

# Building Data Systems with iRODS and Golang

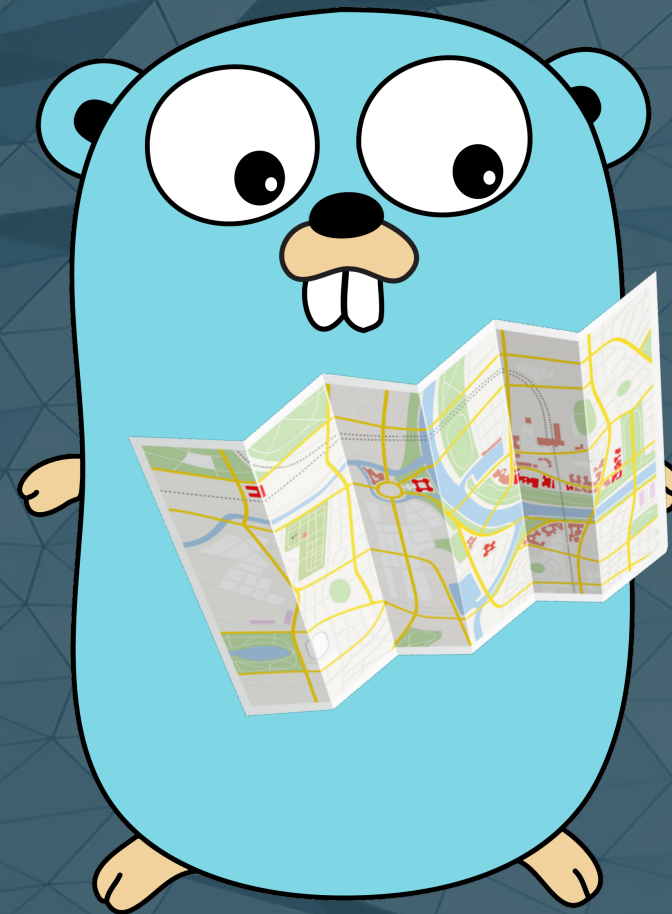


John Jacquay | [john@bioteam.net](mailto:john@bioteam.net)



# Presentation Road Map

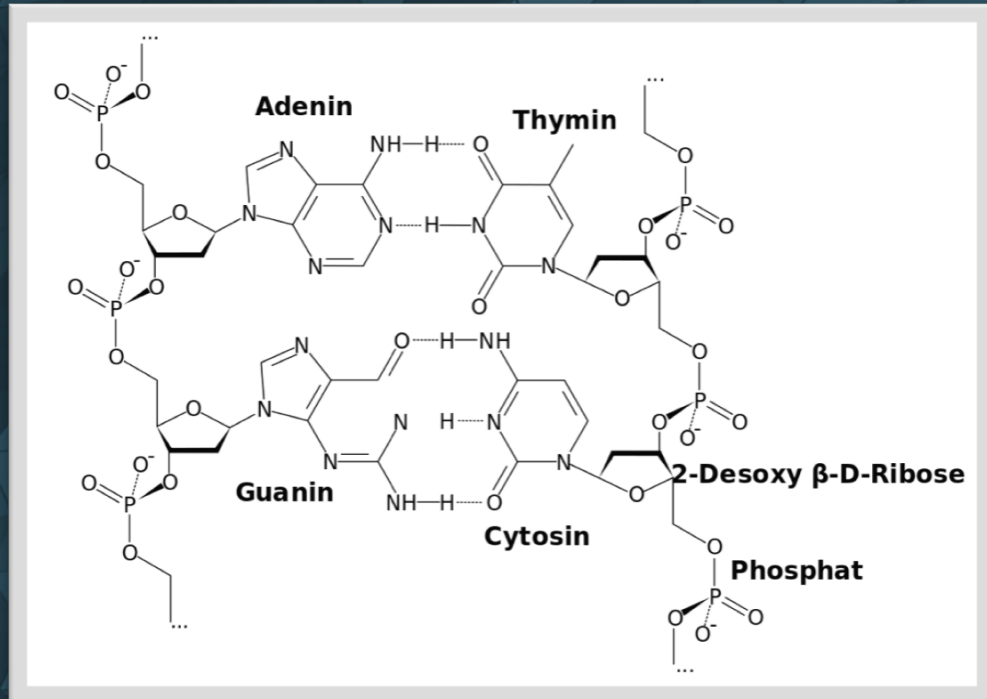
- My vision of an iRODS data system in the domain of life sciences
  - Data formats & automated bioinformatics pipelines
  - Tiered storage
  - GUI for automation design, data discovery, delivery
  - Open-source
- Why Golang?
- Integration Points and Use Cases
  - Client – A deep dive into GoRODS
  - Microservices and iRODS DSL
  - Rule Engine
  - RPC API: Client & Server
- That's it!





# iRODS In The Domain of Life Sciences

## DNA Sequencing





# iRODS In The Domain of Life Sciences

## File formats for High-Throughput sequencing

- Instrument generates BCL
  - Utility translates raw BCL to FASTQ
  - QC is run
  - Cleaning / trimming takes place
  - Assembly, Alignment
  - FASTQ to BAM/SAM
  - BAM/SAM -> diff -> VCF
- 
- Raw: BCL
  - Primary: FASTQ
  - Secondary: SAM, BAM
  - Tertiary: VCF





# iRODS In The Domain of Life Sciences

Tiered Storage of Sequencing Data: Optimizing processing & performance

Raw: BCL

Primary: FASTQ

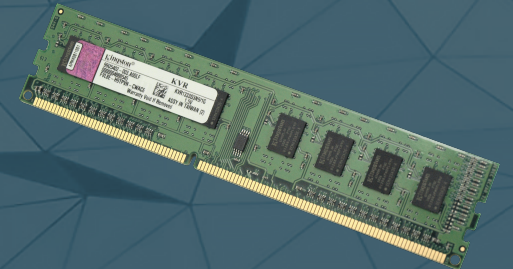
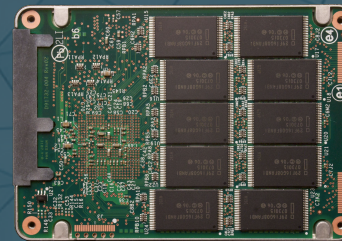
Secondary: SAM, BAM

Tertiary : VCF

iRODS Resource  
(Consumer)



iRODS Resource  
(Consumer)

