# Advanced GATE Embedded

## Module 8, part 3

## Eleventh GATE Training Course
## June 2018

© 2018 The University of Sheffield

This material is licenced under the Creative Commons

Attribution-NonCommercial-ShareAlike Licence

(http://creativecommons.org/licenses/by-nc-sa/3.0/)

## Outline

1 GATE in Multi-threaded/Web Applications
   - Multi-threading and GATE
   - Servlet Example
   - The Spring Framework
   - Making your own PRs duplication-friendly

2 Extending GATE
   - Adding new document formats

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Outline

1 GATE in Multi-threaded/Web Applications
   - Multi-threading and GATE
   - Servlet Example
   - The Spring Framework
   - Making your own PRs duplication-friendly

2 Extending GATE
   - Adding new document formats

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Introduction

- Scenario:
    - Implementing a web application that uses GATE Embedded to process requests.
    - Want to support multiple concurrent requests
    - Long running process - need to be careful to avoid memory leaks, etc.
- Example used is a plain HttpServlet
    - Principles apply to other frameworks (struts, Spring MVC, Grails...)

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Setting up

- GATE libraries in `WEB-INF/lib`
    - dependencies from central via your build tool, or copy the `lib` folder from GATE installation
- Optional:
    - user config file if you need to configure things in there (e.g. "add space on markup unpack")

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

## GATE in a Multi-threaded Environment

- GATE initialization needs to happen once (and only once) before any other GATE APIs are used.
- The `Factory` is synchronized internally, so safe for use in multiple threads.
- Individual PRs/controllers are *not* safe – must not use the same PR instance concurrently in different threads
  - this is due to the design of runtime parameters as Java Beans properties.
- Individual LRs (documents, ontologies, etc.) are only thread-safe when accessed read-only by *all* threads.
  - if you need to share an LR between threads, be sure to synchronize (e.g. using `ReentrantReadWriteLock`)

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Initializing GATE using a ServletContextListener

`ServletContextListener` called by container at startup and shutdown (only startup method shown).

```
1  public void contextInitialized(ServletContextEvent e) {
2    // if you want a user config
3    Gate.runInSandbox(false);
4    ServletContext ctx = e.getServletContext();
5    File userConfig = new File(
6        ctx.getRealPath("/WEB-INF/user.xml"));
7    Gate.setUserConfigFile(userConfig);
8    // otherwise ignore this and use the default (sandbox=true)
9
10   // initialise GATE
11   Gate.init();
12 }
```

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Initializing GATE using a ServletContextListener

You must register the listener in `web.xml`

```
1 <listener>
2   <listener-class>
3     gate.web.example.GateInitListener
4   </listener-class>
5 </listener>
```

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Handling Concurrent Requests

Naïve approach – new PRs for every request

```
1  public void doPost(request, response) {
2    ProcessingResource pr = Factory.createResource(...);
3    try {
4      Document doc = Factory.newDocument(
5          getTextFromRequest(request));
6      try {
7        // do some stuff
8      }
9      finally {
10       Factory.deleteResource(doc);
11     }
12   }
13   finally {
14     Factory.deleteResource(pr);
15   }
16 }
```

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Handling Concurrent Requests

Naïve approach – new PRs for every request

```
1  public void doPost(request, response) {
2    ProcessingResource pr = Factory.createResource(...);
3    try {
4      Document doc = Factory.newDocument(
5          getTextFromRequest(request));
6      try {
7        // do some stuff
8      }
9      finally {
10       Factory.deleteResource(doc);
11     }
12   }
13   finally {
14     Factory.deleteResource(pr);
15   }
16 }
```

Many levels of try/finally
– make sure you clean up
even when errors occur

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Problems with Naïve Approach

- Guarantees no interference between threads
- But inefficient, particularly with complex PRs (large gazetteers, JAPE grammars, etc.)

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Take Two: using ThreadLocal

Store the PR/Controller in a thread-local variable

```
1  private ThreadLocal<CorpusController> controller =
2      new ThreadLocal<CorpusController>() {
3
4    protected CorpusController initialValue() {
5      return loadController();
6    }
7  };
8
9  private CorpusController loadController() { ... }
10
11 public void doPost(request, response) {
12   CorpusController c = controller.get();
13   // do stuff with the controller
14 }
```

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# An Improvement. . .

