

GSRS Development Update

15 May 2019

TYLER PERYEA
DANNY KATZEL
NIH/NCATS

NCATS

Agenda

High Level GSRS Core Status

- 2.3.3 Release
- Ongoing Software Updates
- Data Export Update
- New Image to Structure Feature

GSRS Release Status

- GSRS FDA Production Releases
 - 2.3.1 - October 2018 - previous public release
 - 2.3.2 - December 2018
 - 2.3.3 - February 2019 - most recent public release including public data set
 - 2.3.4 - March 2019
 - 2.3.5 - April 2019
- Upcoming GSRS FDA Releases
 - 2.3.6 - scheduled for June 2019
 - 2.4 development - 3rd quarter 2019
 - Overall goals: Improved data exchange / communication

GSRS 2.3.3 Release Update

Updates available in the 2.3.3 public release

- Structure rehash / reprocessing task to regenerate hashes when structure hash algorithm changes
- Allow scheduled task for remaking backups
- Minor UI Improvements
- Improved speed of facet filtering, sequence searching
- Integrated Image to Structure in registration pages
- Bug fixes, additional validation rules

Ongoing Software Updates

Completed Updates to be Included in Next Public Release

- Initial Beta UI Redesign - View Forms, Global Search (basic), Structure Search, Sequence Search
- Improved Name Display - added Additional Listing name to view of names, added name details to Names card
- Add a button on browse for privileged users to copy the structure/sequence directly to a registration page without bringing the references

Upcoming Updates

- Beta UI Improvements
- Migrating Bookmarklets / Powertools to UI
- Improved Molvec

New Dataset Released

Latest News

April 10, 2019

Newest GSRS Public Data Released

The most recent set of public data, comprised of 105,020 records, has been compiled and is available for download.

[Download](#)

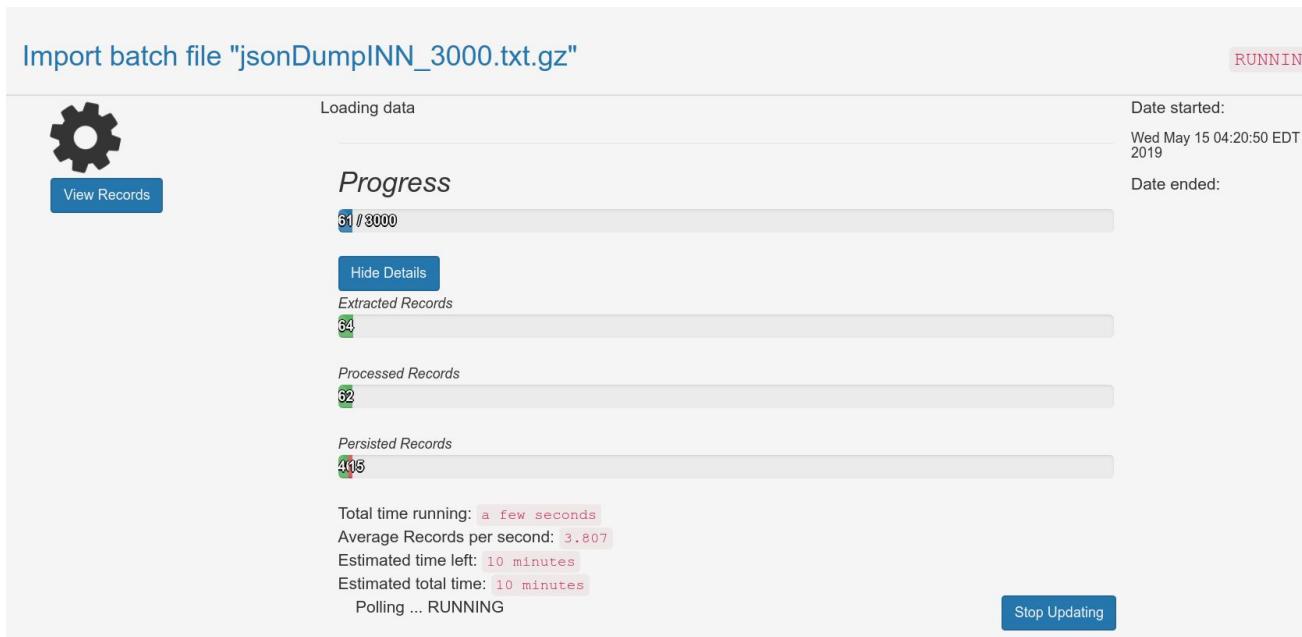
[View All News](#)

Data Type	Count
Chemical	71,098
Structurally Diverse	25,038
Protein	3,547
Mixture	2,479
Polymer	2,080
Nucleic Acid	172
Group 1 Specified Substance	15
Names	735,102
Codes	426,409
Relationships	55,323

New Dataset Released

A few words about the Data Dump format

- Native file format for G-SRS (.gsrs)
 - gzipped, tab-delimited file with 1 G-SRS substance JSON per row
 - loadable directly by the G-SRS application
 - convertible (lossy) to several other formats (txt, sdf)



New Dataset Released

- Browsable on main website,
or via API

Molecular Weight	
Search Molecular Weight...	
<input checked="" type="checkbox"/> 400:600	14,595
<input type="checkbox"/> 0:200	16,616
<input type="checkbox"/> 200:400	33,727
<input type="checkbox"/> 600:800	3,959
<input type="checkbox"/> 800:1000	1,469
<input type="checkbox"/> Exclude Selected	
More ...	Clear
Search	
Substance Type	
<input type="checkbox"/> Chemical	14,352
<input type="checkbox"/> Polymer	242
<input type="checkbox"/> Protein	1
Search	
Source Tag	
Search GInAS Tag...	
<input type="checkbox"/> MI	2,899
<input type="checkbox"/> WHO-DD	2,883
<input type="checkbox"/> INN	2,309
<input type="checkbox"/> USAN	1,668
<input type="checkbox"/> MART.	1,249
More ...	
ATC Level 1	
ATC Level 2	
ATC Level 3	
ATC Level 4	
Code System	
Molecular Weight 400:600	
14,595	< 1 2 3 4 5 6 7 8 ... 912 913 > »
Sort By:	
Sort By:	
Grid View	
Switch	
CINCHONIDINE SULFATE PENTAHYDRATE	
UNII:P424MD870D	
ABSOLUTE	
Names: CINCHONIDINE, SULFATE (1:1), PENTAHYDRATE ✓ CINCHONIDINE SULFATE PENTAHYDRATE ✓ CINCHONAN-9-OL, (8(ALPHA),9R)-, SULFATE (...)	
Codes: CAS: 5907-43-7 🔗 PUBCHEM: 71586828 🔗	
Relationships: 1	
Formula: C19H22N2O ₅ H ₂ O ₄ S	
Mol Weight: 482.55	
Search	
LOBEGLITAZONE SULFATE	
UNII:95C712E83P	
RACEMIC	
Names: 2,4-THIAZOLIDINEDIONE, 5-((4-(2-((6-(4-METH... LOBEGLITAZONE SULFATE ✓ LOBEGLITAZONE SULFATE [WHO-DD] CKD-501 (+/-)-2,4-THIAZOLIDINEDIONE, 5-((4-(2-((6-(4-M...	
Codes: CAS: 763108-62-9 🔗 EVMPD: SUB182020 PUBCHEM: 15951505 🔗	
Relationships: 1	
Formula: C ₂₄ H ₂₄ N ₄ O ₅ S.H ₂ O ₄ S	
Mol Weight: 578.62	
Search	
NEAMINE HYDROCHLORIDE	
UNII:EIU453IDVS	
ABSOLUTE	
Names: D-STREPTAMINE, 2-DEOXY-4-O-(2,6-DIAMINO-... NEOMYCIN A, TETRAHYDROCHLORIDE	

<https://qinas.ncats.nih.gov/qinas/app/substances>

<https://tripod.nih.gov/qinas/#/api>

Links

Full dataset json

<https://tripod.nih.gov/ginias/downloads/dump-public-2019-04-03.gsrs>

Collaborator Slack channel (ask us to add you)

<https://gsrscollaborator.slack.com>

GSRS New User Interface (UI)

Main Goals

- Decouple UI/ client side logic from monolithic application
- Continue to support current application until decoupling finished

Why

- Use newer web technology typescript, angular 6 instead of legacy Play scala templates that are harder to write, slower, and require backend processing
- Better performance
 - Use REST API only remove back end processing
 - SPA loads most assets at once
- Leads to more efficient way of adding new features
- Easier for future development of custom UIs for different organizations

Beta UI

Record Status

■ □ Validated (UNII)	494
■ □ FAILED	6
■ □ pending	1

Substance Type

■ □ chemical	480
■ □ polymer	15
■ □ protein	6

Source Tag

Code System

ATC Level 1

ATC Level 2

ATC Level 3

ATC Level 4

Moiety Type

Stereochemistry

Relationships

Sort By

Items per page: 10

1 - 10 of 501

ACHIRAL

3,4-DICHLOROCINNAMIC ACID

480625A7SY

Names: 3,4-DICHLOROCINNAMIC ACID ✓
2-PROPOENOIC ACID, 3-(3,4-DICHLOROPHENYL)-
3',4'-DICHLOROCINNAMIC ACID
CINNAMIC ACID, 3,4-DICHLORO-
NSC-518800