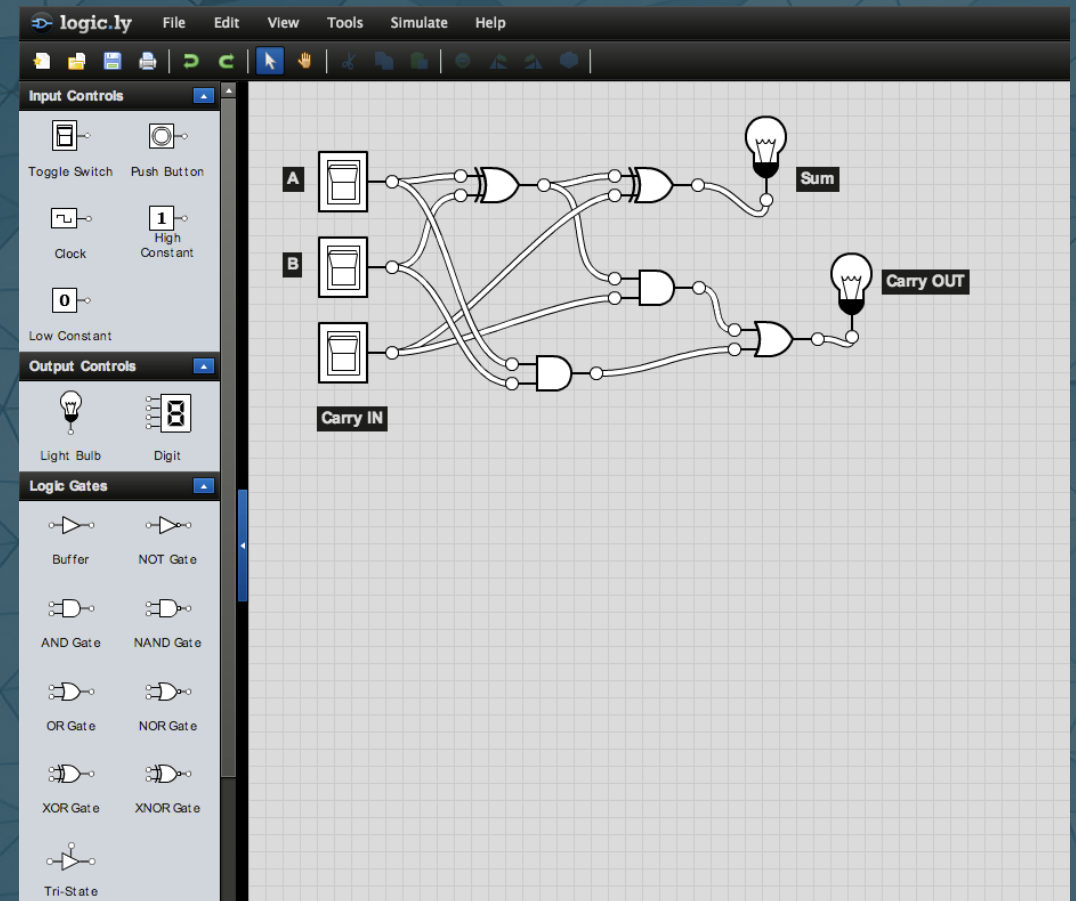




# iRODS In The Domain of Life Sciences

The full picture we're working towards

- Web UI: Primary Interface
  - Invoke workflows of rules / microservices / processes
  - View status and output
  - Search and discover in-flight & processed data
- Rule Engine: Data "Brain"
  - Visually compose workflow automation
  - Code for automated workflows
  - Invoke Multi-step Data Processing Pipelines
  - Distributed computation
  - Tiered Storage for data at-rest and processing
  - Geographically distributed Backups





# Open-source

- Add to the sum of human knowledge and computational ability
- We can work faster and achieve more together
- No sense in reinventing the wheel, unless you're inventing a better wheel
- If an idea is worth creating, it's worth sharing





# But Why Golang?

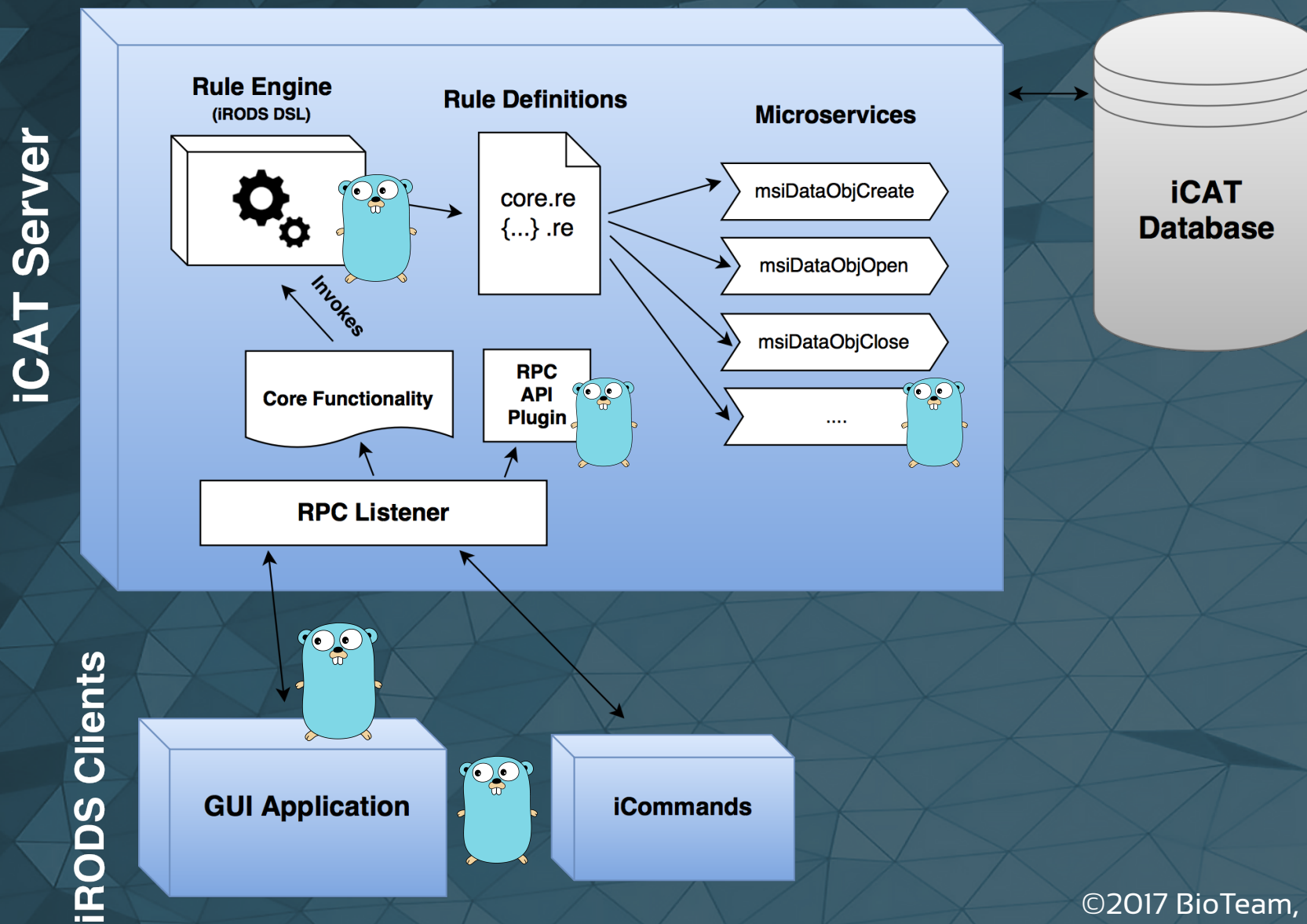
## A biased explanation of the Golang programming language

- Golang is John's favorite programming language
- It plays nicely with C and C++ (via cgo!)
- It's fast
- Simple design
- Encourages consistency (idiomatic go code), modularity (via packages), sturdy code (forced error handling)
- Standard library
- Standard toolchain – profiling, testing, etc...
- Community
- Concurrency model for dummies
- Deployment
- End goal: Commoditize iRODS development for the Golang community





# iRODS Integration Points: Visual



# iRODS Integration Points

## Common use cases

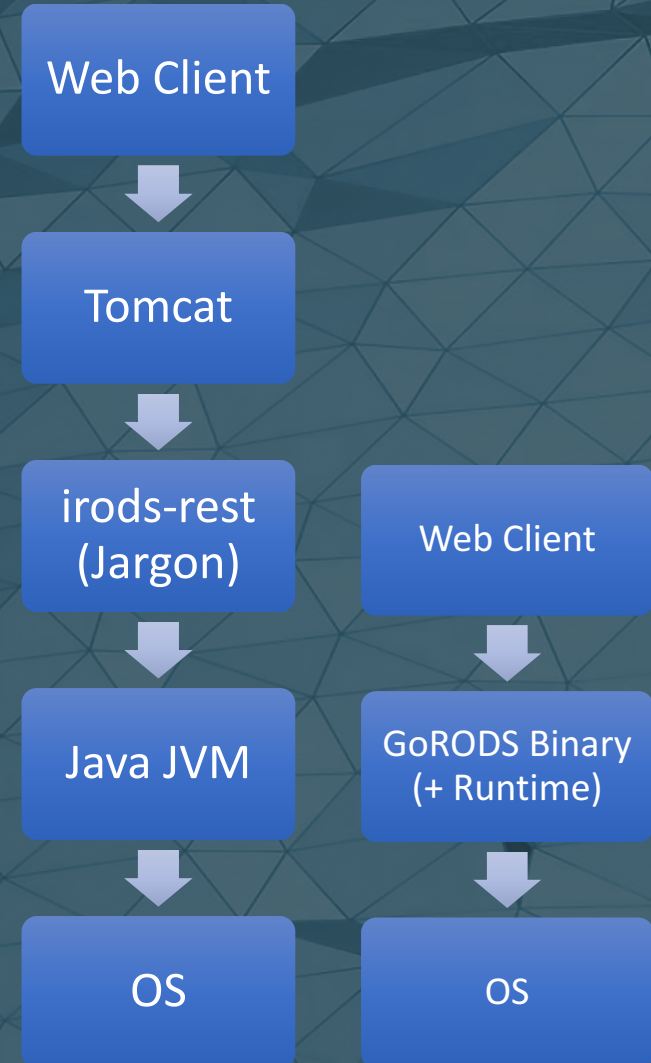
- iRODS Client
  - RESTful APIs
  - WebDAV
  - Web applications
  - Remote integration (decoupled distributed compute)
- Microservices & iRODS Rule Language DSL
  - Automated data processing
  - Automated business logic and data management
  - Replica based distributed computation
- RPC API
  - Deep integration
  - Side-loaded business logic
  - Full DB access  
(good for additional features for stateless clients)
- Rule Engine
  - Full integration into iRODS rule engine component, a step above microservices
  - DRY event notifications -> WebSocket server -> Web UI



