anyone else has a better name and/or acronym, please feel free to offer suggestions)
3) a specification of the DRP Web Service interface
4) some details on possible implementations
5) usage scenarios
6) conclusions we can draw about dynamic resource provisioning

## 2.    DRP Architecture

This section covers the description of the proposed architecture DRP: Dynamic Resource Provisioning.  Figure 1 has a pictorial representation of DRP; similarly as we outlined in the overview section, the basic architectural components are with more concrete functionality as seen in Figure 1:

- **User:** User(s) to produce work

- **DRP Utilizing Web Service, Work Dispatch:** Work Dispatcher

- **DRP Web Service, Resource Monitor:** observes state information at the work dispatcher via Resource Properties exposed by the DRP Utilizing Web Service

- **DRP Web Service, Resource Manager:** Resource Provision resource via GRAM based on information from the Resource Monitor

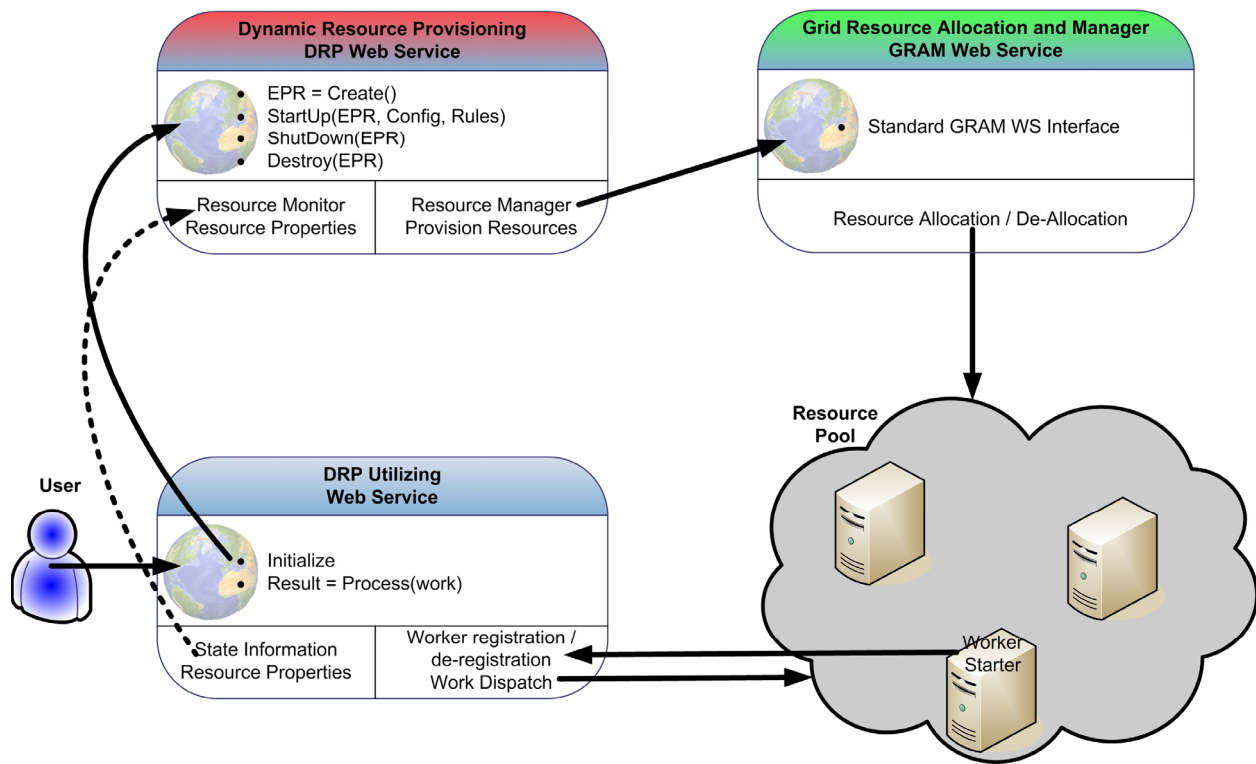- **Resource Pool, Worker Starter:** used to bootstrap worker code by the Resource Manager via GRAM