Codes: CAS: [1202-39-7](#)
EC (EINECS): 214-866-4
BDNUM: 0000113AB

Mol. Weight: 217.049
Formula: C9H6Cl2O2

OXANAMIDE

050271194T

MIXED

Names: OXANAMIDE ✓
oxanamida
oxanamidum
оксанамид
أوكساناميد

Codes: CAS: [126-93-2](#)
WIKIPEDIA: [OXANAMIDE](#)
EVMPD: SUB09495MIG
INN: [812](#)
BDNUM: 0000114AB

Beta UI - Substance View

Overview > **3,4-DICHLOROCINNAMIC ACID**

Structure >

Names 5 >

Identifiers 3 >

Notes 1 >

Audit Info >

References 7 >

Moieties 1 >

Overview

Substance Class **chemical**

Record UNII 480625ATSY

Record Status **Validated (UNII)**

Record Version 1

Show Definitional References ▾

Structure

Stereochemistry **ACHIRAL**

Molecular Formula **C₉H₆Cl₂O₂**

Molecular Weight **217.049**

Optical Activity **UNSPECIFIED**

Defined Stereocenters 0

E/Z Centers 1

Charge 0

Show References ▾

Show SMILES / InChi ▾

Names

Search

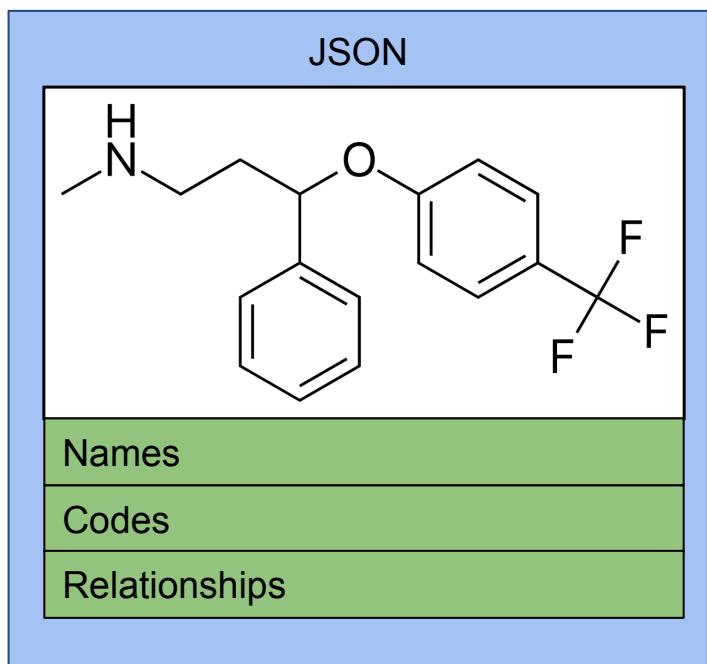
Name	Type	Language	References
3,4-DICHLOROCINNAMIC ACID	Systematic Name	English	View
2-PROPENOIC ACID, 3-(3,4-DICHLOROPHENYL)-	Common Name	English	View
3',4'-DICHLOROCINNAMIC ACID	Common Name	English	View
CINNAMIC ACID, 3,4-DICHLORO-	Common Name	English	View
NSC-518800	Code	English	View

Items per page: 5 ▾ 1 - 5 of 5 | < < > > |

Questions

Molecular Structure Image Recognition

- G-SRS is intended to produce, consume, curate and exchange highly-structured substance data



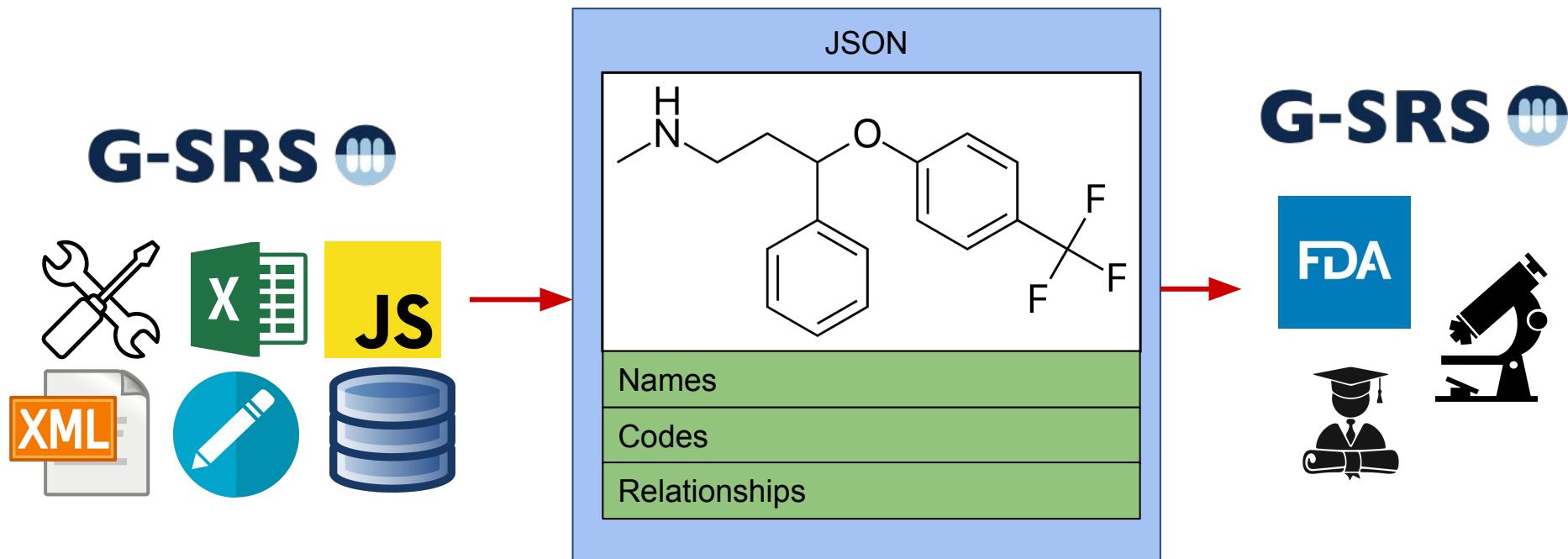
G-SRS



- Tracking
- Safety
- Research
- Analysis

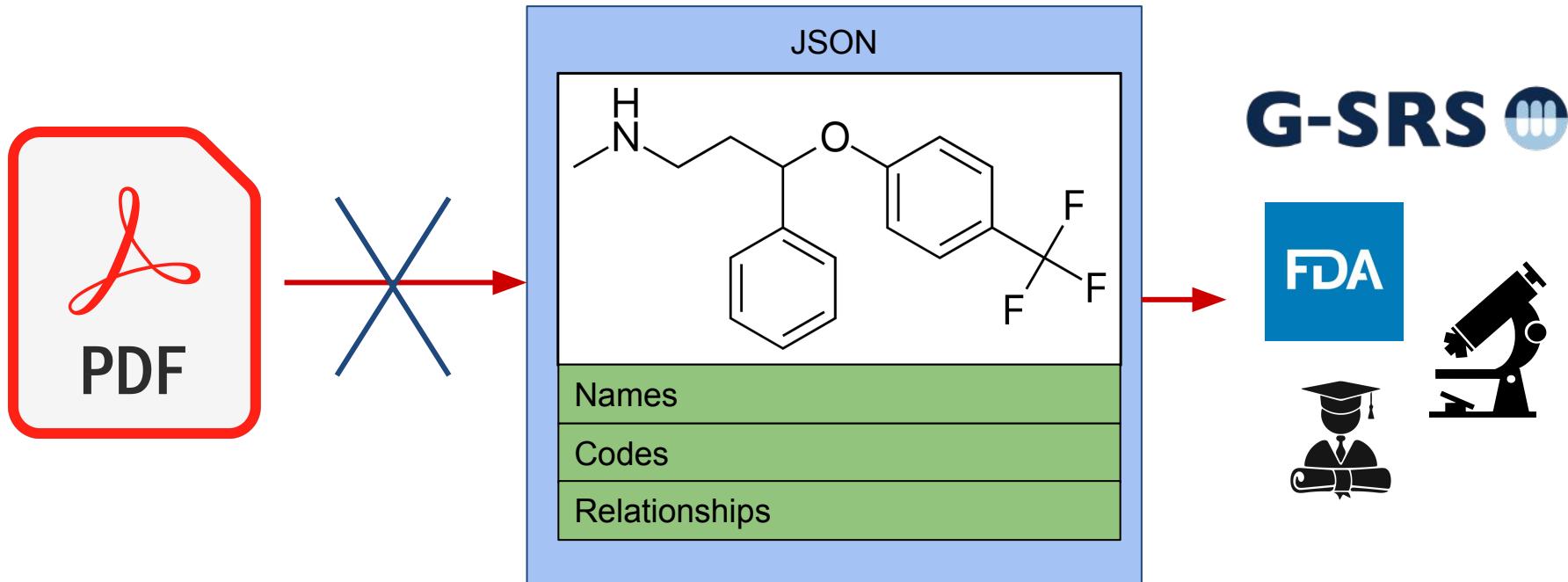
Molecular Structure Image Recognition

- For structured formats and systems transformation, exchange of data is pretty easy



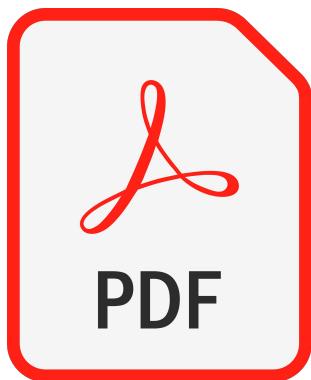
Molecular Structure Image Recognition

- For structured formats and systems transformation, exchange of data is pretty easy
- Unstructured data is harder to process

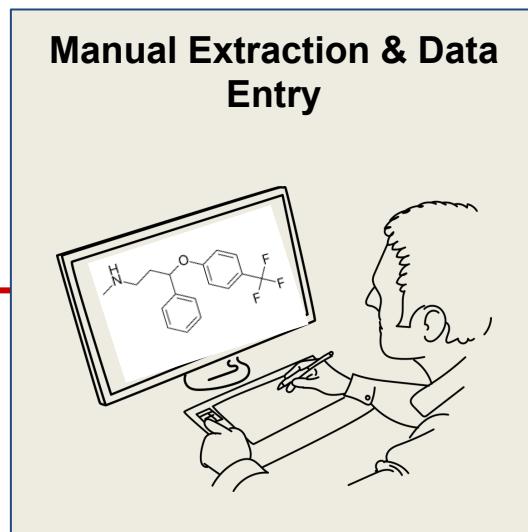


Molecular Structure Image Recognition

- For structured formats and systems transformation, exchange of data is pretty easy
- Unstructured data is harder to process
- A large portion of legacy data **is** unstructured