# Integration Point: iRODS Client

GoRODS: <https://github.com/jjacquay712/GoRODS>

- iRODS C API Binding via cgo
- Supports latest 4.2.1
- Efficient options for memory management
- Working example for HTTP User Interface
  - Byte range support (concurrency for large downloads, "chunked downloads")
- Fully documented at [godoc.org](http://godoc.org)



# Integration Point: iRODS Client

GoRODS: <https://github.com/jjacquay712/GoRODS>

Creating a Connection  
to iRODS

```
1 package main
2
3 import (
4     "fmt"
5     "github.com/jjacquay712/GoRODS"
6     "log"
7 )
8
9 func main() {
10
11     client, conErr := gorods.New(gorods.ConnectionOptions{
12         Type: gorods.UserDefined,
13
14         Host: "localhost",
15         Port: 1247,
16         Zone: "tempZone",
17
18         Username: "rods",
19         Password: "password",
20     })
21
22     // Ensure the client initialized successfully and connected to the iCAT server
23     if conErr != nil {
24         log.Fatal(conErr)
25     }
26 }
```



# Integration Point: iRODS Client

GoRODS: <https://github.com/jjacquay712/GoRODS>

```
1 // Open a collection reference for /tempZone/home/rods
2 if openErr := client.OpenCollection(gorods.CollectionOptions{
3     Path: "/tempZone/home/rods",
4 }, func(col *gorods.Collection, con *gorods.Connection) {
5
6     // Output collection's string representation
7     fmt.Printf("String(): %v \n", col)
8
9     // Loop over the data objects in the collection, print the file name
10    col.ForEachDataObj(func(obj *gorods.DataObj) {
11        fmt.Printf("%v \n", obj.Name())
12    })
13
14    // Loop over the subcollections in the collection, print the name
15    col.ForEachCollection(func(subcol *gorods.Collection) {
16        fmt.Printf("%v \n", subcol.Name())
17    })
18
19 }); openErr != nil {
20     log.Fatal(openErr)
21 }
```

Looping over collection  
objects

Output:

Collection: /tempZone/home/rods  
C: pemtest  
C: source-code  
C: test  
d: hello.txt  
d: mydir1.tar

# Integration Point: iRODS Client

GoRODS: <https://github.com/jjacquay712/GoRODS>

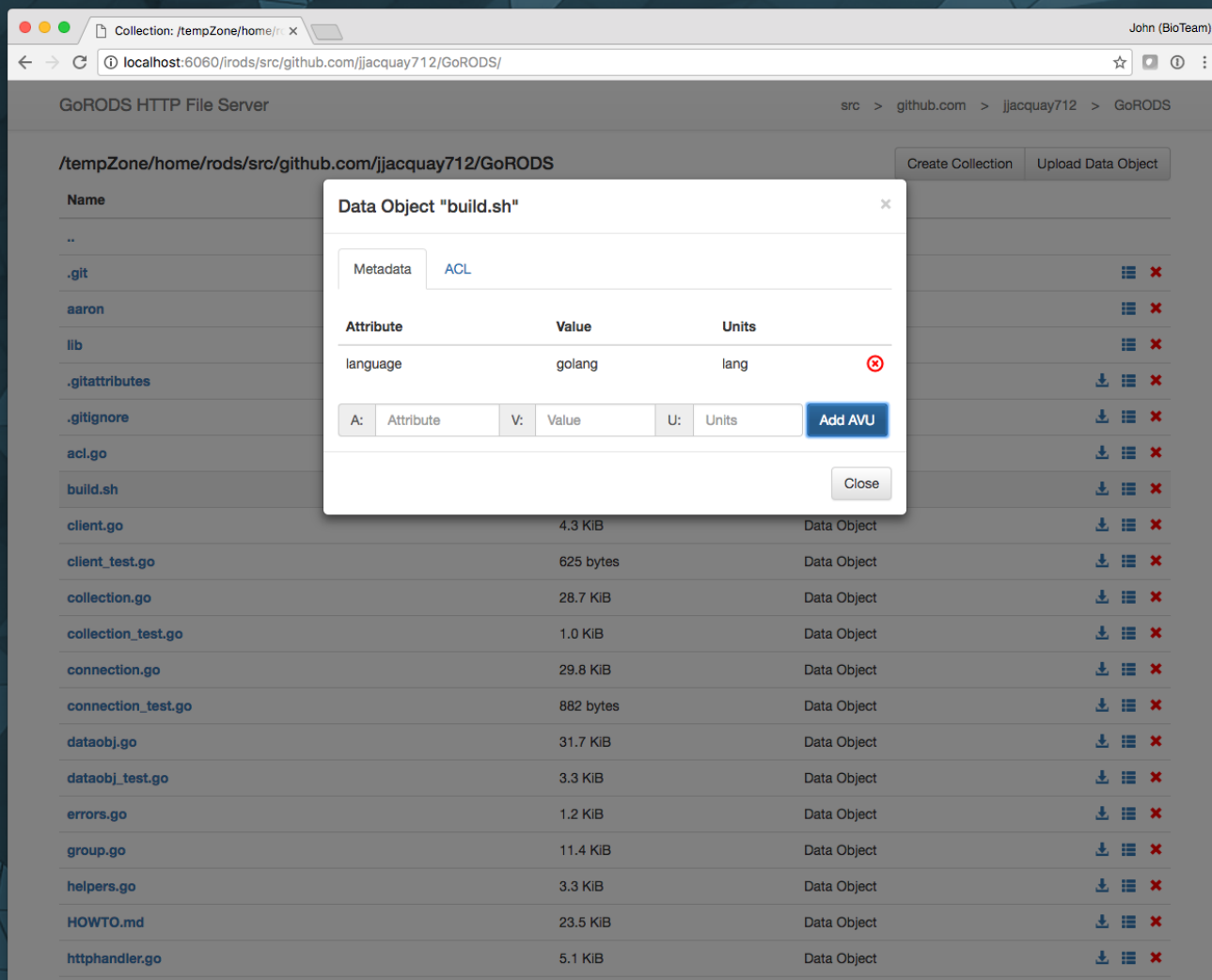
```
1 mountPath := "/irods/"
2
3 // Setup the GoRODS FileServer
4 fs := gorods.FileServer(gorods.FSOptions{
5     Path:  "/tempZone/home/rods",
6     Client: client,
7     Download: true,
8     StripPrefix: mountPath,
9 })
10
11 // Create the URL router
12 mux := http.NewServeMux()
13
14 // Serve the iRODS collection at /irods/
15 mux.Handle(mountPath, http.StripPrefix(mountPath, fs))
16
17 // Start HTTP server on port 8080
18 log.Fatal(http.ListenAndServe(":8080", mux))
```

Serving over HTTP



# Integration Point: iRODS Client

GoRODS: <https://github.com/jjacquay712/GoRODS>



GoRODS HTTP File Server

src > github.com > jjacquay712 > GoRODS

/tempZone/home/rods/src/github.com/jjacquay712/GoRODS

Create Collection Upload Data Object

Name

..

