

The Prudence Manual

Main text written by Tal Liron

June 30, 2010

Copyright 2009-2010 by Three Crickets LLC. This work is licensed under a Creative Commons Attribution-NonCommercial-Share

Contents

Tutorial	5
Run the Prudence Instance	5
Your First Application	5
A Static Web Page	5
A Dynamic Web Page	5
A Resource	9
Advanced Routing	12
Other Features	13
The Prudence Instance	13
Configuration by Script	13
/defaults/	13
/applications/	14
/instance/	14
/bin/	16
/logs/	16
/configuration/	16
/libraries/	16
Prudence Applications	16
Deploy by Zip	16
/web/dynamic/	17
/web/static/	17
/web/fragments/	17
/resources/	17
/libraries/	17
Configuration Scripts	17
Settings	18
Generating HTML	20
Files	20
Scriptlets	21
Mixing Languages	21
Includes	21
Life	21
Caching	22
Streaming	22
HTML forms	22
Mapping and Changing MIME Types	22
Resources	23
Life	23
CRUD	23
Implied Representation	23
Error Codes	23
Explicit Representation	23
Conditional HTTP	23
Concurrency Concerns	23
Mixing Languages	23
Accessing Resources Internally	23
Accessing External REST Services	24
Representations	24

Static Web	24
Mapping MIME Types	24
Replacing Grizzly	24
CacheControlFilter	24
JavaScriptUnifyMinifyFilter	25
CssUnifyMinifyFilter	25
Routing	25
Instance Routing	25
Application Routing	26
Custom Routing	27
URI Templates	31
API	33
application	33
document	34
executable	37
conversation	38
Sharing State	43
Debugging	45
Logging	45
The Debug Page	46
Live Viewing of Source Code	46
Breakpoints?	46
Logging	46
Loggers	47
Sending Messages	47
/configuration/logging.conf	47
Analyzing /logs/web.log	48
Administration	49
Installation	49
Customization	49
Prudence As a Daemon	50
JSW	50
YAJSW	50
HTTP Proxy	50
Perlbal	51
Apache	51
Prudence As a Restlet Container	51
Summary	51
Custom Resources	52
Custom Application	53
How to Choose a Flavor?	53
Python (Succulent!)	53
Ruby (Delectable!)	53
Clojure (Scrumptious!)	53
JavaScript (Savory!)	53
PHP (Ambrosial!)	54
Groovy (Luscious!)	54

The Case for REST	54
Resources	54
Identifiers	54
Delete	55
Read	55
Update	55
Create	55
Aggregate Resources	56
Formats	56
Shared State	56
Summary of Features	56
Let's Do It!	57
The Punchline	57
It's All About Infrastructure	57
Does REST Scale?	57
Prudence	58
Scaling Tips	58
Performance Does Not Equal Scalability	58
Caching	60
Dealing with Lengthy Requests	64
Backend Partitioning	65
Data Backends	67
Under the Hood	69
The JVM	69
Scripturian	69
Jython, JRuby, Clojure, Rhino, Quercus, Groovy	69
Restlet	69
Succinct	70
Jygments	70
H2	70
Hazelcast	70
FAQ	70
REST	70
Languages	71
Concurrency	71
Scalability	71
Performance	71
Licensing	72

Tutorial

Run the Prudence Instance

Prudence comes ready to rumble!

All you need is a Java Virtual Machine (JVM), minimally version 5. Your operating system may already have one. Typing “java” from the command line will usually tell you so. Otherwise, an excellent, open-source JVM is available from the OpenJDK project.

A Java Runtime Environment (JRE) is enough for Prudence. You need a Java Development Kit (JDK) only if you plan to write code in Java. Also, Prudence does not require anything from Java Enterprise Edition (JEE). In fact, you can see Prudence as a RESTful, minimal alternative to developing web applications under JEE.

Try Prudence! Run `/bin/run.sh` for Unix-like systems (Linux, *BSD and OS X), or `/bin/run.bat` for Windows. Prudence should declare its version and list the installed demo applications. When it announces that it is listening on port 8080, it’s ready to go. Open your web browser to `http://localhost:8080/`. You should see the Prudence Administration application, where you can access the demos.

The `/bin/run` scripts are there for getting you quickly up and running, however it is strongly recommended that you run Prudence as a daemon in production environments.

Your First Application

Your Prudence instance hosts all applications under the `/applications/` directory. To install a new application, simply create a new subdirectory there. The subdirectory name will be used as a default for various things: the base URL, logging, etc. We can change those later.

So, let’s create `/applications/wackywiki/`

If you restart Prudence, you’ll see wackywiki listed in the admin application. There’s nothing to see there quite yet, though.

By the way: It’s perfectly fine to use symbolic links or mounts to put your application subdirectories elsewhere.

A Static Web Page

Create `/applications/wackywiki/web/static/`. Files you put there will served just like from any static web server. You can put images, HTML files, CSS or anything else. Let’s start with a default web page.

`/web/static/index.html`:

```
<html>
<head>
  <title>Wacky Wiki</title>
</head>
<body>
  <div>
    Nothing to see here, for now. Carry on.
  </div>
</body>
</html>
```

Restart Prudence, and browse to `http://localhost:8080/wackiwiki/` to see the page.

A Dynamic Web Page

Important! While Prudence does dynamic web pages well, it really stands out from other web frameworks in its support for REST resources, which we’ll see in the next section. We decided to start this tutorial with web pages, because the topic would likely be more familiar to most newcomers to Prudence.

Create `/applications/wackywiki/web/dynamic/`. Unlike `/web/static/`, files in this directory must be text files (HTML, XML, plain text, etc.) They are specially processed so that they can include “scriptlets” of programming code.

Let’s move our `index.html` from `/web/static/` to `/web/dynamic/` and edit it to add some Python scriptlets.
`/web/dynamic/index.html`:

```
<html>
<head>
  <title>Wacky Wiki</title>
</head>
<body>
  <div>
    It is currently
    <%python
from datetime import datetime
now = datetime.now()
print now.strftime('%H:%M')
%>
    o’ clock , and
    <%
print now.strftime('%S')
%>
    seconds
  </div>
</body>
</html>
```

Restart Prudence, and browse to `http://localhost:8080/wackiwiki/` to see the page.

You’ll notice the `<%` and `%>` are used to delimit scriptlets, and that we’ve added “python” to the opening delimiter of the first scriptlet. Subsequent scriptlets on the page will automatically use the language of the previous scriptlets, so we don’t have to repeat “python”.

Note, too, that Python’s “print” works by adding output to the stream of text where the scriptlet is located.

Final note: Though they are conventionally called “scriptlets,” they are not “scripts” in the common meaning of interpreted code. All scriptlets are, in fact, compiled and run as JVM bytecode. That’s fast.

Expressions

Let’s make our dynamic page more readable by using special “expression scriptlets.” They’re really just a shorthand for “print.”

`/web/dynamic/index.html`:

```
<html>
<head>
  <title>Wacky Wiki</title>
</head>
<body>
  <div>
    <%python
from datetime import datetime
now = datetime.now()
%>
    It is currently <%= now.strftime('%H:%M') %> o’ clock , and <%= now.strftime('%
    </div>
  </body>
</html>
```

What happens if you have a coding error in your scriptlets? Let’s try and introduce an error on purpose.

`/web/dynamic/index.html`:

```
... from dratetime import datetime ...
```

If you try to browse to the page, Prudence will show you a detailed debug page, from where you can also access the source code with the troublesome code line highlighted. Once you deploy your application, you can cancel this debug feature. Also, you can show your own, custom, user-friendly error pages. We'll get to that below.

Fragments

Prudence lets you include others documents in-place using a special scriptlet. This makes it very easy to reuse fragments of documents (which can include scriptlets) throughout your dynamic pages. Common reusable fragments are page headers, footers, and navigation menus.

By convention in Prudence, fragments are put in `/web/fragments/`, though you can include them from anywhere, even from your regular `/web/dynamic/` or `/web/static/` subdirectories. The advantage of keeping them out of those subdirectories is that users won't be able to access the individual fragments directly. Instead, we must explicitly include them in our `/web/dynamic/` files.

Let's create a fragment that displays the current time.

`/web/fragments/time.html`:

```
<%python
from datetime import datetime
now = datetime.now()
%>
It is currently <%= now.strftime('%H:%M') %> o' clock , and <%= now.strftime('%S') %> seconds

/web/dynamic/index.html:

<html>
<head>
    <title>Wacky Wiki</title>
</head>
<body>
    <div>
        <%& 'time/' %>
    </div>
</body>
</html>
```

Note that the included scriptlet, which begins with the `<%&` delimiter, accepts a real Python expression. For example, you can do Pythony things like string interpolation: `<%& '%s-%s' % (language, encoding) %>`, allowing you to include different fragments according to changing circumstances.

Caching

Dynamically generating HTML and other text is very powerful. You may be wondering, however, how well this scales. Even if the Python code is compiled, it would still need to be run for every user request, right? Prudence supports straightforward per-page caching to greatly help you scale. By changing `document.cacheDuration` within a scriptlet, you can dynamically change how long Prudence will look for updates in the document. Note that even very short caching times can be a great boost to scalability. For example, if your cache duration is 1 second, it means that your scriptlets would only be run once every second. In computing terms, that's a very long time! Imagine that your popular web site is getting 1000 hits per second. . .

Finally, you'll be happy to know that by default Prudence also enables client-side caching. It does this by setting a standard HTTP header that tells the client how long the content should be cached. All popular desktop and mobile web browsers recognize it. This means bandwidth savings for you, because the client will not download content it already has in its cache.

Let's enable a simple 60 second cache.

`/web/dynamic/index.html`:

```
<% document.cacheDuration = 60000 %>
<html>
<head>
    <title>Wacky Wiki</title>
</head>
```

```

<body>
    <div>
        <%& 'time/' %>
    </div>
</body>
</html>

```

To test the above, go to our page and keep refreshing it in your web browser. You will see that the printed time will only be updated once every 10 seconds.

To verify that client-side caching works correctly, we recommend using the free Firefox browser with the Firebug add-on. In the Firebug network panel, you can see which requests are sent to the server. Within those 60 seconds between updates, you will see that the request for the page returns a 304 status code: “the document has not been modified.”

For an even more elaborate experiment, try running two browsers at the same time on different machines, and keep refreshing them. You’ll see each client caching individually, and yet the server still returning a different page only after the 60 second cache duration has passed.

Caching is very powerful, but obviously not useful for all pages. For example, if it’s important for a page to always show the up-to-date to the moment information, you might not be able to afford even a one second cache. It’s useful, then, to know that Prudence’s `document.cacheDuration` works at the level of the actual document, so that each included fragment can set its own cache duration. Just remember that the top document or fragment requests will determine the caching for included fragments, too. If your `index.html` has a cache duration of 60 seconds, then even if an included fragment has a cache duration of 1 second, it would only appear to update every 60 seconds.

Scriptlets can also dynamically change the cache duration according to changing circumstances. For example, under heavy load, you might want your application to cache for longer periods of time.

HTML Forms

TODO

Let’s allow users to edit our wacky wiki.

```

form = prudence.resource.request.resourceRef.queryAsForm
fresh = form.getFirstValue('fresh') == 'true'

```

/web/fragments/wiki.html:

```

<%python
wiki = 'Nothing in the wiki yet.'
%>
Wiki content is: <%= wiki %>

```

/web/dynamic/index.html:

```

<html>
<head>
    <title>Wacky Wiki</title>
</head>
<body>
    <div>
        <%& 'time/' %>
    </div>
    <div>
        <%& 'wiki/' %>
    </div>
</body>
</html>

```


More HTTP

Prudence does a lot of the HTTP work automatically, but you can manipulate the response yourself. For example, to redirect the client:

```
prudence.resource.response.locationRef = prudence.resource.request.resourceRef.baseRef
prudence.statusCode = 303
raise
```

Note the use of “raise”: it’s a handy trick to end processing of our page. Prudence recognizes that we explicitly changed the status code, so it doesn’t consider this to be a real error.

You can use also `prudence.statusCode` to explicitly set an error, for example to send 403 (“Forbidden”) if the user is not logged in.

Need more HTTP? Prudence supports HTTP cookies and HTTP authentication challenges. We won’t get into that in this tutorial, though.

Templating

You might wonder what other options are available to you other than “python” for scriptlets. You can choose other languages—Ruby, JavaScript, Clojure, PHP, etc.—if you have them installed. But, more importantly, you can choose to use one of the pre-included templating languages, Succinct and Velocity.

Templating languages have far less features than any of the above programming languages, but this can be an advantage if all you need is simple templating. They tend to be more readable for templates. Best of all, you can mix scriptlets from various languages in one page. Here’s an example of using Velocity and Python:

/web/fragments/time.html:

```
<%python
from datetime import datetime
now = datetime.now()
clock = now.strftime('%H:%M')
seconds = now.strftime('%S')
%>
<%velocity It is currently $clock o' clock, and $seconds seconds %>
```

A Resource

As stated earlier, where Prudence truly shines is in its support for REST resources. A resource is a piece of data associated with a URL, which clients can optionally create, read, update or delete. An HTML web page is an example of a resource which is usually read, and sometimes supports writing (for handling HTML forms). Prudence lets you create resources that support one or more formats, all create/read/update/delete (CRUD) operations, and advanced caching.

A resource-oriented architecture is perfect for “AJAX” and other rich client platforms (Flash, Silverlight, etc.). Another use is for a scalable backend conduit between servers. Depending on how you look at it, resources are either an online database, an API for your service, or both.

You might want to take a look at Making the Case for REST for a better understanding of REST and what it can do for you.

Let’s start with a resource that can be read for some plain text. Resources are put in the /resources/ subdirectory, and, like /web/dynamic/ and /web/static/, the filename and directory tree define the URL:

/resources/wiki.py:

```
def handle_init(conversation):
    conversation.addMediaTypeByName('text/plain')
def handle_get(conversation):
    return 'Nothing in the wiki yet.'
```

handleInit() always gets called first, before any of the other handler functions.

Try browsing to <http://localhost:8080/wackiwiki/wiki/> to see your resource.

Plain text is good for us humans to read. But, we might want to support formats that are more useful for client applications, such as JSON for browser-based JavaScript or Flash's ActiveScript.

Let's add support for JSON to our resource:

/resources/wiki.py:

```
wiki = 'Nothing in the wiki yet.'
```

```
def handle_init(conversation):
    conversation.addMediaTypeByName('text/plain')
    conversation.addMediaTypeByName('application/json')
```

```
def handle_get(conversation):
    if conversation.mediaTypeByName == 'text/plain':
        return wiki
    else:
        return '{"content": "%s"}' % wiki
```

[Add jQuery to dynamic page to test in browser.]

Just as /web/fragments/ was a useful place for us to put reusable web fragments, we can use /libraries/ for reusable code for resources, as well as dynamic web scriptlets. Let's create a library to handle our wiki:

/libraries/wiki.py:

```
wiki = 'Nothing in the wiki yet.'
```

```
def get_wiki():
    return wikiwiki = 'Nothing in the wiki yet.'
```

```
def handle_init(conversation):
    conversation.addMediaTypeByName('text/plain')
    conversation.addMediaTypeByName('application/json')
```

```
def handle_get(conversation):
    if conversation.mediaTypeByName == 'text/plain':
        return wiki
    else:
        return '{"content": "%s"}' % wiki
```

Let's make use of our library in our resource:

/resources/wiki.py:

```
from wiki import get_wiki
```

```
def handle_init(conversation):
    conversation.addMediaTypeByName('text/plain')
    conversation.addMediaTypeByName('application/json')
```

```
def handle_get(conversation):
    if conversation.mediaTypeByName == 'text/plain':
        return get_wiki()
    else:
        return '{"content": "%s"}' % get_wiki()
```

[Put, Post, Delete.]

/libraries/wiki.py:

```
from threading import RLock
```

```
wiki = 'Nothing in the wiki yet.'
```

```
wiki_lock = RLock()

def get_wiki():
    wiki_lock.acquire()
    try:
        return wiki
    finally:
        wiki_lock.release()

def set_wiki(content):
    wiki_lock.acquire()
    try:
        wiki = content
    finally:
        wiki_lock.release()
```

[Fixing our form handling from earlier using doPost.]

Remember how we supported client-side caching for our dynamic web pages? We have much more control over this process in resources. We can set explicit modification dates and ETags.

Let's give our wiki resources a modification date, so that clients can properly cache it:

/libraries/wiki.py:

```
from threading import RLock
from datetime import datetime

def to_milliseconds(dt):
    return long(mktime(dt.timetuple()) * 1000)
    wiki = {'content': 'Nothing in the wiki yet.', 'timestamp': datetime.now()}
    wiki_lock = RLock()

def get_wiki():
    wiki_lock.acquire()
    try:
        # Note that we are returning a copy of our wiki dictionary!
        return {'content': wiki['content'], 'timestamp': wiki['timestamp']}
    finally:
        wiki_lock.release()

def set_wiki(content):
    wiki_lock.acquire()
    try:
        wiki['content'] = content
        wiki['timestamp'] = datetime.now()
    finally:
        wiki_lock.release()
```

/resources/wiki.py:

```
from wiki import get_wiki, set_wiki, to_milliseconds

def handle_init(conversation):
    conversation.addMediaTypeByName('text/plain')
    conversation.addMediaTypeByName('application/json')

def handle_get(conversation):
    wiki = get_wiki()
    conversation.modificationTimestamp = to_milliseconds(wiki['timestamp'])
    if conversation.mediaTypeByName == 'text/plain':
        return wiki['content']
```

```

else:
    return '{"content": "%s"}' % wiki['content']

```

Support for modification date can go a long way towards saving us bandwidth. However, note that we still have to fetch the entire wiki in order to get its modification date. This could potentially be a costly operation, say, if the wiki is stored in a database.

Prudence lets you optimize this by letting you check only the modification date (or ETag) before actually fetching the resource.

Let's add `handleGetInfo()`:
`/resources/wiki.py`:

```

from wiki import get_wiki, set_wiki, to_milliseconds

def handle_init(conversation):
    conversation.addMediaTypeByName('text/plain')
    conversation.addMediaTypeByName('application/json')

def handle_get(conversation):
    wiki = get_wiki()
    conversation.modificationTimestamp = to_milliseconds(wiki['timestamp'])
    if conversation.mediaTypeByName == 'text/plain':
        return wiki['content']
    else:
        return '{"content": "%s"}' % wiki['content']

def handle_get_info(conversation):
    wiki = get_wiki()
    return to_milliseconds(wiki['timestamp'])

```

Our trivial example didn't offer any significant savings. But, suppose our resource involved fetching lots of wiki pages at once? We could save computing power if you we only had to fetch a modification date.

As we said earlier, our resources can be seen as an API. Indeed, there's no reason to maintain two sets of APIs, one for "the world" to access via your resources, and one for your application to use internally. You can use easily use your resources internally:

`/fragments/wiki.html`:

```

<%python
wiki = document.internal('wiki', 'application/json')
wiki = eval('(' + wiki.text + ')')
%>
Wiki content is: <%= wiki['content'] %>

```

The above is wasteful because we are turning the wiki data into JSON and then back again. For internal API uses, it's possible to support direct access to the object:

[ObjectRepresentation]

For completion, here's an example of how to access our resources from another Prudence-based application:

`/fragments/wiki.html`:

```

<%python
wiki = document.external('http://192.168.1.2/wackywiki/wiki', 'application/json')
wiki = eval('(' + wiki.text + ')')
%>
Wiki content is: <%= wiki['content'] %>

```

Advanced Routing

Capturing URL patterns.

Improve our resource above over the query params.

By the way: Difference between URL and URI.

Virtual hosting, multiple servers.
Status handlers—custom 404.
Just for fun: robots.txt, sitemap.xml.

Other Features

Logging.
Debug page.
Defrosting and pre-heating.
SSL. <http://www.restlet.org/documentation/snapshot/jse/ext/index.html?org/restlet/ext/grizzly/HttpsServerHelper.html>
Advanced representations. (Restlet extensions.)

The Prudence Instance

Configuration by Script

In Prudence, most instance and application configuration files are written in programming language code, rather than XML or other configuration formats. This method is sometimes called “bootstrapping” or just “scripting.”

Scripting is powerful. It lets your configuration be dynamic, such that the same configuration script can do different things depending on the actual deployment environment and runtime circumstances. For example, you can deploy the same instance to a development environment, in which case your instance would start various debugging processes and logs, and a production environment, in which case it would optimize for performance and scale. You can even have different optimizations for different deployments. For example, a weak cloud-based instance with a single virtual CPU could use a different HTTP engine from a dedicated box with 8 cores and a lot of RAM. With scripts, you can dynamically test for the presence of installed optional components, etc.

The default scripts are designed to fall back on other default scripts, which are all in the `/defaults/` subdirectory. The common way to override any default script is to create your own version, execute the default script first, and then apply your overrides. You can also *not* execute the default script, and instead handle things your own way. Or, you can edit the defaults directly to apply changes across the board.

The disadvantage of configuration by script is more limited portability. For example, if you write your configuration scripts in Ruby, they will not work in Prudence for Python—*unless* you manually install Python support. If you want to support all Prudence flavors, you would have to include scripts for all languages.

For changes in your configuration scripts to take effect, you will need to restart the Prudence instance.
This is true for Prudence 1.0: we plan to support runtime configuration in a future version of Prudence.

For the purposes of this manual, we’ll use “filename.*” for script names. Replace the “*” with the filename extension appropriate for your Prudence flavor: “js” for JavaScript, “rb” for Ruby, “py” for Python, etc.

`/defaults/`

Here you’ll find the default scripts used in `/instance/` and `/applications/`. Generally, you won’t want to edit the `/defaults/` scripts directly, but instead override them there. Prudence tries to make configuration easy by implementing sensible default behavior with a straightforward override mechanism: simply execute the `/defaults/` script and then apply your own code. It’s always a good idea to at least look at the `/defaults/` scripts so you can clearly understand what you’re overriding.

For an example in the JavaScript flavor, to override the default cache installed in the component, we would create an “`/instance/component.js`” script:

```
// Implement defaults
document.execute('defaults/instance/component/');

// Replace default cache with a Hazelcast-based cache
importClass(com.threecrickets.prudence.cache.HazelcastCache);
component.context.attributes.put('com.threecrickets.prudence.cache',
    new HazelcastCache());
```

Rather than explain the role of individual `/defaults/` scripts here, we will refer to them as they are used in the sections below.