- Only initialise resources once per thread
- Interacts nicely with typical web server thread pooling
- But if a thread dies (e.g. with an exception), no way to clean up its controller

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# One Solution: Object Pooling

- Manage your own pool of Controller instances
- Take a controller from the pool at the start of a request, return it (in a finally!) at the end
- Number of instances in the pool determines maximum concurrency level

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Simple Example of Pooling

Setting up and cleaning up:

```
1  private BlockingQueue<CorpusController> pool;
2
3  public void init() {
4    pool = new LinkedBlockingQueue<CorpusController>();
5    for(int i = 0; i < POOL_SIZE; i++) {
6      pool.add(loadController());
7    }
8  }
9
10 public void destroy() {
11   for(CorpusController c : pool) {
12     Factory.deleteResource(c);
13   }
14 }
```

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Simple Example of Pooling

Processing requests:

```
15  public void doPost(request, response) {
16     CorpusController c = pool.take();
17     try {
18        // do stuff
19     }
20     finally {
21        pool.add(c);
22     }
23  }
```

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Simple Example of Pooling

Processing requests:

```
15  public void doPost(request, response) {
16    CorpusController c = pool.take();
17    try {
18      // do stuff
19    }
20    finally {
21      pool.add(c);
22    }
23  }
```

This blocks when the pool is empty. Use `poll` for non-blocking check.

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Creating the pool

- Typically to create the pool you would use `PersistenceManager` to load a saved application several times.
- But this is not always optimal, e.g. large gazetteers consume lots of memory.
- GATE provides API to *duplicate* an existing instance of a resource: `Factory.duplicate(existingResource)`.
- By default, this simply calls `Factory.createResource` with the same class name, parameters, features and name.
- But individual Resource classes can override this by implementing the `CustomDuplication` interface (more later).
  - e.g. `DefaultGazetteer` uses a `SharedDefaultGazetteer` — same behaviour, but shares the in-memory representation of the lists.

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

## Other Caveats

- With most PRs it is safe to create lots of identical instances
- But *not all*!
    - e.g. training a machine learning model with the learning framework
    - but it is safe to have several instances *applying* an existing model.
- When using `Factory.duplicate`, be careful not to duplicate a PR that is being used by another thread
    - i.e. either create all your duplicates up-front or else keep the original prototype "pristine".

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
**The Spring Framework**
Making your own PRs duplication-friendly

# Exporting the Grunt Work: Spring

- http://www.springsource.org/
- "Inversion of Control"
- Configure your business objects and connections between them using XML, Groovy or Java annotations.
- Handles application startup and shutdown
- GATE provides helper library to initialise GATE, load saved applications, etc.
- Built-in support for object pooling
- Web application framework (Spring MVC)
- Used by other frameworks (Grails, etc.)

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
**The Spring Framework**
Making your own PRs duplication-friendly

# Initializing GATE via Spring XML

```
1  <beans
2      xmlns="http://www.springframework.org/schema/beans"
3      xmlns:gate="http://gate.ac.uk/ns/spring">
4    <gate:init run-in-sandbox="false"
5        user-config-file="gate-files/user-gate.xml" />
6
7    <gate:extra-plugin group-id="uk.ac.gate.plugins"
8                       artifact-id="annie"
9                       version="8.5" />
10 </beans>
```

- Paths can be full URLs (`file:/...`) or *resource* paths that are resolved appropriately by Spring

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
**The Spring Framework**
Making your own PRs duplication-friendly

## Loading a Saved Application

To load an application state saved from GATE Developer:

```
1 <gate:saved-application
2         id="myApp"
3         location="gate-files/application.xgapp"
4         scope="prototype" />
```

- scope="prototype" means create a new instance each time we ask for it
- Default scope is "singleton" — one instance is created at startup and shared.

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
**The Spring Framework**
Making your own PRs duplication-friendly

# Duplicating an Application

- Alternatively, load the application once and then duplicate it

```
1 <gate:duplicate id="myApp" return-template="true">
2   <gate:saved-application location="..." />
3 </gate:duplicate>
```

- `<gate:duplicate>` creates a new duplicate each time we ask for the bean.
- `return-template` means the original controller (from the `saved-application`) will be returned the first time, then duplicates thereafter.
- Without this the original is kept pristine and only used as a source for duplicates.