- APIs
- Impurities
- Metabolites
- Legacy Data



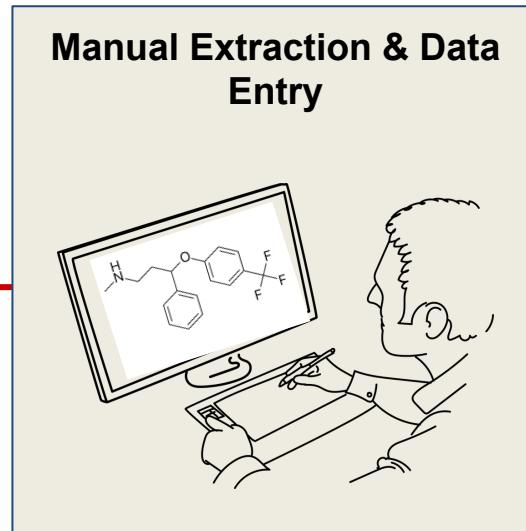
G-SRS



Molecular Structure Image Recognition



- APIs
- Impurities
- Metabolites
- Legacy Data

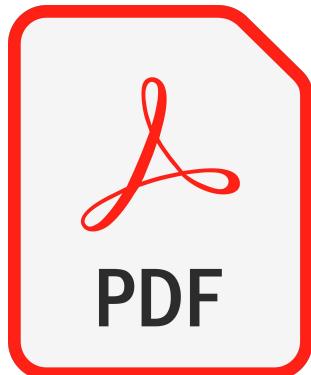


- Tedious and slow
- Error-prone

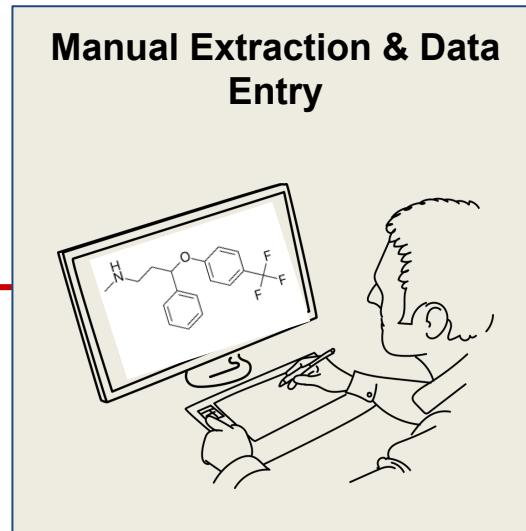
G-SRS



Molecular Structure Image Recognition



- APIs
- Impurities
- Metabolites
- Legacy Data



- Tedious and slow
- Error-prone
- Carpal Tunnel Syndrome

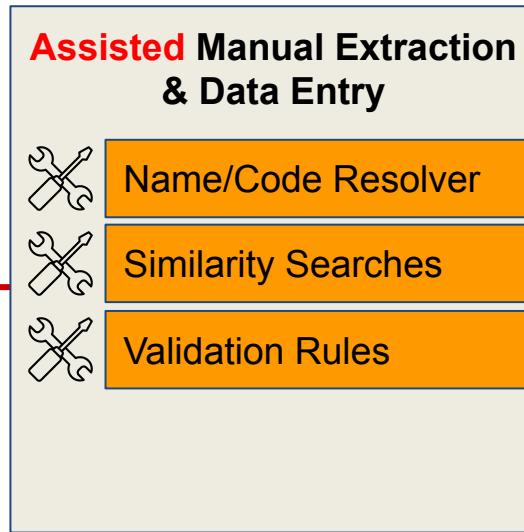
G-SRS



Molecular Structure Image Recognition



- APIs
- Impurities
- Metabolites
- Legacy Data



G-SRS

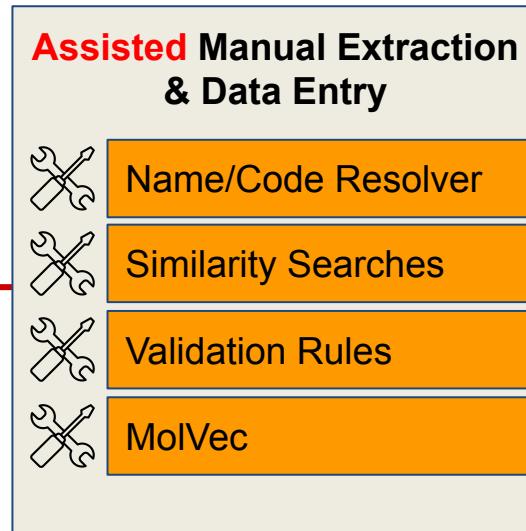


- **Less** Tedious and slow
- **Less** Error-prone
- Carpal Tunnel Syndrome

Molecular Structure Image Recognition



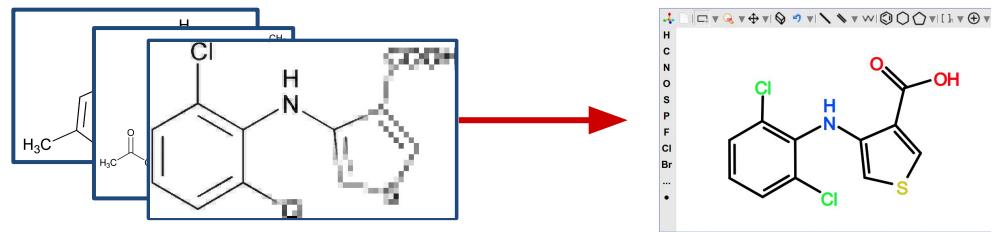
- APIs
- Impurities
- Metabolites
- Legacy Data



- **Less** Tedious and slow
- **Less** Error-prone
- **Less** Carpal Tunnel Syndrome*

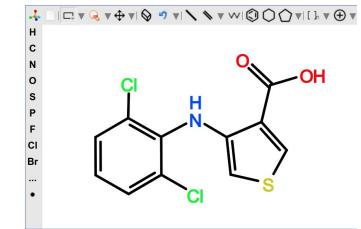
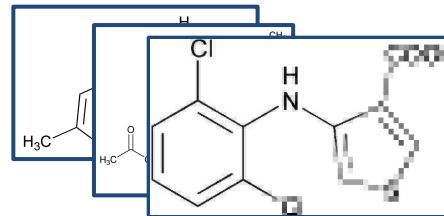
G-SRS: MolVec

- **MolVec** takes chemical structure *images* and turns them into chemical structure **data**

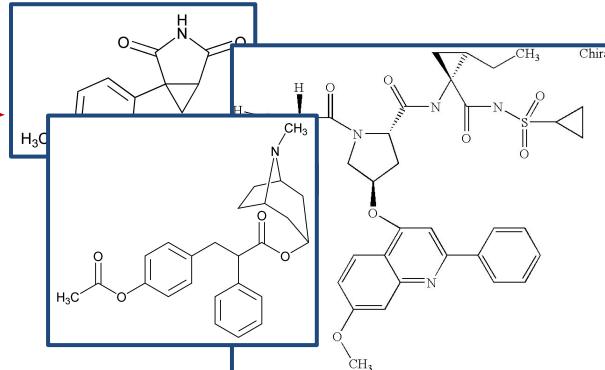


G-SRS: MolVec

- **MolVec** takes chemical structure *images* and turns them into chemical structure **data**