/applications/

This is where you deploy your applications, each occupying a subdirectory. Your application's subdirectory name serves as a useful default name for your application. It is also the default base URL, and the default logging name. You can change any of these defaults easily. For example:

Here's an example for configuring an application's routing:

```
// Implement defaults
document.execute('defaults/application/routing/');

// Capture URLs
router.capture('profile/{id}/'), 'profile/');
router.capture('data/schedule/{id}/'), 'data/schedule/');
```

See [configuring applications \(page 16\)](#).

/instance/

This is where the Prudence instance gets initialized. The instance is the overall container which manages the servers and the applications.

/instance/default.*

This is the only required instance configuration script. It simply executes `/defaults/instance/`, which bootstraps the Prudence instance.

You'll usually prefer to override the other configuration scripts. Override `default.*` if you want to do things after everything else in the instance has been initialized.

The default script does the following:

1. Prints the Prudence welcome message to the console
2. Sets up logging, including applying `/configuration/logging.conf`
3. Executes `/instance/component/` or `/defaults/instance/component/`
4. Executes `/instance/clients/` or `/defaults/instance/clients/`
5. Executes `/instance/routing/` or `/defaults/instance/routing/`
6. Executes `/instance/servers/` or `/defaults/instance/servers/`
7. Starts the component
8. Submits bootstrap tasks for execution

/instance/component.*

The “component” in REST terminology refers to the highest-level communication abstraction: it's all “components” communicating with other “components.” Every Prudence instance is in essence a single component, which can include multiple servers and clients.

Override this script to change the way the component is created. In particular you might want to change the default cache backend or executor.

The default script does the following:

1. Creates a Restlet component
2. Configures its logService
3. Sets its statusService to be a DelegatedStatusService
4. Creates a global, fixed-thread-pool executor (accessible as “com.threecrickets.prudence.executor” in the component's context)
5. Creates an InProcessMemoryCache backend (accessible as “com.threecrickets.prudence.cache” in the component's context); see above for an example on how to configure it

/instance/clients.*

Override this script to add more clients to your component.

The default script adds file and HTTP clients.

The file client is required for static web support, and also for you to access files on your filesystem as REST resources.

The HTTP client is required for you to access external REST resources on the web.

/instance/routing.*

Override this script to change the way the component initializes its routing. Because this script delegates to `/instance/hosts/` and to your `/applications/`, it's more likely that you'll want to override those.

The default script does the following:

1. Executes `/instance/hosts/` or `/defaults/instances/hosts/`
2. Initializes all applications in `/applications/` (this list is accessible as `"com.threecrickets.prudence.applications"` in the component's context)
3. If there are no applications, Prudence quits

/instance/hosts.*

Virtual hosting allows you to serve multiple web sites from a single component. You can configure individual applications to attach one or more hosts, and to have different base URLs on each of these hosts. Prudence allows you to create virtual hosts for multiple domain names and HTTP ports, using simple expressions and wildcards. See [virtual hosts \(page 25\)](#).

By default, an "allHosts" virtual host is created, and set as the component's default host, which in turn all applications attach to by default. "allHosts" accepts all incoming requests, to any domain name, to any port.

Override this script to create additional virtual hosts for your component, and to change the default host.

For example, you might want your Prudence component to only serve requests for `"www.mysite.org"` instead of the permissive "allHosts." Or, you might want to serve multiple web sites with different sets of applications on each.

The default script does the following:

1. Creates the "allHosts" virtual host
2. Sets it as the component's default host

/instance/servers.*

Servers do the low-level work of listening to a port, accepting requests, and returning responses according to supported protocols. Prudence currently supports your choice among various HTTP server technologies: Grizzly, Jetty, Netty and an "internal" connector. To change the HTTP server technology, you have to add it to `/libraries/`. See the Restlet documentation for more information.

You can create as many servers as you need to listen on multiple ports. Remember that routing is actually handled by virtual hosts, not the servers. So, if you have servers on ports 8081, 8082, and 8083, and your applications all attach to "allHosts," then the applications will be available on all ports. To limit applications to specific ports, you will need to create additional virtual hosts. See [virtual hosts \(page 25\)](#).

Override this script to change the default port (8080) or add additional servers on other ports.

The default script does the following:

1. Creates an HTTP server on port 8080, with support for the X-FORWARDED-FOR used by proxies
2. Prints out information about this server

/bin/

These shell scripts start up the Prudence instance. Use `run.sh` for Unix-like operating systems, such as Linux, *BSD and Mac OS X, and `run.bat` for Windows.

In production environments, it's best to run [Prudence as a daemon \(page 50\)](#) in Unix-like systems or a service in Windows, via a lightweight wrapper.

What the `/bin/` scripts do is:

1. Set the JVM classpath to include all the JARs in `/libraries/`
2. Start the JVM
3. Delegate to your `/instance/` script

Note that if you're using JVM version 6 or above, you can use wildcards for specifying the classpath in `run.sh` and `run.bat`, instead of having to list individual JAR files. So, you can just use `"/libraries/*.jar"`!

/logs/

This is where your rolling logs will appear. Prudence logs are highly configurable and powerful. In particular, your `web.log` will show all requests hitting your Prudence server, using a standard format that can be consumed and analyzed by many monitoring tools.

See [logging \(page 46\)](#).

/configuration/

There are a few essential configuration files here. Most important is [logging.conf \(page 47\)](#), but you will also find some optional files to help you configure Prudence to run as a daemon, configure Velocity, Hazelcast, etc.

/libraries/

Here you will find Prudence's main libraries as well as support libraries. The main subdirectory is for Java archives (JARs), but you may find subdirectories for libraries in other languages, such as Python.

If you add your own JARs, make sure to edit the scripts in `/bin/` to accommodate your additions.

We've named or renamed all Prudence libraries according to their main Java package prefix, but you do not have to follow this convention.

Prudence Applications

Prudence applications live in the `/applications/` directory, with one subdirectory per application. The subdirectories mentioned below should all be considered as subdirectories of an application subdirectory.

Deploy by Zip

The `/applications/` directory is "zip-aware," meaning that any file placed here with the `.zip` extension will automatically be unzipped when Prudence starts up. Prudence will only unzip the file again if its modification date changes (it creates an "applications.properties" file to keep track of that). Zips can contain pre-packaged single applications, groups of applications, and even patches for individual applications. Zip-awareness makes it very easy to package, deploy and upgrade applications, and even entire sites containing many applications.

This above is true for Prudence 1.0. A future version of Prudence will feature "live" zip-awareness, so that you will not have to restart the Prudence instance in order to take new or updates zips into account.

We encourage you to follow our convention and name your deployable zips with the `".prudence.zip"` extension, for example `"wackywiki-1.1.prudence.zip"`.

/web/dynamic/

This is where you'll put your dynamic HTML files. By “dynamic” is meant that they are generated on-demand, such that each request can potentially produce a different, cacheable result. Prudence has a powerful framework for embedding programming language code into HTML as “scriptlets.” (Actually, it supports textual formats other than HTML, such as XML, RSS, JSON, etc.) See [generating HTML \(page 20\)](#).

The names of the files and subdirectories under /web/dynamic/ attach to URLs, with simple intelligence to make it easy for you to create sensible, pretty URL schemes. See [pretty URLs \(page 26\)](#) for more information.

/web/static/

This subdirectory works like a standard “static” web server. Files put here are attached to URLs and accessible to clients as is. Prudence uses non-blocking I/O for high performance, scalable delivery of data to clients.

Like many web servers, MIME types for HTTP headers are automatically assigned according to the filename extension. For example, “.html” files will be sent as “application/html”, and “.png” files will be sent as “image/png”.

The static web subdirectory will likely “just work” for you as is. See [static web \(page 24\)](#) for extras, such as adding filename-extension-to-MIME-type mappings and compressing JavaScript/CSS.

/web/fragments/

Your dynamic pages in /web/dynamic/ can include any page fragments from here. The advantage of putting them here rather than there is that here they will not be attached to URLs and be available to users.

Fragments allow you to compose complex pages out of basic building blocks. Another important use is for fine-grained caching: each fragment has its own caching behavior. See [generating HTML \(page 20\)](#).

/resources/

Whereas the /web/dynamic/ subdirectory has HTML files with embedded programming language code, /resources/ is pure code. This is Prudence's most flexible development feature: files here are attached as REST resources, capable of handling all HTTP verbs and responding with appropriate representations in any format.

From the perspective of web development, consider that if /web/dynamic/ lets you write HTML-based front ends for “thin” clients, such as simple web browsers, /resources/ lets you handle “rich” clients, such as AJAX, Flash and other dynamic platforms.

See [resources \(page 23\)](#).

/libraries/

All your code, whether it's in /resources/ or in scriptlets embedded in /web/dynamic/ and /fragments/, can include code from /libraries/.

Use whatever mechanism is appropriate for your language: “import” for Python or Ruby, “use” for Clojure, etc. For languages that don't have inclusion mechanisms—Groovy, JavaScript—you can use Prudence's inclusion mechanism, [document.execute \(page 36\)](#).

Configuration Scripts

The application's configuration scripts are in its base subdirectory.

Make sure you read [configuration by script \(page 13\)](#) for general guidelines.

/default.*

You'll rarely need to do it, but you can also override Prudence's default application bootstrap.

Here you can modify an application's media type mappings.

TODO

/settings.*

Here you can override some of Prudence’s defaults for your application, such as the subdirectory structure detailed here, the default URLs, include some distribution information, configure the logging name, etc.

You can also add your own runtime settings for your code to use, such as database usernames and passwords.

See the Settings section, below, for full detail on overrides and default settings.

/routing.*

The settings file gives you some control over the default URLs, but here you can manipulate them extensively. Your routing tools are very powerful, including redirection based on URL patterns, regular expressions, and route scoring.

In particular, this is where you install URL patterns for your resources. For example, you can attach `/item/{id}/` to your item resource, and have “id” automatically extracted from the URL.

This is also where you can attach your own custom (non-Prudence) resources to URLs. Actually, anything that’s a “restlet” will do, because Prudence uses Restlet for its resource routing. See the section on [integrating custom restlets \(page 51\)](#).

See [routing \(page 25\)](#).

/application.*

Use this to install non-Prudence Restlet applications into Prudence. By default, Prudence creates a basic Restlet application, but you can override that creation in this file. See the section on [integrating Restlets \(page 51\)](#).

Settings

Each application’s subdirectory name under the `/applications/` directory is used as a default for many settings.

In this manual, we’ve used setting names in camel case, as is used by JavaScript and Groovy. For Python, Ruby and PHP use lowercase underscore notation, for example: “`application_home_url`.” For Clojure, use lowercase hyphenated notation, for example “`application-home-url`.”

Information

These are for administrative purposes, such as the Prudence administration application, and are also used for the default error message pages.

They are directly available at runtime via [application.application \(page 34\)](#).

applicationName Defaults to the application’s subdirectory name.

applicationDescription

applicationAuthor

applicationOwner

applicationHomeURL

applicationContactEmail

Debugging

showDebugOnError Set to true to show debug information on error. See [debugging \(page 45\)](#).

showSourceCodeURL The base URL for showing source code (only relevant when `showDebugOnError` is true).

Logging

applicationLoggerName Defaults to the application’s subdirectory name.

Hosts

hosts This is a vector of vectors of two elements: the first is the virtual host to which our application will be attached, the second is the base URLs on the hosts. Specify null for the URL to default to the application's directory name.

TODO

Resources

Sets up a directory under which you can place script files that implement RESTful resources. The directory structure underneath the base directory is directly linked to the base URL.

resourcesBaseURL

resourcesBasePath

resourcesDefaultName If the URL points to a directory rather than a file, and that directory contains a file with this name, then it will be used. This allows you to use the directory structure to create nice URLs without relying on filenames.

resourcesDefrost Set this to true if you want to start to load and compile your resources as soon as Prudence starts.

resourcesSourceViewable This is so we can see the source code for scripts by adding ?source=true to the URL. You probably wouldn't want this for most applications.

resourcesMinimumTimeBetweenValidityChecks This is the time (in milliseconds) allowed to pass until a script file is tested to see if it was changed. During development, you'd want this to be low, but during production, it should be high in order to avoid unnecessary hits on the filesystem.

Dynamic Web

Sets up a directory under which you can place text files that support embedded scriptlets. Note that the generated result can be cached for better performance.

dynamicWebBaseURL

dynamicWebBasePath

dynamicWebDefaultDocument If the URL points to a directory rather than a file, and that directory contains a file with this name, then it will be used. This allows you to use the directory structure to create nice URLs that do not contain filenames.

dynamicWebDefrost Set this to true if you want to compile your scriptlets as soon as Prudence starts.

dynamicWebPreheat Set this to true if you want to load all your dynamic web documents as soon as Prudence starts.

dynamicWebSourceViewable This is so we can see the source code for scripts by adding ?source=true to the URL. You probably wouldn't want this for most applications.

dynamicWebMinimumTimeBetweenValidityChecks This is the time (in milliseconds) allowed to pass until a script file is tested to see if it was changed. During development, you'd want this to be low, but during production, it should be high in order to avoid unnecessary hits on the filesystem.

dynamicWebClientCachingMode If you set server-side caching with [document.cacheDuration \(page 34\)](#), then you can use this setting to define if client-side caching should be enabled, too. The default mode is 1:

- 0: Disabled. Client caching headers are not sent.
- 1: Conditional. The client is asked to use conditional mode HTTP to verify that the cache has not changed, via the `modificationDate` and `expirationDate` headers. This is a good default, because it generally offers most of the advantages of caching with no risks.
- 2: Offline. The client is asked to cache without verification, via the `maxAge` header. This involves some risk: if you ask to cache a page for one week, but then find out that you have a mistake in your application, then users will not see any fix you publish until their local cache expires, which can take up to a week! It's important that you understand the implications before using this mode.

It's generally safer to apply offline caching for some of your `/web/static/` resources, such as graphics and other resources. See [CacheControlFilter \(page 24\)](#). For `/resources/`, you have precise control over each header; see [conversation.modificationDate](#), [conversation.expirationDate](#) and [conversation.tag \(page 41\)](#).

Static Web

Sets up a directory under which you can place static files of any type.

Servers like Grizzly and Jetty can use non-blocking I/O to stream static files efficiently to clients.

staticWebBaseURL

staticWebBasePath

staticWebDirectoryListingAllowed If the URL points to a directory rather than a file, then this will allow automatic creation of an HTML page with a directory listing.

Preheater

preheatResources List resources here that you want heated up as soon as Prudence starts.

URL Manipulation

urlAddTrailingSlash The URLs in this array will automatically be redirected to have a trailing slash added to them if it's missing.

Predefined Globals

Values set here will be available in [application.globals \(page 33\)](#) at runtime.

Generating HTML

Prudence has good support for generating HTML by allowing you to embed programming language code in it. Despite some unique and useful features, Prudence is not so much different in this than other platforms that support dynamic HTML generation. Where Prudence stands out is in its support for REST resources, which we'll deal with in [resources \(page 23\)](#). Go ahead and skip to there if HTML generation bores you!

Files

Files under `/web/dynamic/` are all assumed to be text files. While we're titling this section "Generating HTML," these files can, in fact, be in any textual format. The filename's extension will be used to map the default MIME type for the file, though you can easily change it in code.

Create a file named `index.html` under `/web/dynamic/`. Open your web browser to `http://localhost:8080/hello/`. You should see your file rendered in the browser. Note that you did not have to enter `index.html`, though it would also work. Additionally, just `index` would work: in the `/web/dynamic/` and `/resources/` subdirectories, Prudence considers filename extensions to be optional for URLs.

Scriptlets

Programming language code can be embedded within your `/web/dynamic/` files using either `<% %>` or `<? ?>`. Note, however, that you can only use one of either in the same file.

For historical reasons, these embedded bits of code are called “scriptlets.” However, they are usually compiled, not interpreted.

Blocks can span several scriptlets (as long as you are using the same programming language: see below). For example, this will work as expected:

You can achieve the same result by printing out the text in Python code:

In fact, behind the scenes, all non-scriptlet text is turned into code. It’s a simple print of the non-scriptlet text.

A common idiom is to print out expressions interwoven with non-scriptlet text. The expression scriptlet, marked by an equal sign, can help you reduce clutter. For example:

Just remember that behind the scenes the expression scriptlet is turned into a regular scriptlet. So, you need to follow all the rules of Python expressions.

Mixing Languages

By default, Prudence will assume your scriptlets to be Python code. However, Prudence lets you scriptlets of various languages by adding the language name after the first scriptlet delimiter. For example:

You need the appropriate language JARs under your `/libs/` directory for this to work.

Once you change the language, all subsequent scriptlets will default to that language. Make sure to change back if you need to.

You might think that mixing programming languages is a bad idea in general, and only necessary for making use of legacy code. However, it can be a great idea if you consider that Prudence comes with high-performance templating languages, both Succinct and Velocity. This lets you write all the more straightforward templating as Succinct scriptlets, switching to Python scriptlets only when you need Python features. For example:

Not only is the Succinct code more readable and easier to manage, but it also performs better, is less prone to errors, and more secure. Prudence lets you use the right tool for the right job.

One tiny hiccup to be aware of is that blocks can no longer trivially span several scriptlets. For example, this won’t immediately work:

The reason is that, by necessity, Prudence must run scriptlets of different languages separately, in sequence. The language switch thus represents a boundary. But, fear not, Prudence lets you solve this problem via the “in-flow” scriptlet, marked by a colon:

This works! But how, you might wonder? Behind the scenes, the in-flow scriptlet is run from within the enclosing language. We thus do never leave the enclosing language for the purposes of running through the file. Don’t worry about performance here: in-flow scriptlets are compiled only once per file. Also, in-flow scriptlets are treated as regular scriptlets if there is no language switch.

Prudence’s terrific handling of scriptlets in multiple languages is internally handled by the Scripturian library. Scripturian was originally developed for Prudence, but its use is generic and useful enough that it merited separating it into a separate project. With Scripturian, you can easily add scriptlet power to your Java applications.

Includes

Life

As mentioned, files are only compiled once, when they are first requested by a client. From then on, each request is handled by the compiled version.

What if you edit the file? Prudence can automatically pick up your changes and recompile. This happens on the fly, while the application is running. Are you worried that this check would happen for every client request? You can easily disable this feature for production deployments, or control the minimum time interval in which Prudence assumes the file is unchanged.

Your compiled file can be run by multiple threads concurrently. It is up to you to avoid race conditions and guarantee thread safety.

Caching

You’ve chosen to use dynamically-generated HTML because you want request to be able to have different results. However, sometimes you do not expect results to change very often. For example, a home web page might display the local temperature, but it would probably be good enough to update it every hour, instead of per request.

Depending on what your scriptlets are doing, dynamically generating a web page can be very costly, and could be a performance and scalability bottleneck for your application under heavy load. You don’t want to waste energy and resources to generate results that do not change.

The solution is to cache results. Sometimes even caching for tiny time intervals can make a huge different in the ability of your application to handle load. For example, let’s say that in order to fetch the local temperature for our web page we need to query a service on the network. Without caching, every client request would result in a service query. Let’s say our web page gets 100 hits per second. Even just caching our home page for 5 seconds would throttle our service queries down to 1 every 5 seconds, vs. 500 every 5 seconds without caching. And users would get temperature readings up to 5 seconds old. It’s a negligible usability with enormous savings.

This was a trivial example. Truly scalable software requires smart caching everywhere, from the level of file and database access all the way to the generated results.

Prudence handles caching of results. Every file has its own caching interval, which you can access and change via scriptlet:

Note that you can change the cache interval dynamically. For example, you might want to increase it if you see that your machine’s CPU load is too high, or you might even want to decrease it during “rush hours” where users expect to see especially up-to-date results.

By using includes, you can fine-tune your caching strategy even more. Since each file has its own caching interval, you can fragment your page into parts with different caching sensitivity. For example, you can put the temperature reading in its own fragment, with a high cache interval, with another fragment being a “number of visitors today” counter with no cache interval, always up-to-date. The home page itself could have its own interval, or none. This setup can help you think of caching problems independently, while allowing for subtle overall schemes.

TODO: default caching time in application context

Streaming

Sometimes the client expects a slow, and long response from your application. For example, a list of bank account transactions for the past year might be hundreds of rows in length. The client does not want to wait until you’ve produced the entire result, but can start consuming results as they come.

Streaming mode lets send results to your client as you produce them.

There are risks to streaming mode. If your scriptlet code fails along the way, the stream will stop and the client might get broken results, For example, an HTML table might not get its closing `</table>` tag, leading to a broken rendering on a web browser. Of course, this is a general risk, and can happen due to connection failures along the way. It’s just something you need to be extra careful about in streaming mode. So, make sure you catch all exceptions and gracefully finish the request for the client in case of error.

HTML forms

The files in `/web/dynamic/` are meant to be sent to the clients, in response to an HTTP GET. Later on, we’ll look at files under `/resources/`, which can respond to all HTTP verbs. However, there is one case in which you might want to respond to HTTP POST and PUT directly in your `/web/dynamic/` files: HTML form submissions. Though hardly ideal, HTML forms are the most universally supported technique of getting input from clients, and are easy to set up.

Prudence’s solution is trivial: it attaches your file to HTTP POST and PUT in addition to HTTP GET. This leaves it up to you to handle each method accordingly. Here is an example:

Handling form responses via scriptlets can make your files hard to read, as they turn into a patchwork of various results all overlapped in one file. This can be cleaned up via includes:

Still, it might make more sense to remove all this form handling logic from `/web/dynamic/` to `/resources/`. We’ll show you an example of it in the next section.

Mapping and Changing MIME Types

See [prudence applications \(page 17\)](#).

By extension
Programmatically (text, XML, JSON)

Resources

Resources as API.

There's no need to maintain a separate "internal API," with REST resources serving as an "external API." Your resources are available internally, without having to use HTTP. A single, uniform API is more maintainable, and less prone to bugs. Additionally, the internal connectivity is an easy way to test your resources. See below for internal access to your resources.

Life

CRUD

An example of `handleGet` in Python:

```
def handleGet():
    id = int(prudence.resource.request.attributes.get('id'))
    session = get_session()
    try:
        note = session.query(Note).filter_by(id=id).one()
    except NoResultFound:
        return 404 finally: session.close()
    prudence.modificationTimestamp = note.timestamp
    if prudence.mediaTypeName == 'application/json':
        return json.write(note.to_dict())
    elif prudence.mediaTypeName == 'application/xml':
        return dict_to_xml(note.to_dict())
    else:
        return 415
```

Implied Representation

Error Codes

Explicit Representation

Conditional HTTP

`getInfo()`, modification dates and ETags

Concurrency Concerns

Mixing Languages

Accessing Resources Internally

Prudence has very easy-to-use REST client, and a rich API for creating and manipulation data representations. You can use it to access your own resources, as well as resources from other applications in your Prudence installation. Examples:

```
importClass(org.restlet.resource.ClientResource);

// A resource in our application
var r = new ClientResource('riap://application/hosts');
print(r.get().text);

// A resource from another application on our host
```

```
var id = 9;
var r = new ClientResource('riap://host/stickstick/note/' + id + '/');
var note = JSON.parse(String(r.get(MediaType.APPLICATION_JSON).text));
```

Accessing External REST Services

The client API mentioned above can also be used to access external resources from anywhere on the web.