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
**The Spring Framework**
Making your own PRs duplication-friendly

# Worked example – a Spring Boot webapp

- *Spring Boot* is a framework to get a Spring-based application up and running with a few lines of code.
- "Convention over configuration" approach, providing sensible defaults that you can override.
- Builds using Maven 3 (or Gradle)
- Example is a Spring MVC web application
  - see hands-on materials for the code.

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
**The Spring Framework**
Making your own PRs duplication-friendly

# Setting up a Spring Boot app

- One "parent" and one "plugin" to add to Maven POM
- Add dependencies on the relevant "starter" modules, which themselves depend on the required libraries
    - `spring-boot-starter-web` for a basic Spring MVC application
    - your preferred view technology, in our case `spring-boot-starter-thymeleaf`
- And whatever other libraries you require
    - in this case we need `gate-spring`, which in turn depends on `gate-core`

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
**The Spring Framework**
Making your own PRs duplication-friendly

## Setting up a Spring Boot app

- Entry point is a simple boilerplate class with annotations.
- Automatically scans sub-packages for other annotated classes (controllers, etc.).

```
1  package gatetutorial;
2  import org.springframework.boot.SpringApplication;
3  import org.springframework.boot.autoconfigure.
       SpringBootApplication;
4
5  @SpringBootApplication
6  public class TutorialApp {
7
8    public static void main(String... args) {
9      SpringApplication.run(TutorialApp.class, args);
10   }
11 }
```

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
**The Spring Framework**
Making your own PRs duplication-friendly

## A simple example controller

- Our example is a single controller that presents an HTML form to enter text.
- When form is submitted, process the document with a GATE application and show the document features as a table.

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
**The Spring Framework**
Making your own PRs duplication-friendly

# Spring pooling support

- Spring's built-in AOP features offer support for object pooling.
- Given a bean definition, we can expose a *proxy* object with the same behaviour backed by a pool of instances
- *Each method call* on the proxy is dispatched to one of the objects in the pool.
- Each target bean is guaranteed to be accessed by no more than one thread at a time.
- When the pool is empty, can configure further requests to block or fail.

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
**The Spring Framework**
Making your own PRs duplication-friendly

# The `gate:pooled-proxy` helper

- The machinery is all Spring but complex to configure.
- `gate-spring` provides a helper to automate this in Spring XML configuration.
- Don't pool GATE applications directly, instead pool a helper class that calls GATE.

```
1 <bean id="gateService"
2       class="gatetutorial.service.GateService">
3   <gate:pooled-proxy max-size="3"
4              initial-size="3" />
5 </bean>
```

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
**The Spring Framework**
Making your own PRs duplication-friendly

# More advanced pooling options

- Many more options to control the pool, e.g. for a pool that grows as required and shuts down instances that have been idle for too long, and where excess requests fail rather than blocking:

```
1 <gate:pooled-proxy
2     max-size="10"
3     max-idle="3"
4     time-between-eviction-runs-millis="180000"
5     min-evictable-idle-time-millis="90000"
6     when-exhausted-action-name="WHEN_EXHAUSTED_FAIL"
7 />
```

- Under the covers, `<gate:pooled-proxy>` creates a Spring `CommonsPoolTargetSource`, attributes correspond to properties of this class.
- See the Spring documentation for full details.

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
**The Spring Framework**
Making your own PRs duplication-friendly

# The GateService

- The GateService is written assuming single-threaded access.

```java
public class GateService {
  // will be injected automatically
  @Autowired private CorpusController application;

  private Corpus corpus;

  @PostConstruct
  public void init() throws GateException {
    corpus = Factory.newCorpus("GateService");
    application.setCorpus(corpus);
  }
```

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# The GateService

```
13    @PreDestroy
14    public void destroy() {
15      Factory.deleteResource(corpus);
16      // not strictly necessary as gate:duplicate will handle this
17      Factory.deleteResource(application);
18    }
```