- APIs
- Impurities
- Metabolites
- Legacy Data

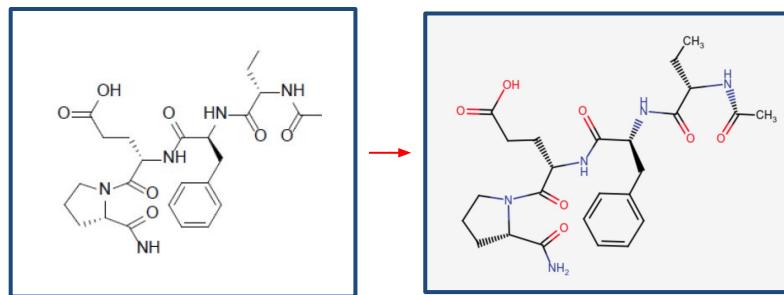


→ **G-SRS** 

+ manual extraction /
curation

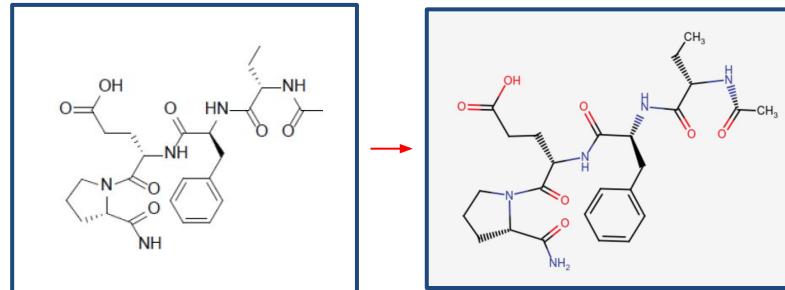
G-SRS: MolVec

How it Works

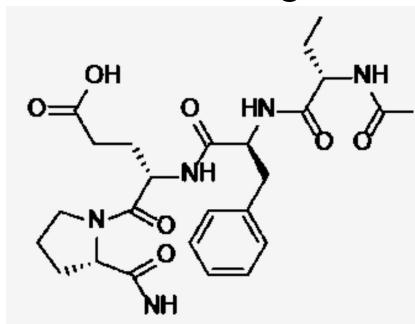


G-SRS: MolVec

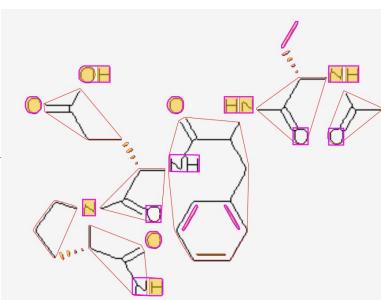
How it Works



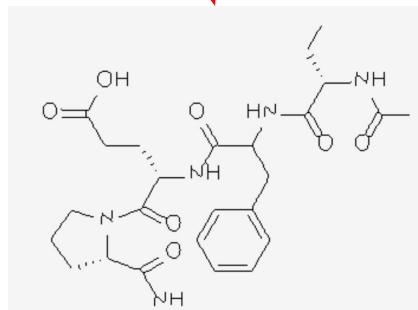
Thresholding



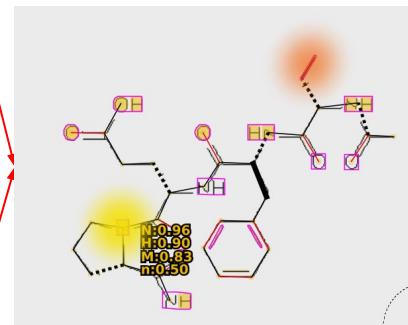
Shape and Feature Detection



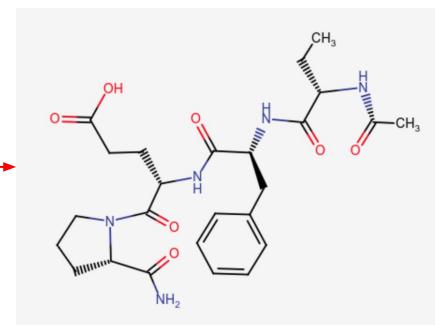
Thinning



Node and Edge Detection



Heuristics and Adjustments

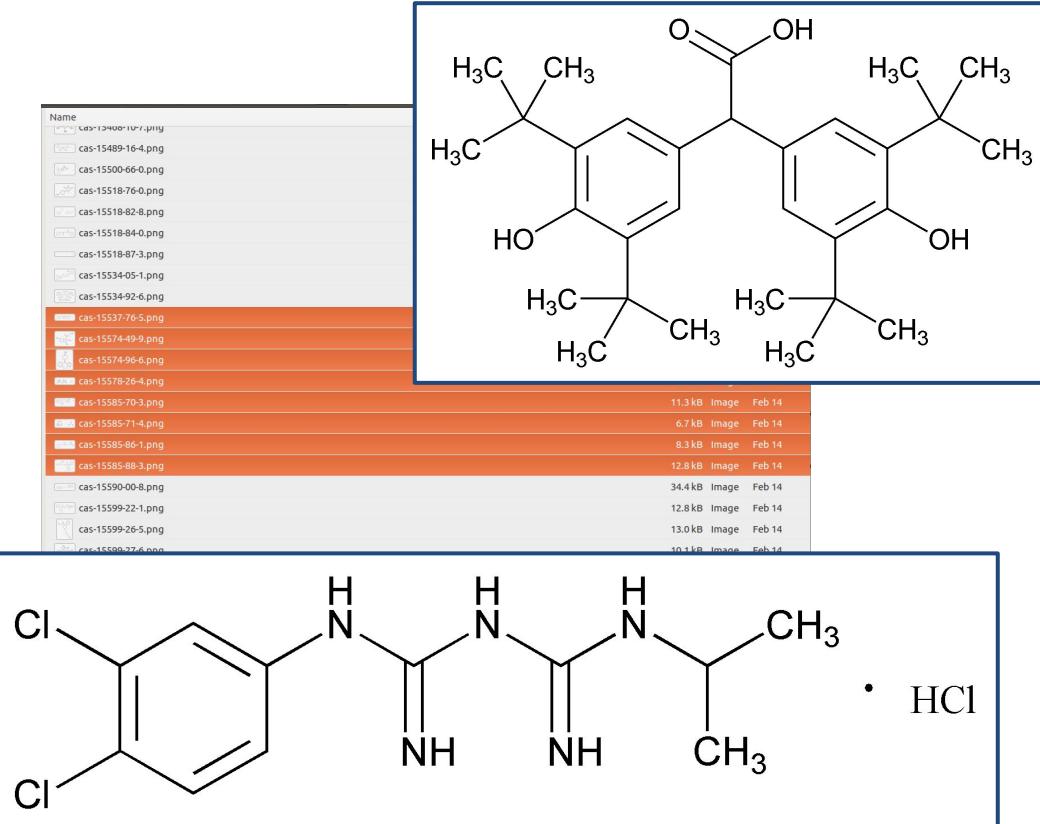


Assembly

G-SRS: MolVec

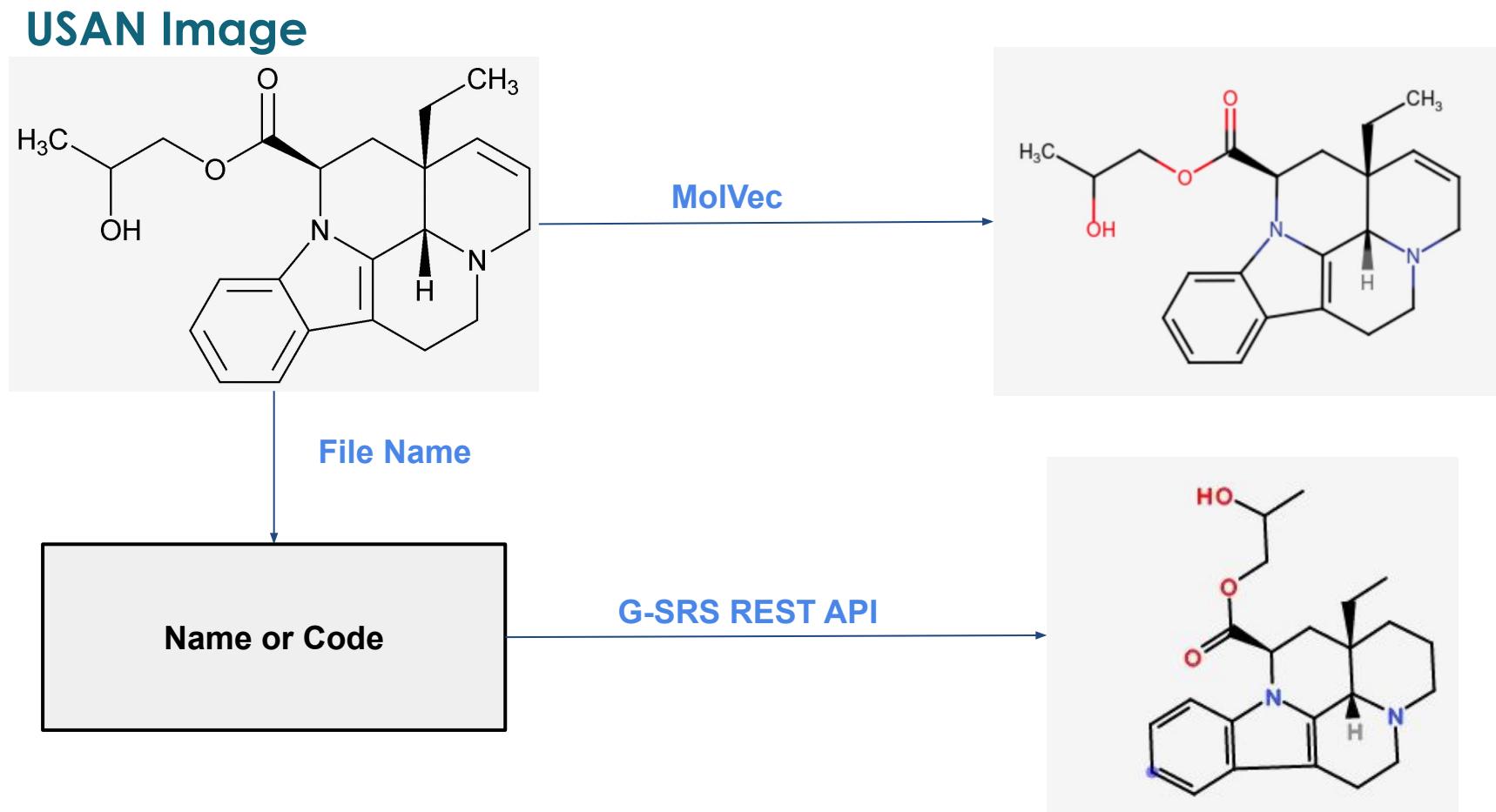
Evaluation Use Case

- **USAN** chemical structure image files
 - ~7500 png files
 - Well-drawn
 - Highly-curated
 - Real and meaningful
 - Definitions already known to FDA G-SRS



G-SRS: MoIVec

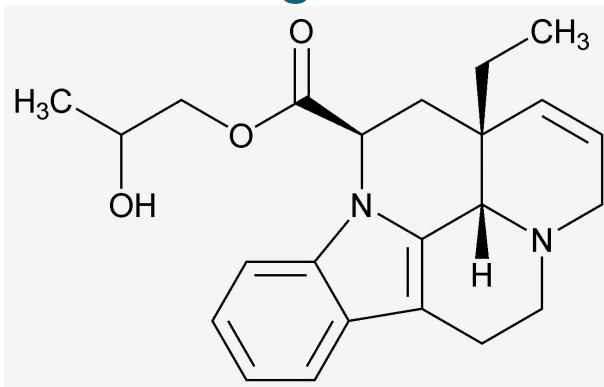
Evaluation Use Case



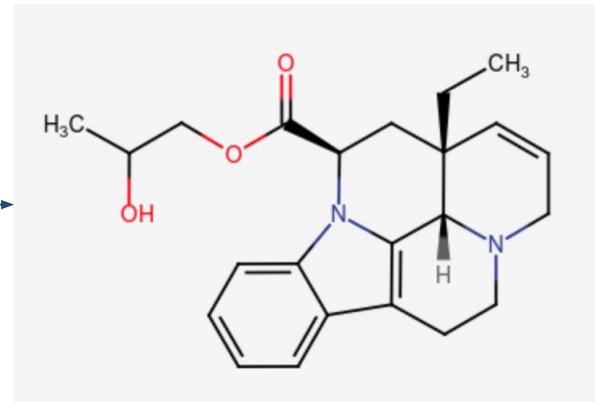
G-SRS: MolVec

Evaluation Use Case

USAN Image



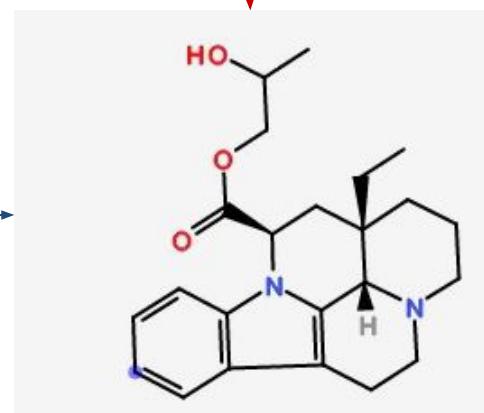
MolVec



File Name

Name or Code

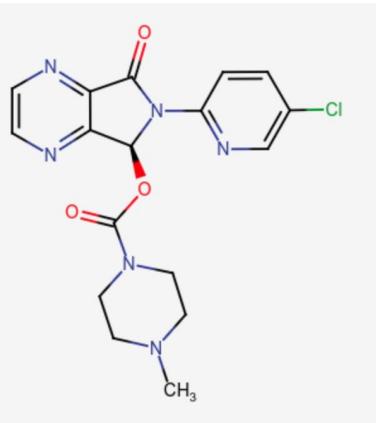
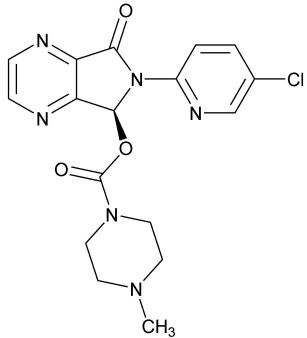
G-SRS REST API



