# NCI SEER Edits Engine

## Interoperable Approach to Data Validation

**Fabian Depry**
**Information Management Services**
**NAACCR Conference, June 2011**

# Today's Presentation

- Introduction to the SEER Edits Engine
  - What is it?
  - Why do we need it?
  - How was it implemented?
  - How is it used?
  - What are our plans for the future?

# What is the SEER Edits Engine?

- Java framework to validate incidence data
  - Used by Java applications that cannot use SEER and other edit sets via GenEdits
  - A library of Java source code that is used by Java programs maintained by the SEER program (SEER*DMS, SEER*Abs, SEER*Edits)

- Executes edits against data in any format
  - Text files, e.g., NAACCR Abstract files
  - Values entered on data entry forms
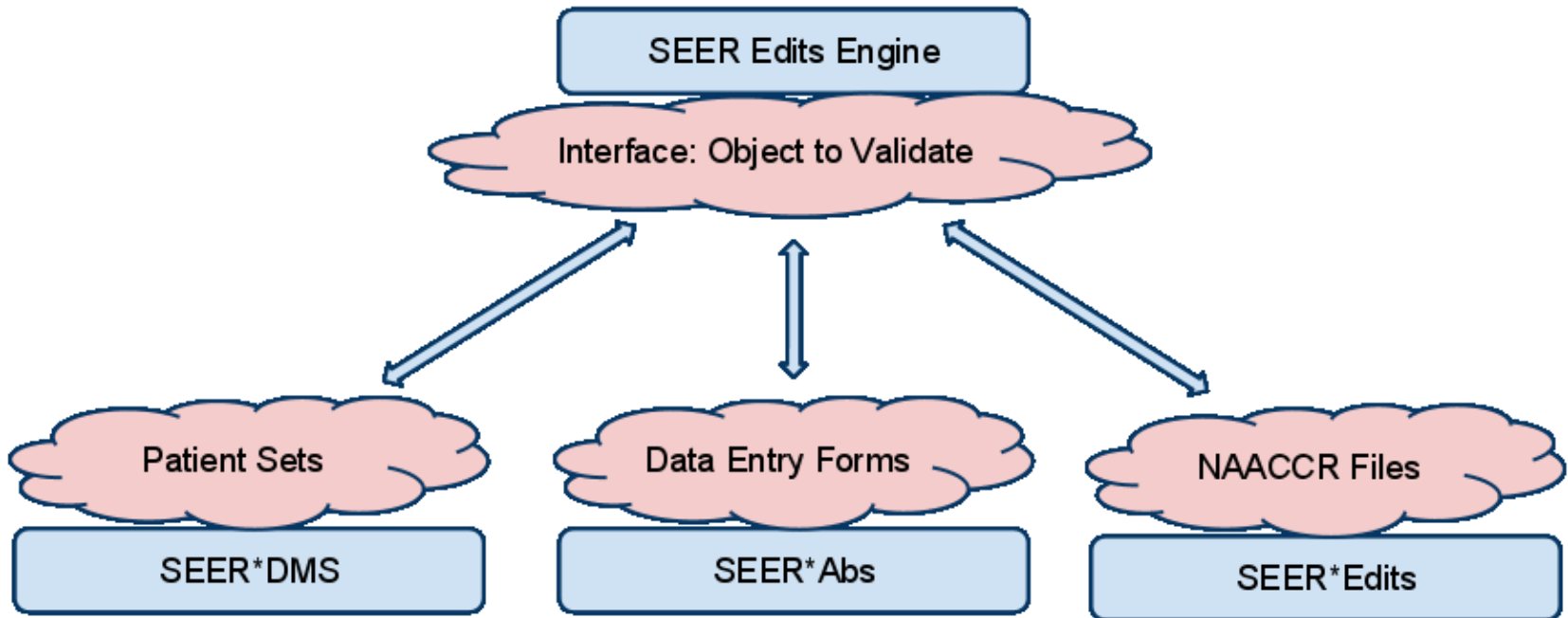  - Values stored in proprietary database structures

# Why build a new engine?

- ## SEER*DMS has specific needs:
  - Validate data stored in a relational database
  - Allow registry-specific edits to be maintained through the application
  - Support a mechanism to easily test the edits logic

- ## SEER*Abs and SEER*Edits:
  - Validate data in data entry forms (SEER*Abs)
  - Validate data in text files (SEER*Edits)
  - Import edits written in other software (both)

National Cancer Institute

IMS

# System Requirements

- **For validating cancer data:**
  - Reliability
  - Speed

- **For maintaining edits:**
  - Simple but powerful syntax
  - Testing framework to verify the edits logic
  - Graphical tools to assist writing and updating edits