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# The GateService

```
20   public FeatureMap processWithGate(Document doc)
21                       throws GateException {
22     try {
23       corpus.add(doc);
24       application.execute();
25       return doc.getFeatures();
26     } finally {
27       corpus.clear();
28     }
29   }
30 }
```

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Tying it together

- Initialize GATE and configure service pool in XML (in `src/main/resources`)

```
1 <gate:init />
2 <gate:duplicate id="gateApplication"
3                 return-template="true">
4   <gate:saved-application
5       location="gate-files/application.xgapp" />
6 </gate:duplicate>
7
8 <bean id="gateService"
9     class="gatetutorial.service.GateService">
10  <gate:pooled-proxy max-size="3"
11              initial-size="3" />
12 </bean>
```

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Tying it together

- Add an annotation to `TutorialApp` to load the XML.

```
1 import org.springframework.context.annotation.
       ImportResource;
2 // ...
3
4 @SpringBootApplication
5 @ImportResource("/gate-beans.xml")
6 public class TutorialApp {
```

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
**The Spring Framework**
Making your own PRs duplication-friendly

# Tying it together

- And finally, autowire the `GateService` into controller, and call its methods without having to worry about threading.

```
1  // imports as required
2
3  @Controller
4  public class GateController {
5
6    @Autowired
7    private GateService gateService;
```

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
**The Spring Framework**
Making your own PRs duplication-friendly

# Tying it together

```
9    @RequestMapping(value="/",
10                   method = RequestMethod.POST)
11   public String process(
12       @ModelAttribute("params")
13       AnnotationRequest params,
14       Map<String, Object> model) throws GateException {
15     Document doc = // extract text from request
16     try {
17       FeatureMap features =
18             gateService.processWithGate(doc);
19       model.put("features", features);
20       return "index";
21     } finally { Factory.deleteResource(doc); }
22   }
23 }
```

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
**The Spring Framework**
Making your own PRs duplication-friendly

## Exercise 1: The Spring Boot example

- In hands-on/webapps you have the source code for the Spring Boot example we've been discussing.
- What's provided:
  - the pom.xml with the necessary dependencies,
  - source code for the controller and GateService (in src/main/java),
  - the Thymeleaf view with the text entry form and results table,
  - configuration in src/main/resources, including GATE config files and the bean definition XML.

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
**The Spring Framework**
Making your own PRs duplication-friendly

## Exercise 1: The Spring Boot example

- What's missing:
    - the GATE application iself. . .
- Use the document statistics PR from earlier.
- In GATE Developer, create a "corpus pipeline" application containing a tokeniser and your statistics PR.
- Right-click on the application and "Export for GATE Cloud".
    - This will save the application state along with all the plugins it depends on in a single zip file.
- Unpack the zip file under
  `src/main/resources/gate-files`
    - don't create any extra directories – you need
      `application.xgapp` and `maven-cache.gate` to end up
      in `gate-files`.

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
**The Spring Framework**
Making your own PRs duplication-friendly

## Exercise 1: The Spring Boot example

- Now you can run the application – in `hands-on/webapps` run `mvn spring-boot:run`
- Browse to `http://localhost:8080/`, enter some text and submit
- Watch the log messages. . .
- Notice the result page includes a feature "handledBy" – each service instance in the pool has a unique ID.
- Multiple submissions go to different instances in the pool.
- Try editing `src/main/resources/gate-beans.xml` and change the pooling configuration.
- Test concurrent requests – the service has a built-in delay to simulate a slow application.

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Not Just for Webapps

- Spring Boot (and Spring in general) isn't just for web applications
- You can use the same tricks in other embedded apps
- GATE provides a `DocumentProcessor` interface suitable for use with Spring pooling, which exposes one `void` method `processDocument`

```
1 <bean id="processor"
2   class="gate.util.LanguageAnalyserDocumentProcessor">
3   <property name="analyser" ref="gateApplication"/>
4   <gate:pooled-proxy max-size="3" initial-size="3" />
5 </bean>
```

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
**The Spring Framework**
Making your own PRs duplication-friendly

# A simple command-line app

```
1  @Component
2  public class GateCommand implements CommandLineRunner {
3    @Autowired private DocumentProcessor proc;
4
5    public void run(String... args) throws Exception {
6      Document doc = Factory.newDocument(args[0]);
7      try {
8        proc.processDocument(doc);
9        // ...
10     } finally { Factory.deleteResource(doc); }
11   }
12 }
```

The main entry point `TutorialApp` is unchanged from the web example.

