

An open-source, hackable, enterprise-ready chemical registration system

Oleksandra Serhiienko, Datagrok

https://github.com/datagrok-ai/mol-track



Do we really need another registration system?

Yes, we do!

What's out there?

Commercial offerings:

- powerful, but expensive
- not self-hosted
- not easily customizable

Lwreg:

- great start!
- self-hosted and hackable

Who is left behind?



Academic labs focused on synthesis



Small biotechs just starting out



Organizations with legacy tools



Groups looking to explore new modalities



MolTrack: A hackable and flexible system

I have a dream...

What if we could combine the ease and intelligence of commercial systems with the flexibility, affordability, and control of open-source?



Fast



Open-source

Compound/batch/assay

Hackable but enterprise-ready

Dynamic properties

Chemical intelligence

Postgresbacked

Searches, aggregations



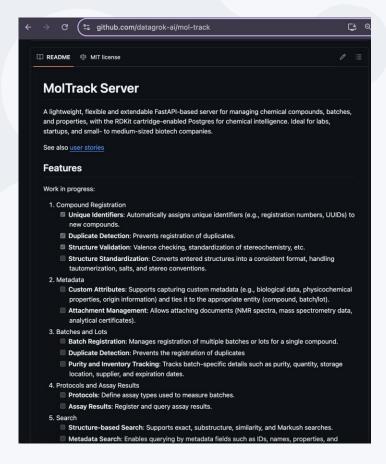


MIT license



- Free as in beer
- Community can adapt the code to their own workflows





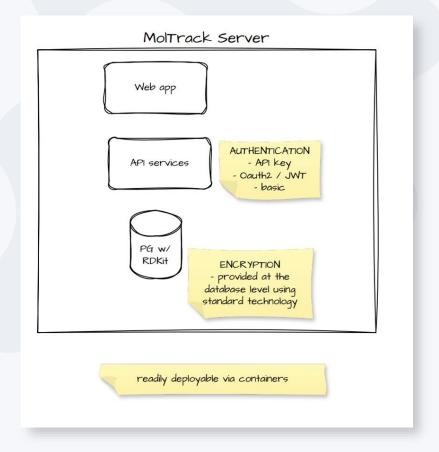


Postgres-backed



Leverages the power of relational databases:

- Transactions
- Backups and replication
- Querying and indexing
- Scalability
- Cloud RDS





Hackable

- MIT license
- Python plugin system
- Events and web hooks
- Modular architecture

...but enterprise-ready

- Postgres: managed cloud databases, transactions...
- Authentication, API keys, tokens
- Users, roles, privileges
- Configuration
- Molecule-level and property-level permissions
- Audit trail
- Integration-ready

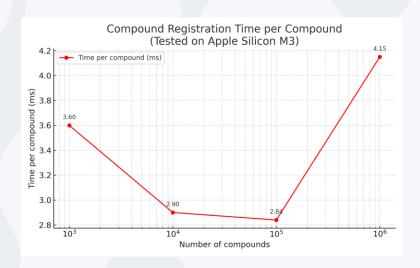
Get it up and running for your lab in a minute - or use as a backbone for the enterprise solution



Fast



1000 compounds	3.57s	3.6 ms/compound
10,000 compounds	29.4s	2.9 ms/compound
100,000 compounds	284s	2.84 ms/compound
1,000,000 compounds	4149s	4.15 ms/compound

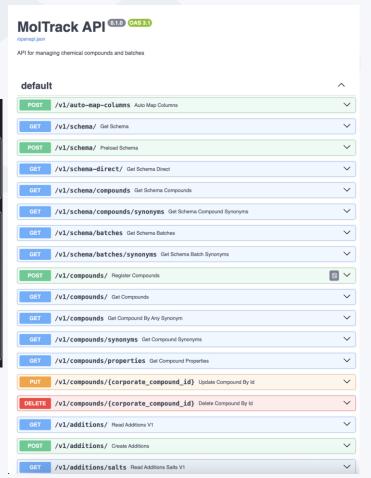




CLI

```
Usage: mtcli.py [OPTIONS] COMMAND [ARGS]...
 --install-completion
                               Install completion for the current shell.
                               Show completion for the current shell, to copy it or
 --show-completion
                               customize the installation.
 --help
                               Show this message and exit.
  Commands
 schema
              Schema management commands
              Compound management commands
 compounds
 batches
              Batch management commands
 properties
              Property management commands
 additions
              Addition management commands
              Assays management commands
 assays
 database
              Database management commands
 directory
              Directory loading commands
              CSV utility commands
 csv
              Search functionality
 search
              Administrative functions
 admin
```

CLI provides full access to Moltrack API capabilities or you can interact directly via the OpenAPI page.