**Figure 1: DRP: Dynamic Resource Provisioning Architecture**

It is worthwhile to note that there is only one core component to DRP, specifically the DRP Web Service (DRP WS).  The DRP WS implements the necessary interface (see the next section for more details) for arbitrary web services to be able to take advantage of dynamic resource provisioning mechanisms without exposing all the details if the arbitrary web service wanted to implement the same functionality from scratch.  To make the DRP WS functional, there are two more pieces that must be in place: the worker starter (which could be absorbed by the worker code which is a component of the arbitrary web service) and the state information that is being exposed for the DRP WS Resource Monitor to read.  Therefore, if an arbitrary web service has the worker code implemented that can execute on a remote resource, then the only things that need to be modified to the arbitrary web service are:

- Initialize the DRP WS with the appropriate information necessary for the DRP WS to make the appropriate resource provisioning decisions; see the next section on what this information is

- Expose state information via Resource Properties

# 3.     DRP Web Service Interface

The DRP WS Interface has been outlined in Figure 2.

**Config**
- URL of resource allocator (i.e. GRAM)
- Host type (i.e. ia32, ia64, etc)
- Minimum number of resources to allocate
- Maximum number of resources to allocate
- Resource allocation granularity (i.e. one, all, linear, exponential, etc)
- Minimum time to allocate a resource
- Maximum time to allocate a resource
- Stage in files (files needed by the executable)
- Executable to be launched by the resource allocator (i.e. GRAM)
- Arguments t be passed to executable
- Stage out files (files that need to be transferred after the completion of the executable)

**WS Interface**
- EPR = Create()
- Boolean = StartUp(EPR, Config, Rules)
- Boolean = ShutDown(EPR)
- Boolean = Destroy(EPR)

External monitoring via resource properties

**Rules**
- List of simple resource allocation rules (value or change based)
- List of simple resource de-allocation rules (value or change based)
- List of complex resource provisioning rules (function based)

**Resource Properties**
- List of Resources Allocated
- Number of Allocated Resources
- Number of Queued Resources
- Number of De-allocated Resource
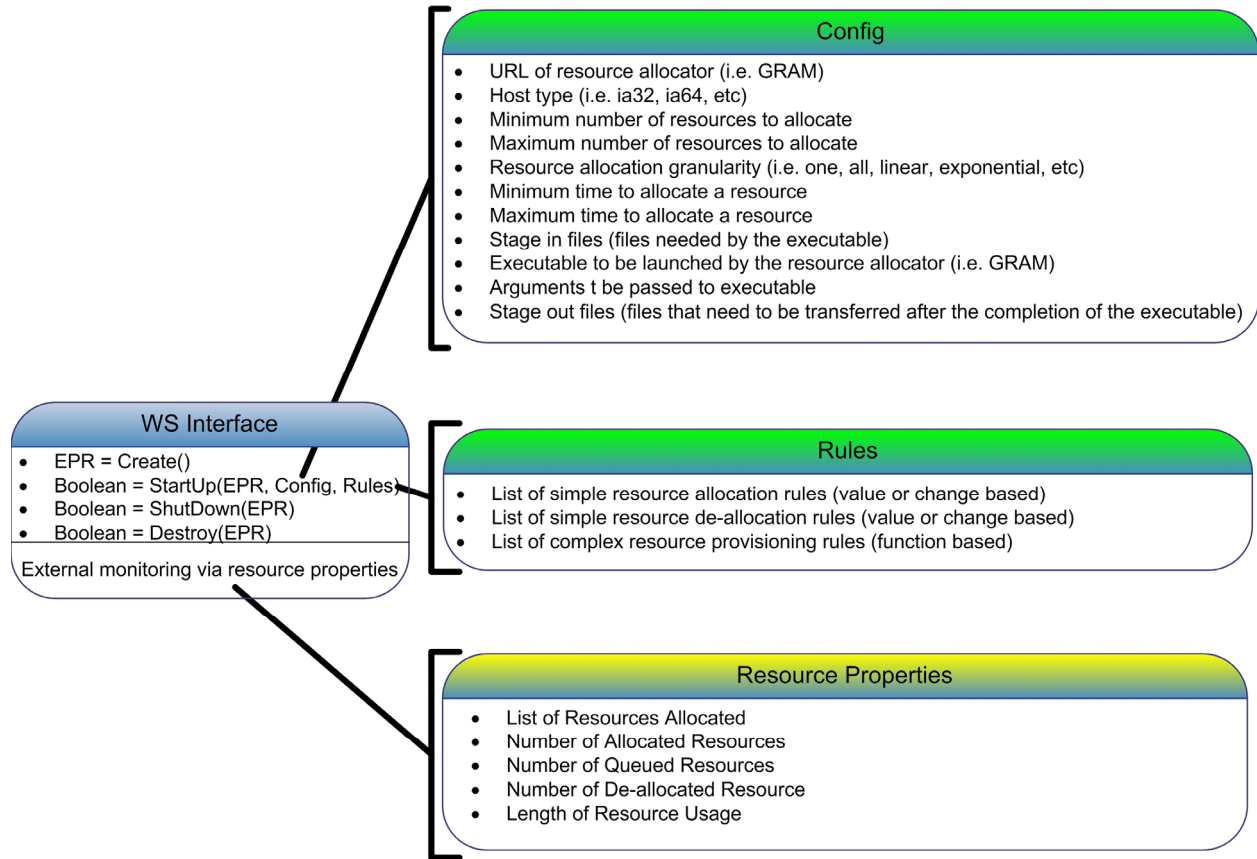- Length of Resource Usage