.git

aaron

lib

.gitattributes

.gitignore

acl.go

build.sh

client.go

client\_test.go

collection.go

collection\_test.go

connection.go

connection\_test.go

dataobj.go

dataobj\_test.go

errors.go

group.go

helpers.go

HOWTO.md

httphandler.go

Data Object "build.sh"

Metadata ACL

Attribute	Value	Units
language	golang	lang

A: Attribute V: Value U: Units Add AVU

Close

Serving over HTTP: Output

# Integration Point: iRODS Client

GoRODS: <https://github.com/jjacquay712/GoRODS>

```
1117 outBuff := make(chan *ByteArr, 100)
1118
1119 go func() {
1120     if readEr := obj.ReadChunkFree(10240000, func(chunk *ByteArr) {
1121         outBuff <- chunk
1122     }); readEr != nil {
1123         log.Print(readEr)
1124
1125         handler.response.WriteHeader(http.StatusInternalServerError)
1126         handler.response.Write([]byte("Error: " + readEr.Error()))
1127     }
1128
1129     close(outBuff)
1130 }()
1131
1132 for b := range outBuff {
1133     handler.response.Write(b.Contents)
1134     b.Free()
1135 }
1136
1137 }
1138
```

Reading Data / Concurrency

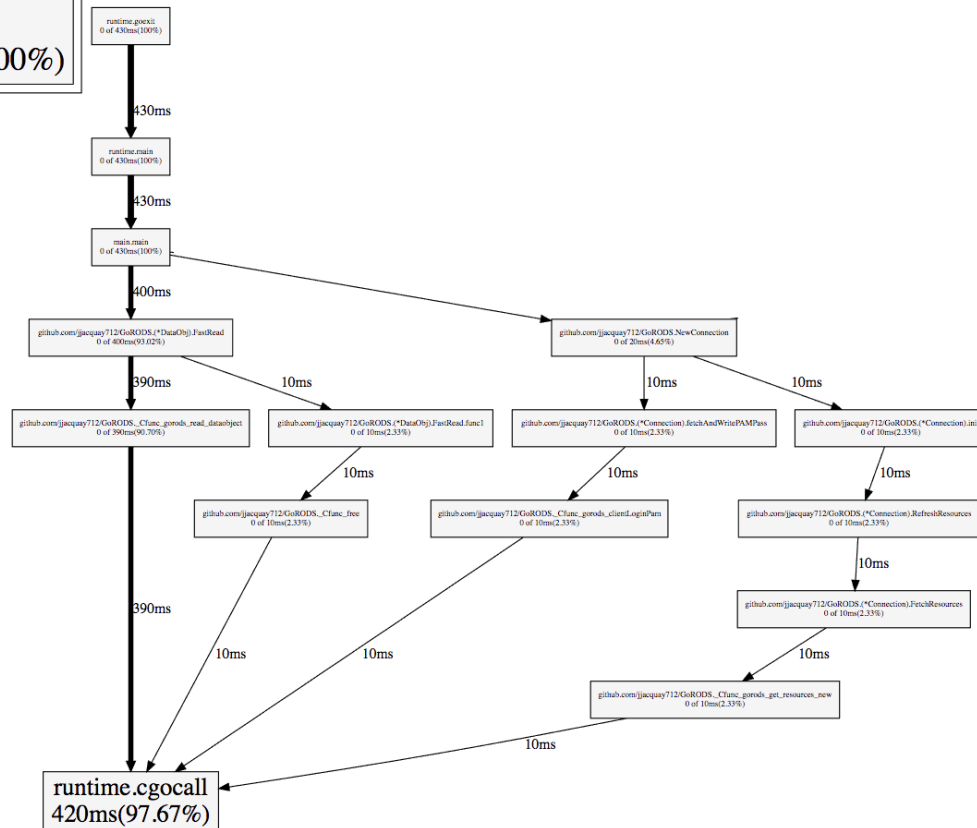


# Integration Point: iRODS Client

GoRODS: <https://github.com/jjacquay712/GoRODS>

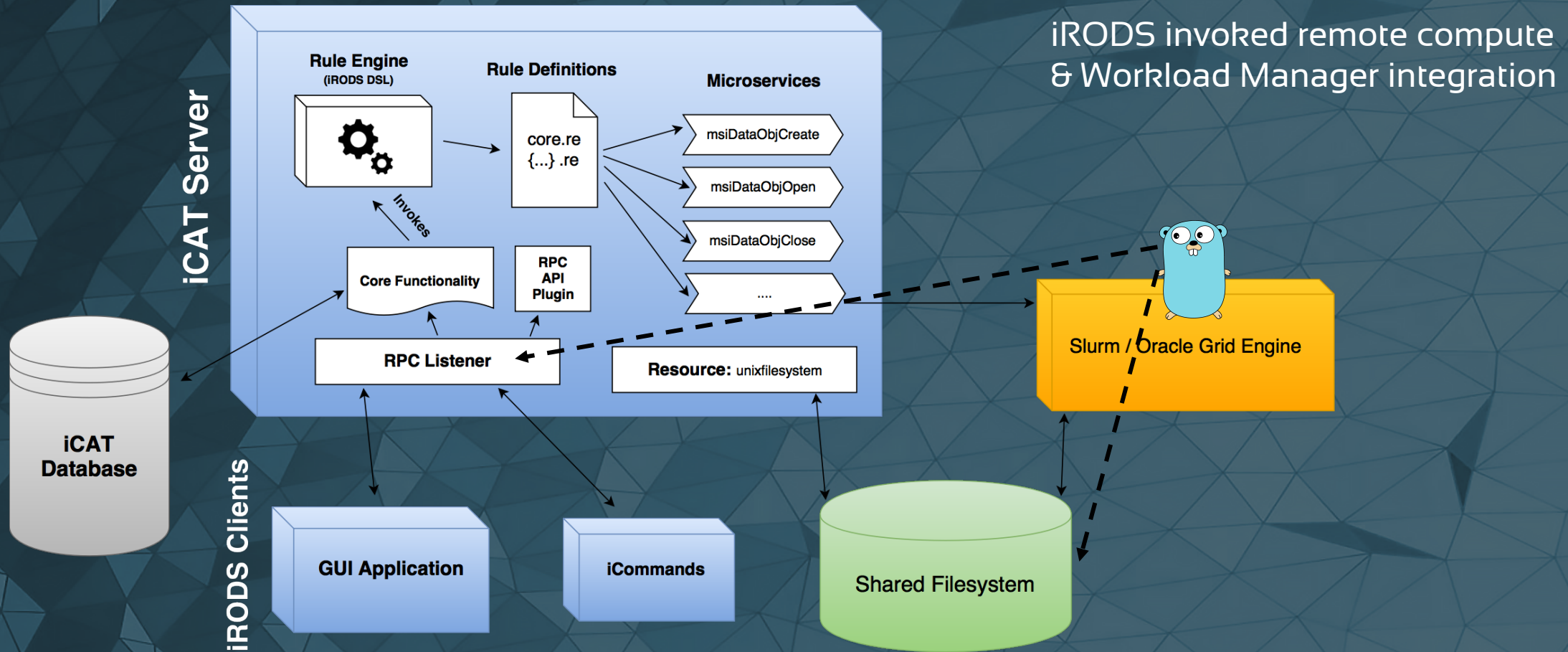
## Golang CPU Profiling

File: tester  
Type: cpu  
430ms of 430ms total ( 100%)