Examples:

```
importClass(org.restlet.resource.ClientResource);

// An external resource
var r = new ClientResource('http://threecrickets/prudence/data/rhino');
print(r.get().text);
```

Note that for the above to work, you need HTTP configured in your clients file:

```
component.clients.add(Protocol.HTTP);
```

Representations

MIME Types

Static Web

Prudence attaches your application's `/web/static/` subdirectory to URLs, making files there available via HTTP GET. In other words, `/static/web/` is a “web server.” It's a convenient place for storing immutable resources that your clients will need to run your application: images, styles, scripts, etc.

You're likely, though, wondering if `/web/static/` is merely a convenience, and if you'd be better off using Apache or other dedicated web servers instead of Prudence to serve your the static files.

Our recommendation is to take that route only if those web servers offer important features that you won't find in Prudence. Remember that Prudence has many powerful features, including URL redirecting, rewriting and routing, and that Prudence's non-blocking I/O actually does a great job at serving files. You will likely not see many performance or scalability gains replacing `/web/static/` with standard Apache.

Mapping MIME Types

```
application.metadataService.addExtension('php', MediaType.TEXT_HTML)
```

Replacing Grizzly

Grizzly can be replaced with Netty, a non-blocking I/O HTTP server with slightly different performance characteristics, and Jetty, which offers similar performance, but with many more features. Or, use any other connector supported by Restlet.

CacheControlFilter

```
importClass(
    org.restlet.data.MediaType,
    com.threecrickets.prudence.util.CacheControlFilter);

// Cache static file types
var staticWebCacheControlFilter = new CacheControlFilter(router.context, staticWeb);
staticWebCacheControlFilter.maxAgeForMediaType.put(MediaType.TEXT_CSS, -1);
staticWebCacheControlFilter.maxAgeForMediaType.put(MediaType.APPLICATION_JAVASCRIPT, -1);
```


JavaScriptUnifyMinifyFilter

CssUnifyMinifyFilter

Routing

“Routing” refers to the decision-making process by which an incoming client request reaches its server handler. Usually, it’s information in the request itself that is used to make the decision, such as the URI, cookies, the client type, capabilities and geolocation. But routing can also take server-side and other circumstances into account. For example, a round-robin load-balancing router might send each incoming request to a different destination in sequence.

A request normally goes through many routers before reaching its handler. Filters along the way can change information in the request, which could also affect routing. Filtering is thus also a routing tool.

This abstract, flexible routing mechanism is one of Prudence’s most powerful features, but it’s important to understand its basic structure. A common misconception is that routing is based on the hierarchical structure of URIs. While it’s possible to route by parent URI in Prudence, routing is primarily to be understood in terms of the order of routers and filters along the way.

In writing applications for Prudence, you will mostly be interested in application-level routing, which we will cover in-depth below. However, to give you a better understanding of how Prudence routing works, let’s follow the journey of a request, starting with routing at the instance level.

Instance Routing

Before a request reaches your application, it is routed by your Prudence instance.

Step 1: Servers

Requests come in from servers. Prudence instances have at the minimum one server, but can have more than one. Each server listens at a particular HTTP port, and multiple servers may in turn be restricted to particular network interfaces on your machine. By default, Prudence has a single server that listen to HTTP requests on port 8080.

You can configure your servers in [/instance/servers.* \(page 20\)](#).

Step 2: The Component

There is only one component per Prudence instance, and *all* servers route to it. This allows Prudence a unified mechanism to deal with all incoming requests.

Step 3: Virtual Hosts

The component’s router decides which virtual host should receive the request. The decision is often made according the domain name in the URL, but can also take into account which server it came from. Nominally, virtual hosting is a tool to let you host multiple sites on the same Prudence instance, but it can do all kinds of subtle routing.

At the minimum you must have one virtual host. By default, Prudence has one that routes all incoming requests, no matter which server they come from. If you have multiple servers and want to treat them differently, you can create virtual hosts for each.

You can configure your virtual hosts in [/instance/hosts.* \(page 15\)](#).

Step 4: Applications

Using [“hosts” your application’s /settings.* \(page 19\)](#), you can configure which virtual hosts your application will be attached to, and the base URI for each. An application can thus be accepting requests from several virtual hosts at once.

Note that you can create a “nested” URI scheme for your applications. For example, one application might be attached at the root URI at a certain virtual host, “/”, while other applications might be at different URIs beneath the root, “/wackywiki” and “/wackywiki/support/forum”. The root application will not “steal” requests from the other applications, because the request is routed to the right application by the virtual host. The fact that the latter URI is the hierarchical descendent of the former makes no difference to the virtual host router.

A Complete Route

Let's assume a client from the internet send a request to URI "http://www.wacky.org/wackywiki/support/forum/thread/12/."

Our machine has two network interfaces, one facing the internet and one facing the intranet, and we have two servers to listen on each. This particular request has come in through the external server. The request reaches the component's router.

We have a few virtual hosts: one to handle "www.wacky.org", our organization's main site, and another to handle "support.wacky.org", a secure site where paying customers can open support tickets.

Our forum application (in the /applications/forum/ subdirectory) is attached to both virtual hosts, but at different URIs. It's at "www.wacky.org/wackywiki/support/forum" and at "support.wacky.org/forum". In this case, our request is routed to the first virtual host. Though there are a few applications installed at this virtual host, our request follows the route to the forum application.

The remaining part of the URI, "thread/12/" will be further routed inside the forum application, as detailed below.

Application Routing

Step 5: Application Handlers

Prudence applications come with default support for three kinds of handlers: resources (in the /resources/ subdirectory), dynamic web pages (in the /web/dynamic/ subdirectory) and static web resources (in the /web/static/ subdirectory). By default, all three handlers are attached at the root URI, "/", of the application (which may vary per virtual host). However, it is possible to change this in your application's /settings.*, see [resourcesBaseUrl](#), [dynamicWebBaseUrl](#) and [staticWebBaseUrl](#). For example, you may want your resources to be routed under "/rest-interface/".

You may ask, for any given request, how can the application's router know which handler to send it to, if all handlers by default assume the same "/" base URL? The answer is that it doesn't. It tries each handler in sequence, and if one handler cannot handle the request, it falls back to the next one. For example, a "style/main.css" URI will be tried first as a resource. If that resource doesn't exist, it will be tried as a dynamic web page. Finally, it will be sent to the static web handler.

Be careful: this system allows for previous handlers in the sequence to supersede later handlers. For example, if you decide to remove scriptlets from a file named /dynamic/web/about.html and move it to /static/web/about.html, make sure to delete the former, or else the latter will not be reached.

Subdirectories and Filenames As URI Segments

The three application handlers—resources, dynamic web pages and static web resources—are all routed by mapping the filesystem structure to a URI. Each subdirectory path or filename is directly translated into a URI segment.

This is exactly the scheme used by most static web servers, and it has the benefit of using a readily-available, easy-to-use hierarchical structure—the filesystem—as a straightforward way for creating URIs.

There's one deficiency to this scheme: by directly mapping filenames, it can allow for "ugly" URIs that include filename extensions. For example, you're probably used to seeing many web sites with URLs that end in ".html", ".php" and ".jsp". While these extensions are meaningful to the site developer, they complicate the URIs and expose internal implementation details to outsiders.

Pretty URIs

To allow for prettier URIs, Prudence does a few things:

Filename Extension Hiding Prudence ignores filename extensions for the purpose of mapping to URIs. For example, the file "/web/dynamic/wiki/table-of-contents.html" will be mapped to the URI "wiki/table-of-contents/".

Default Files In /web/dynamic/ and /web/static/, if a file "index.html" exists in a subdirectory, it is mapped to the URI for the subdirectory itself. In /resources/, the default file is "default.*".

For example, the file "/web/dynamic/wiki/contributors/index.html" is mapped to the URI "wiki/contributors/". Note that another way to create the same URI is to use "/web/dynamic/wiki/contributors.html" (the ".html" is hidden by Prudence).

Default files are simply another organizational option for you. In very large applications, it can help keep your files well-organized without having any effect on the URIs that exposed externally. You can even change your scheme as your application evolves.

The name “index.html” is used by web servers for archaic reasons: it was conceived of as a place where you could list the contents, or “index,” the subdirectory. These days, however, we tend to have a more general understanding of URIs. By default, we’ve followed this convention in Prudence. You can change this with the “com.threecrickets.prudence.GeneratedTextResource.defaultName” application attribute.

Trailing Slash Requirement You’ll note that we used “wiki/contributors/” for the above, rather than “wiki/-contributors” (the difference is the trailing slash). This is because Prudence requires trailing slashes by default: trying to access “wiki/table-of-contents” would permanently redirect to “wiki/table-of-contents/”.

Prudence does this for two main reasons:

1. To keep the URI space consistent, whether you use subdirectories or filenames to create the URI segments.
2. This is Prudence’s way of fighting to the “trailing slash” problem, which plagues the use of relative URIs in HTML and CSS.

You can turn off the trailing slash requirement via the “com.threecrickets.prudence.GeneratedTextResource.trailingSlashRequired” and “com.threecrickets.prudence.DelegatedResource.trailingSlashRequired” application globals.

In Prudence 1.0, the filename extension hiding and trailing slash requirement work only for resources and dynamic web pages. Static web resources still map full filenames. A future version of Prudence may extend these features to all resources.

Filename Extensions

Though Prudence hides the filename extensions from the URIs, they do have two important functions:

Filename Extensions and MIME Types Filename extensions define the default media type for “GET” requests to pages in /web/dynamic/ and /web/static/. For example, a “.xml” file will be mapped to the “application/xml” MIME type.

Every application has its own [filename extension mapping table \(page 17\)](#).

Filename Extensions and Programming Languages In /resources/, and for configuration scripts, the filename extension tells Prudence the programming language of the source code. Prudence supports “.py”, “.rb”, “.js”, “.php”, “.clj” and “.php” files.

Filename Extensions Preference In /resources/, the Prudence flavor you are using will determine which file to use in the case of ambiguity. For example, the Groovy flavor will prefer “process.gv” file even if “process.js” is present. This allows you to write and deploy applications that can run in multiple Prudence flavors. The same is true for configuration scripts.

However, in /web/dynamic/ there is no explicit preference, and behavior is undefined if you have more than one file with the same name but different extension in the same subdirectory. For example, contents.html and contents.xml.

Custom Routing

The straightforward URI routing based on directories and filenames is useful, but you’ll likely need other kinds of routing, too. That’s what [your application’s /routing.* is for \(page 18\)](#).

Hiding

The following will make sure that the given URI always returns a 404, even if the page or resource exists. For example, in /routing.* (JavaScript flavor):

```
router.hide('administration/');
```

The resource will still be available internally, via document.internal, and for capturing (see below).

Capturing

Prudence lets you “capture” arbitrary URI patterns into any internal URI you wish, in `/resources/`, `/web/dynamic/` or `/web/static/`.

You can use curly-bracket-delimited tags in the URI to parse URI segments and store them in `conversation.locals`.

For example, in `/routing.*` (JavaScript flavor):

```
// Implement defaults
document.execute('defaults/application/routing/');

router.capture('forum/help/{topic}', 'forum/help/');
```

And then, in `/web/dynamic/forum/help.html`:

```
<html>
<body>
<p>You are viewing help topic <%= conversation.locals.get('topic') %>.</p>
</body>
</html>
```

Capturing might look like redirection, but in fact it’s an *internal* redirection, similar to how the [document.internal API](#) works (The client remains entirely ignorant as to what internal URIs you might be using).

It’s important to understand this distinction: the client might be seeing an entirely different URI than your internal one. Thus, if you’re using HTML and CSS, you need to make sure that your relative references reach the right place. This is easy in Prudence with the [conversation.pathToBase API](#) (page 39), which will use the captive client URI. For example:

```
<html>
<body>
<p>You are viewing help topic <%= conversation.locals.get('topic') %>.</p>

</body>
</html>
```

Capturing in this manner is often used in conjunction with URI hiding (see above). By hiding `/forum/help/`, users would be able to access `/forum/help/faq/` via HTTP, but not `/forum/help/`. This lets you effectively control the exposed URI space.

Another feature of internal redirection is that your code you can check for internal access and enforce it. As an example, let’s implement URI hiding via a scriptlet:

```
<%
if (!conversation.internal){
    conversation.statusCode = 404; // resource not found
    conversation.stop();
}
%>
<html>
<body>
<p>You are viewing help topic <%= conversation.locals.get('topic') %>.</p>
</body>
</html>
```

Prudence also lets you capture into a different application in the instance. Refer to applications using their subdirectory name. For example:

```
// Implement defaults
document.execute('defaults/application/routing/');

router.captureOther('/forum/help/{topic}', 'wackyhelp', '/help/');
```

In addition to parsing URI segments, you can use pre-defined curly-bracket-delimited tags to route based on attributes of the request (see URI Templates, below). For example, we can route with the `{m}` tag (HTTP method) to elegantly handle HTTP forms:

```
router.capture('message/', 'message/{m}/');
```

We can then have a dynamic web page named `/message/GET.html`:

```
<html>
<body>
<form method="post">
  Enter your message:
  <p>
  <textarea name="message"></textarea>
  <p>
  <input type="submit" />
</form>
</body>
</html>
```

And one named `/message/POST.html`:

```
<html>
<body>
  <%= conversation.form.get('message') %>
</body>
</html>
```

Another example of capturing with URI templates:

```
router.capture('data/client/', 'data/client/{emt}/');
```

Because “`{emt}`” is the MIME type of data sent by the client, we can have files such as these:

- `/resources/data/client/text/html/default.js`
- `/resources/data/client/application/json/default.js`
- `/resources/data/client/default.js` (this will be used if no client data is sent, for example in an HTTP GET)

Of course, you might prefer to check the MIME type in code your code via [conversation.mediaTypeName \(page 40\)](#), but the above is a good example of how powerful routing can be with URI templates.

Capturing Errors

A special case of capturing is for errors. “Errors” can be set explicitly by you: for example, we can set `conversation.statusCode` to 404, as we did above. However, a 500 error occurs automatically for uncaught exceptions in your code.

If you have debug mode enabled, the user would see the special debug page for 500 errors; see [debugging \(page 45\)](#). On a production site, you may instead prefer to capture the 500 error and provide a friendlier page. (You’d also want to test carefully and make sure your code never throws exceptions. . .)

You can capture errors both at the application level or at the instance level, which the former taking precedence. It may be a good idea to always capture at the instance level, in case applications don’t do their own custom error pages.

Examples from `/instance/routing.*` (JavaScript flavor):

```
// Implement defaults
document.execute('defaults/instance/routing/');

// 404 errors
component.statusService.capture(404, 'wackyhelp', 'help/main/',
  component.context);

// 500 errors
component.statusService.captureHostRoot(500, 'wackywiki', 'oops/',
  component.context);
```

Examples from an application's `/routing.*`:

```
// Implement defaults
document.execute('defaults/application/routing/');

// 404 errors
applicationInstance.statusService.capture(404, 'wackyhelp', 'help/main/',
    applicationInstance.context);
```

Notes for error capturing:

- You always need to specify the application name (like `router.captureOther`).
- As with regular capturing, you can hide these pages if you don't want users to be able to access them directly. (Though, if you hide your 404 page, users would still get to it because it's what would be displayed for hidden pages! The difference is that it would appear to users with a 404 status, not a 200 success status.)
- The difference between `statusService.capture` and `statusService.captureHostRoot` is in how the base URI is set, which affects the `conversation.pathToBase` API. `statusService.capture` uses to the application's base URI on the current virtual host, while `statusService.captureHostRoot` uses to the virtual host root itself.
- For 500 error capturing, you should prefer a `/web/static/` page, which has the least chance of generating an exception and causing a 500 error again, resulting in a loop.

Dynamic Redirection

You can handle redirection at either `/resources/` or `/web/dynamic/` with the following API:

- `conversation.response.redirectSeeOther`, for HTTP status 305
- `conversation.response.redirectTemporary`, for HTTP status 307
- `conversation.response.redirectPermanent`, for HTTP status 301

For example:

```
<%
conversation.response.redirectSeeOther('../contact-us/');
conversation.stop();
%>
```

Note that most clients will support relative URI paths for redirection. Also, they will likely ignore any data in the redirected response, so you can just return null in `/resources/` or empty text in `/web/dynamic/` (which is why we call `conversation.stop` in the example above).

For other redirections in the 300 family, follow this example (JavaScript):

```
conversation.response.locationRef = 'http://wackywiki.org/report-a-bug/';
conversation.statusCode = 302;
conversation.stop();
```

Filtering

TODO

Attaching

This is the lowest-level routing API. It allows you to route URI patterns to any custom “restlet” (a REST conversation handler), as well as Restlet resources (which internally use a Finder restlet). Do your attachments in the application's `/routing.*` configuration script. Prudence offers some attachment “sugar” in addition the standard Restlet API.

Here are a few examples (JavaScript):

```
// Implement defaults
document.execute('defaults/application/routing/');

// Attach a directory instance
importClass(org.restlet.resource.Directory);
router.attach('forum/info/',
    new Directory(applicationInstance.context,
        'file:///user/info/files/'))
    .matchingMode = Template.MODE_STARTS_WITH;

// Or, you can use Prudence sugar to use MODE_STARTS_WITH
router.attachBase('forum/help/',
    new Directory(applicationInstance.context,
        'file:///user/data/files/'));

// More Prudence sugar: attach a resource via its classname
router.attach('forum/help/{topic}/', 'org.wackywiki.HelpAccessResource');

// Prudence also lets you detach restlet instances
router.detach(staticWeb);
```

Static Redirection

Static redirection is handled by attaching (see above) a Redirector restlet. By “static” here is meant that redirection is configured into the application’s `/routing.*` script. However, because attachment supports URI templates (see below), this method can actually handle complex “dynamic” redirection.

Here are a few examples (JavaScript):

```
// Implement defaults
document.execute('defaults/application/routing/');

importClass(org.restlet.routing.Redirector);

// Simple redirection
router.attach('bug/',
    new Redirector(applicationInstance.context,
        'http://wackywiki.org/contact-us/bug/'));

// With a URI template
router.attach('forum/',
    new Redirector(applicationInstance.context,
        '{rp}?debug=true'))
    .matchingMode = Template.MODE_STARTS_WITH;

// A URI template with URI segments
router.attach('contact/{reason}/',
    new Redirector(applicationInstance.context,
        'http://wackywiki.org/contact-us/{reason}'));
```

See the Restlet API documentation of Redirector for more information.

URI Templates

The custom routing techniques described above support URI templates, which are URIs with optional curly-bracket-delimited tags that are to be replaced by conversation values at runtime.

Note that the same templates are used to generate [cache keys \(page 34\)](#), and that there are some tags that can only be used for cache keys. In fact, many of these tags might be more useful for cache keys than for routing.

For more information, see the Restlet documentation.

Data Attributes

All these refer to the data (“entity”) sent by the client or that you are returning to the client. Lowercase is used for request values, uppercase for response values. We’ll note these as pairs:

- {es} or {ES}: entity size in bytes
- {emt} or {EMT}: entity media type
- {ecs} or {ECS}: entity character set
- {el} or {EL}: entity language
- {ee} or {EE}: entity encoding
- {et} or {ET}: entity tag (HTTP ETag)
- {eed} or {EED}: entity expiration date
- {emd} or {EMD}: entity modification date

Request Attributes

- {d}: date (Unix timestamp)
- {m}: the method (in HTTP, it would be “GET,” “POST,” “PUT,” “DELETE,” etc.)
- {cia}: client IP address
- {ciua}: client upstream IP address (if the request reached us through an upstream load balancer)
- {cig}: client agent name

Response Attributes

- {S}: the status code
- {SIA}: server IP address
- {SIP}: server port number
- {SIG}: server agent name

URIs

We’ll use a hyphen to show that you need to add one of the modifiers detailed after this list. For example, “{ri}” is the complete actual URI.

- {p}: the protocol (“http,” “https,” “ftp,” etc.)
- {r-} or {R-}: actual URI (the capital “R” here refers to the response, which may be different from the request if you’re redirecting)
- {h-}: virtual host URI
- {o-}: the application’s root URI on the virtual host
- {f-}: the referring URI (usually means that the client clicked a hyperlink or was redirected here)

Add the following modifiers to URI values above in order to access the various parts of the URI:

- {-i}: the complete URI
- {-a}: the authority (for URLs, this is the host or IP address)
- {-p}: the path

- `{-q}`: the query
- `{-f}`: the fragment
- `{-r}`: the remaining part of the path
- `{-e}`: the part of the path relative to the application's root URI

Every URI also has a “base” URI, accessed via the “b” modifier. Usually, this is the application's root URI on the virtual host. You can add the URI modifiers above to this URI. For example: `{rbi}`.

Conversation Locals

Tags that aren't any of the above tags will be assumed to be `conversation.locals` (page 42). You can thus inject any data you wish.

Note that URI capturing (see above) parses URIs *into* `conversation.locals`.

Cache Key Tags

The following tags are only available for generating cache keys, not for routing:

- `{dn}`: the document name (full path from the Prudence instance root)

API

Prudence exposes its API as a set of services to your source code. These services are available in scriptlets in `/web/dynamic/`, complete source code in `/resources/`, and in your configuration scripts.

Note for Python flavor: If you're using the Jepp engine, not the default Jython engine, you will need to use `get-` and `set-` notation to access attributes. For example, use `application.getArguments()` to access `application.arguments`.

Note for Ruby flavor: Our Ruby engine, JRuby, conveniently lets you use the Ruby code convention. For example, you can use `$application.get_global` rather than `$application.getGlobal`.

Note for JavaScript flavor: Our JavaScript engine, Rhino, does not provide dictionary access to maps, so you must use `get-` and `put-` notation to access map attributes. For example, use `application.globals.get('myapp.data.name')` rather than `application.globals['myapp.data.name']`.

Note for Clojure flavor: Clojure does not support bean attributes, so you will need to use `get-` and `set-` notation to access them. For example, use `(.getArguments application)` to access `application.arguments`. You can, though, use Clojure's bean function to create a read-only representation of Prudence services.

application

The same “application” service is shared between all code in a single application. Note that there is always a single application instance per application per component, even if the application is attached to several virtual hosts and servers.