**Figure 2: DRP Web Service Interface**

The full WSDL definitions have been included in the following 4 pages.  There are two WSDL files, one denoting the definition of the DRP_Factory while the other containing the definition of the DRP_Service.

The DRP_Factory contains the create() interface, which returns back an EPR specific to the newly created resource.  The EPR can be used later for the rest of the WS calls to both the factory and the service.  The DRP_Factory also contains the Destroy() interface, which can be used to destroy the specified resource, which completely relinquishes all allocated resources, both remote and local pertaining to the specified EPR.

The DRP_Service contains the StartUp() interface, which is primarily used to initialize the DRP resource so it can be ready to perform the resource provisioning.  The arguments to the StartUp() call are the EPR (denoting the specific DRP resource to use, the Config object (shown in the Config section in Figure 2), and the Rules object (shown in the Rules section in Figure 2).  It might be worth elaborating more on the Rules object, and how rules can be easily derived, expressed, and effectively represented so that the desired resource provisioning takes place based on the observed state information.  For administrative purposes, or even for the DRP utilizing web service to gain knowledge about the resource provisioning state, several resource properties should be exposed by the DRP WS, as shown in the Resource Properties in Figure 2.

Please see the following pages for the specific WSDL definitions to help clarify the DRP WS Interface, the Config and Rules objects, and the resource properties exposed.

```
1    <?xml version="1.0" encoding="UTF-8"?>
2    <definitions name="DRP Factory"
3        targetNamespace="http://www.globus.org/namespaces/DRP/DRP Factory"
4        xmlns="http://schemas.xmlsoap.org/wsdl/"
5        xmlns:tns="http://www.globus.org/namespaces/DRP/DRP Factory"
6        xmlns:wsa="http://schemas.xmlsoap.org/ws/2004/03/addressing"
7        xmlns:wsdl="http://schemas.xmlsoap.org/wsdl/"
8        xmlns:xsd="http://www.w3.org/2001/XMLSchema">
9        <!--=========================================================
10
11                         T Y P E S
12
13        =========================================================-->
14   <types>
15       <xsd:schema targetNamespace="http://www.globus.org/namespaces/DRP/DRP Factory"
16           xmlns:tns="http://www.globus.org/namespaces/DRP/DRP Factory"
17           xmlns:xsd="http://www.w3.org/2001/XMLSchema">
18           <xsd:import
19               namespace="http://schemas.xmlsoap.org/ws/2004/03/addressing"
20               schemaLocation="../../ws/addressing/WS-Addressing.xsd" />
21           <!-- REQUESTS AND RESPONSES -->
22
23           <xsd:element name="createResource">
24               <xsd:complexType/>
25           </xsd:element>
26           <xsd:element name="createResourceResponse">
27               <xsd:complexType>
28                   <xsd:sequence>
29                       <xsd:element ref="wsa:EndpointReference"/>
30                   </xsd:sequence>
31               </xsd:complexType>
32           </xsd:element>
33
34       </xsd:schema>
35   </types>
36   <!--=========================================================
37
38                     M E S S A G E S
39
40        =========================================================-->
41   <message name="CreateResourceRequest">
42       <part name="request" element="tns:createResource"/>
43   </message>
44   <message name="CreateResourceResponse">
45       <part name="response" element="tns:createResourceResponse"/>
46   </message>
47   <!--=========================================================
48
49                     P O R T T Y P E
50
51        =========================================================-->
52   <portType name="FactoryPortType">
53       <operation name="createResource">
54           <input message="tns:CreateResourceRequest"/>
55           <output message="tns:CreateResourceResponse"/>
56       </operation>
57   </portType>
58   </definitions>
```

```xml
1   <?xml version="1.0" encoding="UTF-8"?>
2   <definitions name="APService"
3       targetNamespace="http://www.globus.org/namespaces/DRP/DRP Service"
4       xmlns="http://schemas.xmlsoap.org/wsdl/"
5       xmlns:tns="http://www.globus.org/namespaces/DRP/DRP Service"
6       xmlns:wsdl="http://schemas.xmlsoap.org/wsdl/"
7       xmlns:wsrlw="http://docs.oasis-open.org/wsrf/2004/06/wsrf-WS-ResourceLifetime-1.2-draft-01.wsdl"
