# G-SRS software for substance registration

**TYLER PERYEA**
NATIONAL CENTER FOR ADVANCING TRANSLATIONAL SCIENCES
NATIONAL INSTITUTES OF HEALTH
7 SEPTEMBER 2015

NCATS

NIH National Center for Advancing Translational Sciences

# Outline

- Why we are here
- How we got here
- Where we were
- Where we are now
- Where we are going

# Outline

- **Why we are here**
- How we got here
- Where we were
- Where we are now
- Where we are going

# Why we are here:
# NCATS Mission

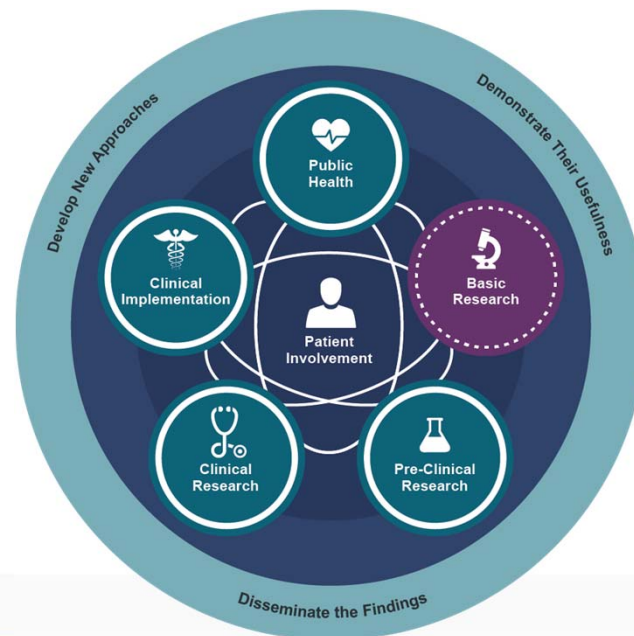- **N**ational
- **C**enter for
- **A**dvancing
- **T**ranslational
- **S**ciences

Translational science is a multidisciplinary form of science that bridges the recalcitrant gaps that sometimes exist between fundamental science and applied science, necessitating something in between **to translate knowledge into applications**.

# Why we are here:
# NCATS Mission

Translational science is a multidisciplinary form of science that bridges the recalcitrant gaps that sometimes exist between fundamental science and applied science, necessitating something in between **to translate knowledge into applications**.

## Three D's:

1. Demonstrate
2. Develop
3. Distribute

# Why we are here:
# NCATS Interest

- ➢ **NCATS need to define substances**
  - ▪ Informatics Discovery / Exploration
  - ▪ Drug Repurposing
  - ▪ Regulatory inquiries
- ➢ **Need for:**
  - ▪ Good identifiers
  - ▪ Curated data
  - ▪ Common data elements
  - ▪ Common data strategy

# Why we are here:
# Substance Problem

➢ **Fundamental Science**

- Groundwork set by International Standard (ISO IDMP)
- Great Subject Matter Experts with detailed knowledge of substances
- Supporting software libraries are available

➢ **Application Need**

- Global identifier
- Common strategy
- Common understanding
- Common data system

# Outline

- Why we are here
- **How we got here**
  - ➢ **Familiarization with the Standard**
  - ➢ Familiarization with the Data
  - ➢ Familiarization with the Procedures
- Where we were
- Where we are now
- Where we are going

# Familiarization with the Standard



Figure 1 –

Figure 6 — Information model for modifications

High-level specified substance information model

National Center
for Advancing
Translational Sciences

NCATS

# Familiarization with the Data

**FDA SRS XML**

**FDA SRS Chemical Structures**

**FDA SRS AUX Data**

**External Data**



```
single_substance  ..
  structure  ..
  element_type  ..
    protein  ..
      sequence_type   COMPLETE
      number_of_subunits   1
      subunit_group  ..
        subunit   1
        length   36
        sequence   APLEPVYPGDNATPEQMAQYAADLRRYINMLTRPRY
      disulfide_linkage
      glycosylation  ..
        glycosylation_type
        n_glycosylation
```

# Familiarize with procedures

- Ginas working group, mockups for how-to at user-experience level
    - ➢ Feedback on each mockup proposal for every substance class
    - ➢ Implementation to functional version
    - ➢ Feedback on functional version

# Familiarize with the procedures

Mockups

Live Front End

Seed data

# Putting it all together

- REST API
- Distributable Messages and Data dumps
- User's guide
- Conformance to standard
- User Acceptance Testing
- Rollout

# Where we were

- **February 2015**
  - ➢ Version 0.9
    - ▪ feature complete
      - ○ Registration for all substance classes, up to group 1 specified substance
      - ○ Alternative definitions
      - ○ Approval process
      - ○ Basic search available for textual, structural, sequence-based data
    - ▪ Self-contained and distributable
    - ▪ Programmatically Accessible
  - ➢ Launched on Health Canada site