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
**The Spring Framework**
Making your own PRs duplication-friendly

# A JMS message consumer

```
1  @Component
2  public class Receiver {
3    @Autowired private DocumentProcessor proc;
4
5    @JmsListener(destination = "someQueue",
6                 concurrency = "3")
7    public void receive(String stringMessage) {
8      Document doc = Factory.newDocument(stringMessage);
9      try {
10       proc.processDocument(doc);
11       doStuffWithResults(doc);
12     } finally { Factory.deleteResource(doc); }
13   }
14 }
```

In this case we need to add @EnableJms to the entry point class,
and relevant dependencies to the POM.

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Conclusions

Two golden rules:

- Only use a GATE Resource in one thread at a time
- Always clean up after yourself, even if things go wrong
  (deleteResource in a finally block).

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Duplication and Custom PRs

- Recap: by default, `Factory.duplicate` calls `createResource` passing the same type, parameters, features and name
- This can be sub-optimal for resources that rely on large read-only data structures that could be shared
- If this applies to your custom PR you can take steps to make it handle duplication more intelligently
- For simple cases: *sharable properties*, for complex cases: *custom duplication*.

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

## Sharable properties

- A way to share object references between a PR and its duplicates
- A JavaBean setter/getter pair with the setter annotated (same as for `@CreoleParameter`)

```
1 private Map dataTable;
2
3 public Map getDataTable() { return dataTable; }
4
5 @Sharable
6 public void setDataTable(Map m) {
7   dataTable = m;
8 }
```

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

## Sharable properties

- Default duplication algorithm will get property value from original and set it on the duplicate before calling `init()`
- `init()` must detect when sharable properties have been set and react appropriately.

```java
1  public Resource init() throws /* ... */ {
2    if(dataTable == null) {
3      // only need to build the data table if we weren't given a shared one
4      buildDataTable();
5    }
6  }
7
8  public void reInit() throws /* ... */ {
9    // clear sharables on reInit
10   dataTable = null;
11   super.reInit();
12 }
```

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Sharable properties – Caveats

- Anything shared between PRs *must* be thread-safe
  - use appropriate synchronization if any of the threads modifies the shared object (e.g. a ReentrantReadWriteLock which is itself @Sharable).
  - or (for the dataTable example), use an inherently safe class such as ConcurrentHashMap
  - for shared counter, use AtomicInteger
- If you use sharable properties, take care not to break reInit

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

## Exercise 2: Multi-threaded cumulative statistics

- `hands-on/shared-stats` contains a variation on the `DocStats` PR that keeps a running total of the number of Tokens it has seen.
- Build this (using the Maven pom), load the plugin, create an application containing a tokeniser and a "Shared document statistics" PR, export for GATE Cloud and unzip into your webapp as before.
- Try posting some requests to the webapp.
- You will see a `running_total` feature, but this is per handler, not global across handlers.

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

## Exercise 2: Multi-threaded cumulative statistics

- Your task: make the running total global.
- Make the totalCount field into a sharable property
    - it's already a thread-safe AtomicInteger
    - add a getter and setter, with the right annotation
    - init() logic to handle the shared/non-shared cases
    - implement a sensible reInit()
- You will need to re-build your PR and re-export (or just copy the compiled plugin JAR file to the appropriate place under gate-files/maven-cache.gate)

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

## Exercise 2: Solution

Getter and setter:

```
1 private AtomicInteger totalCount;
2
3 public AtomicInteger getTotalCount() {
4   return totalCount;
5 }
6
7 @Sharable
8 public void setTotalCount(AtomicInteger tc) {
9   this.totalCount = tc;
10 }
```

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

## Exercise 2: Solution

`init()` and `reInit()`:

```
1  public Resource init() throws
2          ResourceInstantiationException {
3    if(totalCount == null) {
4      totalCount = new AtomicInteger(0);
5    }
6    return this;
7  }
8
9  public void reInit() throws
10          ResourceInstantiationException {
11    totalCount = null;
12    super.reInit();
13  }
```

`execute()` is unchanged.