# Integration Point: iRODS Client

GoRODS: <https://github.com/jjacquay712/GoRODS>



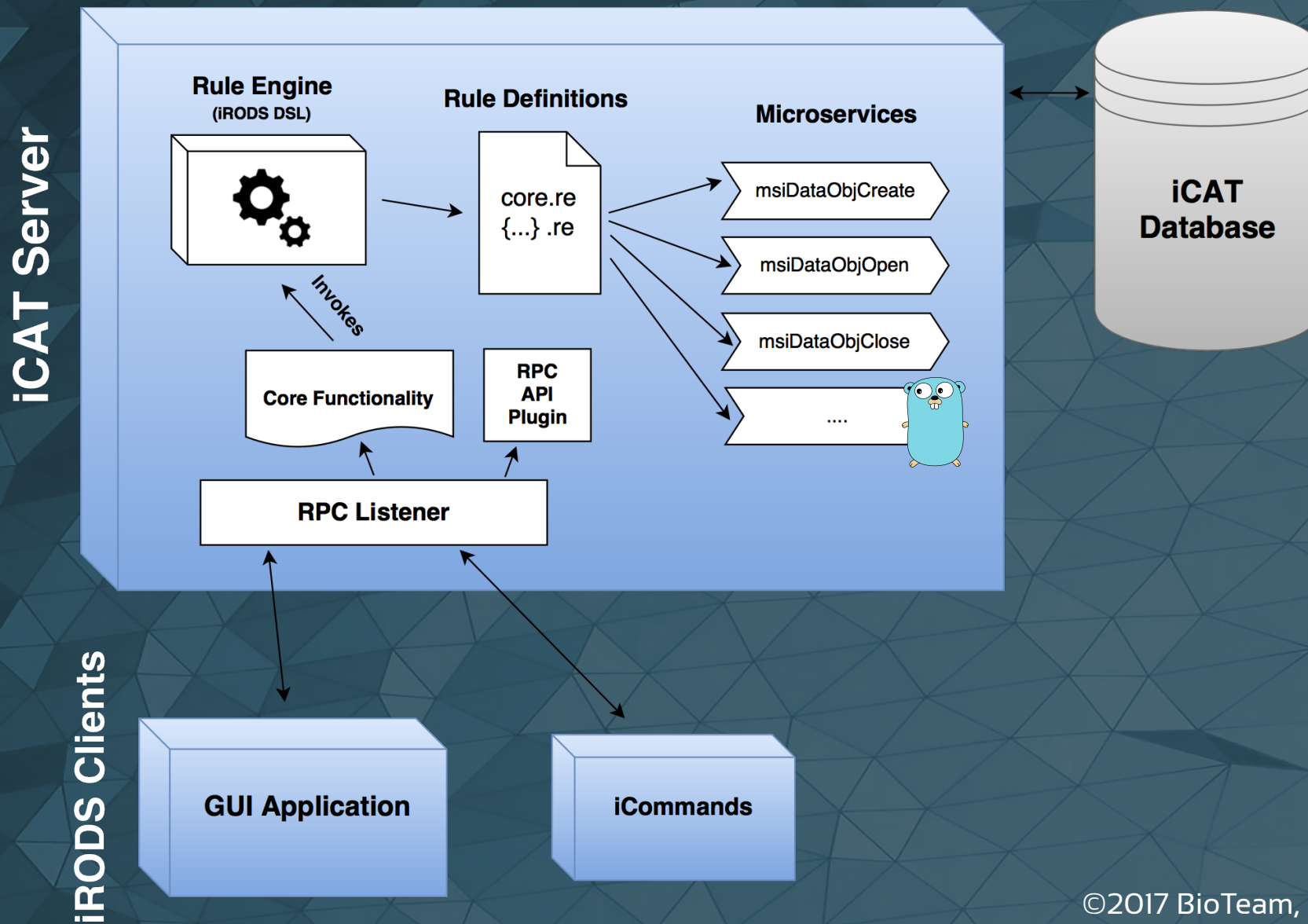


# iRODS Integration Points

## & Common use cases

- iRODS Client
  - RESTful APIs
  - WebDAV
  - Web applications
  - Remote integration (decoupled distributed compute)
- RPC API
  - Deep integration
  - Side-loaded business logic
  - Full DB access  
(good for additional features for stateless clients)
- Microservices & iRODS Rule Language DSL
  - Automated data processing
  - Automated business logic and data management
  - Replica based distributed computation
- Rule Engine
  - Full integration into iRODS rule engine component, a step above microservices
  - DRY event notifications -> WebSocket server -> Web UI

# iRODS Integration Points: Visual





## Introducing: GoRODS/msi

- Everything is on GitHub, ready to go. Go get it!
- Just spin up the Vagrant machine (2.5GB box file)
- Uses new GoRODS “msi” subpackage, written just for you!

```
82 // Associate metadata to data object
83 if err := msi.Call("msiAssociateKeyValuePairsToObj", labelsKVP, imagePath, "-d"); err != nil {
84     log.Print(err)
85     return msi.SYS_INTERNAL_ERR
86 }
87
```

- Simple example “msibasic\_example” returns CSV metadata as KVP struct
  - Testing example
- Advanced example implements Google Cloud Vision and Translation APIs to extract meaningful metadata, optionally compresses data objects



# Integration Point: Microservices & iRODS Rule Language DSL

<https://github.com/jjacquay712/irods-ugm-2017>

```
1 package main
2
3 import (
4     "encoding/csv"
5     "github.com/jjacquay712/GoRODS/msi"
6     "io"
7     "log"
8     "strings"
9     "unsafe"
10 )
11
12 // #cgo CFLAGS: -I/usr/include/irods
13 // #cgo LDFLAGS: -lirods_server -lirods_common -lpthread
14 /*
15 #include "msParam.h"
16 #include "re_structs.h"
17 */
18 import "C"
19
20 //export BasicExample
21 func BasicExample(inputParam *C.msParam_t, outputParam *C.msParam_t, rei *C.ruleExecInfo_t) int {
22     // Setup GoRODS/msi
23     msi.Configure(unsafe.Pointer(rei))
24
25     // Convert *C.msParam_t to go lang types
26     inputCSV := msi.ToParam(unsafe.Pointer(inputParam)).String()
27     outputKVP := msi.ToParam(unsafe.Pointer(outputParam)).ConvertTo(msi.KeyValPair_MS_T)
28
29     // Set output KVP
30     outputKVP.SetKVP(GetKVMap(inputCSV))
31
32     return msi.SUCCESS
33 }
34 }
```