Chemical Intelligence...powered by RDKit

- Cleanup
- Standardization
- De-duplication
- Validation
- Sketching
- Search via hash, substructure, similarity
- Configurable
- Extensible

- RDKit toolkit
- RDKit postgresql cartridge
- Standardization & Filtering methods
- Molecular Hashes

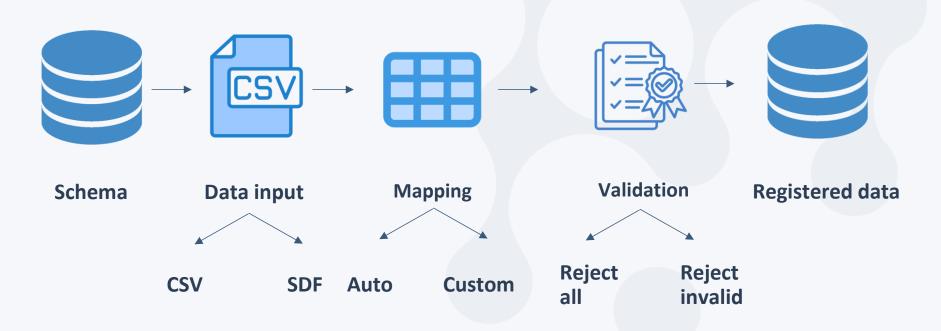
Administrators can define, extend, and enforce chemistry data quality

Schemas and properties

- Structured metadata for compounds, batches, assays
- Used for validation, field mapping
- Extensible with new properties

	Properties		
Name	Entity Type	Value Type	
CAS CHEMBL ID Common Name corporate_compound_id IUPAC Name MolLogP Source Source Compound Code USAN Acquired Date corporate_batch_id ELN Reference EPA Batch ID Project Purity Responsible Party	Entity Type COMPOUND COMPOUND COMPOUND COMPOUND COMPOUND COMPOUND COMPOUND COMPOUND BATCH BATCH BATCH BATCH BATCH BATCH BATCH BATCH BATCH	string string string string double string string string string string string datetime string string string	
Source Source Batch Code Synthesized Date assay format biological system Assay Run Date Assayer Cell Concentration Cell Lot Cell Species pIC50	BATCH BATCH BATCH BATCH ASSAY ASSAY ASSAY_RUN ASSAY_RUN ASSAY_RUN ASSAY_RUN ASSAY_RUN ASSAY_RUN ASSAY_RUN ASSAY_RUN	string string string datetime string string datetime string double string string double	

Registration: Bringing data to life



MolTrack UI: a Datagrok plugin



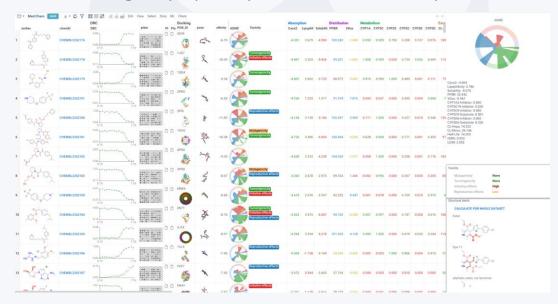
Why Datagrok?

- Swiss-Army knife for cheminformatics
- First-class cheminformatics support
- Rendering, sketching
- Chemical space, activity cliffs
- Contextual info on molecules
- Connecting to databases
- ADME, Boltz-2
- Docking, folding, retrosynthesis
- Hit Design, Assay Plates

MolTrack service: MIT licensed

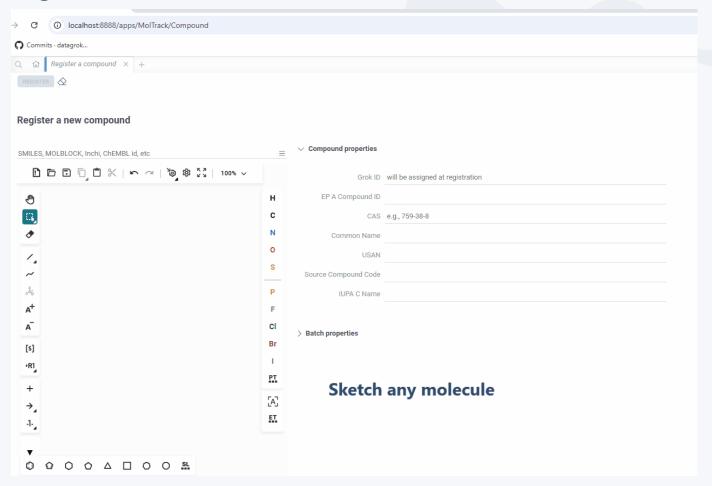
MolTrack plugin: MIT licensed; depends on Datagrok

Datagrok: proprietary; free for academic and personal use



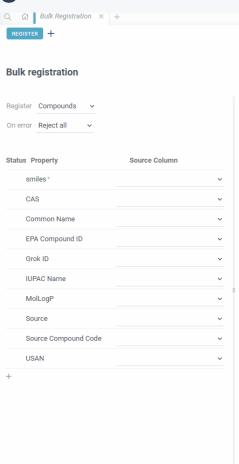
Single registration





Bulk registration





Easy file import

Drag and drop a CSV file here, or click to select a file.

Search



Future plans

Hi datagrok

- Additional modalities: no-structures, peptides, oligos, helm, biomolecules
- Synchronization with commercial registration systems
- Multiuser improvements (groups, roles, etc)
- Staging area to handle large files
- Monomer management
- Mixtures
- Computations and predictions

Acknowledgements



We want to thank RDKit community for the fantastic toolkit, and more specifically lwreg authors (Greg Landrum, Jessica Braun, Paul Katzberger, Marc Lehner, Sereina Riniker).

This has been a joint project between Datagrok and a company whose identity we cannot disclose yet. No further details at this time, but watch this space :)

Datagrok

- Oleksandra Serhiienko
- Ed Jaeger
- Maria Dolotova
- Andrew Skalkin

The TBA Company

- Danica Prodanovic Djapa
- Aleksei Khariukov
- Valentin Anoprenko



Thank You!