The “application” service is a good place to store shared state for the application.

application.globals, application.getGlobal Application globals are general purpose attributes accessible by any code in the application.

Names can be any string, but the convention is to use “.” paths to allow for unique “namespaces” that would not overlap with future extensions, expansions or third-party libraries. For example, use “myapp.data.sourceName” rather than “dataSourceName” to avoid conflict.

Though `application.globals` is thread safe, it's important to understand how to use it properly. Make sure you read the section on concurrency in Sharing State, below.

application.arguments Available only in configuration scripts. This is a list of command-line arguments provided to the Prudence instance script.

application.application This is a reference to the underlying Restlet application instance. Here you can access some information defined in `/settings.*`, such as `application.application.owner`, `application.application.author`, `application.application.statusService.contactEmail`, etc.

For more information, refer to the Restlet API documentation.

application.logger Use the logger to print messages in the log files. The messages are prefixed by the `applicationLoggerName` service attribute, which defaults to the application's subdirectory name.

By default, all log messages from all applications will be sent to `prudence.log`, but you can change this in `/configuration/logging.conf`.

See [logging.conf \(page 47\)](#).

application.getSubLogger Uses a logger that inherits the `application.logger` configuration by default. The name you use will be appended to your base logger name with a “.”.

See [logging.conf \(page 47\)](#).

application.getMediaType Utility to get a media type from a MIME type name or filename extension.

Note that each application has its own [mappings \(page 17\)](#).

document

The “document” service has two distinct uses: first, it represents the file you are in: *this* document. This is where you can access the document's attributes and possibly change them. The second use of this service is for accessing *other* documents. Prudence combines these two uses into one service, but they functionally separate.

In the case of `/resources/` and the configuration scripts, “this document” is simply the source code file. In the case `/web/dynamic/`, it's the whole “text-with-scriptlets” page, so it is shared by all scriptlets on the page, even if they are written in different languages.

This Document

Many of these attributes have to do with caching. Caching is your best tool to make sure your application can scale well. Read more about it in 58.

document.cacheDuration (only available in `/web/dynamic/` and `/web/fragments/`)

The duration in milliseconds for which the output of this document will be cached. If this value is zero, the default, then caching is disabled. So, you must explicitly set this to a greater than zero value to enable caching. The key used to store and retrieve the cached output is determined by `document.cacheKey` (see below).

document.cacheKey (only available in `/web/dynamic/` and `/web/fragments/`)

This lets you control the key that is used to store and retrieve the cached output of the current document. *Not that this is not necessarily the key itself*, but instead a template that can contain variables that are set dynamically, and can also include values from `conversation.locals`. See “Routing” 25 for a complete list of Prudence's URI templating variables.

`document.cacheKey` is ignored if `document.cacheDuration` is zero.

Cache key templating is a very powerful feature that lets you easily create different cached versions of documents for different kinds of users and requests, but it's not always trivial to determine the best cache key for every situation. It depends strongly on how you use and cache your fragments.

The default cache key template is “`{ri}{dn}`”, which is a string containing the request identifier (the URI), a pipe character (“|”), and then the document name. An actual key could thus be: “`http://mysite.org/wackywiki/main|common/header`”. This default is sensible, because it makes sure that included fragments are cached individually. For example, only using “`{ri}`” would have each included fragment use the same key and override others.

However, though sensible, the default cache key template may not be the most efficient. For example, if the header fragment used in the example above never changes per page, then it's wasteful to cache it separately per URI. It would make more sense to set `document.cacheKey = “{dn}”` in a scriptlet in `header.html`.

Of course, you do not have to use Prudence's templating system, and can create your cacheKeys explicitly using your own scheme. Just make sure not to use “`{`” and “`}`”, which are reserved for template tags.

Rule of thumb: Set `document.cacheKey` to be as short as you possibly can. See [“Scaling Tips” \(page 58\)](#) for more information on how caching can help you scale.

Note the default cache key template can be changed by setting the `“com.threecrickets.prudence.GeneratedTextResource.defaultCacheKey”` application global.

document.cacheTags (only available in `/web/dynamic/` and `/web/fragments/`)

This is a list of one or more strings that will be attached to the cached output of the page. Any number of tags can be associated with a cache entry. Cache tags are used for `document.cache.invalidate` (see below). Note that you can set cache tags to hardcoded strings (for example: `“browse-pages”`) or dynamically generate them using code (for example: `“blog-comments-for-entry-” + blogId`).

`document.cacheTags` is ignored if `document.cacheDuration` is zero.

document.cache (note that the same cache is accessible in both `/web/dynamic/` and `/resources/`)

Provides access to the cache implementation used by this document. Prudence supports a pluggable cache backend mechanism, allowing you to use RAM, disk, database, distributed and other cache systems. It also allows for chaining of various backends together for improved performance and reliability.

Though Prudence automatically caches the output of dynamic HTML and fragments, you can use the cache as you please. Cache entries are instances of `CacheEntry`, which embeds various formatting attributes that you are free to ignore if you don’t need them.

Prudence’s default cache backend is an in-process memory cache chained before an H2 database cache. The memory cache ensures extremely fast retrieval times, while the database cache makes sure entries will persist even if you restart Prudence.

We have an [example for configuring the cache backend \(page 14\)](#).

document.cache.invalidate This lets you remove entries, zero or many, from your cache at once. It is useful for when your application state changes in such a way that certain pages must be regenerated. The argument is a cache tag, as defined by `document.cacheTags`, above.

A common use case is to invalidate display pages when a user posts new data. For a detailed example, consider a forum hosting site. The home page has a section showing “recent posts in our forums” and additionally each forum has its own front page showing “forum highlights”. Both of these query the data backend in order to generate the last, and have a 24-hour cache. You can associate each forum with cache tag “forum-X”, where X is the forum number, and associate the home page with all these cache tags. When a user posts a new forum thread in forum X, you just need to call `document.cache.invalidate(“forum-X”)` to make sure all associated pages will be regenerated on the next user request.

See [“Scaling Tips” \(page 58\)](#) for more information on how caching can help you scale.

document.source This provides access to the underlying Scripturian `DocumentSource` instance used for this document.

Scripturian is Prudence’s mechanism for loading, caching and executing source code in various languages. By default, Prudence uses a `FileDocumentSource`. From here, you can access attributes of it, for example: `document.source.basePath` and `document.source.minimumTimeBetweenValidityChecks`.

Other Documents

document.include (only available in `/web/dynamic/` and `/web/fragments/`)

Executes another document “in place,” meaning that its output is appended at the location in the document where you call `document.include`. Global variables, function definitions, class definitions, etc., in the other document would be made available locally.

The included path is an internal URI, not the external URL visible to the world. The URI can be relative to either `/web/dynamic/` or to `/web/fragments/`. The difference is that the former is open to the world, while the latter isn’t.

Prudence might compile or otherwise prepare and cache scriptlets in the included document. This means that the first time it is included it would be delayed, but subsequent includes would be much faster. To avoid that first-time wait, Prudence supports “defrosting” of your documents when it starts. This is enabled by default for all documents in `/web/dynamic/`, and would affect included documents from `/web/fragments/`. What “preparation” actually involves depends on the language used by the scriptlets.

Calling `document.include` is equivalent to using the include scriptlet, `<%& ... %>`. Internally, the include is turned into a regular scriptlet that called `document.include`.

There are three main use cases for inclusion:

1. This mechanism allows you to divide your documents into fragments that you can re-use in many documents, helping you manage large applications and keep them consistent. Fragments can include other fragments, those can include others, etc. A common strategy is to separate the document header and footer into fragments and include these in all pages.
2. Because each document fragment can have its own caching properties, fragmentation is also an important strategy for fine-grained caching. It's important to keep in mind, though, that the outermost document's `cacheDuration` will override all others. If a cached version of a document is used, then it is not executed, which means that `document.include` calls in it are not executed, too.
3. Scriptlets in the fragments can include re-usable code, such as function and class definitions. You can thus use `document.include` to include code libraries. You might want to consider, though, using `document.execute` instead (see below), as it will let you use regular source code documents and not have to use scriptlets.

document.execute Executes a program defined by source code in another document. Global variables, function definitions, class definitions, etc., in the other document would be made available locally.

The included path is an internal URI, not the external URL visible to the world. The URI can be relative to either `/web/dynamic/` in the case of scriptlets, `/resources/` in the case of resources, or to `/libraries/` in either case. The `/libraries/` subdirectory is indeed the best place to put code usable by all parts of your application.

Prudence might compile or otherwise prepare and cache code in the executed document. This means that the first time it is executed it would be delayed, but subsequent execution would be much faster. To avoid that first-time wait, Prudence supports “defrosting” of your documents when it starts. This is enabled by default for all documents in `/resources/`, and would affect documents they execute from `/libraries/`. What “preparation” actually involves depends on the language of the source code.

The main difference between `document.include` and `document.execute` is that the former expects “text-with-scriptlets” documents, while the latter uses plain source code.

Common use cases:

1. The executed code can be re-usable, such as function and class definitions. This allows you to treat it as a code library. Notes:
 - (a) In most cases, you would probably want to use your language's code inclusion mechanism instead of `document.execute`. For example, use “import” in Python, Ruby and PHP, and “require” in Clojure. The native inclusion mechanism would do a better job at caching code, managing namespaces, avoiding duplication, etc. For example, if you use `document.execute` in Clojure, then you would have to use `defonce` rather than `def` to avoid duplication in case you execute the same document multiple times in the same context.
 - (b) For JavaScript and Groovy flavors: Prudence's `document.execute` is your only option for code inclusion in Prudence, because both JavaScript and Groovy do not have a code inclusion mechanism.
2. The executed code does not have to be in the same language as the calling code. This lets you use multiple languages in your applications, using the strengths of each. Note that languages cannot normally share function and class definitions, but can share state, via mechanisms such as `application.globals`, if the languages use compatible structures.
3. Use `document.execute` as an alternative to `document.include` if you prefer not to use scriptlets. For example, executing `mylibrary.js` might be more readable than including a fragment called `mylibrary.html` that is all just one JavaScript scriptlet.

document.internal Creates a proxy for a resource in the current application, or in other applications in your component. You can perform all the usual REST verbs via the proxy: GET, PUT, POST, DELETE, etc. For more information, refer to the Restlet API documentation for `ClientResource`.

The URIs used in `document.internal` are not relative to any server or virtual host, but to the application root. Common use cases:

1. By letting you use your RESTful resources both internally and externally, `document.internal` lets you create unified resource API. `document.internal` directly accesses the resource (not via HTTP), so it is just about as fast as a simple function call, and is definitely scalable. Thus, there may be no need to create a separate set of functions for you to use internally and HTTP resources for external clients to use. A unified API would minimize the possibility of bugs and add coherence to your code, by enforcing a RESTful architecture all around. A second advantage is that you could trivially make your API remote by running it on a different Prudence instance, and using `document.external` instead (see below). This could allow for an easy way to run your application in a cluster, behind a load balancer.
2. You can create unit tests for your resources without having to start an HTTP server.
3. The defrosting mechanism uses `document.internal` to load resources.

document.external Creates a proxy for a resource. You can perform all the usual REST verbs via the proxy: GET, PUT, POST, DELETE, etc. For more information, refer to the Restlet API documentation for `ClientResource`.

Common use cases:

1. You can add scale and redundancy to your internal REST API by running several Prudence instances behind a load balancer, and using `document.external` instead of `document.internal`.
2. You can use Prudence to create an internal communication backbone for your enterprise, with various backend services exposing resources to each other. You can expose the same resources to business partners, allowing for “B2B” (business-to-business) services.
3. There are many, many use cases for `document.external`, and they keep growing as REST is adopted by service providers on the internet. There are online storage, publishing and content delivery systems, public databases, archives, geolocation services, social networking applications, etc. Perhaps, with Prudence, you will create the next one.

This is a good place to remind you that REST is not just HTTP. By default, Prudence supports `http:` and `file:` scheme URIs for `document.external`, and you can add more protocols via your [/instance/clients.* configuration script \(page 15\)](#).

Here’s an example of using `document.external` to read a JSON representation from the filesystem:

```
fixture =
    document.external('file:///myfiles/fixture.json',
        'application/json').get();
```

document.preferredExtension (only available in configuration scripts)

If multiple files with the same name but different extension exist in the same directory, then this extension will be preferred.

This value is set automatically according the Prudence flavor you are using.

executable

The “executable” is the low-level equivalent of “this document” (see above). Here you can explore which languages are installed in your Prudence instance, and gain access to their implementation mechanism. You’ll likely never need to do any of this in your Prudence application. For more information on executables, see *Scripturian*, the library that handles code execution for Prudence.

A feature that you might want to use here is the executable globals. These are similar to the application globals (see above), except that they are global to the entire Prudence instance (in fact, to the virtual machine). It’s a good place to store state that you want shared between Prudence applications.

For more information, refer to the *Scripturian* API documentation.

executable.globals, executable.getGlobal There are similar in use to `application.globals`, but are shared by all applications on the VM. See *Sharing State*, below, for more information.

executable.context The context is used for communication between the Prudence container and the executable.

- `executable.context.writer`: direct access to the output writer (writes to a memory buffer in `/web/dynamic/`, and to standard output in `/resources/` and configuration scripts)
- `executable.context.exposedVariables`: Prudence services are here (application, document, executable, conversation)
- `executable.context.attributes`: for use by the language engines

Note: The Python flavor, when using Jython, does not redirect its standard output, `sys.stdout`, to `executable.context.writer`. You must use `executable.context.writer` directly.

executable.manager This is the language manager used to create executable in many languages. Here you can query which languages are supported by the current Prudence instance (`executable.manager.adapters`).

executable.container This is identical to the “document” service detailed above.

For example, “`document.include`” is the same as “`executable.container.include`”. Internally, Prudence uses this equivalence to hook the include scriptlet, a Scripturian feature, into `document.include`, Prudence’s implementation of this feature.

conversation

The “conversation” represents the request received from the user as well your response to it, hence it’s a “conversation.” Because Prudence is RESTful, conversations encapsulate only a single request and its response. Higher level “session” management is up to you.

Here you can access various aspects of the request: the URI, formatting preferences, client information, and actual data sent with the request (the “entity”). You can likewise set response characteristics.

Note that in `/web/dynamic/` and `/web/fragments/`, “conversation” is available as a global variable. In `/resources/`, it is sent to the handling entry points (functions, closures, etc.) as an argument. Usage is otherwise identical.

Request Attributes

conversation.reference This is the URI used by the client.

A few useful attributes:

- `conversation.reference.identifier`: the complete URI
- `conversation.reference.path`: the URI, not including the domain name and the query matrix
- `conversation.reference.segments`: a list of URI segments in the path
- `conversation.reference.lastSegment`: the last segment in the URI path
- `conversation.reference.fragment`: the URI fragment (whatever follows “#”)
- `conversation.reference.query`: the URI query (whatever follows “?”); you might prefer to use `conversation.query`, instead (see below)
- `conversation.reference.relativeRef`: a new reference relative to the base URI (usually the application root URI on the current virtual host)

Refer to the Restlet API documentation for more details.

conversation.pathToBase This is a URI path relative to the base URI, which is usually the application root URI on the current virtual host.

This is very useful to allow for relative references in HTML. It's especially useful in fragments that you might include at various parts of your application, and for [captured URIs \(page 28\)](#). For example, let's say you have a contact page at `/dynamic/web/contact/index.html`. The following HTML snippet can be used anywhere in your application:

Click `<a href="<%= conversation.pathToBase %>/contact/">here` to contact us.

conversation.form, conversation.formAll Use these to access data sent via HTML forms, or any other request with data of the "application/x-www-form-urlencoded" MIME type.

`conversation.form` is a map of form names to values. In case the form has multiple values for the same name, only the last one is mapped.

`conversation.formAll` is a list of parameters. Use this if you need to support multiple values for the same name.

conversation.query, conversation.queryAll Use these to access data sent via the URI query.

`conversation.query` is a map of query parameter names to values. In case the URI has multiple values for the same name, only the last one is mapped.

`conversation.queryAll` is a list of parameters. Use this if you need to support multiple values for the same name.

conversation.entity This is data sent by the client. For more information, refer to the Restlet API documentation for Representation.

In the case of an HTML form, it would be more convenient `conversation.form` to access the parsed entity data. For other kinds of data, you have to parse the entity data yourself. Before attempting to parse entity data on your own, make sure to look through the Restlet API and its extensive set of plugins for tools to help you parse representations. Plugins exist for many common internet formats.

A few useful representation attributes:

- `conversation.entity.size`: the size of the data in bytes, or -1 if unknown
- * `conversation.entity.text`: the data as text (only useful if the data is textual)
- * `conversation.entity.reader`: an open JVM Reader to the data (only useful if the data is textual)
- * `conversation.entity.stream`: an open JVM InputStream to the data (useful for binary data)

Note: Client data is provided as a stream that can only be "consumed" once. Attributes that cause consumption are marked with a "*" above. *Note that `conversation.entity.text` is one of them!* If you want to access `conversation.entity.text` more than once, save it to a variable first.

conversation.variant This allows low-level access to the Restlet Variant instance sent by the client. In negotiated HTTP mode, the "variant" encapsulates a set of preferences the client might have for returned representations. You usually don't have to access the variant: Prudence expands these via `conversation.mediaType`, `conversation.characterSet`, `conversation.encoding`, and `conversation.language`. Refer to the Restlet API documentation for details.

Note that in non-negotiated HTTP mode, Prudence "passes through" the client entity as a variant.

Request and Response Attributes

The attributes are initially set according to the client request, and by default are passed through as is to the response. For example, if a client sends you data in `mediaType` "application/json," then the default would be for the response to be in the same format. Obviously, this is not always the desired outcome, and you can explicitly examine and change these attributes.

The one exception to this rule is HTML forms. The data comes from the client as MIME type "application/x-www-form-urlencoded," but it's highly unlikely that you'd want to return an HTML form back to the client. For this case only, Prudence sets the `mediaType` to "application/html". You can explicitly change it something else if you need to.

Note that in `/resources/` you always have the option of explicitly returning a Representation instance to the client, in which case none of these attributes will be used for the response.

conversation.cookies This is initialized as a list (not a map) of cookies sent from the client.

If you want to ask the client to change any of them, be sure to call `save()` on the cookie in order to send it in the response. You can also call `delete()` to ask the client to delete the cookie (no need to call `save` in that case; internally sets `maxAge` to zero). Note that you can call `save()` and `delete()` as many times as you like, and that only the last changes will be sent in the response.

Note that you can only *ask* a client to change, store cookies, or for them to be used in various. It's up to the client to decide what to do with your requirements. For example, many web browsers allow users to turn off cookie support or filter out certain cookies.

Cookies have the following attributes:

- `cookie.name`: (read only)
- `cookie.version`: (integer) per a specific `cookie.name`
- `cookie.value`: textual, or text-encoded binary data (note that most clients have strict limits on how much total data is allowed to be stored in all cookies per domain)
- `cookie.domain`: the client should only use the cookie with this domain and its subdomains (web browsers will not let you set a cookie for a domain which is not the domain of the request or a subdomain of it)
- `cookie.path`: the client should only use the cookie with URIs that begin with this path (“/” would mean to use it with all URIs)

The following attributes are not received from the client, but you can set them for sending to the client:

- `cookie.maxAge`: age in seconds, after which the client should delete the cookie. `maxAge=0` deletes the cookie immediately, while `maxAge=-1` (the default) asks the client to keep the cookie only for the duration of the “session” (this is defined by the client; for most web browsers this means that the cookie will be deleted when the browser is closed).
- `cookie.secure`: true if the cookie is meant to be used only in secure connections (defaults to false)
- `cookie.accessRestricted`: true if the cookie is meant to be used only in authenticated connections (defaults to false)
- `cookie.comment`: some clients store this, some discard it

conversation.createCookie You must provide the cookie name as an argument. Returns a new cookie instance if the cookie doesn't exist yet, or the existing cookie if it does.

For new cookies, be sure to call `save()` on the cookie in order to send it in the response, thus asking the client to create it, or `delete()` if you want to cancel the creation (in which case nothing will be sent in the response).

conversation.mediaTypeName, conversation.mediaTypeExtension, conversation.mediaType These three variants all represent the same value, letting you access the value in different ways.

`conversation.mediaTypeName` is the MIME representing the media type. MIME (Multipurpose Internet Mail Extensions) is an established web standard for specifying media types. Examples include: “text/plain,” “text/html,” “application/json,” and “application/x-www-form-urlencoded.” The exact list of supported MIME types depends on the underlying Restlet implementation.

`conversation.mediaTypeExtension` is the media type as the default filename extension for the media type. For example, “txt” is equivalent to MIME “text/plain,” and “xml” is equivalent to “application/xml.” Each application has its own mappings of filename extensions to media types. See also `application.getMediaType`, above, and how to change the [mappings for your application \(page 17\)](#).

`conversation.mediaType` is the underlying Restlet Media Type instance. Refer to the Restlet API documentation for details.

To find out which media types the client accepts and prefers, use `conversation.request.clientInfo.acceptedMediaTypes` and `conversation.request.clientInfo.getPreferredMediaType()`.

In `/web/dynamic/` and `/web/static/`, if the client request does not contain any data representation, for example in a “GET” request, then the media type's default value will be set according to the filename extension (page 27).

conversation.characterSetName, conversation.characterSetShortName, conversation.characterSet These three variants all represent the same value, letting you access the value in different ways.

`conversation.characterSetName` is ISO's UTC (Universal Character Set) name of the character set. For example, "ISO-8859-1" is the "Latin 1" character set and "UTF-8" is the 8-bit Unicode Transformation Format, "US-ASCII" is ASCII, etc. The exact list of supported ISO names depends on the underlying Restlet implementation.

`conversation.characterSetShortName` is a shortcut name for the character set. Shortcuts include "ascii," "utf8," and "win" (for the Windows 1252 character set). Restlet handles shortcuts names together with [filename extension mappings](#). (page 20)

`conversation.characterSet` is the underlying Restlet `CharacterSet` instance. Refer to the Restlet API documentation for details.

To find out which character sets the client accepts and prefers, use `conversation.request.clientInfo.acceptedCharacterSets` and `conversation.request.clientInfo.getPreferredCharacterSet()`.

If the client request does not specify a character set, then the character set will fall back to a default determined by the `"com.threecrickets.prudence.GeneratedTextResource.defaultCharacterSet"` application global. If not explicitly set, it will be UTF-8.

Response Attributes

conversation.statusCode, conversation.status These two variants represent the same value, letting you access the value in different ways.