### Index of /pub/ginasISO

| Name | Last modified | Size | Description |
|---|---|---|---|
| Parent Directory | | - | |
| legacy/ | 05-Oct-2014 00:48 | - | |
| v0.9501/ | 05-Oct-2014 00:50 | - | |
| v0.9511/ | 13-Jan-2015 20:35 | - | |
| v0.9512/ | 29-Jan-2015 11:30 | - | |
| v0.9513/ | 01-Feb-2015 18:05 | - | |
| v0.9516/ | 16-Apr-2015 19:02 | - | |
| v0.9517/ | 08-May-2015 11:36 | - | |

# Where we were
## What we learned

### The Good

➢ Lots of features

➢ Lots of details

➢ Lots of discussion

➢ Physical instantiation of data

➢ Utility functions for registration

### The Bad

➢ Lots of features

➢ Lots of details

➢ Lots of discussion

➢ Complex installation

➢ Difficult to customize, beyond substances

➢ Lack of mobile support

➢ Persistence layer not as flexible / accessible

➢ **Browsing / discoverability lacking**

# Where we were
## What we learned



FDA SRS Data

External Data

Transformation

Standardization

Validation

Seed Data

Data Model
JSON Schema

Backend Software

REST API

Persistence

Utility Functions

Frontend Software

Procedures

Registration Forms

Layout

NCATS

# Where we were
## What we learned

# Where we were
## What we learned

Data Model JSON Schema

- **JSON Schema Overview**
- JSON Schema versioning
- Schema Review

# Where we were
## What we learned

Data Model JSON Schema

- JSON Schema Overview
- **JSON Schema versioning**
- Schema Review

| | | | |
|---|---|---|---|
| 📄 ginasSchema106.json | 3 months ago | | 👤 Tyle |
| 📄 ginasSchema107.json | 3 months ago | | 👤 Tyle |
| 📄 ginasSchema108.json | 3 months ago | | 👤 Tyle |
| 📄 ginasSchema109.json | about a month ago | | 👤 Tyle |
| 📄 ginasSchema110.json | 18 days ago | | 👤 Tyle |

# Where we were
## What we learned



Data Model JSON Schema

- JSON Schema Overview
- JSON Schema versioning
- **Schema Review**



National Center for Advancing Translational Sciences

NCATS

# Where we were
## What we learned

FDA SRS Data

External Data

Transformation

Standardization

Validation

Data Model JSON Schema

Seed Data

Backend Software

REST API

Persistence

Utility Functions

Frontend Software

Procedures

Registration Forms

Layout

National Center for Advancing Translational Sciences

NIH

NCATS

# Where we were
## What we learned



FDA SRS Data

External Data

**New Components**

Transformation

Standardization

Validation

Seed Data

**Data Model
JSON Schema**

Backend Software

REST API

Persistence

Utility Functions

Procedures

Layout

Registration Forms

NIH National Center for Advancing Translational Sciences

NCATS

# Outline

- Why we are here
- How we got here
- Where we were
- **Where we are now**
- Where we are going

# What we are now

- Frontend / Backend rewrite
  - ➢ Discoverability first
  - ➢ Mobility first
  - ➢ Focus on browsing data

# Where we are now

Browse all substances in seed data

# Where we are now

Quick filtering to entities of interest *("racemic INN-named chemical substances with EVMPD codes")*