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
**Making your own PRs duplication-friendly**

# Custom Duplication

- For more complex cases, a resource can take complete control of its own duplication by implementing `CustomDuplication`
- This tells `Factory.duplicate` to call the resource's own `duplicate` method instead of the default algorithm.

```
1 public Resource duplicate(DuplicationContext ctx)
     throws ResourceInstantiationException;
```

- `duplicate` should create and return a duplicate, which need not be the same concrete class but must "behave the same"
    - Defined in terms of implemented interfaces.
    - Exact specification can be found in the `Factory.duplicate` JavaDoc.

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Custom Duplication

- If you need to duplicate other resources, use the two-argument `Factory.duplicate`, passing the `ctx` as the second parameter, to preserve object graph
  - two calls to `Factory.duplicate(r, ctx)` for the same resource `r` in the same context `ctx` will return the same duplicate.
  - calls to the single argument `Factory.duplicate(r)` or to the two-argument version with different contexts will return different duplicates.
- Can call the default duplicate algorithm (bypassing the `CustomDuplication` check) via `Factory.defaultDuplicate`
  - it is safe to call `defaultDuplicate(this, ctx)`, but calling `duplicate(this, ctx)` from within its own custom `duplicate` will cause infinite recursion!

GATE in Multi-threaded/Web Applications
Extending GATE

Multi-threading and GATE
Servlet Example
The Spring Framework
Making your own PRs duplication-friendly

# Custom Duplication Example (SerialController)

```
1  public Resource duplicate(DuplicationContext ctx)
2          throws ResourceInstantiationException {
3    // duplicate this controller in the default way - this handles subclasses nicely
4    Controller c = (Controller)Factory.defaultDuplicate(
5                  this, ctx);
6
7    // duplicate each of our PRs
8    List<ProcessingResource> newPRs =
9          new ArrayList<ProcessingResource>();
10   for(ProcessingResource pr : prList) {
11     newPRs.add((ProcessingResource)Factory.duplicate(
12           pr, ctx));
13   }
14   // and set this duplicated list as the PRs of the copy
15   c.setPRs(newPRs);
16
17   return c;
18 }
```

# Outline

## Adding new document formats

- GATE provides default support for reading many source document formats, including plain text, HTML, XML, PDF, DOC, ...
- The mechanism is extensible – the format parsers are themselves resources, which can be provided via CREOLE plugins.
- GATE chooses the format to use for a document based on *MIME type*, deduced from
  - explicit `mimeType` parameter
  - file extension (for documents loaded from a URL)
  - web server supplied Content-Type (for documents loaded from an `http:` URL)
  - "magic numbers", i.e. signature content at or near the beginning of the document

## The `DocumentFormat` resource type

- A GATE document format parser is a resource that extends the `DocumentFormat` abstract class or one of its subclasses.
- Override `unpackMarkup` method to do the actual format parsing, creating annotations in the `Original markups` annotation set and optionally modifying the document content.
- Override `init` to register with the format detection mechanism.
- In theory, can take parameters like any other resource ...
- ... but in practice most formats are singletons, created as *autoinstances* when their defining plugin is loaded.

# Repositioning info

- Some formats are able to record *repositioning info*
- Associates the offsets in the extracted text with their corresponding offsets in the original content.
- Allows you to save annotations as markup inserted into the original content.
- Of the default formats, only HTML can do this reliably.
    - If you're interested, see the `NekoHtmlDocumentFormat`

# Implementing a `DocumentFormat`

- Define a class that extends `DocumentFormat`, with CREOLE metadata

```
1 import gate.*;
2 import gate.creole.metadata.*;
3 import gate.corpora.*;
4
5 @CreoleResource(name = "Example DocumentFormat",
6     autoinstances = {@AutoInstance})
7 public class MyDocumentFormat
8         extends TextualDocumentFormat {
9     // ...
10 }
```

- `autoinstances` causes GATE to create an instance of this resource automatically when the plugin is loaded.

# `DocumentFormat` methods

- Most formats need to override three or four methods.