## Basic Example

**Input:**

2 Column CSV

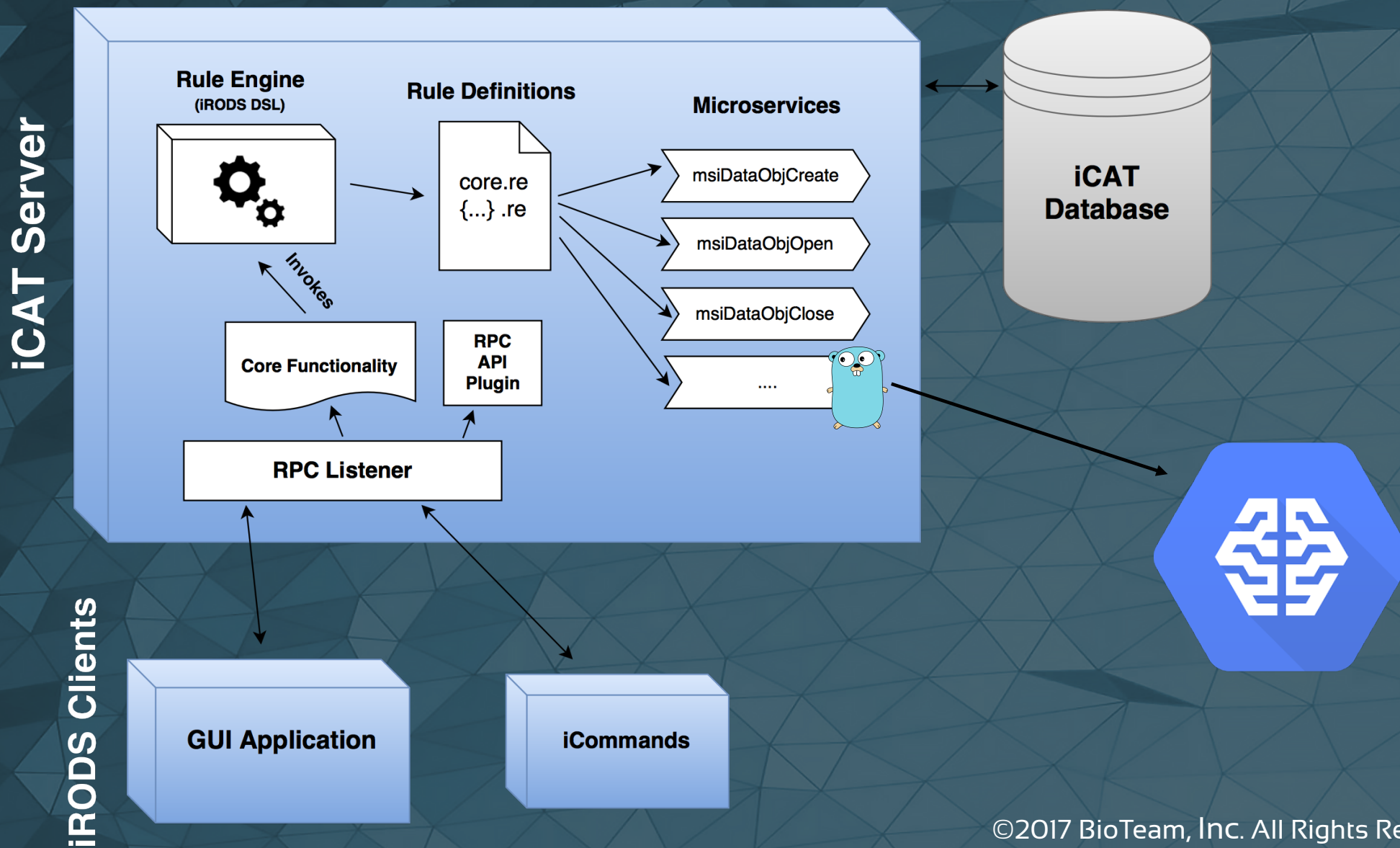
**Output:**

Key value pair structure (for  
metadata assignment)

```
1 TestBasicExample {
2     msibasic_example("keytest, valuetest", *outKVP);
3     msiPrintKeyValPair("stderr", *outKVP)
4 }
5
6 INPUT null
7 OUTPUT ruleExecOut
```



# iRODS Integration Points: Visual



# Go's Magical Interfaces

```

56     if enableGzip {
57         gzipDataObjPath := imagePath + ".gz"
58         gzipDesc := msi.NewParam(msi.INT_MS_T)
59
60         if err := msi.Call("msiDataObjCreate", gzipDataObjPath, "", gzipDesc); err != nil {
61             log.Print(err)
62             return msi.SYS_INTERNAL_ERR
63         }
64
65         gzipDataObj := msi.NewObjReaderFromDesc(gzipDesc)
66         defer gzipDataObj.Close()
67
68         gzWriter := gzip.NewWriter(gzipDataObj)
69         defer gzWriter.Close()
70
71         origDataObj, err := msi.NewObjReader(imageFilePath)
72         if err != nil {
73             log.Fatalf("Failed to read file: %v", err)
74         }
75         defer origDataObj.Close()
76
77         io.Copy(gzWriter, origDataObj)
78
79         imagePath = gzipDataObjPath
80     }

```

Assigns metadata from APIs to data object, and optionally compresses the object.



## GoRODS/msi

### Go's Magical Interfaces

```
196 file, err := msi.NewObjReader(filepath)
197 if err != nil {
198     log.Fatalf("Failed to read file: %v", err)
199 }
200 defer file.Close()
201
202 image, err := vision.NewImageFromReader(file)
203 if err != nil {
204     log.Fatalf("Failed to create image: %v", err)
205 }
206
207 labels, err := visionClient.DetectLabels(ctx, image, nil, 10)
208 if err != nil {
209     log.Fatalf("Failed to detect labels: %v", err)
210 }
```

### Advanced Example

#### Input:

iRODS Data Object Path

#### Output:

Assigns metadata from APIs to data object, and optionally compresses the object.

## GoRODS/msi

Memory Management, take #3. This time, it's for real

```
21 // NewParam creates a new *Param, with the provided type string
22 func NewParam(paramType ParamType) *Param {
23     p := new(Param)
24
25     p.rodsType = paramType
26
27     cTypeStr := C.CString(string(paramType))
28     defer C.free(unsafe.Pointer(cTypeStr))
29
30     p.ptr = C.NewParam(cTypeStr)
31
32     runtime.SetFinalizer(p, func(param *Param) {
33         C.FreeMsParam(param.ptr)
34     })
35
36     return p
37 }
```



# Integration Point: Microservices & iRODS Rule Language DSL

<https://github.com/jjacquay712/irods-ugm-2017>

msiextract\_image\_metadata: Let's run some data through!



```
[vagrant@irods-icat irods-ugm-2017]$ imeta ls -d gopher.jpg.gz
AVUs defined for dataObj gopher.jpg.gz:
attribute: tags_english
value: mammal,vertebrate,wildlife,squirrel,fauna,whiskers,prairie dog,marmot,rodent,prairie
units:
----
attribute: tags_dutch
value: zoogdier,gewerveld,dieren in het wild,eekhoorn,fauna,bakkebaarden,prairiehond,marmot,knaagdier,prairie
units:
```

```
attribute: exif_PhotometricInterpretation
value: 2
units:
----
attribute: exif_XResolution
value: "720000/10000"
units:
----
attribute: exif_YResolution
value: "720000/10000"
units:
----
attribute: exif_PixelYDimension
value: 630
units:
----
attribute: exif_ThumbJPEGInterchangeFormat
value: 450
units:
----
attribute: exif_ImageWidth
value: 840
units:
```



# Integration Point: Microservices & iRODS Rule Language DSL

<https://github.com/jjacquay712/irods-ugm-2017>



```
[vagrant@irods-icat irods-ugm-2017]$ imeta ls -d pretzel.gif.gz
AVUs defined for dataObj pretzel.gif.gz:
attribute: tags_dutch
value: hond,zoogdier,gewerveld,Hond zoals zoogdier,Pit bull,Amerikaanse pit bull terrier
units:
----
attribute: tags_english
value: dog,mammal,vertebrate,dog like mammal,pit bull,american pit bull terrier
units:
```





# Integration Point: Microservices & iRODS Rule Language DSL

<https://github.com/jjacquay712/irods-ugm-2017>



```
[vagrant@irods-icat irods-ugm-2017]$ imeta ls -d jason.jpg
AVUs defined for dataObj jason.jpg:
attribute: tags_english
value: long luscious hair, facial hair, face, nose, chin, beard, forehead, eyebrow, hairstyle, head, mouth, hair metal band
units:
----
attribute: tags_dutch
value: lang lekker haar, gezichtshaar, gezicht, neus, kin, baard, voorhoofd, wenkbrauw, kapsel, hoofd, mond, haar metalen band
units:
```