`conversation.statusCode` is an HTTP status code as an integer.

`conversation.status` is the underlying Restlet Status instance. Refer to the Restlet API documentation for details.

conversation.languageName, conversation.language These two variants represent the same value, letting you access the value in different ways.

`conversation.languageName` is the IETF locale name for the language. Examples include "en" for English, "en-us" for USA English, "fr" for French, etc. The exact list of supported IETF names depends on the underlying Restlet implementation.

`conversation.language` is the underlying Restlet Language instance. Because IETF names are hierarchical, you might prefer to use this as a way to test for containment. For example, `conversation.language.includes` will tell you that "en-us" is included in "en." Refer to the Restlet API documentation for details.

To find out which languages the client accepts and prefers, use `conversation.request.clientInfo.acceptedLanguages` and `conversation.request.clientInfo.getPreferredLanguage()`.

Note that the language can be a null value. Responses do not have to specify a language.

conversation.encodingName, conversation.encoding

This value is currently not being used.

These two variants represent the same value, letting you access the value in different ways.

`conversation.encodingName` is an internal name used for the encoding. Examples include "zip," "gzip," "compres," and "deflate." The "*" name represents all possible encodings. The exact list of supported names depends on the underlying Restlet implementation.

`conversation.encoding` is the underlying Restlet Language instance. Refer to the Restlet API documentation for details.

To find out which encodings the client accepts and prefers, use `conversation.request.clientInfo.acceptedEncodings` and `conversation.request.clientInfo.getPreferredEncoding()`.

Response Attributes for Resources

Note that `modificationTimestamp`, `expirationTimestamp` and `maxAge` are set indirectly for `/web/dynamic/`, via `document.cacheDuration`. See the [dynamicWebClientCachingMode setting](#) (page 20).

conversation.modificationTimestamp, conversation.modificationDate These two variants represent the same value, letting you access the value in different ways.

`conversation.modificationTimestamp` is a long integer value representing the number of milliseconds since January 1, 1970, 00:00:00 GMT ("Unix time"). While not too useful in itself, it's easy to compare timestamps to see which moment in time is before the other.

`conversation.modificationDate` is the underlying JVM Date instance. Refer to the Java API documentation for details.

conversation.expirationTimestamp, conversation.expirationDate These two variants represent the same value, letting you access the value in different ways.

`conversation.expirationTimestamp` is a long integer value representing the number of milliseconds since January 1, 1970, 00:00:00 GMT (“Unix time”). While not too useful in itself, it’s easy to compare timestamps to see which moment in time is before the other.

`conversation.expirationDate` is the underlying JVM Date instance. Refer to the Java API documentation for details.

conversation.httpTag, conversation.tag These two variants represent the same value, letting you access the value in different ways.

`conversation.httpTag` is an HTTP ETag string.

`conversation.tag` is the underlying Restlet Tag instance. Refer to the Restlet API documentation for details.

conversation.maxAge The maximum number of seconds for which clients would cache the response, if they support caching. Defaults to zero.

See the [dynamicWebClientCachingMode](#) setting (page 20).

conversation.addMediaTypeByName, conversation.addMediaTypeByExtension, conversation.addMediaType TODO

Conversation Flow

conversation.stop Throws an exception, thereby ending execution of your code, and hence the conversation (unless you have deferred it: see `conversation.defer`, below). Note that the client will still get a response, so you can set attributes (`conversation.statusCode`, `conversation.expirationTimestamp`, etc.) before calling `conversation.stop`.

conversation.internal True if the client’s request was internal, false if it was external.

Internal requests are usually created in one of two ways:

1. The `document.internal` API (see above)
2. [URI capturing](#) (page 28)

conversation.locals A map of general purpose attributes that is destroyed at the end of the conversation. Importantly, `conversation.locals` are maintained even if the conversation has been deferred via `conversation.defer` (see below).

See [Sharing State](#), below, for more information.

conversation.defer (only available in `/web/dynamic/` and `/web/fragments/`)

Releases the current conversation thread, and queues handling of this conversation on a separate thread pool. When the conversation turn comes to be handled, it will cause the page to be executed again, but with `conversation.deferred` (see below) set to true. Use `conversation.locals` if you want to pass state for the deferred execution.

Returns true if indeed the conversation has been successfully deferred. Will return false if the conversation is already deferred.

Note that calling `conversation.defer` does not stop the current execution. You’d likely follow a successful call to `conversation.defer` with a call to `conversation.stop`. For example:

```
if (conversation.defer()) {
    conversation.stop();
}
```

This is an experimental feature in Prudence 1.0. The use of a separate thread pool is only supported when using the internal Restlet connector. For other connectors (such as Grizzly, the default), a successful call to defer will cause the page to be executed again in the same thread.

conversation.deferred (only available in `/web/dynamic/` and `/web/fragments/`)

True if the conversation has been deferred via a call to `conversation.defer` (see above).

Low-level Access

Use these to access the Restlet instances underlying the conversation. This is useful for features not covered by Prudence's standard API.

For more information, refer to Prudence's Java API documentation and also Restlet's API documentation.

conversation.resource The Restlet resource. In the case of `/web/dynamic/`, this will be Prudence's `GeneratedTextResource`. In the case of `/resources/`, this will be `DelegatedResource`.

conversation.request Equivalent to `conversation.resource.request`.

conversation.response Equivalent to `conversation.resource.response`.

Sharing State

Prudence is designed to allow massive concurrency and scalability while at the same time shielding you from the gorier details. However, when it comes to sharing state between different parts of your code, it's critical that you understand Prudence's state services.

Global Variables

You know how local variables work in your programming language: they exist only for the duration of a function call, after which their state is discarded. If you want state to persist beyond the function call, you use a global variable (or a "static" local, which is really a global).

But in Prudence, you cannot expect global variables to persist beyond a user request. To put it another way, you should consider every single user request as a separate "program" with its own global variables. If you need global variables to persist, you must use `application.globals` or `executable.globals`.

Why does Prudence discard your language's globals? This has to do with allowing for concurrency while shielding you from the complexity of having to guarantee the thread-safety of your code. By making each user request a separate "program," you don't have to worry about overlapping shared state, coordinating thread access, etc., for every use of a variable.

The exception to this is code in `/resources/`, in which language globals *might* persist. To improve performance, Prudence caches the global context for these in memory, with the side effect that your language globals persist beyond a single user request. For various reasons, however, Prudence may reset this global context. You should not rely on this side effect, and instead always use `application.globals` or `executable.globals`.

`application.globals` vs. `executable.globals`

You should prefer `application.globals`. By doing so, you'll minimize interdependencies between applications, and help make each application deployable on its own.

It's best to use `executable.globals` as an *optional* bridge between applications. Examples:

1. To save resources. For example, if an application detects that a database connection has already been opened by another application in the Prudence instance, and stored in `executable.globals`, then it could use that connection rather than create a new one. This would only work, of course, if a few applications share the same database.
2. To send messages between applications. This would be necessary if operations in one application could affect another. For example, you could place a task queue in `executable.globals`, where application could queue required operations. A thread in another application would consume these and act accordingly. Of course, you will have to plan for asynchronous behavior, and especially allow for failure. What happens if the consumer application is down?

Generally, if you find yourself having to rely on `executable.globals`, ask yourself if your code would be better off encapsulated as a single application. Remember that Prudence has powerful URL routing, support for virtual hosting, etc., letting you easily have one application work in several sites simultaneously

Note for Clojure flavor: All Clojure vars are VM-wide globals equivalent in scope to `executable.globals`. You usually work with namespaces that Prudence creates on-the-fly, so they do not persist beyond the execution. However, if you explicitly define a namespace, then you use it as a place for shared state. It will then be up to you to make sure that your namespace doesn't collide with that of another application installed in the Prudence instance. Though this approach might seem to break our rule of thumb here, of preferring `application.globals` to `executable.globals`, it is more idiomatic to Clojure and Lisps more generally.

Concurrency

Though `application.globals` and `executable.globals` are thread safe, it's important to understand how to use them properly.

Note for Clojure flavor: Though Clojure goes a long way towards simplifying concurrent programming, it does not solve the problem of concurrent access to global state. You still need to read this section!

For example, this code is broken:

```
def get_connection()
    data_source = application.globals['myapp.data.source']
    if data_source is None:
        data_source = data_source_factory.create()
        application.globals['myapp.data.source'] = data_source
    return data_source.get_connection()
```

The problem is that in the short interval between comparing the value in the “if” statement and setting the global value in the “then” statement, another thread may have already set the value. Thus, the “data_source” instance you are referring to in the current thread would be different from the “myapp.data.source” global used by other threads.

This may seem like a very rare occurrence to you: another thread would have to set the value *exactly* between our comparison and our set. If your application has many concurrent users, and your machine has many CPU cores, it can actually happen quite frequently. And, even if rare, your application has a chance of breaking if just two users use it at the same time!

Use this code instead:

```
def get_connection()
    data_source = application.globals['myapp.data.source']
    if data_source is None:
        data_source = data_source_factory.create()
        data_source = application.getGlobal('myapp.data.source', data_source)
    return data_source.get_connection()
```

The `getGlobal` call is an atomic compare-and-set operation. It guarantees that the returned value is the unique one.

Optimizing for Performance You may have noticed, in the code above, that if another thread had already set the global value, then our created data source would be discarded. If data source creation is heavy and slow, then this could significantly affect our performance. The only way to guarantee that this would not happen would be to make the entire operation atomic, by synchronizing it with a lock:

Here's an example:

```
def get_connection()
    lock = application.getGlobal('myapp.data.source.lock', RLock())
    lock.acquire()
    try:
```

```

        data_source = application.globals['myapp.data.source']
        if data_source is None:
            data_source = data_source_factory.create()
            application.globals['myapp.data.source'] = data_source
        return data_source.get_connection()
    finally:
        lock.release()

```

Note that we have to store our RLock as a unique global, too.

Not only is the code above complicated, but synchronization has its own performance penalties, which *might* make this apparent optimization actually perform worse. It's definitely not a good idea to blindly apply this technique: attempt it only if you are experiencing a problem with resource use or performance, and then make sure that you're not making things worse with synchronization.

If all else fails, then Prudence's globals may not be the best solution for your problem. Look into creating an external service (possibly written in Java) to manage global connections for you.

Here's a final version of our `get_connection` function that lets you control whether to lock access:

```

def get_connection(lock_access=False):
    if lock_access:
        lock = application.getGlobal('myapp.data.source.lock', RLock())
        lock.acquire()

    try:
        data_source = application.globals['myapp.data.source']
        if data_source is None:
            data_source = data_source_factory.create()
            if lock_access:
                application.globals['myapp.data.source'] = data_source
            else:
                data_source = application.getGlobal('myapp.data.source', data_source)
        return data_source.get_connection()
    finally:
        if lock_access:
            lock.release()

```

Complicated, isn't it? Unfortunately, complicated code and fine-tuning is the price you must pay in order to support concurrent access, which is the key to Prudence's scalability.

But, don't be discouraged. The standard protocol for using Prudence's globals will likely be good enough for the vast majority of your state-sharing needs.

conversation.locals

These are not "local" in the same way that function locals are. The term "local" here should be read as "local to the conversation" (compare with the term "thread locals"). They are "global" in the sense that they can be accessed by any function in your code, but are "local" in the sense that they persist only for the duration of the user request.

You may ask, then, why you wouldn't want to just use your language globals, which have the same scope and life. `conversation.locals` exist for three reasons:

1. To easily share conversation-scope state between scriptlets written in different languages.
2. To share state for deferred conversations (see `conversation.defer`, above).
3. URI segments are stored here (page 28).

Debugging

Logging

Logging (page 46) is the developer's best friend. Use it wisely, and you'll be able to clearly analyze the flows of control, code and data.

The Debug Page

The debug page is returned as a response whenever an un-handled exception is thrown in your code (or, if there is a bug in Prudence!).

Though very useful during development, you'd probably want to turn it off for production systems. Simply set the [showDebugOnError setting \(page 18\)](#) to false. Note that in that case, you might want to capture errors, sending them to a friendly, apologetic page for the users. See [error capturing \(page 29\)](#).

The debug page shows you the following information about the conversation:

- A stack trace of the exception, with a link to a view of the source code (see below).
- The reference (URI) used by the client, including the application's root reference, the virtual host reference, the original version of the reference (it might have been altered by filters along the route), and whether it was captured.
- The query of the URI.
- Cookies included in the request.
- Breakdown of the request metadata: media types, character sets, encodings, languages, etc.
- Request conditions, if included.
- Entity data, if available and not yet consumed by your code.
- Caching directives, if included.
- Information about the client: IP address, browser, operating system, user-agent, etc.
- Request attributes, including captured attributes.
- Warnings, if included.
- Underlying JVM stack trace.

Live Viewing of Source Code

Prudence supports live viewing of source code in `/web/dynamic/` and `/resources/`, with syntax coloring provided by Jygments. Links to source code are provided by the debug page, but you can also GET source code directly via:

- Prefixing URIs with `"/sourcecode/"`; see the [showSourceCodeURL setting \(page 18\)](#). For example, use GET on `"/sourcecode/support/forum/"` to see the source code for the resource or dynamic page at `"/support/forum/"`.
- Adding a `"?source=true"` to the query for URIs of `/web/dynamic/` pages.

In both cases, you can add `"?highlight=n"` to the query, when `"n"` is the line number to highlight.

Breakpoints?

As of Prudence version 1.0, you cannot set breakpoints in your code, unless it's in Java. Future versions of Prudence may allow breakpoints for some flavors, as supported languages add more debugging features. However, we do not feel that this is such a great loss, or that it would adversely affect your ability to develop for Prudence. The combination of robust logging and the debug page can go a long way towards helping you diagnose your problems. Debugging highly concurrent applications, such as network servers, does not work very well with breakpoints, anyway.

Logging

Prudence comes pre-configured with robust logging, based on log4j.

You are encouraged to make use of logging in your applications, but even if you don't, you will still find the logs useful. Prudence's servers, routers, programming languages and other components all send messages to the logs, making them an invaluable tool for debugging, monitoring and understanding how Prudence works.

By default, logs are sent to the `/logs/` directory, using configurable rolling log schemes. `/logs/web.log` records all server requests, using Apache's format, while everything else goes to `/logs/prudence.log`.

Loggers

The Base Logger

Use `application.logger` (page 34) to send text messages to the log. The logger name defaults to your application's subdirectory name, but can be configured via the `applicationLoggerName` setting (page 18).

Sub-Loggers

Large applications might benefit from using more than one logger. Use `application.getSubLogger` (page 34) with any name you wish. This name will be appended to your base logger name with a ".", and will inherit the base logger's configuration by default. For example, if your base logger is named "wackywiki", then `application.getSubLogger('backend')` will appear as "wackiwiki.backend" in the log files.

See more on logger inheritance below.

Sending Messages

Whether messages actually are written to the log file depends on your logging configuration. Messages are ranked by levels, and loggers are configured to allow only messages up to a certain level. Smart, consistent use of log levels will increase the debuggability of your application.

The logging methods and their common uses are these (shown for the base logger, but work the same for sub-loggers):

1. `application.logger.severe`: "Severe" messages are used for unrecoverable errors, alerts about unavailable computing resources, network backends, etc. You'd always want to configure logging to include these messages!
2. `application.logger.warning`: While not quite severe, the event could still point out a problem that, if left un-handled, might become severe, either now or in the future. Many applications emit copious warnings that can be safely ignored.
3. `application.logger.info`: These don't report a problem, but instead are used to mark an occurrence of an event. Useful for monitoring and high-level debugging.
4. `application.logger.config`: Treat these as components of an event that together would constitute a single "info" message. They are meant to show how the event was initialized or released. Useful for low-level debugging of events.
5. `application.logger.fine`: General purpose, low-level debugging.
6. `application.logger.finer`: Even lower!
7. `application.logger.finest`: Lowest of the low!

/configuration/logging.conf

We'll cover the basics here. See log4j documentation for more information.

The Prudence defaults are mostly at "info" level. You are encouraged to experiment with lower levels in order to see how Prudence's internals function!

Appenders

An "appender" is the service that actually writes log messages. Appenders are configured with properties prefixed with "log4j.appender.X", where X is the name of the appender.

Prudence by default uses two rolling file appenders, one called "web" for web.log, and one called "prudence" for prudence.log. Additionally, a console appender named "console" and a remote appender named "remote" are configured, though they are not used by default.

Loggers

A “logger” is the service to which you send your log messages. It decides whether to write the message according to its level, and if so sends it to one or more appenders. Loggers are configured with properties prefixed with “log4j.logger.X”, where X is the name of the logger. Logger names appear in the logs for every message, and are thus useful for organization your log.

You do not have to define all your loggers in logging.conf. Any logger name can be used by your application. If it is not found in logging.conf, then default attributes are inherited. Inheritance works by treating logger names hierarchically: if you do not specify a certain attribute for a logger, then its parent logger is used. The nameless root logger, configured with the “log4j.rootLogger” prefix, defines defaults for all loggers.

You can configure your application’s logger and sub-loggers. For example, if your application is named “wacky-wiki,” you can set its maximum logging level thus:

```
log4j.logger.wackywiki=WARN
```

Log levels in logging.conf are named a bit differently from the commands used in application.logger:

1. application.logger.severe: ERROR
2. application.logger.warning: WARN
3. application.logger.info: INFO
4. application.logger.config: DEBUG
5. application.logger.fine: DEBUG
6. application.logger.finer: DEBUG
7. application.logger.finest: TRACE (or ALL)

The differences are due to the preponderance of logging solutions for the JVM, which are used in some of Prudence’s underlying libraries. We hope to streamline this further in a future version of Prudence.

You can also add or change appenders for your loggers. For example, to send wackywiki messages to the console appender:

```
log4j.logger.wackywiki=WARN, console
```

Separate Logs Per Application

We’ll create a new appender for each new log file we need. In this example, we’ll just copy the “prudence” appender with a new name:

```
log4j.appender.wackywiki=org.apache.log4j.RollingFileAppender
log4j.appender.wackywiki.File=logs/wackywiki.log
log4j.appender.wackywiki.MaxFileSize=5MB
log4j.appender.wackywiki.MaxBackupIndex=9
log4j.appender.wackywiki.layout=org.apache.log4j.PatternLayout
log4j.appender.wackywiki.layout.ConversionPattern=%d: %-5p [%c] %m%n
```

Then, we’ll direct our application logger to use this appender:

```
log4j.logger.wackywiki=WARN, wackywiki
```

Analyzing /logs/web.log

You can throw Prudence’s web.log into almost any Apache log file analyzer. Here’s an example using the ubiquitous Analog:


```

analog \
-C'LOGFORMAT (%Y-%m-%d\t%h:%n:%j\t%S\t%u\t%j\t%j\t%j\t%r\t%q\t%c\t%b\t%j\t%T\t%v\t%B\t%f)' \
-C'LOCALCHARTDIR local/images/' \
-C'CHARTDIR images/' \
-C'HOSTNAME "mysite.org" \
logs/web.log \
-Oapplications/myapp/web/static/analog/index.html

```

Administration

Prudence comes with the “Prudence Administration Application.” As of Prudence 1.0, this is a simple application that lets you see the servers, virtual hosts and applications in the Prudence instance, and lets you start and stop them. Future versions of Prudence will build on this foundation, allowing for more runtime control, configuration and monitoring.

Installation

Prudence comes with this application installed as “prudence-admin,” bound to the default virtual host at the root URL.

Customization

Changing the Root URL

Let’s say we want prudence-admin at the “/admin/” URL on the default host, and the “/prudence/admin/” URL on myHost. We’ll edit its settings.*. Here’s an example for the JavaScript flavor:

```
hosts = [[component.defaultHost, '/admin/'], [myHost, '/prudence/admin/']];
```

Requiring a Password

As an example, we’ll route the application’s router through an HTTP authentication filter (page 25).

We’ll create a routing.js file (for the JavaScript flavor):

```

// Implement defaults
document.execute('defaults/application/routing/');

importClass(
    org.restlet.security.ChallengeAuthenticator,
    org.restlet.security.MapVerifier,
    org.restlet.data.ChallengeScheme);

// Create an authenticator
var verifier = new MapVerifier();
verifier.localSecrets.put('admin', new java.lang.String('opensesame').toCharArray());
var authenticator = new ChallengeAuthenticator(applicationInstance.context,
    ChallengeScheme.HTTP_BASIC, 'Prudence Administration');
authenticator.verifier = verifier;

// Put authenticator before root
authenticator.next = applicationInstance.inboundRoot;
applicationInstance.inboundRoot = authenticator;

```

Prudence As a Daemon

In production environments, it's best to run Prudence as a daemon in Unix-like systems or a service in Windows, via a lightweight wrapper. This will provide proper process monitoring and control. For example, if Prudence's JVM crashes for some reason, hangs, grabs too much CPU or RAM, the wrapper can automatically shut down and even restart it.

While Prudence comes with the necessary configuration files for this, it does not include an actual wrapper, which you will need to install manually. Each of the two available wrappers is problematic: Tanuki's JSW uses a restrictive license (GPL) that will not allow us to distribute it with Prudence, and YAJSW is far too big a download to include in Prudence. You can, however, easily download and install either of them yourself. Just be sure to follow the distribution restrictions in the case of JSW.

History lesson: The copyright for JSW is held by its original developer, Tanuki Software. For the first few versions Tanuki released JSW under a very permissive license, making it popular in many open source projects. However, since version 3.2 it has been distributed under the GPL 2.0 (and also via a commercial license). We applaud Tanuki's commitment to open source, and are big fans of the GPL. However, the GPL makes it impossible to distribute JSW with open source projects using less restrictive licenses, such as Prudence. Many projects have kept distributing version 3.2 of JSW, which is now out of date and missing bug fixes. For less restrictive open source projects, the gap has been filled by YAJSW, which seeks not only to be a drop-in replacement for JSW, but to also go beyond it with many additional features.

JSW

JSW is written in C rather than Java, making it much more lightweight than YAJSW. It supports many operating systems.

See `/configuration/wrapper.conf` for a sample JSW configuration. The wrapper will log to `/logs/wrapper.log`. In particular, you'll want to change:

- `wrapper.working.dir`: Set this to your Prudence installation's base directory
- `wrapper.java.command`: Set this to the JVM runtime (the "java" command) you'll want to use for Prudence
- `wrapper.java.maxmemory`: Set this according to your deployment environment. More memory translates to better performance, and more room for the in-process memory cache.
- `wrapper.java.initmemory`: A reasonably high value here can help speed up Prudence's startup time.

YAJSW

YAJSW is written in 100% Java, using JNA to handle the native operating-system-dependent features.