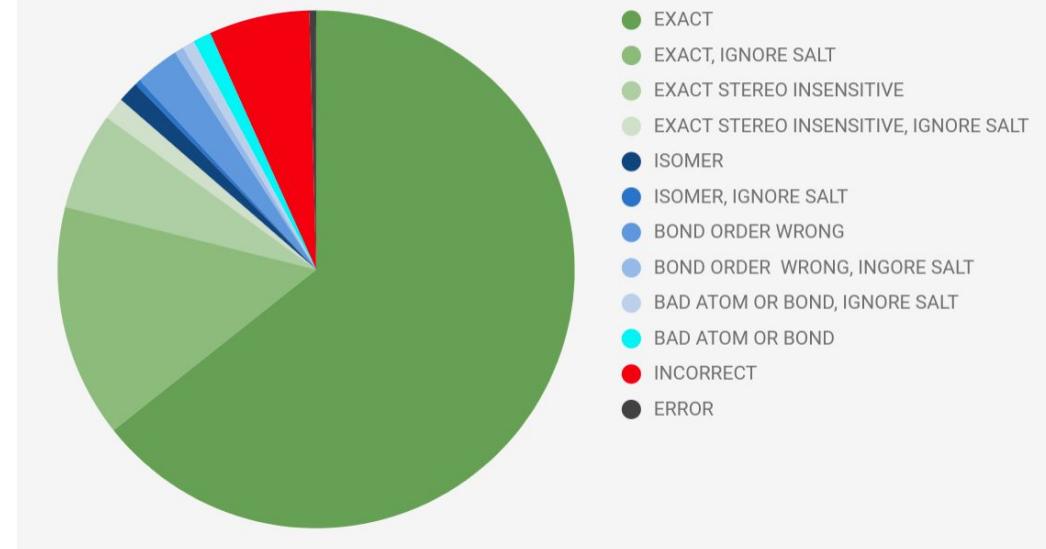
G-SRS Structure

G-SRS: MolVec

Evaluation Use Case

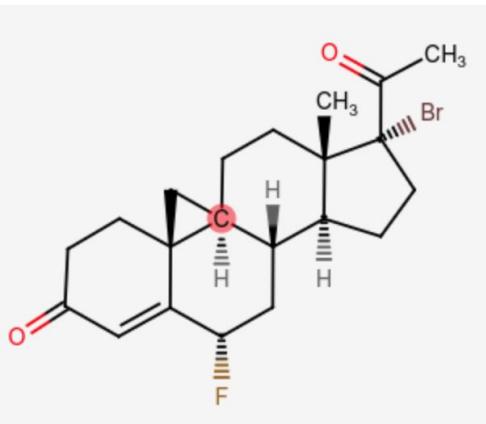
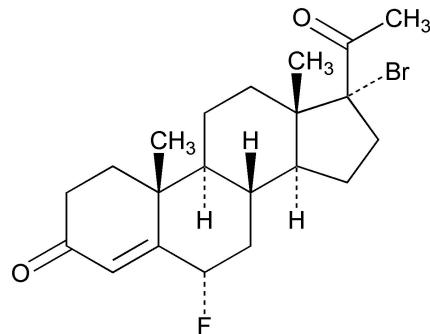