8       xmlns:wsrp="http://docs.oasis-open.org/wsrf/2004/06/wsrf-WS-ResourceProperties-1.2-draft-01.xsd"
9       xmlns:wsrpw="http://docs.oasis-open.org/wsrf/2004/06/wsrf-WS-ResourceProperties-1.2-draft-01.wsdl
    "
10      xmlns:wsdlpp="http://www.globus.org/namespaces/2004/10/WSDLPreprocessor"
11      xmlns:xsd="http://www.w3.org/2001/XMLSchema">
12      <wsdl:import
13          namespace=
14          "http://docs.oasis-open.org/wsrf/2004/06/wsrf-WS-ResourceProperties-1.2-draft-01.wsdl"
15          location="../../wsrf/properties/WS-ResourceProperties.wsdl" />
16      <wsdl:import
17          namespace=
18          "http://docs.oasis-open.org/wsrf/2004/06/wsrf-WS-ResourceLifetime-1.2-draft-01.wsdl"
19          location="../../wsrf/lifetime/WS-ResourceLifetime.wsdl" />
20      <types>
21          <xsd:schema targetNamespace="http://www.globus.org/namespaces/DRP/DRP Service"
22              xmlns:tns="http://www.globus.org/namespaces/DRP/DRP Service"
23              xmlns:xsd="http://www.w3.org/2001/XMLSchema">
24              <!-- Requests and responses -->
25              <!-- determined number of separate allocation calls (i.e. GRAM calls) to get the desired
    number of resources -->
26              <xsd:complexType name="Granularity">
27                  <xsd:sequence>
28                      <!-- to allocate 16 resources, it would allocate them 1, 1, 1, 1, 1, 1, 1, 1, 1,
    1, 1, 1, 1, 1, 1, 1  in each allocation call until all desired resources have been allocated-->
29                      <xsd:element name="one" type="xsd:boolean"/>
30                      <!-- to allocate 16 resources, it would allocate them 16 in each allocation call
    until all desired resources have been allocated-->
31                      <xsd:element name="all" type="xsd:boolean"/>
32                      <!-- to allocate 16 resources, it would allocate them 1, 2, 3, 4, 5, 1  in each
    allocation call until all desired resources have been allocated-->
33                      <xsd:element name="linear" type="xsd:boolean"/>
34                      <!-- to allocate 16 resources, it would allocate them 1, 2, 4, 8, 1  in each
    allocation call until all desired resources have been allocated-->
35                      <xsd:element name="exponential" type="xsd:boolean"/>
36                  </xsd:sequence>
37              </xsd:complexType>
38              <!-- configuration of the values needed to communicate with the resource allocator (i.e.
    GRAM), as well as setting upper and lower bounds on the number of resources and time to allocate the
    resources -->
39              <xsd:complexType name="Config">
40                  <xsd:sequence>
41                      <xsd:element name="gramURL" type="xsd:string"/>
42                      <xsd:element name="hostType" type="xsd:string"/>
43                      <xsd:element name="minNumResourceAllocate" type="xsd:int"/>
44                      <xsd:element name="maxNumResourceAllocate" type="xsd:int"/>
45                      <xsd:element name="resourceAllocateGranularity" type="tns:Granularity"/>
46                      <xsd:element name="minTimeResourceAllocate" type="xsd:int"/>
47                      <xsd:element name="maxTimeResourceAllocate" type="xsd:int"/>
48                      <xsd:element name="stageIn" type="xsd:string"/>
49                      <xsd:element name="executable" type="xsd:string"/>
50                      <xsd:element name="arguments" type="xsd:string"
51                       minOccurs="0" maxOccurs="unbounded"/>
52                      <xsd:element name="stageOut" type="xsd:string"/>
53                  </xsd:sequence>
54              </xsd:complexType>
55              <!-- needed in the SimpleRuleAllocate to determine what operation to apply to the
    MarkValue; one element from below must be true, whie the rest must be false  -->
56              <xsd:complexType name="Operand">
57                  <xsd:sequence>
58                      <xsd:element name="lessThan" type="xsd:boolean"/>
59                      <xsd:element name="lessThanOrEqual" type="xsd:boolean"/>
60                      <xsd:element name="equal" type="xsd:boolean"/>
61                      <xsd:element name="notEqual" type="xsd:boolean"/>
62                      <xsd:element name="greaterThanOrEqual" type="xsd:boolean"/>
63                      <xsd:element name="greaterThan" type="xsd:boolean"/>
64                      <xsd:element name="difference" type="xsd:boolean"/>