See `/configuration/yajsw.conf` for a sample JSW configuration. The wrapper will log to `/logs/yajsw.log`. Change `yajsw.conf` in the same way as described for `wrapper.conf`, above.

Monitoring with JMX

YAJSW allows for remote monitoring and control via JMX.

TODO

HTTP Proxy

There's nothing special about how Prudence handles HTTP, and it can work easily behind any reverse proxy. This lets you easily unite Prudence with other web servers or run it behind a load balancer. Though it's not unique to Prudence, we thought to add this section to the manual in order to get you up and running quickly with this useful scenario.

Perlbal

You can run many instances of Prudence behind a load balancer. This offers fault tolerance, maintenance options, and the possibility to dramatically scale up the number of requests you can support. Your application can tolerate failure of any number of instances, as long as you have one running, because load balancers will automatically route to working instances. Similarly, load balancing allows you to bring some instances down for maintenance while keeping your application up and running.

Scaling up can be straightforward: simply add more and more instances behind the load balancer, which will make sure to distribute requests among them, while monitoring their responsiveness to accommodate for how well they handle their load. More complex systems can involve different kinds of instances, with the load balancer being in charge of routing requests to the appropriate pool of instances. This “partitioning” can be according to features (one pool handles chat room, one pool handles file downloads), geography (one pool handles England, one pool handles France), or other clever ways to keep the whole system efficient and responsive.

There are many great load balancers out there, but we especially like Perlbal.

Apache

Apache is often called the “Swiss army knife of the Internet” for how well it manipulates URLs and routes HTTP requests. Prudence already does powerful URI-based routing, including virtual hosting, meaning that you probably won’t need Apache for that feature.

Where you might want to use Apache is as a container environment for other application platforms, such as `mod_php` and `mod_wsgi`. If you have no choice but to run Apache as your front end, it is straightforward to set it to route to Prudence as a reverse proxy.

Prudence As a Restlet Container

Prudence brings the power of REST and the JVM to programmers in other languages, but has a lot to offer to Java/Restlet developers.

Though applications can be written in Prudence without a single line of Java code, Prudence also acts as a useful container for existing Restlet applications, Restlet resources or just restlets, written in Java or other JVM languages.

Why use Prudence for a Restlet application that already works?

Prudence makes it easy to handle the bootstrapping and [routing \(page 18\)](#) of your applications, and the [Administration Application \(page 49\)](#), [debugging \(page 45\)](#) and [logging \(page 46\)](#) features make it easier to deploy and manage multiple applications together. These issues have more to do with your application’s configuration, rather than its functionality, and it can be useful to handle them outside of the Java build process, using live, dynamic languages in simple text source that you can modify on-the-fly. Deploying your Restlet application to a Prudence instance can be as simple as plopping in your jar.

This need is also fulfilled by servlet and Java Enterprise Edition (JEE) containers, such as Tomcat, Resin and JBoss. Indeed, Restlet has a JEE edition, and good support for servlets. However, if all you need is a deployment container, Prudence can serve as a straightforward, pure REST alternative.

Some people also look to JEE containers for their support of Java Server Pages (JSP). We urge you to take a good look at Prudence’s [dynamic web \(page 20\)](#) support. It may likely surpass JSP for your purposes. In particular, it is based on Restlet, which you already know and love, giving you the entire Restlet API at your fingertips. It also lets you use many wonderful languages other than Java for scriptlets. For example, we at Three Crickets are mad about Clojure. For simpler templating, Velocity and Succinct are built in and immediately ready to use.

Summary

A 100% Restlet-based alternative to servlet/JEE containers. (Requires only Restlet JSE.)

1. Easy deployment
 - (a) Boot scripts: avoid weird configuration formats and start your Restlet component exactly as you want (your choice among 6 languages)
 - (b) The default scripts already handle virtual hosting, multiple servers and internal routing
 - (c) Designed from the ground-up to handle multiple apps on the same component

- (d) Admin app for live management of components
 - (e) Logging is pre-configured and “just works,” including an Apache-compatible web log
 - (f) Single zip-file application deployment (like WAR files in JEE)
2. Easy prototyping of REST resources
 - (a) Your choice among 6 languages
 - (b) Code is compiled, cached and loaded on-the-fly
 - (c) Rich debug page shows errors and source code
 - (d) When you’re happy with it, you can rewrite it as a `ServerResource` in Java
 3. Powerful HTML generation platform, like JSP/ASP/PHP (again, 100% Restlet-based)
 - (a) Your choice among 6 languages, including mixing languages and template engines (Velocity, Succinct) on one page
 - (b) Code is compiled, cached and loaded on-the-fly
 - (c) RAM/database/Hazelcast/memcached server-side caching (uses Restlet’s URI template language for cache key generation)
 - (d) Straightforward support for client-side caching
 - (e) Asynchronous output
 - (f) Easily accept uploaded files
 - (g) Rich debug page shows errors and source code
 4. Restlet sugar (also available as a standalone JAR)
 - (a) Fallback router (attach multiple `MODE_STARTS_WITH` restlets to the same base URI)
 - (b) URI “capturing” (internal redirection)
 - (c) JavaScript/CSS unify-and-minify filter
 - (d) Delegated status service for diverting to custom pages (404, 500 errors, etc.)
 - (e) Rich `DebugRepresentation`
 - (f) Cache backend abstraction, designed for storing `StringRepresentations`
 - (g) Easier file uploads (slightly higher-level than the Restlet `FileUpload` extension)
 - (h) `ConversationCookie` (combines `Cookie` and `CookieSetting`)
 - (i) Filter to selectively add `CacheControl` (to `Directory`, for example)

Custom Resources

Use your application’s routing file to attach your resources, or otherwise manage routing. Example:

```
// Prudence defaults
document.execute('defaults/application/routing/');

// MyOrg resources
router.attach('data/item/{id}', classLoader.loadClass('org.myorg.ItemResource'));
router.attach('data/items', classLoader.loadClass('org.myorg.ItemsResource'));
```

You can also change Prudence’s default routing by detaching and re-attaching routes:

```
// Wrap static web in JavaScript minifying filter
importClass(org.myorg.JavaScriptMinifyFilter);
router.detach(staticWeb);
router.attach(
    staticWebBaseURL,
    new JavaScriptMinifyFilter(application.context, staticWeb, File(applicationBasePath +
        .matchingMode = Template.MODE_STARTS_WITH;
```

Custom Application

By default, Prudence creates an instance of the standard Restlet Application class. Use your application's application file to override this, and create and configure your own application. Example:

```
// MyOrgApplication importClass(org.myorg.MyOrgApplication);
var application = new MyOrgApplication();

// Install our custom tunnel service
importClass(org.myorg.MyOrgTunnelService);
application.tunnelService = new MyOrgTunnelService(MyOrgTunnelService.MODE_QUERY);

// These attributes are specific to MyOrgApplication
application.databaseURI = 'mysql://localhost/myorg';
application.useTransactions = true;
```

How to Choose a Flavor?

Python (Succulent!)

Python is a powerful object-oriented language, far richer in features than JavaScript and PHP, with many high-quality core and 3rd party libraries. Python has already proven itself as a web programming language with many excellent platforms.

Python presents a unique challenge in a scriptlet environment, due to its reliance on indentation. However, because HTML is loose with whitespace, it's possible to force the whole file to adhere to Python's scheme. In fact, as many Python enthusiasts would argue, forcing your code to adhere to Python's indentation requirements can go a long way towards making it more readable and manageable.

In the included example application we show how to use SQLAlchemy as a data backend for Prudence.

Note: Prudence for Python was built primarily around Jython, but also offers limited support for Jepp if it's installed. For those cases where you need access to a natively-built Python library that won't work on Jython, Jepp lets you run code on the CPython platform.

Ruby (Delectable!)

Ruby can do most of what Python can do and more. A true chameleon, it can adapt to many styles of code and programming. If something can be possible, Ruby allows it and supports it. Unlike Python, it has a very loose and forgiving syntax, which is perfect for scriptlets.

Ruby's Rails platform has revolutionized web programming by offering elegant, powerful alternatives to working directly with HTTP. We hope Ruby web programmers will find in Prudence a refreshing alternative to Rails: elegantly embracing HTTP, instead of avoiding it.

Clojure (Scrumptious!)

Prudence's only functional flavor is a Lisp designed from the ground up for high concurrency. If you're serious about scaling, Clojure is the way to go. Though new, Clojure is based on one of the oldest programming languages around, and enjoys a rich tradition of elegant solutions for tough programming challenges.

Clojure embraces the JVM, but also has a growing collection of nifty "contrib" libraries—all included in Prudence. In the included example application, we show how to use Clojure's SQL library to access a data backend.

JavaScript (Savory!)

JavaScript (a dialect of ECMAScript) is a sensible choice for "AJAX" and other rich web client applications, because it's the same language used by web browsers and other client platforms. Web developers are already proficient in JavaScript, and can quickly be brought on board to a server-side project. JavaScript lets you to write server-side and client-side code in the same language, and have both sides share code. Couple it with JSON, and you're on solid ground for rapid development. Of all the web programming languages, it's the one most widely deployed and with the most secure future.

JavaScript is an under-appreciated language. Though not as feature-rich as Python or Ruby, it's still very powerful. Its straightforward closure/prototype mechanisms allow it to support object-orientation, namespaces and other paradigms. Unfortunately, it has been the target of much undeserved angst due to the fickleness of client-side development. Working with the browser DOM, testing with cross-browser HTML rendering quirks—these are the not the fault of the language itself. They are also not relevant to server-side development with Prudence. If you've been burned by JavaScript on the client, try it on the server via Prudence: you might be pleasantly surprised.

JavaScript does not have its own core libraries, making it the most minimal Prudence flavor. Instead, it relies on the excellent JVM core.

PHP (Ambrosial!)

PHP is, of course, ubiquitous. It's a simple language with the most mature libraries of any web programming language, and programmers are available with years of experience. It's also designed from the ground up as a programming language for the web.

Prudence allows a smooth transition from traditional PHP HTML generation to REST resources. It supports PHP “superglobals” such as `$_GET`, `$_POST`, `$_COOKIE` and `$_FILE` (but not `$_SESSION`) to make you feel right at home. It also adds many new features to conventional HTML generation: fine-grained caching, high-performance templating languages, and more.

Note: Prudence PHP was built around the open source edition of Quercus, which does not feature JVM bytecode compilation as is available in the non-free professional edition. Nevertheless, we found the “non-pro” Quercus to be an excellent performer!

Groovy (Luscious!)

In some ways, Groovy is the best of this bunch. It has all the flexibility of Ruby, but is designed from the ground up to enhance and extend Java. Java programmers would immediately feel at home, while gaining access to far less restrictive programming paradigms. Groovy makes Java... groovy!

All the other Prudence flavors offer JVM interaction, but Groovy does it best. If you know your project will require a lot of interaction with Java libraries, Groovy is a terrific—and fun!—choice.

The Case for REST

There's a lot of buzz about REST, but also a lot confusion about what it is and what it's good for. The essay attempts to convey its simple essence.

Let's start, then, not at REST, but at an attempt to create a new architecture for building scalable applications. Our goals are for it to be minimal, straightforward, and still have enough features to be productive. We want to learn some lessons from the failures of other, more elaborate and “complete” architectures.

Let's call ours a “resource-oriented architecture.”

Resources

Our base unit is a “resource,” which, like an object in object-oriented architectures, encapsulates data with some functionality. However, we've learned from object-orientation that implementing arbitrary interfaces is a recipe for enormous complexity. Instead, then, we'll keep it simple and define a limited interface that would still be useful enough.

From our experience with relational databases, we've learned that a tremendous amount of power can be found in “CRUD”: Create, Read, Update and Delete. If we support just these operations, our resources will already be very powerful, enjoying the accumulated wisdom and design patterns from the database world.

Identifiers

First, let's start with a way of discriminating our resources. We'll define a name-based address space where our resources live. Each resource is “attached” to one or more addresses. We'll allow for “/” as a customary separator to allow for hierarchical addressing schemes. For example:

```
/animal/cat/12/image  
/animal/cat/12/image/large  
/animal/cat/12/specs
```

In the above, we’ve allowed for different kinds of animals, a way of referencing individual animals, and a way of referencing specific aspects of these animals. Let’s now go over CRUD operations in order of complexity.

Delete

“Delete” is the most trivial operation. After sending “delete” to an identifier, we expect it to not exist anymore. Whether or not sub-resources in our hierarchy can exist or not, we’d leave up to individual implementations. For example, deleting “/animal/cat/12/image” may or may not delete “/animal/cat/12/image/large”.

Note that we don’t care about atomicity here, because we don’t expect anything after our “delete” operation. A million changes can happen to our cat before our command is processed, but they’re all forgotten after “delete.” (See “update,” below, for a small caveat.)

Read

“Read” is a bit more complicated than “delete.” Since we expect our resource to change, we want to make sure that there’s some kind of way to mark which version we are reading. This will allow us to avoid unnecessary reads if there hasn’t been any change.

We’ll need our resource-oriented architecture to support some kind of version tagging feature.

Update

The problem with “update” is that it always references a certain version that we have “read” before. In some cases, though not all, we need some way to make sure that the data we expect to be there hasn’t changed since we’ve last “read” it. Let’s call this a “conditional update.”

Actually, we’ve oversimplified our earlier definition of “delete.” In some cases, we’d want a “conditional delete” to depend on certain expectations about the data. We might not want the resource deleted in some cases.

We’ll need our resource-oriented architecture to support a “conditional” feature.

Create

This is our most complex operation. Our first problem is that our identifier might not exist yet. One approach could be to try identifiers in sequence:

```
Create: /animal/cat/13 -> Error, already exists  
Create: /animal/cat/14 -> Error, already exists  
Create: /animal/cat/15 -> Error, already exists  
...  
Create: /animal/cat/302041 -> Success!
```

Obviously, this is not a scalable solution. Another approach could be to have a helper resource which provides us with the necessary ID:

```
Read: /animal/cat/next -> 14  
Create: /animal/cat/14 -> Oops, someone else beat us to 14!  
Read: /animal/cat/next -> 15  
Create: /animal/cat/15 -> Success!
```

Of course, we can also have “/animal/cat/next” return IDs that are never used and avoid duplications. If we never create our cat, they will be wasted, though. The main problem with this approach is that it requires two calls per creation: a “read,” and then a “create.” We can handle this in one call by allowing for “partial” creation:

```
Create: /animal/cat -> We send the data for the cat without the ID, and get back the ID
```

Other solutions exist, too. The point of this discussion is to show you that “create” is not trivial, but also that solutions to “create” exist within the resource-oriented architecture we’ve defined. “Create,” though complex, does not demand any new features.

Aggregate Resources

At first glance, handling the problem of getting lots of resources at the same time, thus saving on the number of calls, can trivially be handled by the features we've listed so far. A common solution is to define a “plural” version of the “singular” resource:

```
/animal/cats
```

A “read” would give us all cats. But what if there are ten million cats? We can support paging. Again, we have a solution within our current feature set, using identifiers for each subset of cats:

```
/animal/cats/100/200
```

The above would return no more than 100 cats: from the 100th, to the 200th. There's a slight problem in this solution: the burden is on whatever component in our system handles mapping identifiers to resources. This is not terrible, but if we want our system to be more generic, it could help if things like “100 to 200” could be handled by our resource more directly. For this convenience, let's implement a simple parameter system for all commands:

```
Read(100, 200): /animal/cats
```

In the above, our mapping component only needs to know “/animal/cats”. It can be very dumb, and easy to implement.

Formats

The problem of supporting multiple formats seems similar, at first glance, to that of aggregate resources. Again, we could potentially solve it with command parameters:

```
Read(UTF-8, Russian): /animal/cat/13
```

This would give us a Russian, Unicode UTF-8 encoded version of our cat. Looks good, except that there is a potential problem. The client might prefer certain formats, but actually be able to handle others. We would not want a series of wasteful operations to happen until one succeeds. Of course, we can have another resource where all available formats are listed, but this would require an extra call, and also introduce the problem of atomicity. A better solution would be to have the client associate certain preferences per command, have our resource emit its capabilities, with mapping component in between “negotiating” these two lists, via a simple algorithm, and choose the best mutually preferable format.

This would be a simple feature to add to our resource-oriented architecture, which could greatly help to decouple its support for multiple formats from its addressing scheme.

Shared State

Shared state between the client and server is very useful for managing sessions, and implementing basic security. Of course, it's quite easy to abuse shared state, too, by treating it as a cache for data. We don't want to encourage that. Instead, we just want a very simple shared state system.

We'll allow for this by attaching small, named, shared state objects to every request and response to a command. Nothing fancy or elaborate. There is a potential security breach here, so we have to trust that all components along the way honor the relationship between client and server, and don't allow other servers access to our shared state.

Summary of Features

So, what do we need?

We need a way to map identifiers to resources. We need support for the four CRUD operations. We need support for “conditional” updates and deletes. We need all operations to support “parameters.” We need “negotiation” of formats. And, we need a simple shared state attachment feature.

This list is very easy to implement. It requires very little computing power, and no support for generic, arbitrary additions.

Before we go on, it's worth mentioning one important feature which we did not require: transactions. Transactions are optional, and sometimes core features in many databases and distributed object systems. They can be extremely powerful, as they allow atomicity across an arbitrary number of commands. They are also, however, heavy to implement, as they require considerable shared state between client and server. Powerful as they are, it is possible to live without them. It's possible, for example, to implement this atomicity within a single resource. This

would require us to define special resources per type of transaction which we want to support, but it does remove the heavy burden of supporting arbitrary transactions from our architecture. With some small reluctance, then, we'll do without transactions.

Let's Do It!

OK, so now we know what we need, let's go ahead and implement the infrastructure of components to handle our requirements. All we need is stacks for all supported clients, backend stacks for all our potential server platforms, middleware components to handle all the identifier routing, content negotiation, caching of data. . .

. . . And thousands of man hours to develop, test, deploy, and integrate. Like any large-scale, enterprise architecture, even trivial requirements have to jump through the usual hoops set up by the sheer scale of the task. Behind every great architecture are the nuts and bolts of the infrastructure.

Wouldn't it be great if the infrastructure already existed?

The Punchline

Well, duh. All the requirements for our resource-oriented architecture are already supported by HTTP:

Our resource identifiers are simple URLs. The CRUD operations are in the four HTTP verbs: PUT, GET, POST and DELETE. "Conditional" and "negotiated" modes are handled by headers, as are "cookies" for shared state. Version stamps are e-tags. Command parameters are query matrixes appended to URLs. It's all there.

Most importantly, the infrastructure for HTTP is already fully deployed world-wide. TCP/IP stacks are part of practically every operating system; wiring, switching and routing are part and parcel; HTTP gateways, firewalls, load balancers, proxies, caches, filters, etc., are stable consumer components; certificate authorities, national laws, international agreements are already in place to support the complex inter-business interaction. Importantly, this infrastructure is successfully maintained, with minimal down-time, by highly-skilled independent technicians, organizations and component vendors across the world.

It's important to note a dependency and possible limitation of HTTP: it is bound to TCP/IP. Indeed, all identifiers are URLs: Uniform Resource Locators. In URLs, the first segment is reserved for the domain, either an IP address or a domain name translatable to an IP address. Compare this with the more general URIs (Uniform Resource Identifiers), which do not have this requirement. Though we'll often be tied to HTTP in REST, you'll see the literature attempting, at least, to be more generic. There are definitely use cases for non-HTTP, and even non-TCP/IP addressing schemes. In Prudence, you'll see that it's possible to address internal resources with non-URL kinds of URIs; see [document.external \(page 37\)](#).

It's All About Infrastructure

The most important lesson to take from this experience is the importance of infrastructure. This is why, I believe, Roy Fielding named Chapter 5 of his 2000 dissertation "Representational State Transfer" rather than, say, "resource-oriented architecture," as we have here. Fielding, one of the authors of the HTTP protocol, was intimately familiar with its challenges, and the name is intended to point out the key characteristic of its infrastructure: it's all about the transfer of lightly annotated data representations. "Resources" are merely logical encapsulations of these representations, depending on a contract between client and server. The infrastructure does not, in itself, do anything in particular to maintain, say, a sensible hierarchy of addresses, the atomicity of CRUD operations, etc. That's up to your implementation. But, representational state transfer—REST—is the mundane, underlying magic that makes it all possible.

To put it succinctly, a resource-oriented architecture requires a REST infrastructure. Practically, the two terms become interchangeable.

The principles of resource-orientation can and are applied in many systems. The word wide web, of course, with its ecology of web browsers, web servers, certificate authorities, etc., is the most obvious model. But other core internet systems, such as email (SMTP, POP, IMAP), file transfer (FTP, WebDAV) also implement some subset of REST. Your application can do this, too, and enjoy the same potential for scalability as the above.

Does REST Scale?

Part of the buzz about REST is that it's an inherently scalable architecture. This is true, but perhaps not in the way that you think.

Consider that there are two uses of the term "scalable":

First, it's the ability to respond to a growing number of user requests without degradation in response time, by "simply" adding hardware (horizontal scaling) or replacing it with more powerful hardware (vertical scaling). This is the aspect of scalability that engineers care about. The simple answer is that REST can help, but it doesn't stand out. SOAP can also do it pretty well. REST aficionados sometimes point out that REST is "stateless," or "sessionless," both characteristics that would definitely help scale. But, this is misleading. Protocols might be stateless, but architectures built on top of them don't have to be. For example, we've specifically talked about sessions here. And, you can easily make poorly scalable REST. The bottom line is that there's nothing in REST that guarantees scalability in *this* respect. Indeed, engineers coming to REST due to this false lure end up wondering what the big deal is.

The second use of "scalability" comes from the realm of enterprise and project management. It's the ability of your project to grow in complexity without degradation in your ability to manage it. And that's REST's beauty—you already have the infrastructure, which is the hardest thing to scale in a project. You don't need to deploy client stacks. You don't need to create and update proxy objects for five different programming languages used in your enterprise. You don't need to deploy incompatible middleware by three different vendors and spend weeks trying to force them to play well together. Why would engineers care about REST? Precisely because they don't have to: they can focus on engineering, rather than get bogged down by infrastructure management.

That said, a "resource-oriented architecture" as we defined here is not a bad start for (engineering-wise) scalable systems. Keep your extras lightweight, minimize or eliminate shared state, and encapsulate your resources according to use cases, and you won't, at least, create any obstacles to scaling.

Prudence

Convinced? The best way to understand REST is to experiment with it. You've come to the right place. Start with the [tutorial \(page 5\)](#), and feel free to skip around the documentation and try things out for yourself. You'll find it easy, fun, and powerful enough for you to create large-scale applications that take full advantage of the inherently scalable infrastructure of REST. Happy RESTing!