MolVec Accuracy on USAN Set

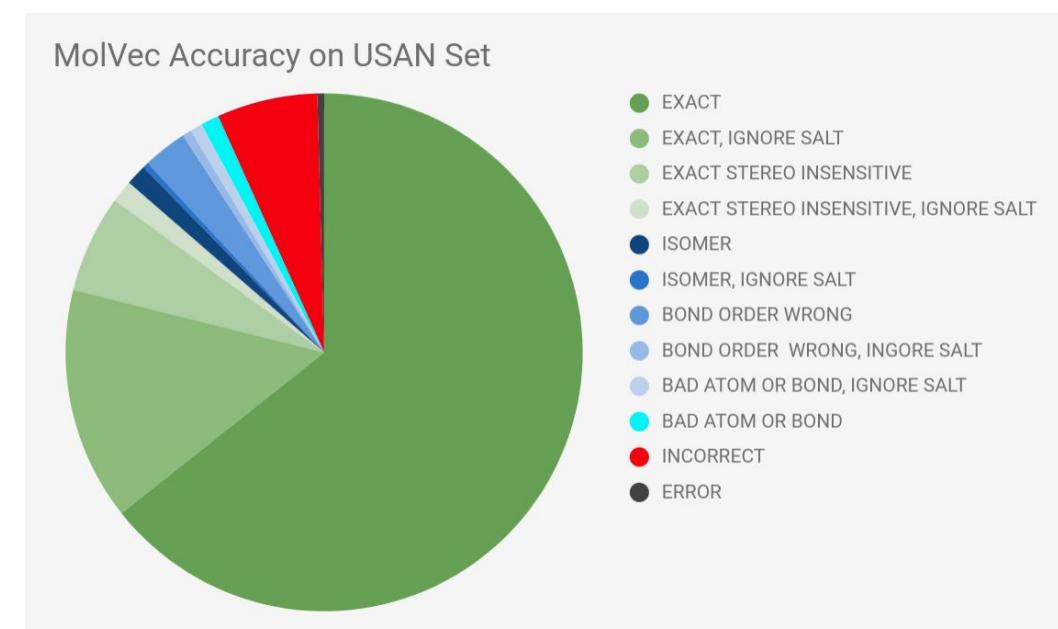


G-SRS: MolVec

Evaluation Use Case

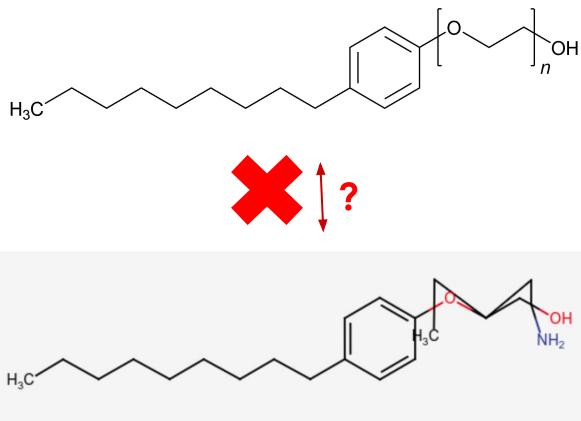


MolVec Accuracy on USAN Set

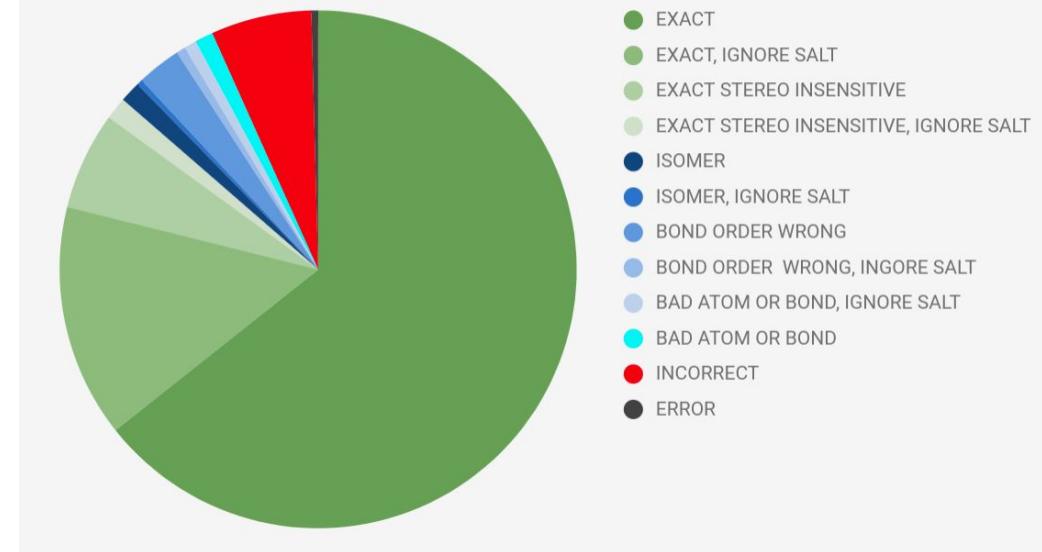


G-SRS: MolVec

Evaluation Use Case



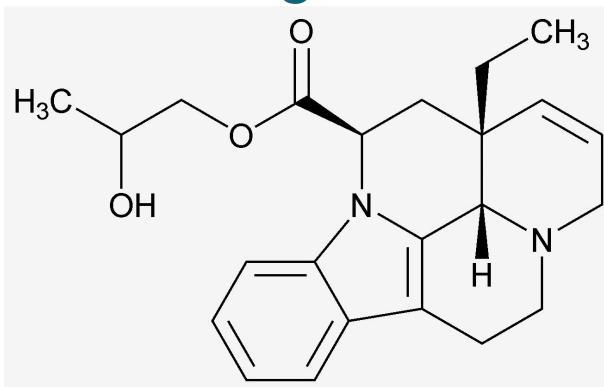
MolVec Accuracy on USAN Set



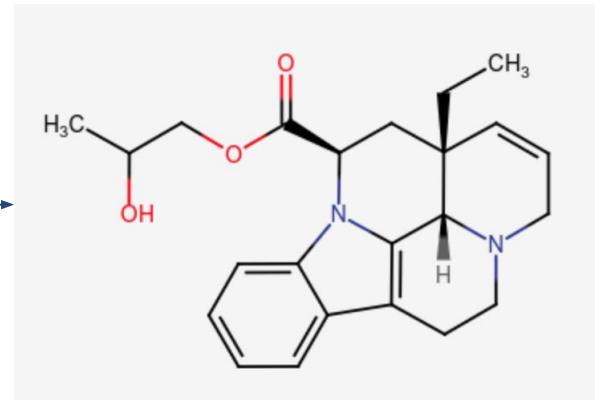
G-SRS: MolVec

Evaluation Use Case

USAN Image



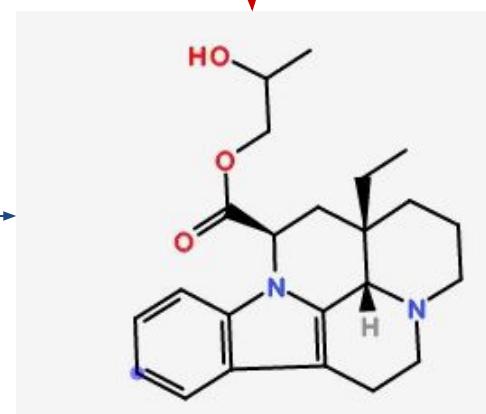
MolVec



File Name

Name or Code

G-SRS REST API

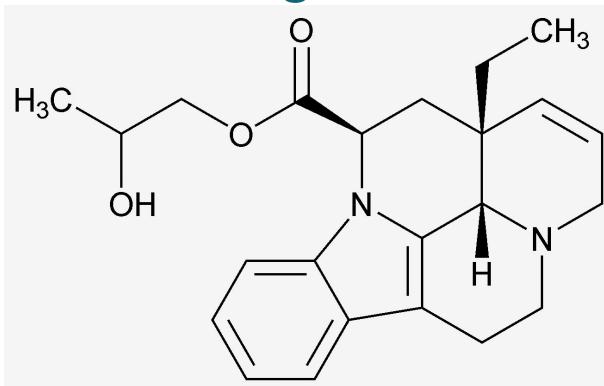


G-SRS Structure

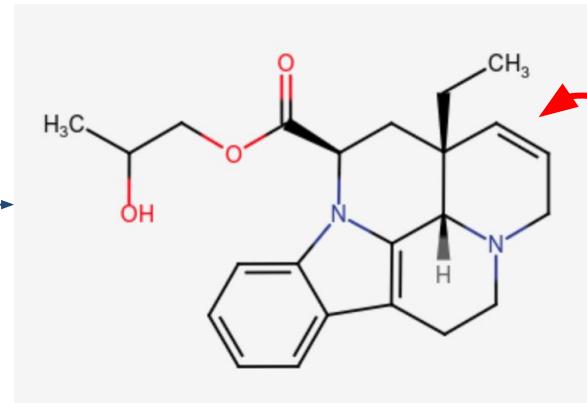
G-SRS: MolVec

Evaluation Use Case

USAN Image



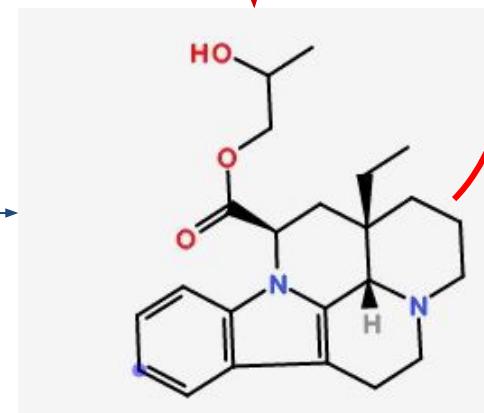
MolVec



File Name

Name or Code

G-SRS REST API

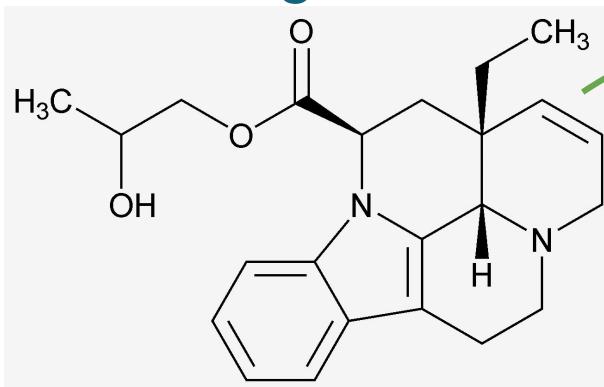


G-SRS Structure

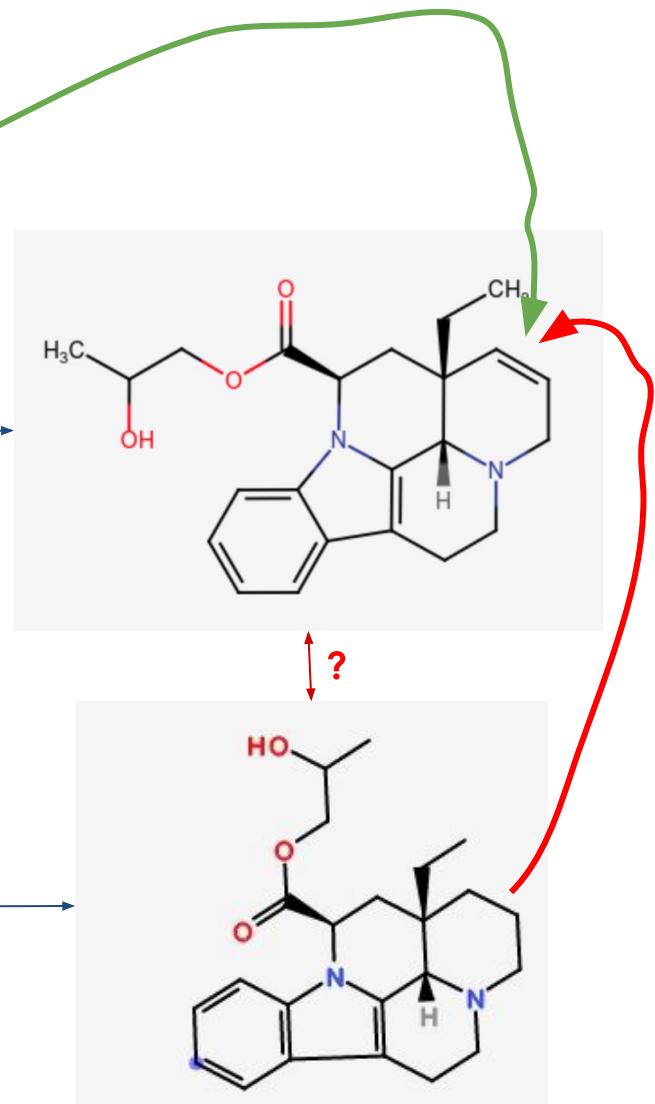
G-SRS: MolVec

Evaluation Use Case

USAN Image



MolVec



File Name

Name or Code

G-SRS REST API

G-SRS Structure

G-SRS: MolVec

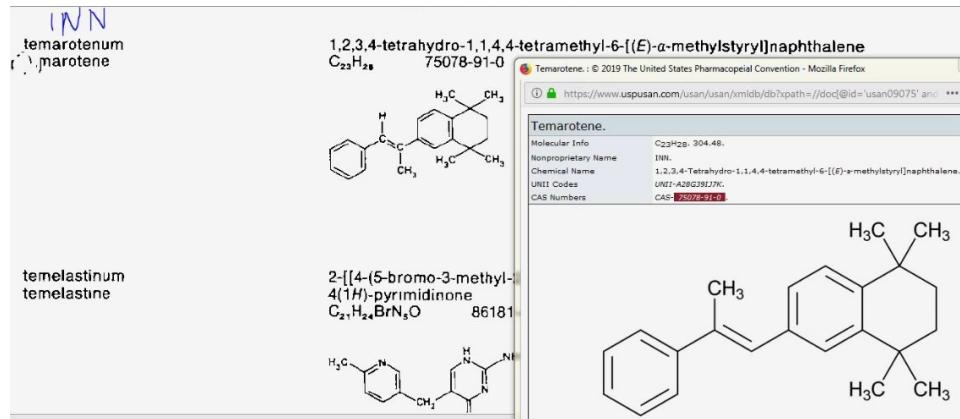
Evaluation Use Case

- 20% of “wrong” **MolVec** images were actually accurate to the image
 - Either the **USAN Image** or **FDA G-SRS** must be incorrect

G-SRS: MolVec

Evaluation Use Case

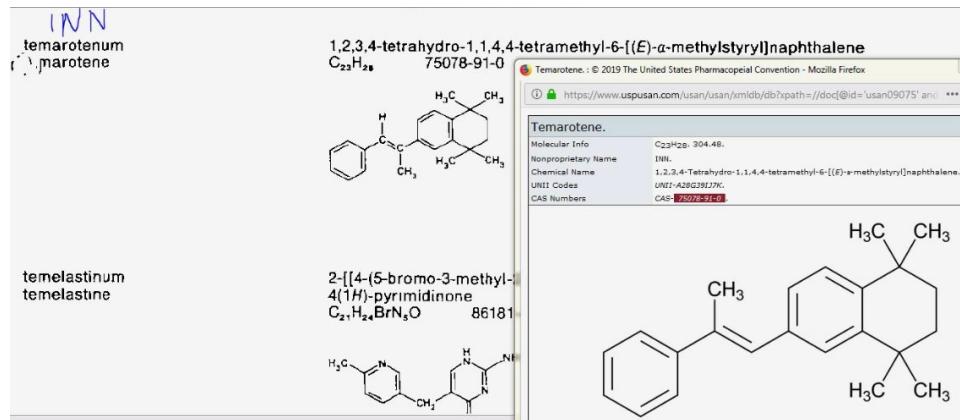
- **1** structure was shown to be wrong in G-SRS, and corrected
- **25** structures were shown to be wrong in the **USAN Images**
 - All but **9** had already been fixed (all changes minor)



G-SRS: MolVec

Evaluation Use Case

- **1** structure was shown to be wrong in G-SRS, and corrected
- **25** structures were shown to be wrong in the **USAN Images**
 - All but **9** had already been fixed (all changes minor)

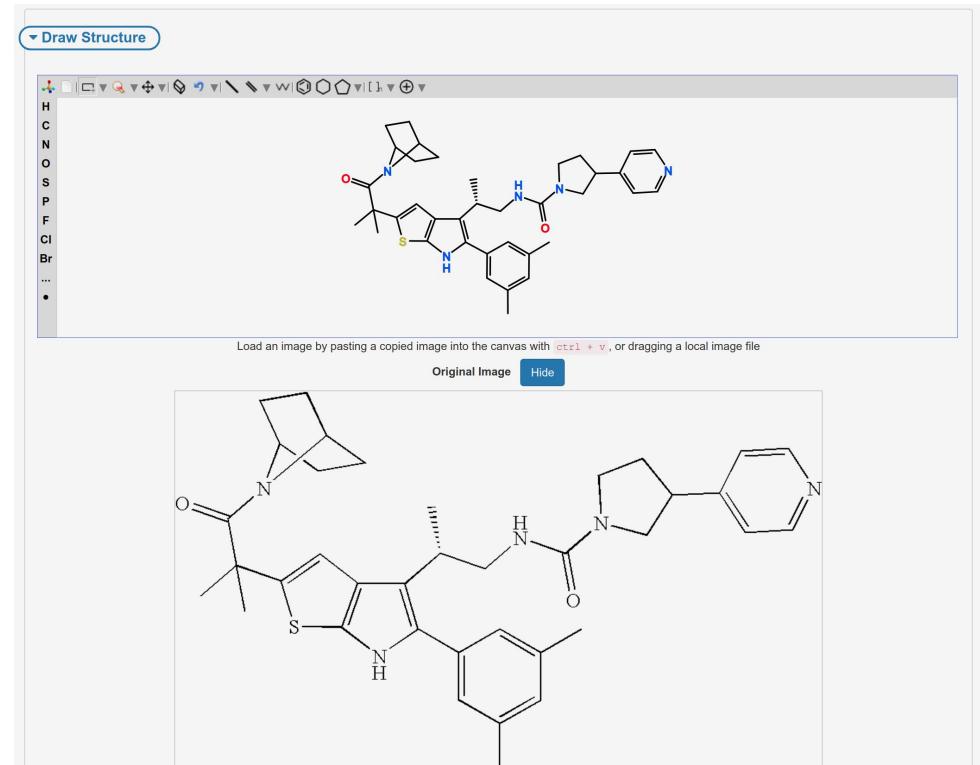


MolVec, while imperfect, can help in bulk curation

G-SRS: MolVec

Evaluation Use Case

- Useful to bootstrap structure searches from *inside* G-SRS
- Useful to bootstrap chemical registrations