65                  </xsd:sequence>
66              </xsd:complexType>
67              <!-- information relevant to access a resource property from a given service -->
68              <xsd:complexType name="RPinfo">
69                  <xsd:sequence>
70                      <xsd:element name="serviceURL" type="xsd:string"/>
71                      <xsd:element name="pollTime" type="xsd:int"/>
72                      <xsd:element name="notificationUse" type="xsd:boolean"/>
73                      <xsd:element name="name" type="xsd:string"/>
```

```
74                        </xsd:sequence>
75                    </xsd:complexType>
76                    <!-- tests between the markValue and the actual value of the resource property identified
         by the name using the markType operand  -->
77                    <xsd:complexType name="SimpleRule">
78                        <xsd:sequence>
79                            <xsd:element name="RP" type="tns:RPinfo"/>
80                            <xsd:element name="markValue" type="xsd:double"/>
81                            <xsd:element name="markType" type="tns:Operand"/>
82                        </xsd:sequence>
83                    </xsd:complexType>
84                    <!-- Complex function support (i.e. y = x^2 + 3*x + 10); this feature needs to be defined
         -->
85                    <xsd:complexType name="ComplexFunction">
86                        <xsd:sequence>
87                            <xsd:element name="undefined" type="xsd:string"/>
88                        </xsd:sequence>
89                    </xsd:complexType>
90                    <!-- Complex Rule Allocate: contains a list of resource property names that can be used
         in the ComplexFunction to determine the amount of resources to allocate  -->
91                    <xsd:complexType name="complexResourceProvision">
92                        <xsd:sequence>
93                            <xsd:element name="RPs" type="tns:RPinfo"
94                             minOccurs="0" maxOccurs="unbounded"/>
95                            <xsd:element name="function" type="tns:ComplexFunction"/>
96                        </xsd:sequence>
97                    </xsd:complexType>
98                    <!-- the collection of rules needed for resource provisioning  -->
99                    <xsd:complexType name="Rules">
100                       <xsd:sequence>
101                           <xsd:element name="simpleResourceAllocate" type="tns:SimpleRule"
102                            minOccurs="0" maxOccurs="unbounded"/>
103                           <xsd:element name="simpleResourceDeAllocate" type="tns:SimpleRule"
104                            minOccurs="0" maxOccurs="unbounded"/>
105                           <xsd:element name="complexResourceProvision" type="tns:ComplexRule"
106                            minOccurs="0" maxOccurs="unbounded"/>
107                       </xsd:sequence>
108                   </xsd:complexType>
109                   <!-- the object arguement to the startUp() WS call which starts up the DRP core logic for
         resource provisioning  -->
110                   <xsd:element name="StartUp">
111                       <xsd:complexType>
112                           <xsd:sequence>
113                               <xsd:element name="config" type="tns:Config"/>
114                               <xsd:element name="rules" type="tns:Rules"/>
115                           </xsd:sequence>
116                       </xsd:complexType>
117                   </xsd:element>
118                   <!-- the object response for the startUp() WS call return true if everything started OK,
         or false if there were any errors (i.e. invalid credentials, incorrect GRAM URL, non-existent
         executable, etc...)  -->
119                   <xsd:element name="StartUpResponse">
120                       <xsd:complexType>
121                           <xsd:sequence>
122                               <xsd:element name="valid" type="xsd:boolean"/>
123                           </xsd:sequence>
124                       </xsd:complexType>
125                   </xsd:element>
126                   <!-- the object arguement to the shutDown() WS call which stops the DRP core logic (i.e.
         de-allocate any allocated resources, etc...)  -->
127                   <xsd:element name="ShutDown">
128                       <xsd:complexType>
129                           <xsd:sequence>