Scaling Tips

If you want your application to handle many concurrent users, then you're fighting this fact: a request will get queued in the best case or discarded in the worst case if there is no thread available to serve it.

Your challenge is to make sure that a thread is always available. And it's not easy, as you'll find out as you read through this article.

Scalability is the ability to respond to a growing number of user requests without degradation in response time. Two variables influence it: 1) your total number of threads and 2) the time it takes each thread to process a request. Increasing the number of threads seems straightforward: you can keep adding more machines behind load balancers. However, the two variables are tied, as there are diminishing returns and even reversals: beyond a certain point, time per request can actually grow longer as you add threads and machines.

Let's ignore the first variable, because the challenge of getting more machines is mostly financial. It's the second that we can do something about as engineers: minimizing the time per request becomes an architectural challenge that encompasses the entire structure of your application.

Meanwhile, feel free to frame these inspirational slogans on your wall:

Requests are hot potatoes: Pass them on!

And:

It's better to have many short requests than one long one.

Performance Does Not Equal Scalability

Performance does not equal scalability. Performance does not equal scalability. Performance does not equal scalability.

Get it? Performance does not equal scalability.

This is an important mantra for two reasons:

1. Opposite Results

Optimizing for performance can adversely affect your scalability. The reason is contextual: when you optimize for performance, you often work in an isolated context, specifically so you can accurately measure response times and fine-tune them. For example, making sure that a specific SQL query is fast would involve just running that query. A full-blown experiment involving millions of users doing various operations on your application would make it very hard to accurately measure and optimize the query. Unfortunately, by working in an isolated context you cannot easily see how your efforts would affect other parts of an application. To do so would require a lot of experience and imagination. To continue our example, in order to optimize your one SQL query you might create an index. That index might need to be synchronized with many servers in your cluster. And that synchronization overhead, in turn, could seriously affect your ability to scale. Congratulations! You've made one query run fast in a situation that never happens in real life, and you've brought your web site to a halt.

One way to try to get around this is to fake scale. Tools such as JMeter, Siege and ApacheBench can create "load." They also create unfounded confidence in engineers. If you simulate 10,000 users bombarding a single web page, then you're, as before, working in an isolated context. All you've done is add concurrency to your performance optimization measurements. Your application pathways might work optimally in these situations, but this might very well be due to the fact that the system is not doing anything else. Add those "other" operations in, and you might get worse site capacity than you did before "optimizing."

2. Wasted Effort

Even if you don't adversely affect your scalability through optimizing for performance, you might be making no gains, either. No harm done? Well, plenty of harm, maybe. Optimizing for performance might waste a lot of development time and money. This effort would be better spent on work that could actually help scalability.

And, perhaps more seriously, it demonstrates a fundamental misunderstanding of the problem field. If you don't know what your problems are, you'll never be able to solve them.

Optimizing for Scalability

So, what can you do?

First, study the problem field carefully. Understand the challenges and potential pitfalls. You don't have to apply every single strategy up-front, but at least make sure you're not making a fatal mistake, such as binding yourself strongly to a technology or product with poor scalability. A poor decision can mean that when you need to scale up in the future, no amount of money and engineering effort will be able to save you in time.

Moreover, be very careful of blindly applying strategies used by other people to your own application. What worked for them might not work for you. In fact, there's a chance that their strategy doesn't even work for them, and they just think it did because of a combination of seemingly unrelated factors. The realm of web scalability is still young, full of guesswork, intuition and magical thinking.

Be especially careful of applying a solution before you know if you even have a problem.

How to identify the problems? You can create simulations and measurements of scalability rather than performance. You need to model actual user behavior patterns, allow for a diversity of such behaviors to happen concurrently, and replicate this diversity on a massive scale.

Creating such a simulation is a difficult and expensive, as is monitoring and interpreting the results and identifying potential bottlenecks. This is the main reason for the lack of good data and good judgment about how to scale. Most of what we know comes from tweaking real live web sites, which either comes at the expense of user experience, or allows for very limited experimentation. Your best bet is to hire a team who's already been through this before.

Generally, be very suspicious of products or technologies being touted as "faster" than others. *"Fast" doesn't say anything about the ability to scale.* Is a certain database engine "fast"? That's important for certain applications, no doubt. But maybe it's missing important clustering features, making it a poor choice for large applications. Does a certain programming language execute faster than another? That's great if you're doing video compression, but speed of execution might not be your bottleneck at all. Web applications mostly do I/O, not computation. The same web application might have very similar performance characteristics whether it's written in C++ or PHP.

Moreover, if the faster language is difficult to work with, has poor debugging tools, limited integration with web technologies, then it would slow down your work and your ability to scale.

Speed of execution can actually help scalability in one respect: If your application servers are constantly at maximum CPU load, then a faster execution platform would let you cram more web threads into each

server. This will help you reduce costs. Because Prudence is built on the fast JVM platform, you're in good hands in this respect. Also see Facebook's HipHop. Note, however, that there's a cost to high performance: more threads per machine would also mean more RAM requirements per machine, which also costs money. Once again, performance does not equal scalability, and you need to optimize specifically for scalability.

In summary, your architectural objective is to increase concurrency, not necessarily performance. Optimizing for concurrency means breaking up tasks into as many pieces as possible, and possibly even breaking requests into smaller pieces. We'll cover numerous strategies here, from frontend to backend.

Project Scalability

That last point about programming languages is worth some elaboration. Beyond how well your chosen technologies perform, it's important to evaluate them in terms to how easy they are to manage. Large web sites are large projects, involving large teams of people and large amounts of money. That's difficult enough to coordinate. You want the technology to present you with as few extra managerial challenges as possible.

Beware especially of languages described as "agile," as if they somehow embody the spirit of the popular Agile Manifesto. Often, this epithet seems to emphasize the following aspects of code: forgiveness for syntax slips, light or no type checking, automatic memory management and concurrency—all features that could just as well be used for sloppy, error-prone, hard-to-debug, and hard-to-fix code. Beware, too, of the cult of "productivity": if you're reading this article, then your goal is likely not to create a quick demo, but a stable application with a long, evolving life span.

Ignore the buzzwords, and instead make sure you're choosing technology that you can control, instead of technology that will control you.

We discuss this topic some more in ["Making the Case for Rest" \(page 54\)](#). By building on the existing web infrastructure, Prudence can make large internet projects easier to manage.

Caching

Retrieving from a cache can be orders of magnitude faster than dynamically processing a request. It's your most powerful tool for increasing concurrency.

Caching, however, is only effective if there's something in the cache. It's pointless to cache fragments that appear only to one user on only one page that they won't return to. On the other hand, there may very well be fragments on the page that will recur often. If you design your page carefully to allow for fragmentation, you will reap the benefits of fine-grained caching. Remember, though, that the outermost fragment's expiration defines the expiration of the included fragments. It's thus good practice to define no caching on the page itself, and only to cache fragments.

In your plan for fine-grained caching, take special care to isolate those fragments that cannot be cached, and cache everything around them.

Make sure to change Prudence's [cache key \(page 34\)](#) to fit the lowest common denominator: you want as many possible requests to use the already-cached data, rather than generating new data. Note that, by default, Prudence includes the request URI in the cache key. Fragments, though, may very well appear identically in many different URIs. You would thus not want the URI as part of their cache key.

Cache aggressively, but also take cache validation seriously. Make good use of Prudence's [cache tags \(page 35\)](#) to allow you to invalidate portions of the cache that should be updated as data changes. Note, though, that every time you invalidate you will lose caching benefits. If possible, make sure that your cache tags don't cover too many pages. Invalidate only those entries that really need to be invalidated.

(It's sad that many popular web sites do cache validation so poorly. Users have come to expect that sometimes they see wrong, outdated data on a page, sometimes mixed with up-to-date data. The problem is usually solved within minutes, or after a few browser refreshes, but please do strive for a better user experience in your web site!)

If you're using a deferred task handler (see below), you might want to invalidate tagged cache entries when tasks are done. Consider creating a special internal API that lets the task handler call back to your application to do this.

How long should you cache? As long as the user can bear! In a perfect world, of limitless computing resources, all pages would always be generated freshly per request. In a great many cases, however, there is no harm at all if users see some data that's a few hours or a few days old.

Note that even very small cache durations can make a big difference in application stability. Consider it the maximum throttle for load. For example, a huge sudden peak of user load, or even a denial-of-service (DOS) attack, might overrun your thread pool. However, a cache duration of just 1 second would mean that your page would never be generated more than once every second. You are instantly protected against a destructive scenario.

Cache Warming

Caches work best when they are “warm,” meaning that they are full of data ready to be retrieved.

A “cold” cache is not only useless, but it can also lead indirectly to a serious problem. If your site has been optimized for a warm cache, starting from cold could significantly strain your performance, as your application servers struggle to generate all pages and fragments from scratch. Users would be getting slow response times until the cache is significantly warm. Worse, your system could crash under the sudden extra load.

There are two strategies to deal with cold caches. The first is to allow your cache to be persistent, so that if you restart the cache system it retains the same warmth it had before. This happens automatically with database-backed caches (see below). The second strategy is to deliberately warm up the cache in preparation for user requests.

Consider creating a special external process or processes to do so. Here are some tips:

1. Consider mechanisms to make sure that your warmer does not overload your system or take too much bandwidth from actual users. The best warmers are adaptive, changing their load according to what the servers can handle. Otherwise, consider shutting down your site for a certain amount of time until the cache is sufficiently warm.
2. If the scope is very large, you will have to pick and choose which pages to warm up. You would want to choose only the most popular pages, in which case you might need a system to record and measure popularity. For example, for a blog, it’s not enough just to warm up, say, the last two weeks of blog posts, because a blog post from a year ago might be very popular at the moment. Effective warming would require you to find out how many times certain blog posts were hit in the past two weeks. It might make sense to embed this auditing ability into the cache backend itself.

Pre-Filling the Cache

If there are thousands of ways in which users can organize a data view, and each of these views is particular to one user, then it may make little sense to cache them individually, because individual schemes would hardly ever be re-used. You’ll just be filling up the cache with useless entries.

Take a closer look, though:

1. It may be that of the thousands of organization schemes only a few are commonly used, so it’s worth caching the output of just those.
2. It could be that these schemes are similar enough to each other that you could generate them all in one operation, and save them each separately in the cache. Even if cache entries will barely be used, if they’re cheap to create, it still might be worth creating them.

This leads us to an important point:

Prudence is a “frontend” platform, in that it does not specify which data backend, if at all, you should use. Its cache, however, is general purpose, and you can store in it anything that you can encode as a string.

Let’s take as a pre-filling example a tree data structure in which branches can be visually opened and closed. Additionally, according to user permissions different parts of the tree may be hidden. Sounds too complicated to cache all the view combinations? Well, consider that you can trigger, upon any change to the tree data structure, a function that loops through all the different iterations of the tree recursively and saves a view of each of them to the cache. The cache keys can be something like “branch1+.branch2-.branch3+”, with “+” signifying “-” whether the branch is visually open or closed. You can use similar +’s and -’s for permissions, and create views per permission combinations. Later, when users with specific permissions request different views of the tree, no problem: all possibilities were already pre-filled. You might end up having to generate and cache thousands of views at once, but the difference between generating one view and generating thousands of views may be quite small, because the majority of the duration is spent communicating with the database backend.

If generating thousands of views takes too long for the duration of a single request, another option is to generate them on a separate thread. Even if it takes a few minutes to generate all the many, many tree views combinations, it might be OK in your application for views to be a few minutes out-of-date. Consider that the scalability benefits can be very significant: you generate views only *once* for the entire system, while millions of concurrent users do a simple retrieval from the cache.

Caching the Data Backend

Pre-filling the cache can take you very far. It is, however, quite complicated to implement, and can be ineffective if data changes too frequently or if the cache has to constantly be updated. Also, it's hard to scale the pre-filling to *millions* of fragments.

If we go back to our tree example above, the problem was that it was too costly to fetch the entire tree from the database. But what if we cache the tree itself? In that case, it would be very quick to generate any view of the tree on-demand. Instead of caching the view, we'd be caching the data, and achieving the same scalability gains.

Easy, right? So why not cache *all* our data structures? The reason is that it's very difficult to do this correctly beyond trivial examples. Data structures tend to have complex interrelationships (one-to-many, many-to-many, foreign keys, recursive tree structures, graphs, etc.) such that a change in data at one point of the structure may alter various others in particular ways. For example, consider a calendar database, and that you're caching individual days with all their events. Weekly calendar views are then generated on-the-fly (and quickly) for users according to what kinds of events they want to see in their personal calendars. What happens if a user adds a recurring event that happens every Monday? You'll need to make sure that all Mondays currently cached would be invalidated, which might mean tagging all these as "monday" using Prudence's cache tags. This requires a specific caching strategy for a specific application.

By all means, cache your data structures if you can't easily cache your output, but be aware of the challenge!

Prudence's sister project, Diligence, is designed specifically to solve this problem. It not only caches your data structures, but it validates them in memory using your coded logic, instead of invalidating them and forcing them to be re-fetched from the database. It supports data structures commonly used with relational databases, pluggable storage technologies, high-performance resource pooling and throttling, and natural integration with Prudence. Together, Diligence and Prudence form a solid platform for building scalable, data-backed web applications. At the time of this writing, Diligence is still under development. We hope to release it as open source soon, so stay tuned!

Cache Backends

Your cache backend can become a bottleneck to scalability if 1) it can't handle the amount of data you are storing, or 2) it can't respond quickly enough to cache fetching.

Before you start worrying about this, consider that it's a rare problem to have. Even if you are caching millions of pages and fragments, a simple relational-database-backed cache, such as Prudence's `SqlCache` implementations, could handle this just fine. A key/value table is the most trivial workload for relational databases, and it's also easy to shard (see backend partitioning, below). Relational database are usually very good at caching these tables in their memory and responding optimally to read requests. Prudence even lets you chain caches together to create tiers: an in-process memory cache in front of a SQL cache would ensure that many requests don't even reach the SQL backend.

High concurrency can also be handled very well by this solution. Despite any limits to the number of concurrent connections you can maintain to the database, each request is handled very quickly, and it would require *very* high loads to saturate. The math is straightforward: with a 10ms average retrieval time (very pessimistic!) and a maximum of 10 concurrent database connections (again, pessimistic!) you can handle 1,000 cache hits per second. A real environment would likely provide results orders of magnitude better.

The nice thing about this solution is that it uses the infrastructure you already have: the database.

But, what if you need to handle *millions* of cache hits per second? First, let us congratulate you for your global popularity. Second, there is a simple solution: distributed memory caches. Prudence comes with Hazelcast and support for memcached, which both offer much better scalability than database backends. Because the cache is in memory, you lose the ability to easily persist your cache and keep it warm: restarting your cache nodes will effectively reset them. There are workarounds—for example, parts of the cache can be persisted to a second database-backed cache tier—but this is a significant feature to lose.

Actually, Hazelcast offers fail-safe, live backups. While it's not quite as permanent as a database, it might be good enough for your needs. And memcached has various plugins that allow for real database persistence, though using them would require you to deal with the scalability challenges of database backends, which we will deal with below.

You'll see many web frameworks out there that support a distributed memory cache (usually memcached) and recommend you use it ("it's fast!" they claim, except that it can be slower per request than optimized databases, and that anyway performance does not equal scalability). We'd urge you to consider that advice carefully: keeping your cache warm is a challenge made much easier if you can store it in a persistent backend, and database backends can take you very far in scale without adding a new infrastructure to your deployment. It's good to know, though, that Prudence's support for Hazelcast and memcached is there to help you in case you reach the popularity levels of LiveJournal, Facebook, YouTube, Twitter, etc.

Client-Side Caching

Modern web browsers support client-side caching, a feature meant to improve the user experience and save bandwidth costs. A site that makes good use of client-side caching will appear to work fast for users, and will also help to increase your site's popularity index with search engines.

Optimizing the user experience is not the topic of this article: our job here is to make sure your site doesn't degrade its performance as load increases. However, client-side caching can indirectly help you scale by reducing the number of hits you have to take in order for your application to work.

Actually, doing a poor job with client-side caching can help you scale: users will hate your site and stop using it—voila, less hits you have to deal with. OK, that was a joke!

Generally, Prudence handles client-side caching automatically. If you cache a page, then headers will be set to ask the client to cache for the same length of time. By default, conditional mode is used: every time the client tries to view a page, it will make a request to make sure that nothing has changed since their last request to the page. In case nothing has changed, no content is returned.

You can also turn on "offline caching" mode, in which the client will avoid even that quick request. Why not enable offline caching by default? Because it involves some risk: if you ask to cache a page for one week, but then find out that you have a mistake in your application, then users will not see any fix you publish until their local cache expires, which can take up to a week! It's important that you understand the implications before using this mode. See the [dynamicWebClientCachingMode application setting](#) (page 20).

It's generally safer to apply offline caching to your static resources, such as graphics and other resources. A general custom is to ask the client to cache these "forever" (10 years), and then, if you need to update a file, you simply create a new one with a new URL, and have all your HTML refer to the new version. Because clients cache according to URL, their cached for the old version will simply not be ignored. See [CacheControlFilter](#) (page 24). There, you'll also see some more tricks Prudence offers you to help optimize the user experience, such as unifying/minimizing client-side JavaScript and CSS.

Upstream Caching

If you need to quickly scale a web site that has not been designed for caching, a band-aid is available: upstream caches, such as Varnish, NCache and even Squid. For archaic reasons, these are called "reverse proxy" caches, but they really work more like filters. According to attributes in the user request (URL, cookies, etc.), they decide whether to fetch and send a cached version of the response, or to allow the request to continue to your application servers.

The crucial use case is archaic, too. If you're using an old web framework in which you cannot implement caching logic yourself, or cannot plug in to a good cache backend, then these upstream caches can do it for you.

They are problematic in two ways:

1. Decoupling caching logic from your application means losing many features. For example, invalidating portions of the cache is difficult if not impossible. It's because of upstream caching, indeed, that so many web sites do a poor job at showing up-to-date information.
2. Filtering actually implements a kind of partitioning, but one that is vertical rather than horizontal. In horizontal partitioning, a "switch" decides to send requests to one cluster of servers or another. Within each cluster, you can control capacity and scale. But in vertical partitioning, the "filter" handles requests internally.

Not only is the “filter” more complex and vulnerable than a “switch” as a frontend connector to the world, but you’ve also complicated your ability to control the capacity of the caching layer. It’s embedded inside your frontend, rather than being another cluster of servers. (More on backend partitioning below.)

Unfortunately, there is a use case relevant for newer web frameworks, too: if you’ve designed your application poorly, and you have many requests that could take a long time to complete, then your thread pools could get saturated when many users are concurrently making those requests. When saturated, you cannot handle even the super-quick cache requests. An upstream cache band-aid could, at least, keep cached pages working, even though your application servers are at full capacity. This creates an illusion of scalability: some users will see your web site behaving fine, while others will see it hanging.

The real solution would be to re-factor your application so that it does not have long request. Below are tips on how to do this.

While upstream caching doesn’t really help you scale, it can provide useful redundancy: in case all your application servers are down, it can continue serving cached pages. Again, non-cached pages will not be available, making this an inadequate “solution,” but it might be better than nothing.

Dealing with Lengthy Requests

One size does not fit all: you will want to use different strategies to deal with different kinds of tasks.

Necessary Tasks

The user can’t continue without the task being resolved.

If the necessary task is deterministically fast, you can do all processing in the request itself.

If not, you should queue the task on a handling service and return a “please wait” page to the user. It would be nice to add a progress bar or some other kind of estimation of how long it would take for the task to be done. The client will poll until the task status is marked “done,” after which they will be redirected somewhere else. Each polling request sent by the client could likely be processed very quickly, so this strategy effectively breaks the task into many small requests (“It’s better to have many short requests than one long one”).

Implementing a handling service is by no means trivial. It adds a new component to your architecture, one that also has to be made to scale. One can also argue that it adversely affects user experience by adding overhead, delaying the time it takes for the task to complete. The bottom line, though, is you’re vastly increasing concurrency and your ability to scale. And, you’re actually improving the user experience: they would get a feedback on what’s going on rather than having their browsers spin, waiting for their requests to complete.

Deferrable Tasks

It’s OK if the task occurs later.

As with necessary tasks, you can queue these with a task handling service, but it’s much simpler because you don’t have to keep track of status or display it to the user. This allows you to use a simpler—and more scalable—task handling service.

Deferring tasks does present a challenge to the user experience: What do you do if the task fails, and the user needs to know about it? One solution can be to send a warning email or other kind of message to the user. Another solution could be to have your client constantly poll in the background (via “AJAX”) to see if there are any error messages, which in turn might require you to keep a queue of such error messages per user.

Before you decide on deferring a task, think carefully of the user experience: for example, users might be constantly refreshing a web page waiting to see the results of their operation. Perhaps the task you thought you can defer should actually be considered necessary?

Prudence comes with Hazelcast, which is nominally used for an optional distributed cache backend (see above), but can also be perfect for handling task queues. Give it a try!

File Uploads

These are potentially very long requests that you cannot break into smaller tasks, because they depend entirely on the client. As such, they present a unique challenge to scalability.

Fortunately, Prudence handles client requests via non-blocking I/O, meaning that large file uploads will not hold on to a single thread for the duration of the upload.

Unfortunately, many concurrent uploads will still saturate your threads. If your application relies on file uploads, you are advised to handle such requests on separate Prudence instances, so that uploads won't stop your application from handling other web requests. You may also consider using a third-party service specializing in file storage and web uploads.

Asynchronous Request Processing

Having the client poll until a task is completed lets you break up a task into multiple requests and increase concurrency. Another strategy is to break an *individual request* into pieces. While you're processing the request and preparing the response, you can free the web thread to handle other requests. When you're ready to deliver content, you raise a signal, and the next available web thread takes care of sending your response to the client. You can continue doing this indefinitely until the response is complete. From the client's perspective it's a single request: a web browser, for example, would spin until the request was completed.

You might be adding some extra time overhead for the thread-switching on your end, but the benefits for scalability are obvious: you are increasing concurrency by shortening the time you are holding on to web threads.

For web services that deliver heavy content, such as images, video, audio, it's absolutely necessary. Without it, a single user could tie up a thread for minutes, if not hours. You would still get degraded performance if you have more concurrent users than you have threads, but at least degradation will be shared among users. Without asynchronous processing, each user would tie up one thread, and when that finite resource is used up, more users won't be able to access your service.