- **For maintaining the edits engine:**
  - Operability: running edits on different data types

# Supporting Different Data Types

# Edits Syntax - Groovy

- Groovy is a scripting language based on Java:
  - Anything that can be done in Java can be done in Groovy
  - Groovy also has unique syntax to allow for small, elegant scripts
  - Groovy scripts can also use readily available Java libraries

- Edits are based on Boolean logic.
  - If an edit returns "true", it passes. Examples:
    - *return line.primarySite ==~ /C\d\d\d/*
    - *return Hist_ICD_O_3_Table.contains(line.histologyIcdO3)*

# Edits storage  - XML

- The engine provides an API to read and write eXtensible Markup Language (XML) files:

```
<rule id="primary-site" name="Primary Site" ruleset="field">
    <expression>return line.primarySite ==~ /C\d\d\d/</expression>
    <message>Primary site is not valid.</message>
</rule>
```

- The edits can actually be stored in any format
  - SEER*DMS persists edit source code in a database table
  - SEER*Abs and SEER*Edits – edits are stored in XML files

# Engine Speed

- The key to super-speed: multi-threading
  - The more resources available, the faster the edits will run
- In SEER*DMS:
  - 1,294 edits ran on data for 269,727 patients using a Linux server
  - 58 minutes → 78 patients/second (may include multiple tumors)
- In SEER*Edits:
  - 582 SEER edits ran on 9,422,096 records using a 64-bit Windows desktop with two dual-core processors
  - 3 hours 37 minutes → 723 records/second
- In SEER*Abs:
  - Edits are executed on one record at a time, not in batches
  - Validates data entry form each time the user exits a field

IMS

# Use case – SEER*DMS

## SEER*DMS Editor:



## Edits Tab of the Editor:

# Use case – SEER*Abs

Data Entry Form:



Edit Documentation:

# Use case – SEER*Edits (Session)

# Use case – SEER*Edits (Results)

# Use case – Edit Writer

**ID:** Primary_Site    **Name:** Primary Site

**Group:** SEER    **Category:** field    **Severity:** 6

**Msg:** Primary site is not valid.

**Code** | Documentation | History | Dependencies

```
if (line.primarySite == null)
    return true
return line.primarySite ==~ /^C\d\d$/
```

Tests
└ Primary_Site
   ├ Line 4
   ├ Line 8
   └ Line 12

Expected pass, got fail.  Line: 8.

Values:
primarySite = C447

```
// no site
line = [:]
line['primarySite'] = null
Testing.assertPass(line)

// a valid site
line['primarySite'] = 'C447'
Testing.assertPass(line)

// an invalid site
line['primarySite'] = 'xxxx'
Testing.assertFail(line)
```

# Using the GenEdits Metafile

- Many SEER Registries also use NPCR or NCDB edit sets
  - Goal:  use these edit sets in SEER*DMS, SEER*Abs,  SEER*Edits

- GenEdits
  - Uses a simple and well-defined language
  - In theory, it is possible to write a compiler to translate GenEdits source code into Groovy code that could be made available to the Java programs via the SEER Edits Engine

- Edits Compiler
  - Work in progress
  - Nearly all edits have been translated, but not all
  - Updates made only in the GenEdits version;  the compiler would be used to create the updated edit sets for the Java programs

# Translating a Simple Edit: Date of Birth

```
If ( EMPTY(#S"Date of Birth"))
    return PASS;

if (VALID_DATE_IOP(#S"Date of Birth"))
    return PASS;
Else
    {
    error_text ("Date of Birth: %DC");
    return FAIL;
    }
```

```
Functions.GEN_RESET_LOCAL_CONTEXT(binding);

if (Functions.GEN_EMPTY(line.birthDate))
    return true

if (Functions.GEN_VALID_DATE_IOP(binding, line.birthDate))
    return true
else {
    Functions.GEN_ERROR_TEXT(binding, "Date of Birth: %DC")
    return false
}

return true
```

# Loading Translated Edits in SEER*Edits

# Next Steps:  SEER*Utils

- Utility programs available in a single Java library
- SEER*Utils currently contains
  - SEER Edits Engine
  - Java bridge to the Collaborative Stage DLL
  - SEER Site Recode mappings
  - SEER*Rx drugs and regimens data
  - Hematopoietic and Lymphoid Database
  - Multiple Primaries Calculator
- NAACCR mappings with an API to facilitate reading and writing of NAACCR data files

# SEER*Utils

- SEER*Utils can be very easily integrated with a Groovy script:

```
import com.imsweb.seerutils.*
import com.imsweb.seerutils.cstage.*

SeerUtils.initializeAll()
println 'CStage DLL version is ' + CollaborativeStage.getVersion()
SeerUtils.uninitializeAll()
```