130                               <xsd:element name="valid" type="xsd:boolean"/>
131                           </xsd:sequence>
132                       </xsd:complexType>
133                   </xsd:element>
134                   <!-- the object response for the shutDown() WS call return true if everything stoped OK,
         or false if there were any errors (i.e. resources that were not de-allocated, etc...)  -->
135                   <xsd:element name="ShutDownResponse">
136                       <xsd:complexType>
137                           <xsd:sequence>
138                               <xsd:element name="valid" type="xsd:boolean"/>
139                           </xsd:sequence>
140                       </xsd:complexType>
141                   </xsd:element>
142                   <!-- Resource properties -->
143                   <!-- Resource description; could include physical resource description as well (i.e. CPU
         type, CPU count, CPU speed, memory size, NIC speed, HD capacity, HD available, etc...) -->
144                   <xsd:element name="ResourceDescription">
145                       <xsd:complexType>
146                           <xsd:sequence>
```

```
147                            <xsd:element name="name" type="xsd:string"/>
148                            <xsd:element name="address" type="xsd:string"/>
149                            <xsd:element name="type" type="xsd:string"/>
150                            <xsd:element name="timeAllocated" type="xsd:string"/>
151                            <xsd:element name="timeExpire" type="xsd:string"/>
152                        </xsd:sequence>
153                    </xsd:complexType>
154                </xsd:element>
155                <xsd:element name="AllocatedResources">
156                    <xsd:complexType>
157                        <xsd:sequence>
158                            <xsd:element name="array" type="tns:ResourceDescription"
159                             minOccurs="0" maxOccurs="unbounded"/>
160                        </xsd:sequence>
161                    </xsd:complexType>
162                </xsd:element>
163                <xsd:element name="NumAllocatedResources" type="xsd:int"/>
164                <xsd:element name="NumQueuedResources" type="xsd:int"/>
165                <xsd:element name="NumDeAllocatedResources" type="xsd:int"/>
166                <xsd:element name="LengthResourceUsage" type="xsd:int"/>
167                <xsd:element name="APResourceProperties">
168                    <xsd:complexType>
169                        <xsd:sequence>
170                            <xsd:element ref="tns:AllocatedResources" minOccurs="1" maxOccurs="1"/>
171                            <xsd:element ref="tns:NumAllocatedResources" minOccurs="1" maxOccurs="1"/>
172                            <xsd:element ref="tns:NumQueuedResources" minOccurs="1" maxOccurs="1"/>
173                            <xsd:element ref="tns:NumDeAllocatedResources" minOccurs="1" maxOccurs="1"/>
174                            <xsd:element ref="tns:LengthResourceUsage" minOccurs="1" maxOccurs="1"/>
175                        </xsd:sequence>
176                    </xsd:complexType>
177                </xsd:element>
178            </xsd:schema>
179        </types>
180        <message name="StartUpInputMessage">
181            <part name="parameters" element="tns:StartUp"/>
182        </message>
183        <message name="StartUpOutputMessage">
184            <part name="parameters" element="tns:StartUpResponse"/>
185        </message>
186        <message name="ShutDownInputMessage">
187            <part name="parameters" element="tns:ShutDown"/>
188        </message>
189        <message name="ShutDownOutputMessage">
190            <part name="parameters" element="tns:ShutDownResponse"/>
191        </message>
192        <portType name="APPortType"
193            wsdlpp:extends="wsrpw:GetResourceProperty
194                     wsrlw:ImmediateResourceTermination
195                     wsrlw:ScheduledResourceTermination"
196            wsrp:ResourceProperties="tns:APResourceProperties">
197            <operation name="startUp">
198                <input message="tns:StartUpInputMessage"/>
199                <output message="tns:StartUpOutputMessage"/>
200            </operation>
201            <operation name="shutDown">
202                <input message="tns:ShutDownInputMessage"/>
203                <output message="tns:ShutDownOutputMessage"/>
204            </operation>
205        </portType>
206    </definitions>
```