Even for lightweight content such as HTML web pages, asynchronous processing can be a good tactic for increasing concurrency. For example, if you need to fetch data from a backend with non-deterministic response time, it's best to free the web thread until you actually have content available for the response.

It's not a good idea to do this for every page. While it's better to have many short requests instead of one long one, it's obviously better to have one short request rather than many short ones. Which web requests are good candidates for asynchronous processing?

1. Requests for which processing is made of independent operations. (They'll likely be required to work in sequence, but if they can be processed in parallel, even better!)
2. Requests that must access backend services with non-deterministic response times.

And, even for #2, if the service can take a *very* long time to respond, consider that it might be better to queue the task on a task handler and give proper feedback to the user.

And so, after this lengthy discussion, it turns out that there aren't that many places where asynchronous processing can help you scale. Caching is far more useful.

As of version 1.0, Prudence has limited support for asynchronous processing, via `conversation.defer`. Better support is planned for a future version.

Backend Partitioning

You can keep adding more nodes behind a load balancer insofar as each request does not have to access shared state. Useful web applications, however, are likely data-driven, requiring considerable state.

If the challenge in handling web requests is cutting down the length of request, then that of backends is the struggle against degraded performance as you add new nodes to your database cluster. These nodes have to synchronize their state with each other, and that synchronization overhead increases exponentially. There's a definite point of diminishing returns.

The backend is one place where high-performance hardware can help. Ten expensive, powerful machines might be equal in total power to forty cheap machines, but they require a quarter of the synchronization overhead, giving you more elbow room to scale up. Fewer nodes is better.

But CPUs can only take you so far.

Partitioning is as useful to backend scaling as caching is to web request scaling. Rather than having one big cluster of identical nodes, you would have several smaller, independent clusters. This lets you add nodes to each cluster without spreading synchronization overhead everywhere. The more partitions you can create, the better you'll be able to scale.

Partitioning can happen in various components of your application, such as application servers, the caching system, task queues, etc. However, it is most effective, and most complicated to implement, for databases. Our

discussion will thus focus on relational (SQL) databases. Other systems would likely require simpler subsets of these strategies.

Reads vs. Writes

This simple partitioning scheme greatly reduces synchronization overhead. Read-only servers will never send data to the writable servers. Also, knowing that they don't have to handle writes means you can optimize their configurations for aggressive caching.

(In fact, some database synchronization systems will only let you create this kind of cluster, providing you with one “master” writable node and several read-only “slaves.” They force you to partition!)

Another nice thing about read/write partitioning is that you can easily add it to all the other strategies. Any cluster can thus be divided into two.

Of course, for web services that are heavily balanced towards writes, this is not an effective strategy. For example, if you are implementing an auditing service that is constantly being bombarded by incoming data, but is only queried once in a while, then an extra read-only node won't help you scale.

Note that one feature you lose is the ability to have a transaction in which a write *might* happen, because a transaction cannot contain both a read-only node and a write-only node. If you must have atomicity, you will have to do your transaction on the writable cluster, or have two transactions: one to lookup and see if you need to change the data, and the second to perform the change—while first checking again that data didn't change since the previous transaction. Too much of this obviously lessens the effectiveness of read/write partitioning.

By Feature

The most obvious and effective partitioning scheme is by feature. Your site might offer different kinds of services that are functionally independent of each other, even though they are displayed to users as united. Behind the scenes, each feature uses a different set of tables. The rule of thumb is trivial: if you can put the tables in separate databases, then you can put these databases in separate clusters.

One concern in feature-based partitioning is that there are a few tables that still need to be shared. For example, even though the features are separate, they all depend on user settings that are stored in one table.

The good news is that it can be cheap to synchronize just this one table between all clusters. Especially if this table doesn't change often—how often do you get new users signing up for your service?—then synchronization overhead will be minimal.

If your database system doesn't let you synchronize individual tables, then you can do it in your code by writing to all clusters at the same time.

Partitioning by feature is terrific in that it lets you partition other parts of the stack, too. For example, you can also use a different set of web servers for each feature.

Also consider that some features might be candidates for using a “NoSQL” database (see below). Choose the best backend per feature.

By Section

Another kind of partitioning is sometimes called “sharding.” It involves splitting up tables into sections that can be placed in different databases. Some databases support sharding as part of their synchronization strategy, but you can also implement it in your code. The great thing about sharding is that it lets you create as many shards (and clusters) as you want. It's the key to the truly large scale.

Unfortunately, like partitioning by feature, sharding is not always possible. You need to also shard all related tables, so that queries can be self-contained within each shard. It's thus most appropriate for one-to-many data hierarchies. For example, if your application is a blog that supports comments, then you put some blogs and their comments on one shard, and others in another shard. However, if, say, you have a feature where blog posts can refer to other arbitrary blog posts, then querying for those would have to cross shard boundaries.

The best way to see where sharding is possible is to draw a diagram of your table relationships. Places in the diagram which look like individual trees—trunks spreading out into branches and twigs—are good candidates for sharding.

How to decide which data goes in which shard?

Sometimes the best strategy is arbitrary. For example, put all the even-numbered IDs in one shard, and the odd-numbered ones in another. This allows for straightforward growth because you can just switch it to division by three if you want three shards.

Another strategy might seem obvious: If you're running a site which shows different sets of data to different users, then why not implement it as essentially separate sites? For example, a social networking site strictly organized around individual cities could have separate database clusters per city.

A "region" can be geographical, but also topical. For example, a site hosting dance-related discussion forums might have one cluster for ballet and one for tango. A "region" can also refer to user types. For example, your social networking site could be partitioned according to age groups.

The only limitation is queries. You can still let users access profiles in other regions, but cross-regional relational queries won't be possible. Depending on what your application does, this could be a reasonable solution.

A great side-benefit to geographical partitioning is that you can host your servers at data centers within the geographical location, leading to better user experiences. Regional partitioning is useful even for "NoSQL" databases.

Coding Tips for Partitioning

If you organize your code well, it would be very easy to implement partitioning. You simply assign different database operations to use different connection pools. If it's by feature, then you can hard code it for those features. If it's sharding, then you add a switch before each operation telling it which connection pool to use.

For example:

```
def get_blogger_profile(user_id):
    connection = blogger_pool.get_connection()
    ...
    connection.close()

def get_blog_post_and_comments(blog_post_id):
    shard_id = object.id % 3
    connection = blog_pools[shard_id].get_connection()
    ...
    connection.close()
```

Unfortunately, some programming practices make such an effective, clean organization difficult.

Some developers prefer to use ORMs (object-relational mappers) rather than access the database directly. Many ORMs do not easily allow for partitioning, either because they support only a single database connection pool, or because they don't allow your objects to be easily shared between connections.

For example, your logic might require you to retrieve an "object" from the database, and only then decide if you need to alter it or not. If you're doing read/write partitioning, then you obviously want to read from the read partition. Some ORMs, though, have the object tied so strongly to an internal connection object that you can't trivially read it from one connection and save it into another. You'd either have to read the object initially from the write partition, minimizing the usefulness of read/write partitioning, or re-read it from the write partition when you realize you need to alter it, causing unnecessary overhead. (Note that you'll need to do this anyway if you need the write to happen in a transaction.)

Object oriented design is also problematic in a more general sense. The first principle of object orientation is "encapsulation," putting your code and data structure in one place: the class. This might make sense for business logic, but, for the purposes of re-factoring your data backend for partitioning or other strategies, you really don't want the data access code to be spread out among dozens of classes in your application. You want it all in one place, preferably even one source code file. It would let you plug in a whole new data backend strategy by replacing this source code file. For data-driven web development, you are better off not being too object oriented.

Even more generally speaking, organizing code together by mechanism or technology, rather than by "object" encapsulation, will let you apply all kinds of re-factorizations more easily, especially if you manage to decouple your application's data structures from any library-specific data structures.

Data Backends

Relational (SQL) databases such as MySQL were, for decades, the backbone of the web. They were originally developed as minimal alternatives to enterprise database servers such as Oracle Database and IBM's DB2. Their modest feature set allowed for better performance, smaller footprints, and low investment costs—perfect for web applications. The free software LAMP stack (Linux, Apache, MySQL and PHP) *was* the web.

Relational databases require a lot of synchronization overhead for clusters, limiting their scalability. Though partitioning can take you far, using a "NoSQL" database could take you even further.

Graph Databases

If your relational data structure contains arbitrary-depth relationships or many “generic” relationships forced into a relational model, then consider using a graph database instead. Not only will traversing your data be faster, but also the database structure will allow for more efficient performance. The implications for scalability can be dramatic.

Social networking applications are often used as examples of graph structures, but there are many others: forums with threaded and cross-referenced discussions, semantic knowledge bases, warehouse and parts management, music “genomes,” user-tagged media sharing sites, and many science and engineering applications.

Though fast, querying a complex graph can be difficult to prototype. Fortunately, the Gremlin and SPARQL languages do for graphs what SQL does for relational databases. Your query becomes coherent and portable.

A popular graph database is Neo4j, and it’s especially easy to use with Prudence. Because it’s JVM-based, you can access it internally from Prudence. It also has embedded bindings for many of Prudence’s supported languages, and supports a network REST interface which you can easily access via Prudence’s `document.external`.

Document Databases

If your data contains mostly “documents”—self-contained records with few relationships to other documents—then consider a document database.

Document databases allow for straightforward distribution and very fine-grained replication, requiring considerably less overhead than relational and graph databases. Document databases are as scalable as data storage gets: variants are used by all the super-massive internet services.

The cost of this scalability is the loss of your ability to do relational queries of your data. Instead, you’ll be using distributed map/reduce, or rely on an external indexing service. These are powerful tools, but they do not match relational queries in sheer speed of complex queries. Implementing something as simple as a many-to-many connection, the bread-and-butter of relational databases, requires some specialization. Document databases shine at listing, sorting and searching through extremely large catalogs of documents.

Candidate applications include online retail, blogs, wikis, archives, newspapers, contact lists, calendars, photo galleries, dating profiles. . . This is a long list, but by no means exhaustive of all that is possible in web applications. Many useful applications cannot be reduced to sets of lightly interconnected “documents” without giving up a lot of useful functionality. For example, merely adding social networking capabilities to a dating site would require complex relations that might be better handled with a graph database.

A popular document database is MongoDB. Though document-based, it has a few basic relational features that might be just good enough for your needs. Another is CouchDB, which is a truly distributed database. With CouchDB it’s trivial to replicate and synchronize data with clients’ desktops or mobile devices, and to distribute it to partners. It also supports a REST interface which you can easily access via Prudence’s `document.external`. And, both MongoDB and CouchDB use JavaScript extensively, making it natural to use with Prudence’s JavaScript flavor.

We’ve started a project to better integrate Prudence JavaScript with MongoDB. Check it out!

Column Databases

These can be considered as subsets of document databases. The “document,” in this case, is required to have an especially simple, one-dimensional structure.

This requirement allows optimization for a truly massive scale.

Column databases occupy the “cloud” market niche: they allow massive companies like Google and Amazon to cheaply offer database storage and services for third parties. See Google’s Datastore (based on Bigtable) and Amazon’s SimpleDB (based on Dynamo; actually, Dynamo is a “key/value” database, which is even more opaque than a column database).

Though you can run your own column database via open source projects like Cassandra (originally developed by/for Facebook) and HBase, the document databases mentioned above offer richer document structures and more features. Consider column databases only if you need truly massive scale, or if you want to make use of the cheap storage offered by “cloud” vendors.

Best of All Worlds

Of course, consider that it’s very possible to use both SQL and “NoSQL” (graph, document, column) databases together for different parts of your application. See backend partitioning, above.

Under the Hood

Prudence brings together many open source libraries, some of which were designed specifically for Prudence.

Consider this as Prudence’s “acknowledgments” page. Hundreds of people have worked on these libraries, and we’re grateful to all of them for sharing their hard work, for embracing open source licensing, and for adhering to design principles that allow reuse of their work in other projects, such as Prudence.

The JVM

How wonderful that the best-performing, most robust, secure and widely ported virtual machine is now open source? How wonderful that you can use it with JavaScript, Python, Ruby, Clojure, PHP and others languages?

We strongly recommend the JVM for enterprise and scalable internet applications, even if you’re not particularly fond of Java-the-programming-language. Treat Java, if you will, as the low machine-level language: Java is to the JVM as C is to Unix. You “drop down” to Java only if you have to do some machine-level work. Otherwise, use the higher-level languages.

Scripturian

Scripturian is Prudence’s “special sauce”: a small, magical library that makes sure that your code runs well on the JVM and can handle the concurrency introduced by HTTP requests and Restlet. It is developed in tandem with Prudence.

Our premise was this: to make Prudence applications easy to deploy, they could not be standard Java applications. The cycle of compilation and packaging required by Java is unnecessarily cumbersome. Although we could have implemented something like the on-the-fly Java compilation done in JSP, we felt that, if that’s the route to go, many exciting choices open up besides Java, and that these languages are more relevant to Prudence’s goals.

Unfortunately, we found that integrating JVM languages into Prudence was anything but trivial. Each implementation had its own idea of what integration could mean. We tried to standardize on JSR-223 (the Java scripting standard), but found that adherence to the specification was inconsistent, and that the specification itself is vague, especially when it comes to threading. We hacked and hacked and hacked. We even submitted patches to fix broken implementations of various languages. All in all, we probably spent more time on this than on any other aspect of Prudence.

The result is an API more abstract than JSR-223, but with a clear threading model. Under the hood, Scripturian contains many tricks and mechanisms to make each language work correctly and well, but you don’t have to worry about any of it. Scripturian just works!

Jython, JRuby, Clojure, Rhino, Quercus, Groovy

Prudence would hardly be as exciting if you had to use Java.

These open source language engines have allowed us to extend the power of Prudence, REST and the JVM to languages outside of Java. Some of these engines are large, complex projects, and are in fact the biggest libraries included in Prudence. We strongly recommend you join in the communities surrounding the language engines corresponding to your favorite flavor of Prudence.

Restlet

Prudence went through many in-house transformations before aligning itself strongly with Restlet. First, we experimented with Facelets, but ended up giving up on JSF, its complex lifecycle, and on the promise of component-based web development in general. Then, we designed REST architectures using servlets and Succinct templates, but found it awkward to force servlets into a REST architecture. Discovering Restlet was a breath of fresh air.

Restlet’s super-powers are three:

1. Clean abstraction of HTTP requests, responses and headers over best-of-breed engines, such as Grizzly, Jetty and Netty. Automatically gain the scalable advantages of non-blocking I/O and even asynchronous request handling (which will be even better supported in Restlet 2.1). Restlet transparently handles conditional requests, content negotiation, and other complicated HTTP labor.

2. Powerful URI routing and manipulation. Expose your service to users and APIs with elegance and coherence. Make sure the URI reaches its destination, with support for virtual hosting, rewriting, templating, and other useful real-world features. Restlet is truly the Swiss army knife of URIs!
3. Straightforward data representation and consumption through a diverse set of extensions. Expose your data using any standard format, and even convert it on-the-fly. Easily parse data received from clients.

Restlet is a great library, with a great ecosystem of extensions. In embracing it, though, we missed some of the advantages of having a servlet container: easy deployment and configuration, centralized logging, etc. We also missed having JSP at our fingertips to quickly push out dynamic HTML.

Prudence is meant to fill in these gaps.

Succinct

Succinct, like Scripturian, started as a part of Prudence, and was in fact one of its earliest components. It has since branched out into an independent library. We created it because we wanted straightforward, scalable templating built in to Prudence, and were unsatisfied by other open source offerings. We think you might like it. (If not, Prudence fully supports Velocity.)

Jygments

Yet another Prudence side-project!

For Prudence's debug mode, we wanted good syntax highlighting for viewing source code. We found nothing adequate enough in the Java world, though we fell in love with Pygments. For a while, we ran Pygments in Prudence via Jython, but found it too heavy for this particular use case. Thus, Jygments was born as a port of Pygments to Java.

H2

We're great fans of this lean and mean database engine! It has allowed us to distribute Prudence with a fully-functioning data-drive demo application, without any external dependencies. We're proud to introduce H2, through Prudence, to more people, and we believe you'll find it fast enough, reliable enough, and flexible enough for many production environments.

Hazelcast

This library is a dream come true: distributed, fault-tolerant implementations of the JVM's standard collection interfaces, with distributed task queues thrown into the box. We are confident that Hazelcast will help many Prudence users scale their applications easily and elegantly.

FAQ

REST

Why are plural URL forms for aggregate resources (`/animal/cats/`) preferred over singular forms (`/animal/cat/`)?

You'll see RESTful implementations that use either convention. The advantage of using the singular form is that you have less addresses, and what some people would call a more elegant scheme:

```
/animal/cat/12 -> Just one cat
/animal/cat/ -> All cats
```

Why add another URL when a single one is enough to do the work? One reason is that you can help the client avoid potential errors. For example, the client probably uses a variable to hold the ID of the cat and then constructs the URL dynamically. But, what if the client forgets to check for null IDs? It might then construct a URL in the form `"/animal/cat/"` which would then successfully access *all* cats. This can cause unintended consequences and be difficult to debug. If, however, we used this scheme:

```
/animal/cat/12 -> Just one cat  
/animal/cats/ -> All cats
```

... then the form “/animal/cat/” would route to our singular cat resource, which would indeed not find the cat and return the expected, debuggable 404 error. From this example, we can extract a good rule of thumb: clearly separate URLs at their base by usage, so that mistakes cannot happen. More addresses means more debuggability.

Languages

Why mix languages?

TODO

What are “in-flow” scriptlets?

TODO

Can different languages share state?

TODO

Can languages call each others’ functions?

TODO

Concurrency

Should I be worried?

TODO

How to make my code thread safe?

TODO

Scalability

I heard REST is very scalable. Is this true? Does this mean Prudence can support many millions of users?

Yes, if you know what you’re doing. See [“The Case for REST” \(page 54\)](#) and [“Scaling Tips” \(page 58\)](#) for in-depth discussions.

The bottom line is that it’s very easy to make your application scale poorly, whatever technology or architecture you use, and that Prudence, in embracing REST and the JVM, can more easily allow for best-practice scalable architectures than most other web platforms.

That’s not very reassuring, but it’s a fact of software and hardware architecture right now. Achieving massive scale is challenging.

Performance

Which Prudence flavor performs best?

TODO

How well does Prudence perform? How well does it scale?

First, recognize that there are two common uses for the term “scale.” REST is often referred to as an inherently scalable architecture, but that has more to do with project management than technical performance. This difference is address in the [“The Case for REST” \(page 54\)](#).

From the perspective of the ability to respond to user requests, there are three aspects to consider:

1. Serving HTTP Prudence comes with Grizzly, an HTTP server based on the JVM's non-blocking I/O API. Grizzly handles concurrent HTTP requests very well, and serves static files at scales comparable to popular HTTP servers. See the [tutorial \(page 5\)](#) for more information.

2. Generating HTML Prudence supports two modes for generating HTML (and other textual formats), each with its own performance characteristics:

Caching mode: First, the entire document is run, with its output sent into a buffer. This buffer is then cached, and *only then* sent to the client. This is the default mode and recommended for most documents. Scriptlets can be used to control the duration of the document's individual cache.

Deferred mode: Output is sent to the client *while* the document runs. This is recommended for documents that need to output a very large amount of text, which might take a long time, or that might otherwise encounter slow-downs while running. In either case, you want the client to receive ongoing output. The output of the document is not cached. Scriptlets can switch between modes according to changing circumstances. For example, to increase caching duration during heavy loads, to decrease it during periods where data changes often, or to stream in the case of an expected large output. See the [tutorial \(page 5\)](#) for more information.

3. Running code There may be a delay when starting up a specific language engine in Prudence for the first time in an application, as it loads and initializes itself. Then, there may be a delay when accessing a dynamic web page or resource for the first time, or after it has been changed, as it might require compilation. Once it's up and running, though, your code performs and scale very well—as well as you've written it. You need to understand concurrency and make sure you make good choices to handle coordination between threads accessing the same data. If all is good, your code will actually perform better throughout the life of the application. The JVM learns and adapts as it runs, and performance can improve the more the application is used.

If you are performing CPU-intensive or time-sensitive tasks, then it's best to profile these code segments precisely. Exact performance characteristics depend on the language and engine used. The Bechmarks Game can give you some comparisons of different language engines running high-computation programs. In any case, if you have a piece of intensive code that really needs to perform well, it's probably best to write it in Java and access it from the your language. You can even write it in C or assembly, and have it linked to Java via JNI.

If you're not doing intensive computation, then don't worry too much about your language being "slow." It's been shown that for the vast majority of web applications, the performance of the web programming language is rarely the bottleneck. The deciding factors are the usually performance of the backend data-driving technologies and architectures.

Licensing

The author is not a lawyer. This is not legal advice, but a personal, and possibly wrong interpretation. The wording of the license itself supersedes anything written here.

Does the LGPL mean I can't use Prudence unless my product is open sourced?

The GPL family of licenses restrict your ability to *redistribute* software, not to use it. You are free to use Prudence as you please within your organization, even if you're using it to serve public web sites (though with no warranty nor an implicit guarantee of support from Three Crickets, the copyright holder).

The GPL would thus only be an issue if you're selling, or even giving away, a product that *would include* Prudence.

Prudence uses the Lesser GPL, which has even less restrictions on redistribution than the regular GPL. As long as you do not alter Prudence in any way, you can include Prudence in any product, free or non-free. (Actually, Prudence uses version 3 of the Lesser GPL, which requires your product, even if it's not free software, to at least not restrict users' ownership of data via schemes such as DRM if you want to include Prudence in its distribution.)

Even if your product does not qualify for including Prudence in it, you always have the option of distributing your product without Prudence, and instructing your customers to download and install Prudence on their own.

Three Crickets, the original developers of Prudence, are not trying to force you to purchase it. Instead, they hope to encourage you 1) to pay Three Crickets for consultation, support and development services for Prudence, and 2) to release your own product as free software, thereby truly sharing your innovation with all of society.

We understand that in some cases open sourcing your product is impossible. As a last resort, we offer you a commercial license. Please contact us for details.

Why the LGPL and not the GPL?

The Lesser GPL used to be called the “Library GPL,” and was originally drafted for glibc. It represents a certain admission of defeat: there are so many alternatives to our library out there, that you might not consider using our library under GPL.

In the case of Linux, the GPL has done a wonderful job in convincing vendors to open source their code in order to ship their products with Linux inside. It just doesn’t seem likely that they would do the same for Prudence.

Note that the LGPL version 3 has a clause allowing you to “upgrade” Prudence to the full GPL for inclusion in your GPL-ed product. This is a terrific feature, and another reason to love this excellent license.