- `supportsRepositioning` to specify whether or not the format is capable of collecting repositioning info – most aren't

```
1 public Boolean supportsRepositioning() {
2   return false;
3 }
```

## `DocumentFormat` methods

- Two variants of `unpackMarkup`
- If you don't support repositioning then best to extend `TextualDocumentFormat` and just override the simple one:

```
1 public void unpackMarkup(Document doc)
2       throws DocumentFormatException {
3   AnnotationSet om = doc.getAnnotations(
4       GateConstants.ORIGINAL_MARKUPS_ANNOT_SET_NAME);
5   // Make changes to the document content, add annotations to om
6 }
```

- Other variant (for repositioning formats) is implemented in terms of this one by `TextualDocumentFormat`

# `DocumentFormat` methods

- Finally, `init` to register the format with GATE
- Mostly boilerplate, using protected `Map` fields defined in `DocumentFormat`

```
1  public Resource init() throws
       ResourceInstantiationException {
2    MimeType mime = new MimeType("text", "x-special");
3    mimeString2ClassHandlerMap.put(
4        mime.getType()+ "/" + mime.getSubtype(), this);
5    mimeString2mimeTypeMap.put(
6        mime.getType() + "/" + mime.getSubtype(), mime);
7    suffixes2mimeTypeMap.put("spec", mime);
8    magic2mimeTypeMap.put("==special==", mime);
9
10   setMimeType(mime);
11   return this;
12 }
```

# Registering a document format

```
2   MimeType mime = new MimeType("text", "x-special");
3   mimeString2ClassHandlerMap.put(
4       mime.getType()+ "/" + mime.getSubtype(), this);
```

- Create a `MimeType` object representing the "primary" MIME type for this format.
- Register this object as the handler for this MIME type.

```
5   mimeString2mimeTypeMap.put(
6       mime.getType() + "/" + mime.getSubtype(), mime);
```

- Establish a mapping between the MIME string "text/x-special" and the primary `MimeType` object.
- To register a format against several different MIME types (e.g. text/json and application/json), add them to the `mimeString2mimeTypeMap`

# Registering a document format

```
7    suffixes2mimeTypeMap.put("spec", mime);
```

- Register the file suffixes (not including the leading dot) that the format will handle, by mapping them to the primary `MimeType`
- Can add several different suffixes for the same type (txt, text, etc.)

```
8    magic2mimeTypeMap.put("==special==", mime);
```

- Add "magic numbers" – strings whose presence within the first 2kB of content will select the format
- E.g. "<?xml" is a strong predictor of XML documents.

# Registering a document format

```
10    setMimeType(mime);
11    return this;
```

- Boilerplate.
- Suffixes and magic numbers are optional – don't use them if they don't make sense for your particular format.
- ...but if neither are specified then only documents created with an explicit `mimeType` parameter will use the format.

# Exercise: Document format registration

- `hands-on/yam-format` contains a simple document format implementation.
- Processes text files in the "YAM" format (the Wiki markup syntax used on `http://gate.ac.uk`).
- `unpackMarkup` has been written for you.
- Annotates `*bold*`, `_italic_` and `^teletype^` text, and section headings (lines starting `%1`, `%2`, etc.).
- For simplicity, does not modify the text or do repositioning, only adds Original markups annotations.

# Exercise: Document format registration

- Your task – write the `init` method registration code
    - Primary MIME type "text/x-yam"
    - File suffixes ".yam" and ".gate"
    - No magic numbers
- To test, `mvn install`, then load the `yam-format` plugin into GATE Developer.
    - Note the auto-instance created when the plugin loads
- Create a document from the `overview.yam` file and inspect the Original markups.

# Solution

```java
@Override
public Resource init() throws
    ResourceInstantiationException {
  // create the primary MIME type
  MimeType mime = new MimeType("text", "x-yam");
  // usual boilerplate
  mimeString2ClassHandlerMap.put(
      mime.getType()+ "/" + mime.getSubtype(), this);
  mimeString2mimeTypeMap.put(
      mime.getType() + "/" + mime.getSubtype(), mime);
  // file suffixes
  suffixes2mimeTypeMap.put("yam", mime);
  suffixes2mimeTypeMap.put("gate", mime);
  // more boilerplate
  setMimeType(mime);
  return this;
}
```