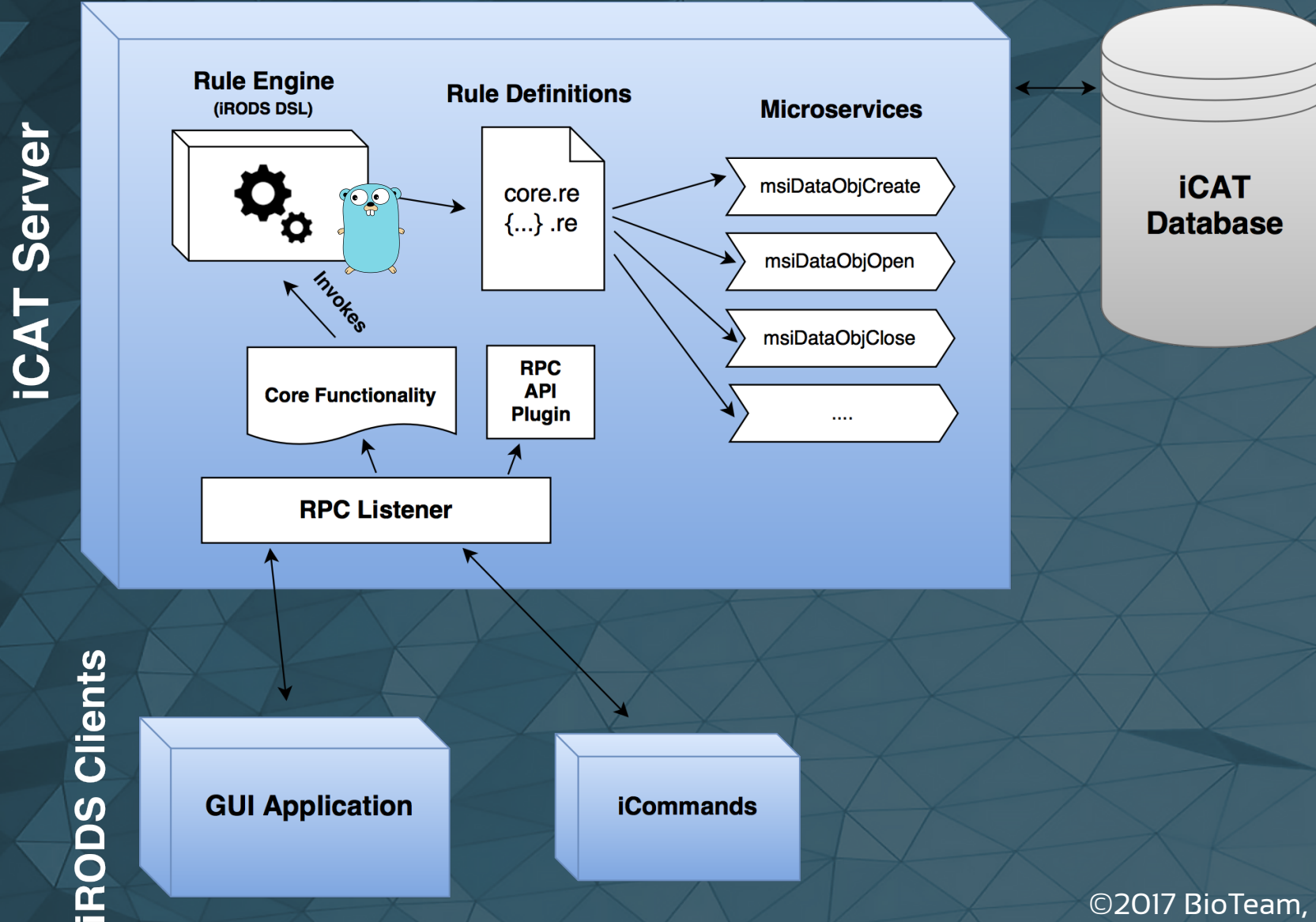
# iRODS Integration Points

## & Common use cases

- iRODS Client
  - RESTful APIs
  - WebDAV
  - Web applications
  - Remote integration (decoupled distributed compute)
- RPC API
  - Deep integration
  - Side-loaded business logic
  - Full DB access  
(good for additional features for stateless clients)
- Microservices & iRODS Rule Language DSL
  - Automated data processing
  - Automated business logic and data management
  - Replica based distributed computation
- Rule Engine
  - Full integration into iRODS rule engine component, a step above microservices
  - DRY event notifications -> AMQP -> WebSocket server -> Web UI

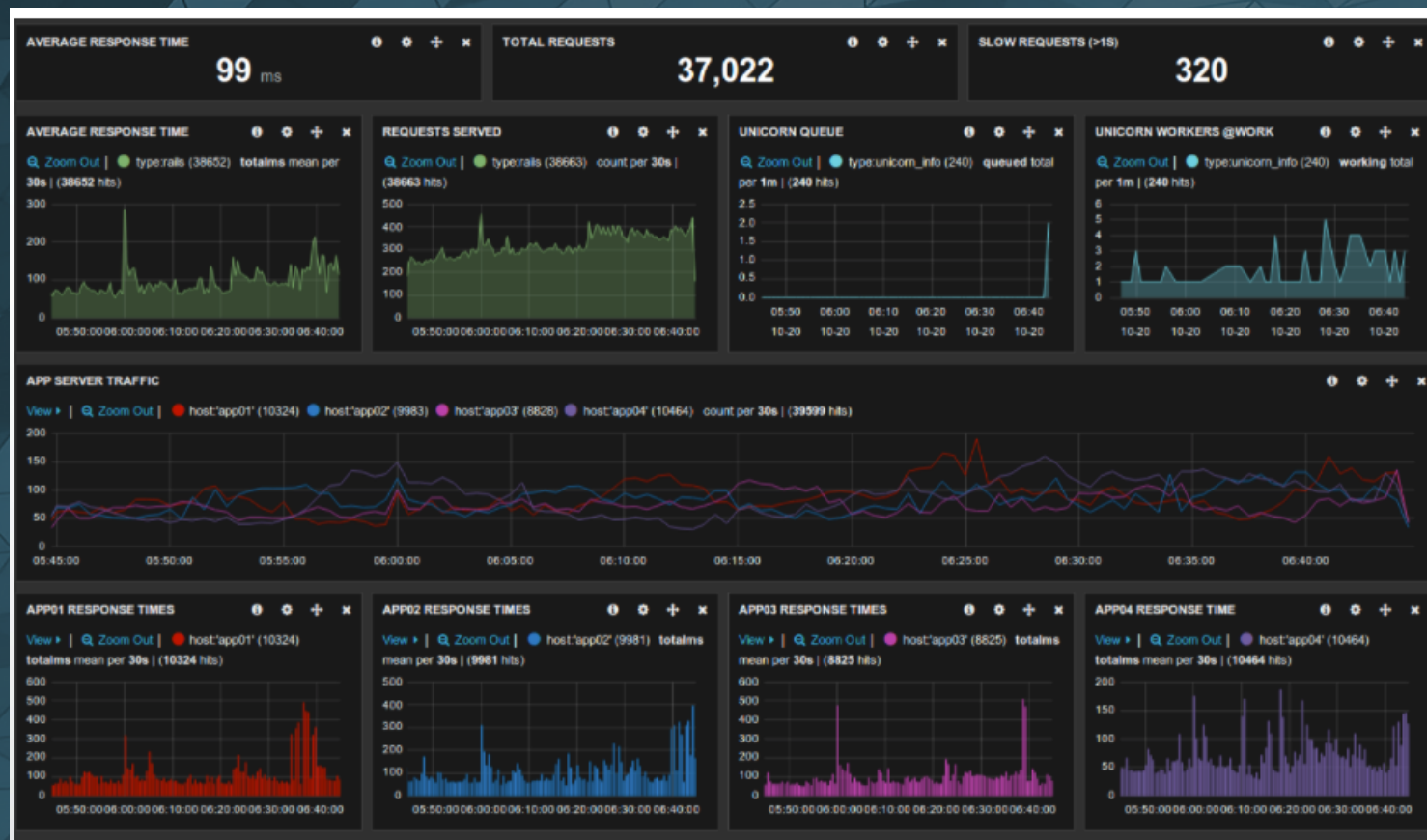


# iRODS Integration Points: Visual



# Integration Point: Rule Engine

## Kibana





```
1 package main
2
3 import (
4     "github.com/xiam/exif"
5     re "gore"
6     "log"
7     "path"
8     "strings"
9 )
10
11 func init() {
12
13     re.Configure(re.Opts{
14         Verbose: true,
15         LoadDefaultRules: true,
16     })
17
18     re.RegisterRules(re.Rules{
19         "acPostProcForPut": func(ps re.Params, cb *re.Callback) *re.Error {
20
21             sess := cb.SessionVars()
22
23             objPath := sess.Get("objPath").String()
24             physPath := sess.Get("filePath").String()
25
26             ext := path.Ext(objPath)
27
28             if ext == ".jpg" {
29                 if data, err := exif.Read(physPath); err == nil {
30
31                     metaStr := ExifToKVPStr(&data)
32
33                     meta := re.NewMsParam(re.KeyValPair_MS_T)
34
35                     if er := cb.Call("msiString2KeyValPair", metaStr, meta); !er.Ok() {
36                         return er
37                     }
38
39                     log.Print(meta.String())
40
41                     if er := cb.Call("msiAssociateKeyValuePairsToObj", meta, objPath, "-d"); !er.Ok() {
42                         return er
43                     }
44                 } else {
```

```
1 package gore
2
3 import (
4     "log"
5 )
6
7 var defaultRules Rules = Rules{
8     "printHello": func(ps Params, cb *Callback) *Error {
9         return cb.Call("print_hello")
10    },
11    "acPreConnect": func(ps Params, cb *Callback) *Error {
12
13        ps.Each(func(p RodsObj) {
14            if p.Type() == StdStringPtr {
15                strPtr := p.(*RodsStdStringPtr)
16
17                strPtr.Set("CS_NEG_DONT_CARE")
18            }
19        })
20
21        return Success()
22    },
23    "acCreateUser": func(ps Params, cb *Callback) *Error {
24
25        var err *Error = cb.Call("acPreProcForCreateUser")
26
27        err = cb.Call("msiCreateUser")
28        if !err.Ok() {
29            cb.Call("msiRollback")
30            return err
31        }
32
33        info := cb.SessionVars()
34
35        otherUserName, _ := info.Get("userNameProxy")
36
37        if otherUserName.String() != "anonymous" {
38            err = cb.Call("acCreateDefaultCollections")
39
40            err = cb.Call("msiAddUserToGroup", "public")
41            if !err.Ok() {
42                cb.Call("msiRollback")
43                return err
44            }
45        }
46
47        cb.Call("msiCommit")
48
49        return cb.Call("acPostProcForCreateUser")
50    }
51 }
```