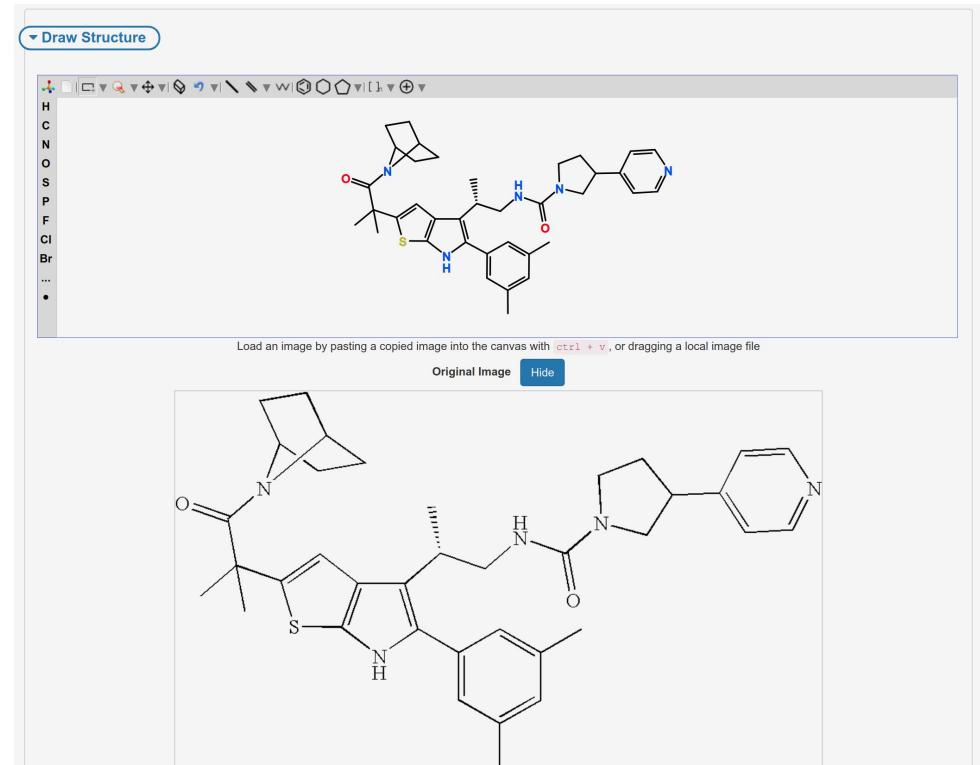


G-SRS: MolVec

DEMO

G-SRS: MolVec

- Currently used within FDA both via API and via G-SRS registration forms / search
- Rapidly improving due to feedback



G-SRS: MolVec

Features

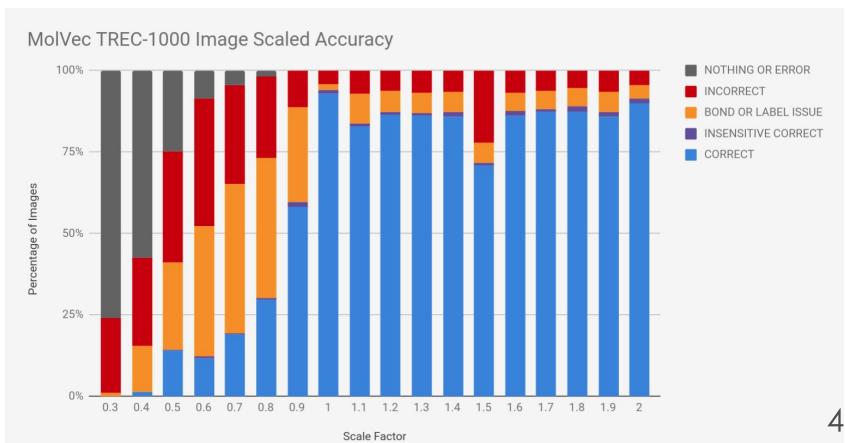
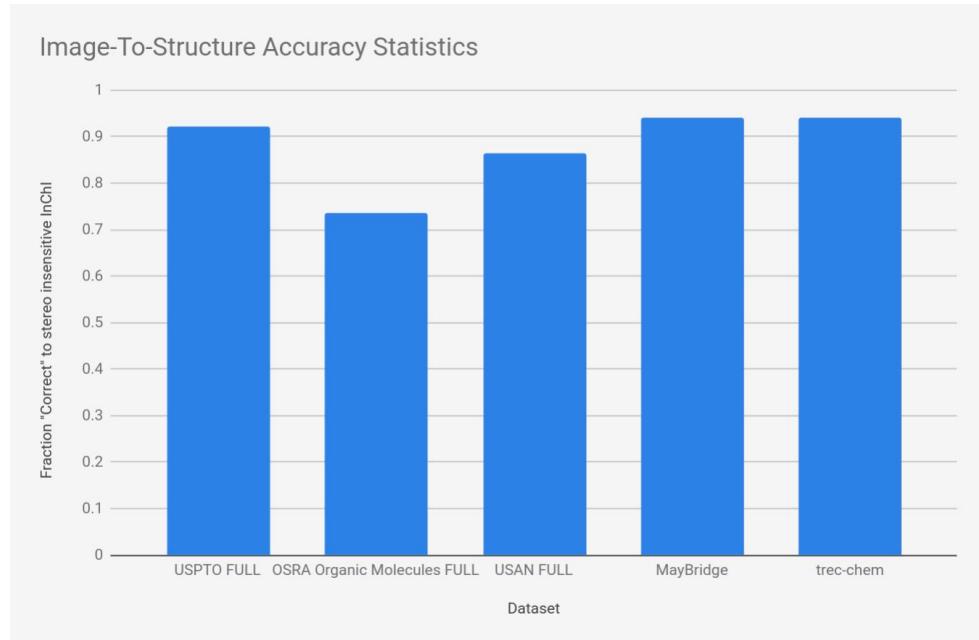
- Free, self-contained Java code, distributable, no server configuration
- Pretty accurate
- Accurately pretty
- Fairly Fast



G-SRS: MolVec

Features

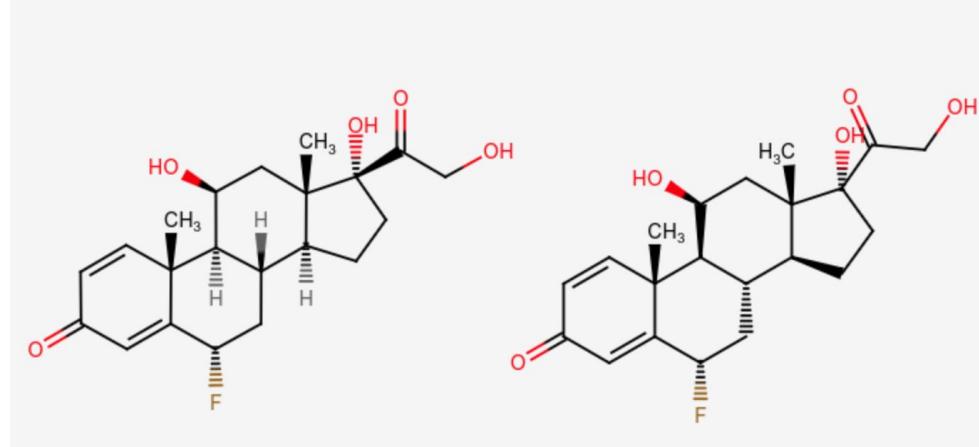
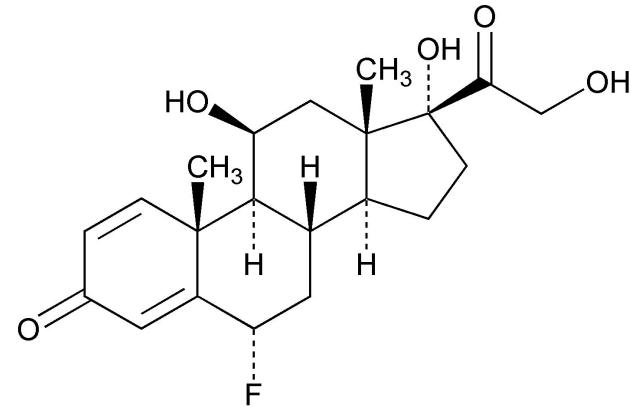
- Free, self-contained Java code, distributable, no server configuration
- **Pretty accurate**
- Accurately pretty
- Fairly Fast



G-SRS: MolVec

Features

- Free, self-contained Java code, distributable, no server configuration
- Pretty accurate
- **Accurately pretty**
- Fairly Fast



From MolVec

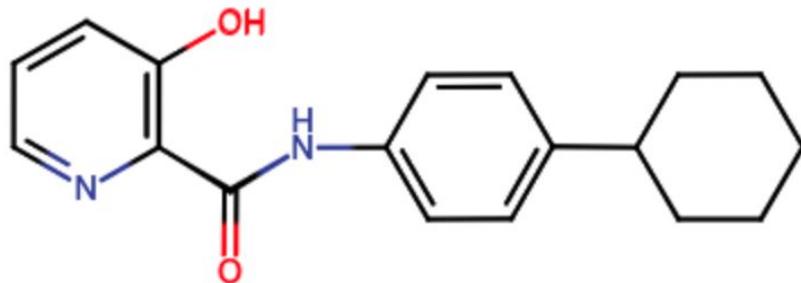
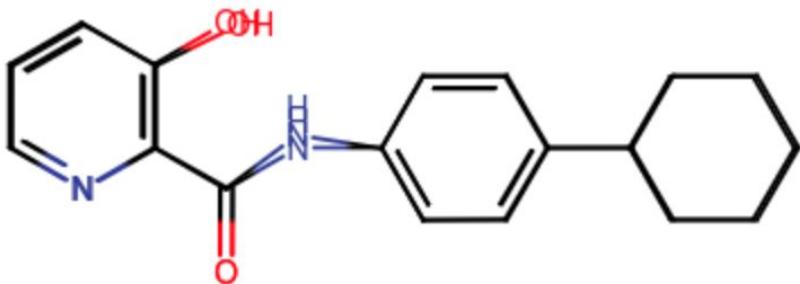
From IUPAC Name

G-SRS: MolVec

Features

- Free, self-contained Java code, distributable, no server configuration
- Pretty accurate
- **Accurately pretty**
- Fairly Fast