## 4.     Usage Scenario

If needed, this can be filled in with how a specific application would use the DRP WS.

## 5.     Implementation Details

Once the DRP WS architecture and interface is accepted, we could start talking about actual implementation details (i.e. programming language, the use of notifications or polling, persistent state information (flat files or databases), transport level security, message level security, or no security, etc).

## 6.     Conclusion

There is not much to say here yet, other than the fact that the consensus is that this kind of functionality (i.e. DRP) is needed!  In the overview section, we covered the motivation for defining such a DRP service; to re-iterate, they were:

- Allows for finer grained resource management, including the control of priorities and usage policies
- Optimize for the grid user's perspective: reduces delays on per job scheduling by utilizing pre-reserved resources
- If a metric of success for the Grid VO is resource utilization, then using a DRP will likely increase the resource utilization; note that the VO does not know of the actual state of the allocated resources (idel vs busy), so once they have been allocated by a user, they are by all practical means busy
- Opens the possibility to customize the resource scheduler per application basis, including the use of both data resource management and compute  resource management information for more efficient scheduling
- Reduced complexity to the application developer as the details of the dynamic resource provisioning are abstracted away behind a web service interface

To be fair, we should also include some of the disadvantages that might be packaged together along with such functionality as DRP; note that theses disadvantages are mostly due from the desired functionality, and not any particular implementation of the DRP:

- All jobs submitted by different members need to map to the same user (this might be an implementation detail, but currently in all the existing prototypes of a DRP like functionality [1, 2, 3], they all map all the users to the same user account from the Grid VO's perspective); adequate logging and accounting at the DRP WS could probably give the Grid VO all the information they need for accounting purposes as long as the grid VO would trust the DRP WS
- Initial startup overhead
- Work could be halted unfinished when the original time lease on a particular resource expires if the time lease not being exposed to the work dispatcher
- From the Grid VO's perspective, if they were to actually look at the allocated resource utilization, it is likely that resources allocated via DRP like functionality would be under utilized, as allocated resources could be idle for long periods of time, yet be unavailable to be reclaimed back into the general pool of available resources; this could be resolved by the user of the DRP functionality properly setting the resource provisioning policy to release idle resources, where idle can be defined by the particular application; the problem with this is that we are then trading off between the resource responsiveness and actual resource utilization, but at least allowing the user to make this trade-off is probably the correct approach

# 7. References

There are probably more relevant references than these, but these are a good starting point as they were the center of our discussion at the Scalable Services break-out session.

[1]     Ioan Raicu, Ian Foster, Alex Szalay.  "Harnessing Grid Resources to Enable the Dynamic Analysis of Large Astronomy Datasets", SuperComputing 2006.

[2]     John Bresnahan, Ian Foster.  "An Architecture for Scavenging Compute Cluster Bandwidth," MS Thesis, University of Chicago, Department of Computer Science, 2006.

[3]     Gaurang Mehta, Carl Kesselman, Ewa Deelman.  "Dynamic Deployment of VO-specific Schedulers on Managed Resources," USC Information Sciences Institute, 2006.