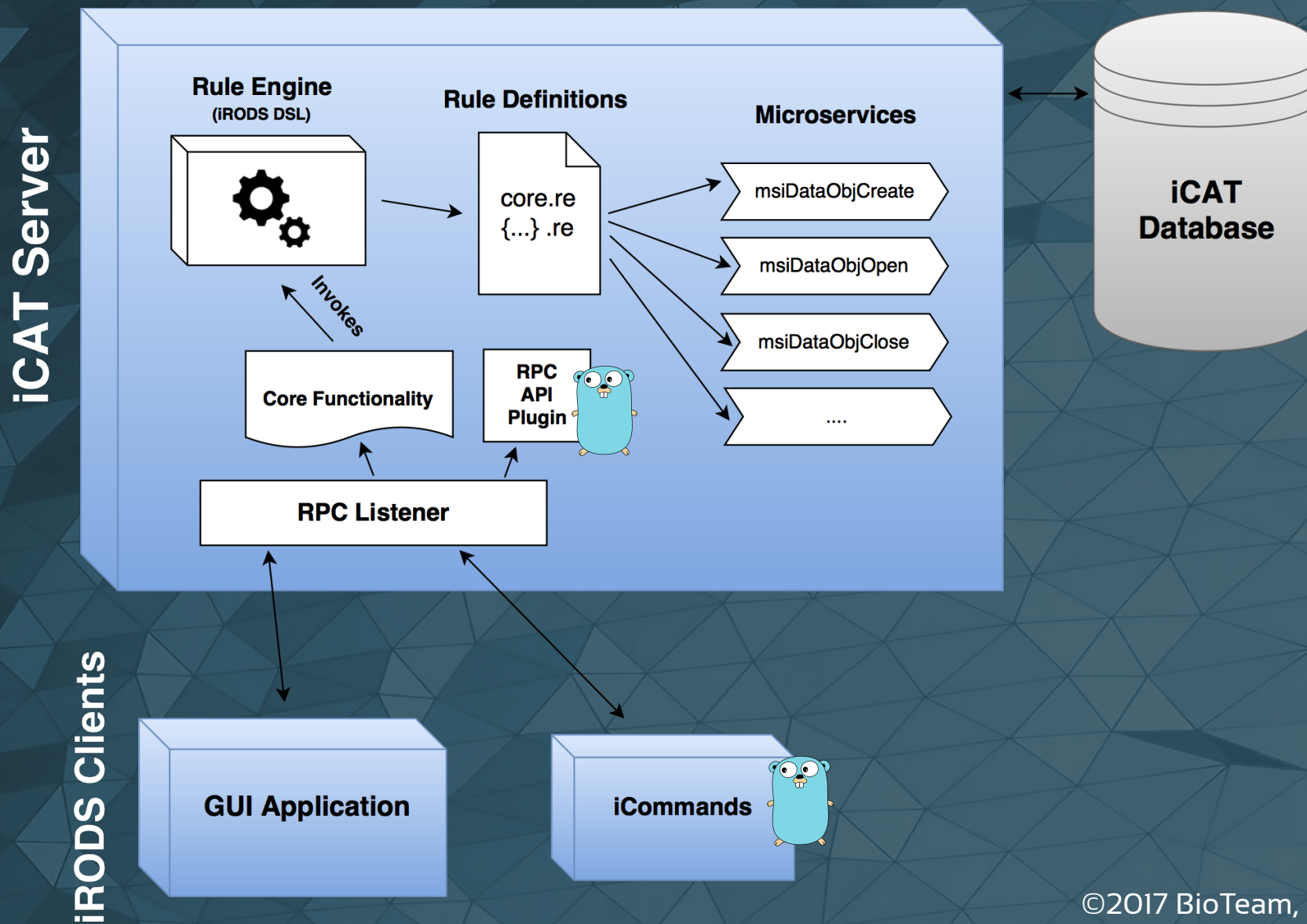


# iRODS Integration Points

## & Common use cases

- iRODS Client
  - RESTful APIs
  - WebDAV
  - Web applications
  - Remote integration (decoupled distributed compute)
- Microservices & iRODS Rule Language DSL
  - Automated data processing
  - Automated business logic and data management
  - Replica based distributed computation
- Rule Engine
  - Full integration into iRODS rule engine component, a step above microservices
  - DRY event notifications -> WebSocket server -> Web UI
- RPC API
  - Deep integration
  - Side-loaded business logic
  - Full DB access  
(good for additional features for stateless clients)

# iRODS Integration Points: Visual

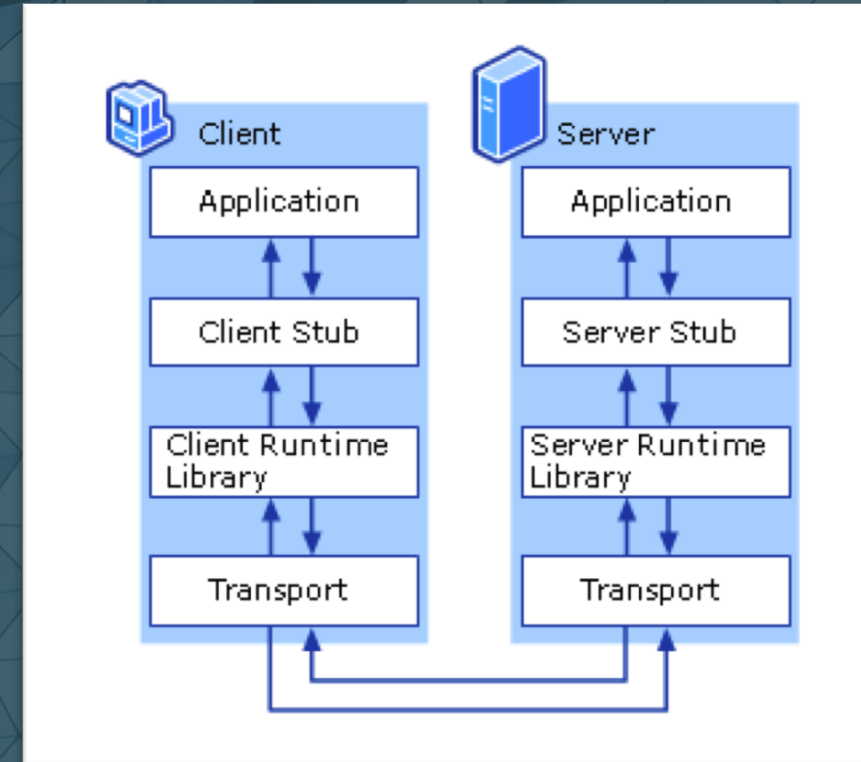




# Integration Point: RPC API

Still figuring this out myself

- Similar binding as rule engine integration
- Low-level binding to C, C++ data structures
  - Not just iRODS C API functions
- Client & Server component



# The End

GoRODS: <https://github.com/jjacquay712/GoRODS>

GoRODS/msi: <https://github.com/jjacquay712/GoRODS/msi>

## UGM 2017 Microservice Demo

<https://github.com/jjacquay712/irods-ugm-2017>

Special Thanks To:  
BioTeam, University of Florida, Product Development  
Team, Aaron Gardner, Alex Oumantsev, iRODS Consortium,  
My Dog Pretzel