Pretty Good Alignment To Image



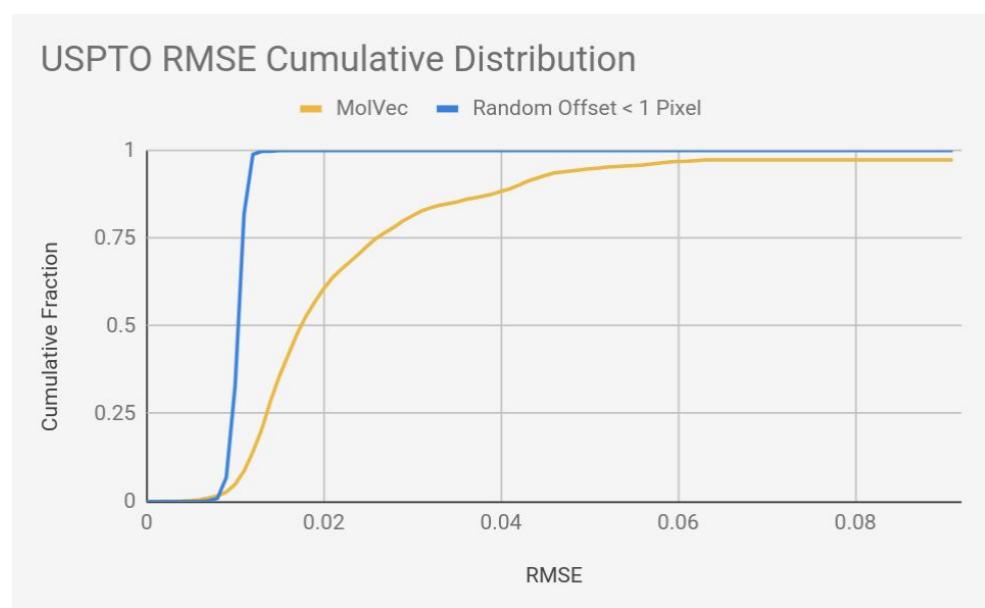
Better Alignment To Image

G-SRS: MolVec

Evaluation Use Case

Features

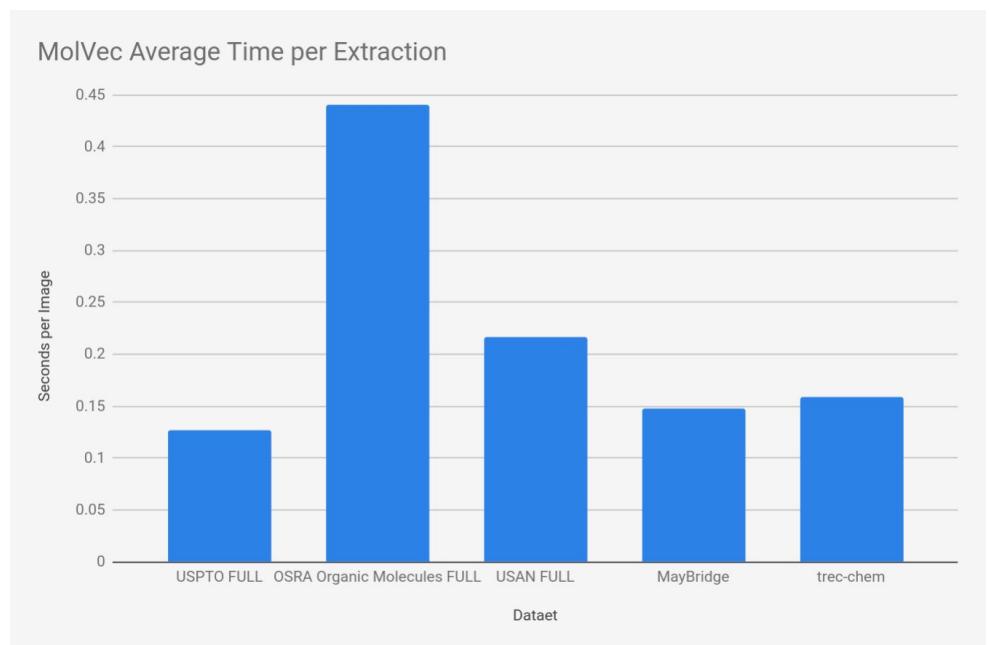
- Free, self-contained Java code, distributable, no server configuration
- Pretty accurate
- **Accurately pretty**
- Fairly Fast



G-SRS: MolVec

Features

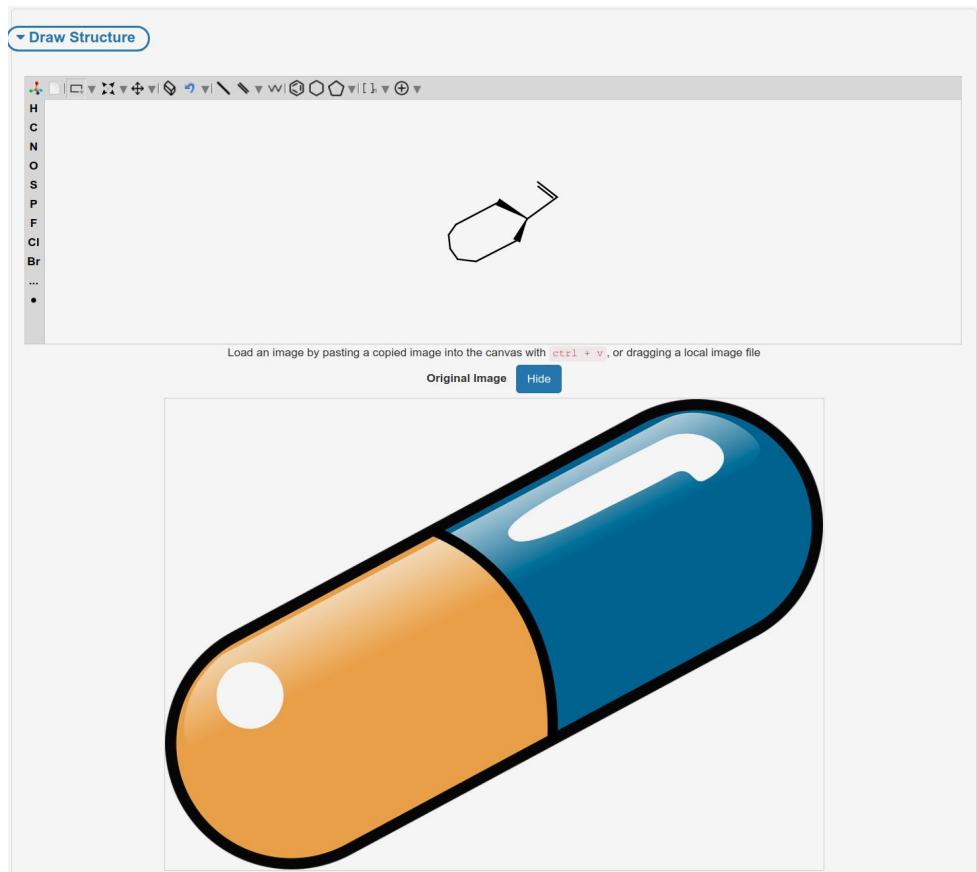
- Free, self-contained Java code, distributable, no server configuration
- Pretty accurate
- Accurately pretty
- **Fairly Fast**



G-SRS: MolVec

Features

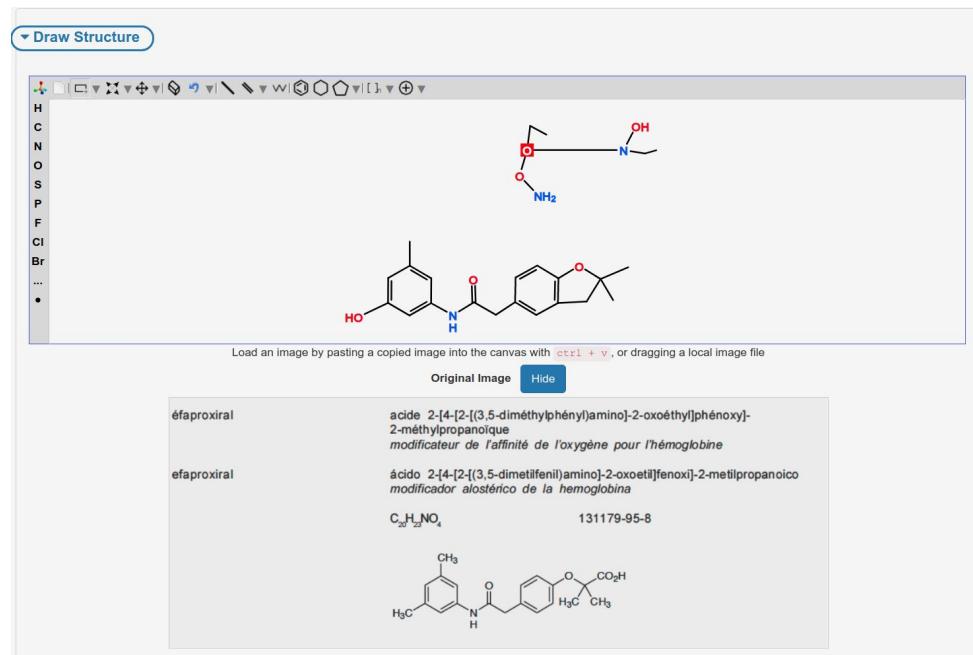
- Free, self-contained Java code, distributable, no server configuration
- Pretty accurate
- Accurately pretty
- Fairly Fast
- **NOT MAGIC**



G-SRS: MolVec

Features

- Free, self-contained Java code, distributable, no server configuration
- Pretty accurate
- Accurately pretty
- Fairly Fast
- **NOT MAGIC**



G-SRS: MolVec

Try it out Online:

<https://ginas.ncats.nih.gov/ginias/app/structure>

or

<https://predictor.ncats.io>

(structure prediction -> use image)

As part of GSRS 2.3.3:

<https://tripod.nih.gov/ginias/>

MolVec Git Repository:

<https://spotlite.nih.gov/ncats/molvec>

Acknowledgements

NCATS

Tyler Peryea
Danny Katzel
Jorge Neyra
Niko Anderson
Mitch Miller
Sarah Stemann
Noel Southall
Dac-Trung Nguyen
Ivan Grishagin
Dammika Amugoda
Chris LeClair
Paul Shinn

FDA

Larry Callahan
Frank Switzer
Archana Newatia
Yulia Borodina
Ramez Ghazzaoui
Elaine Johansen
Ta-Jen Chen
Mary-Ann Slack
Alex Welsch
Yoshiyuki Tokiwa
George Washburn
Sabrina Mosley

Medicines Evaluation Board

Herman Diederik	Ciska Matai
Marcel Hoefnagel	Burt Kroes
Joris Kampmeijer	

US Pharmacopeial Convention

Fouad Atouf	Ahmed Abdulhadi
Andrej Wilk	Steven Emrick
Tina Morris	

Health Canada

Craig Anderson
Daniel Buijs
Ted Kim

Questions

Thank You