# Where we are now

Full substance views for substance class

# Where we are now

Advanced searching

# Where we are now

Registration and Validation

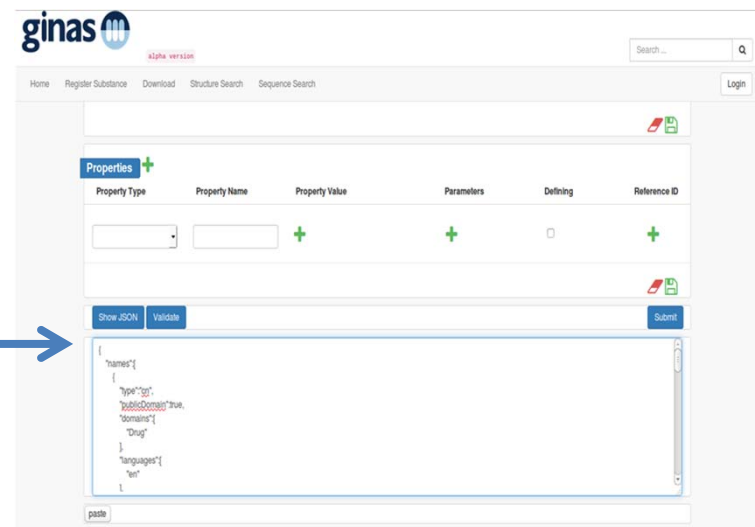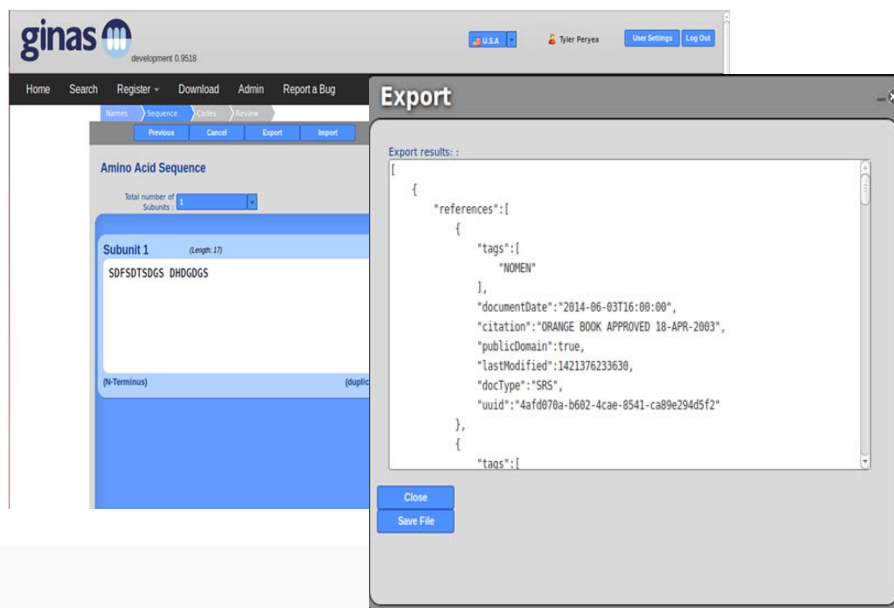# Where we are now

Batch import of records

# Outline

- Why we are here
- How we got here
- Where we were
- Where we are now
- **Where we are going**

# Where we are going

- **Complete migration of Registration forms to new pages**
  - Nucleic Acids
  - Polymers
  - Proteins

# Where we are going

- **Complete migration of Registration forms to new pages**
  - ➢ Nucleic Acids
  - ➢ Polymers
  - ➢ Proteins
- **Complete migration of utility functions**
  - ➢ Name-to-definition resolver
  - ➢ Name analyzers (INN protein format, etc)

# Where we are going

- **Complete migration of Registration forms to new pages**
  - Nucleic Acids
  - Polymers
  - Proteins
- **Complete migration of utility functions**
  - Name-to-definition resolver
  - Name analyzers (INN protein format, etc)
- **Simple Export to variety of formats**
  - **JSON**
  - **SDF**
  - **Excel**

# Where we are going

- **Complete migration of Registration forms to new pages**
  - Nucleic Acids
  - Polymers
  - Proteins
- **Complete migration of utility functions**
  - Name-to-definition resolver
  - Name analyzers (INN protein format, etc)
- **Simple Export to variety of formats**
  - **JSON**
  - **SDF**
  - **Excel**
- **Integrate and validate other data and entities into seed set**
  - Metabolism
  - Monograph information
  - Seed examples for products

# Where we are going

- Demonstration
  - ➤ Later today
  - ➤ Tomorrow during lunch
- Software Details / Architecture
  - ➤ Tim Sheils will give more information tomorrow

# Acknowledgements

## U.S. Food and Drug Administration

| | |
|---|---|
| Yulia Borodina | Archana Newatia |
| Larry Callahan | Vada Perkins |
| Ta-Jen Chen | Mary-Ann Slack |
| Ramez Ghazzaoui | Frank Switzer |
| Elaine Johansen | Alex Welsch |

## U.S. Pharmacopeial Convention

| | |
|---|---|
| Fouad Atouf | Andrej Wilk |
| Tina Morris | |

## Federal Institute for Drugs and Medical Devices (Germany)

Thomas Balzer

## SwissMedic

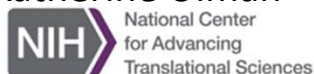Philipp Weyermann

## European Directorate Quality of Medicines

Christopher Jarvis

## Dow Corning

Katherine Ulman

## NIH/NCATS

| | |
|---|---|
| Chris Austin | Tyler Peryea |
| Ajit Jadhav | Tim Sheils |
| Dac-Trung Nguyen | Tongan Zhao |

## Medicines Evaluation Board (Netherlands)

| | |
|---|---|
| Herman Diederik | Burt Kroes |
| Marcel Hoefnagel | Ciska Matai |
| Joris Kampmeijer | |

## Health Canada

Vikesh Srivastava

## Royal Botanic Gardens, Kew (UK)

| | |
|---|---|
| Bob Allkin | Elizabeth Dauncey |

## European Medicines Agency

| | |
|---|---|
| Paolo Alcini | Telonis Pangiotis |
| Sabine Brosch | Ilaria Del Seppia |

## Uppsala Montintoring Centre / WHO

| | |
|---|---|
| Malin Jakobsson | Malin Fladvad |
| Martin Strömberg | |